

Intel® MPI Library for Linux* OS

Reference Manual

Contents

Legal Information	4
1. Introduction.....	6
1.1. Introducing Intel® MPI Library	6
1.2. Intended Audience	6
1.3. What's New	6
1.4. Notational Conventions.....	7
1.5. Related Information	7
2. Command Reference	8
2.1. Compiler Commands.....	8
2.1.1. Compiler Command Options	9
2.1.2. Configuration Files	12
2.1.3. Profiles	12
2.1.4. Environment Variables	12
2.2. Simplified Job Startup Command	16
2.3. Scalable Process Management System (Hydra) Commands	19
2.3.1. Global Options.....	19
2.3.2. Local Options.....	34
2.3.3. Extended Device Control Options.....	35
2.3.4. Environment Variables	37
2.3.5. Cleaning up Utility	47
2.3.6. Checkpoint-Restart Support.....	49
2.4. Hetero Operating System Cluster Support	55
2.5. Intel® Xeon Phi™ Coprocessor Support.....	55
2.5.1. Usage Model	55
2.5.2. Environment Variables	57
2.5.3. Compiler Commands	61
2.6. Multipurpose Daemon Commands	62
2.6.1. Job Startup Commands	69
2.6.2. Configuration Files	86
2.6.3. Environment Variables	87
2.7. Processor Information Utility	90
3. Tuning Reference	93
3.1. Using mpitune Utility	93
3.1.1. Cluster Specific Tuning	97
3.1.2. Application Specific Tuning.....	98
3.1.3. Tuning Utility Output.....	101
3.2. Process Pinning	101
3.2.1. Default Settings for Process Pinning	101
3.2.2. Processor Identification	102
3.2.3. Environment Variables	102
3.2.4. Interoperability with OpenMP* API	110
3.3. Fabrics Control	122
3.3.1. Communication Fabrics Control	123
3.3.2. Shared Memory Control	128
3.3.3. DAPL-capable Network Fabrics Control	135
3.3.4. DAPL UD-capable Network Fabrics Control.....	144
3.3.5. TCP-capable Network Fabrics Control	152
3.3.6. TMI-capable Network Fabrics Control.....	154

3.3.7. OFA-capable Network Fabrics Control	156
3.3.8. OFI*-capable Network Fabrics Control	161
3.4. Collective Operation Control.....	162
3.4.1. I_MPI_ADJUST Family	162
3.4.2. I_MPI_MSG Family	172
3.5. Miscellaneous.....	175
3.5.1. Timer Control	176
3.5.2. Compatibility Control	176
3.5.3. Dynamic Process Support.....	177
3.5.4. Fault Tolerance	177
3.5.5. Statistics Gathering Mode	179
3.5.6. ILP64 Support	194
3.5.7. Unified Memory Management.....	195
3.5.8. File System Support	195
3.5.9. Multi-threaded memcpy Support	196
4. Glossary	199
5. Index	200

Legal Information

No license (express or implied, by estoppel or otherwise) to any intellectual property rights is granted by this document.

Intel disclaims all express and implied warranties, including without limitation, the implied warranties of merchantability, fitness for a particular purpose, and non-infringement, as well as any warranty arising from course of performance, course of dealing, or usage in trade.

This document contains information on products, services and/or processes in development. All information provided here is subject to change without notice. Contact your Intel representative to obtain the latest forecast, schedule, specifications and roadmaps.

The products and services described may contain defects or errors which may cause deviations from published specifications.

MPEG-1, MPEG-2, MPEG-4, H.261, H.263, H.264, MP3, DV, VC-1, MJPEG, AC3, AAC, G.711, G.722, G.722.1, G.722.2, AMRWB, Extended AMRWB (AMRWB+), G.167, G.168, G.169, G.723.1, G.726, G.728, G.729, G.729.1, GSM AMR, GSM FR are international standards promoted by ISO, IEC, ITU, ETSI, 3GPP and other organizations. Implementations of these standards, or the standard enabled platforms may require licenses from various entities, including Intel Corporation.

Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products.

Intel, the Intel logo, BlueMoon, BunnyPeople, Celeron, Celeron Inside, Centrino, Centrino Inside, Cilk, Core Inside, E-GOLD, Flexpipe, i960, Intel, the Intel logo, Intel AppUp, Intel Atom, Intel Atom Inside, Intel Core, Intel Inside, Intel Insider, the Intel Inside logo, Intel NetBurst, Intel NetMerge, Intel NetStructure, Intel SingleDriver, Intel SpeedStep, Intel Sponsors of Tomorrow., the Intel Sponsors of Tomorrow. logo, Intel StrataFlash, Intel vPro, Intel XScale, Intel True Scale Fabric, InTru, the InTru logo, the InTru Inside logo, InTru soundmark, Itanium, Itanium Inside, MCS, MMX, MPSS, Moblin, Pentium, Pentium Inside, Puma, skool, the skool logo, SMARTi, Sound Mark, Stay With It, The Creators Project, The Journey Inside, Thunderbolt, Ultrabook, vPro Inside, VTune, Xeon, Xeon Phi, Xeon Inside, X-GOLD, XMM, X-PMU and XPOSYS are trademarks of Intel Corporation in the U.S. and/or other countries.

*Other names and brands may be claimed as the property of others.

Microsoft, Windows, and the Windows logo are trademarks, or registered trademarks of Microsoft Corporation in the United States and/or other countries.

Java is a registered trademark of Oracle and/or its affiliates.

Bluetooth is a trademark owned by its proprietor and used by Intel Corporation under license.

Intel Corporation uses the Palm OS* Ready mark under license from Palm, Inc.

OpenCL and the OpenCL logo are trademarks of Apple Inc. used by permission by Khronos.

© 2015 Intel Corporation. Portions (PBS Library) are copyrighted by Altair Engineering, Inc. and used with permission. All rights reserved.

Optimization Notice

Intel's compilers may or may not optimize to the same degree for non-Intel microprocessors for optimizations that are not unique to Intel microprocessors. These optimizations include SSE2, SSE3, and SSSE3 instruction sets and other optimizations. Intel does not guarantee the availability, functionality, or effectiveness of any optimization on microprocessors not manufactured by Intel. Microprocessor-dependent optimizations in this product are intended for use with Intel microprocessors. Certain optimizations not specific to Intel microarchitecture are reserved for Intel microprocessors. Please refer to the applicable product User and Reference Guides for more information regarding the specific instruction sets covered by this notice.

Notice revision #20110804

1. Introduction

This *Reference Manual* provides you with command and tuning reference for the Intel® MPI Library. The *Reference Manual* contains the following sections.

Document Organization

Section	Description
Section 1. Introduction	Introduces this document
Section 2. Command Reference	Describes options and environment variables for compiler commands, job startup commands, and MPD daemon commands as well
Section 3. Tuning Reference	Describes environment variables used to influence program behavior and performance at run time
Section 4. Glossary	Explains basic terms used in this document
Section 5. Index	References options and environment variables names

1.1. Introducing Intel® MPI Library

The Intel® MPI Library is a multi-fabric message passing library that implements the Message Passing Interface, v3.0 (MPI-3.0) specification. It provides a standard library across Intel® platforms that enable adoption of MPI-3.0 functions as their needs dictate.

The Intel® MPI Library enables developers to change or to upgrade processors and interconnects as new technology becomes available without changes to the software or to the operating environment.

The library is provided in the following kits:

- *The Intel® MPI Library Runtime Environment* (RTO) has the tools you need to run programs, including scalable process management system (Hydra*), Multipurpose Daemon* (MPD), and supporting utilities, shared (.so) libraries, and documentation.
- *The Intel® MPI Library Software Development Kit* (SDK) includes all of the Runtime Environment components plus compilation tools, including compiler drivers such as `mpicc`, include files and modules, static (.a) libraries, debug libraries, and test codes.

1.2. Intended Audience

This *Reference Manual* helps an experienced user understand the full functionality of the Intel® MPI Library.

1.3. What's New

This document reflects the updates for Intel® MPI Library 5.1 Update 3 for Linux* OS:

The following latest changes in this document were made:

- Added description for the new environment variable `I_MPI_ADJUST_BCAST_SEGMENT` in [I_MPI_ADJUST Family](#).

- Added description for the new non-blocking collective algorithms in [I_MPI_ADJUST Family](#).
- Extension of the target operations list for `I_MPI_STATS_SCOPE` with non-blocking operations. See [Native Statistics Format](#).
- Added description for the new options `-psm2` and `-PSM2` in [Extended Device Control Options](#).

1.4. Notational Conventions

The following conventions are used in this document.

<i>This type style</i>	Document or product names
This type style	Hyperlinks
<code>This type style</code>	Commands, arguments, options, file names
<code>THIS_TYPE_STYLE</code>	Environment variables
<code><this type style></code>	Placeholders for actual values
<code>[items]</code>	Optional items
<code>{ item item }</code>	Selectable items separated by vertical bar(s)
(SDK only)	For Software Development Kit (SDK) users only

1.5. Related Information

The following related documents that might be useful to the user:

[Product Web Site](#)

[Intel® MPI Library Support](#)

[Intel® Cluster Tools Products](#)

[Intel® Software Development Products](#)

2. Command Reference

This section provides information on different command types and how to use these commands:

- [Compiler commands](#)
- [Simplified job startup command](#)
- [Scalable process management system \(Hydra\) commands](#)
- [Hetero Operating System Cluster Support](#)
- [Intel® Xeon Phi™ Coprocessor Support](#)
- [Multipurpose daemon commands](#)
- [Processor information utility](#)

2.1. Compiler Commands

(SDK only)

The following table lists available MPI compiler commands and the underlying compilers, compiler families, languages, and application binary interfaces (ABIs) that they support.

Table 2.1-1 The Intel® MPI Library Compiler Drivers

Compiler Command	Default Compiler	Supported Language(s)	Supported ABI(s)
Generic Compilers			
mpicc	gcc, cc	C	64 bit
mpicxx	g++	C/C++	64 bit
mpifc	gfortran	Fortran77*/Fortran 95*	64 bit
GNU* Compilers Versions 4.3 and Higher			
mpigcc	gcc	C	64 bit
mpigxx	g++	C/C++	64 bit
mpif77	g77	Fortran 77	64 bit
mpif90	gfortran	Fortran 95	64 bit
Intel® Fortran, C++ Compilers Versions 14.0 through 16.0 and Higher			
mpiicc	icc	C	64 bit
mpiicpc	icpc	C++	64 bit
mpiifort	ifort	Fortran77/Fortran 95	64 bit

- Compiler commands are available only in the Intel® MPI Library Development Kit.

- Compiler commands are in the `<installdir>/<arch>/bin` directory. Where `<installdir>` refers to the Intel® MPI Library installation directory and `<arch>` is one of the following architectures:
 - `intel64` - Intel® 64 architecture
 - `mic` - Intel® Xeon Phi™ Coprocessor architecture
- Ensure that the corresponding underlying compiler (64-bit, as appropriate) is already in your `PATH`.
- To port existing MPI-enabled applications to the Intel® MPI Library, recompile all sources.
- To display mini-help of a compiler command, execute it without any parameters.

2.1.1. Compiler Command Options

-static_mpi

Use this option to link the Intel® MPI Library statically. This option does not affect the default linkage method for other libraries.

-static

Use this option to link the Intel® MPI Library statically. This option is passed to the compiler.

-nostrip

Use this option to turn off the debug information stripping while linking the Intel® MPI Library statically.

-config=<name>

Use this option to source the configuration file. See [Configuration Files](#) for details.

-profile=<profile_name>

Use this option to specify an MPI profiling library. The profiling library is selected using one of the following methods:

- Through the configuration file `<profile_name>.conf` located in the `<installdir>/<arch>/etc`. See [Profiles](#) for details.
- In the absence of the respective configuration file, by linking the library `lib<profile_name>.so` or `lib<profile_name>.a` located in the same directory as the Intel® MPI Library.

-t or -trace

Use the `-t` or `-trace` option to link the resulting executable file against the Intel® Trace Collector library. Using this option has the same effect as if you use `-profile=vt` as an argument to `mpiicc` or another compiler script.

To use this option, include the installation path of the Intel® Trace Collector in the `VT_ROOT` environment variable. Set the environment variable `I_MPI_TRACE_PROFILE` to the `<profile_name>` to specify another profiling library. For example, set `I_MPI_TRACE_PROFILE` to `vtfs` to link against the fail-safe version of the Intel® Trace Collector.

-check_mpi

Use this option to link the resulting executable file against the Intel® Trace Collector correctness checking library. Using this option has the same effect as using `-profile=vtmc` as an argument to the `mpiicc` or another compiler script.

To use this option, include the installation path of the Intel® Trace Collector in the `VT_ROOT` environment variable. Set `I_MPI_CHECK_PROFILE` to the `<profile_name>` environment variable to specify another checking library.

-ilp64

Use this option to enable partial ILP64 support. All integer arguments of the Intel MPI Library are treated as 64-bit values in this case.

-no_ilp64

Use this option to disable the ILP64 support explicitly. This option must be used in conjunction with `-i8` option of Intel® Fortran Compiler.

NOTE

If you specify the `-i8` option for the separate compilation with Intel® Fortran Compiler, you still have to use the `i8` or `ilp64` option for linkage. See [ILP64 Support](#) for details.

-dynamic_log

Use this option in combination with the `-t` option to link the Intel® Trace Collector library dynamically. This option does not affect the default linkage method for other libraries.

To run the resulting programs, include `$VT_ROOT/slib` in the `LD_LIBRARY_PATH` environment variable.

-g

Use this option to compile a program in debug mode and link the resulting executable file against the debugging version of the Intel® MPI Library. See [Environment variables](#), `I_MPI_DEBUG` for information on how to use additional debugging features with the `-g` builds.

NOTE

The optimized library is linked with the `-g` option by default.

NOTE

Use `mpivars.{sh|csh} [debug|debug_mt]` during runtime to load particular `libmpi.so` configuration.

-link_mpi=<arg>

Use this option to always link the specified version of the Intel® MPI Library. See the [I_MPI_LINK](#) environment variable for detailed argument descriptions. This option overrides all other options that select a specific library.

NOTE

Use `mpivars.{sh|csh} [debug|debug_mt]` during runtime to load particular `libmpi.so` configuration.

-O

Use this option to enable compiler optimization.

-fast

Use this Intel compiler option to maximize speed across the entire program. This option forces static linkage method for the Intel® MPI Library.

See [xHost](#) for information on implication of this option on non-Intel processors.

NOTE

This option is supported on `mpiicc`, `mpiicpc`, and `mpiifort` Intel compiler drivers.

-echo

Use this option to display everything that the command script does.

-show

Use this option to learn how the underlying compiler is invoked, without actually running it. Use the following command to see the required compiler flags and options:

```
$ mpiicc -show -c test.c
```

Use the following command to see the required link flags, options, and libraries:

```
$ mpiicc -show -o a.out test.o
```

This option is particularly useful for determining the command line for a complex build procedure that directly uses the underlying compilers.

-show_env

Use this option to see the environment settings in effect when the underlying compiler is invoked.

-{cc,cxx,fc,f77,f90}=<compiler>

Use this option to select the underlying compiler.

For example, use the following command to select the Intel® C++ Compiler:

```
$ mpiicc -cc=icc -c test.c
```

Make sure `icc` is in your path. Alternatively, you can specify the full path to the compiler.

-gcc-version=<nnn>

Use this option for compiler drivers `mpicxx` and `mpiicpc` when linking an application running in a particular GNU* C++ environment. The valid `<nnn>` values are:

<code><nnn></code> value	GNU* C++ version
430	4.3.x
440	4.4.x
450	4.5.x
460	4.6.x
470	4.7.x

By default, the library compatible with the detected version of the GNU* C++ compiler is used. Do not use this option if the GNU* C++ version is lower than 4.0.04.1.0.

-compchk

Use this option to enable compiler setup checks. In this case, each compiler driver performs checks to ensure that the appropriate underlying compiler is set up correctly.

-v

Use this option to print the compiler driver script version and its native compiler version.

2.1.2. Configuration Files

You can create Intel® MPI Library compiler configuration files using the following file naming convention:

```
<installdir>/<arch>/etc/mpi<compiler>-<name>.conf
```

where:

<arch> = {intel64em64t,mic} for the Intel® 64 architectures, and Intel® Xeon Phi™ Coprocessor architecture, respectively

<compiler> = {cc,cxx,f77,f90}, depending on the language being compiled

<name> = name of the underlying compiler with spaces replaced by hyphens

For example, the **<name>** value for `cc -64` is `cc--64`

To enable changes to the environment based on the compiler command, you need to source these files, or use the `-config` option before compiling or linking.

2.1.3. Profiles

You can select a profile library through the `-profile` option of the Intel® MPI Library compiler drivers.

The profile files are located in the `<installdir>/<arch>/etc` directory. The Intel® MPI Library comes with several predefined profiles for the Intel® Trace Collector:

`<installdir>/etc/vt.conf` - regular Intel® Trace Collector library

`<installdir>/etc/vtfs.conf` - fail-safe Intel® Trace Collector library

`<installdir>/etc/vtmc.conf` - correctness checking Intel® Trace Collector library

You can also create your own profile as `<profile_name>.conf`

The following environment variables can be defined there:

`PROFILE_PRELIB` - libraries (and paths) to include before the Intel® MPI Library

`PROFILE_POSTLIB` - libraries to include after the Intel® MPI Library

`PROFILE_INCPATHS` - C preprocessor arguments for any include files

For instance, create a file `/myprof.conf` with the following lines:

```
PROFILE_PRELIB="-L<path_to_myprof>/lib -lmyprof"
```

```
PROFILE_INCPATHS="-I<paths_to_myprof>/include"
```

Use the command-line argument `-profile=myprof` for the relevant compile driver to select this new profile.

2.1.4. Environment Variables

`I_MPI_{CC,CXX,FC,F77,F90}_PROFILE` (`MPI{CC,CXX,FC,F77,F90}_PROFILE`)

Specify a default profiling library.

Syntax

```
I_MPI_{CC,CXX,FC,F77,F90}_PROFILE=<profile_name>
```

Deprecated Syntax

```
MPI{CC,CXX,FC,F77,F90}_PROFILE=<profile_name>
```

Arguments

<code><profile_name></code>	Specify a default profiling library.
-----------------------------------	--------------------------------------

Description

Set this environment variable to select a specific MPI profiling library to be used by default. This has the same effect as using `-profile=<profile_name>` as an argument to the `mpiicc` or another Intel® MPI Library compiler driver.

I_MPI_TRACE_PROFILE

Specify a default profile for the `-trace` option.

Syntax

`I_MPI_TRACE_PROFILE=<profile_name>`

Arguments

<code><profile_name></code>	Specify a tracing profile name. The default value is <code>vt</code> .
-----------------------------------	--

Description

Set this environment variable to select a specific MPI profiling library to be used with the `-trace` option to `mpiicc` or another Intel® MPI Library compiler driver.

The `I_MPI_{CC,CXX,F77,F90}_PROFILE` environment variable overrides `I_MPI_TRACE_PROFILE`.

I_MPI_CHECK_PROFILE

Specify a default profile for the `-check_mpi` option.

Syntax

`I_MPI_CHECK_PROFILE=<profile_name>`

Arguments

<code><profile_name></code>	Specify a checking profile name. The default value is <code>vtmc</code> .
-----------------------------------	---

Description

Set this environment variable to select a specific MPI checking library to be used with the `-check_mpi` option to `mpiicc` or another Intel® MPI Library compiler driver.

The `I_MPI_{CC,CXX,F77,F90}_PROFILE` environment variable overrides the `I_MPI_CHECK_PROFILE`.

I_MPI_CHECK_COMPILER

Turn on/off compiler compatibility check.

Syntax

`I_MPI_CHECK_COMPILER=<arg>`

Arguments

<code><arg></code>	Binary indicator.
<code>enable yes on 1</code>	Enable checking the compiler.
<code>disable no off 0</code>	Disable checking the compiler. This is the default value.

Description

If `I_MPI_CHECK_COMPILER` is set to `enable`, the Intel MPI compiler drivers check the underlying compiler for compatibility. Normal compilation requires using a known version of the underlying compiler.

I_MPI_{CC,CXX,FC,F77,F90} **(MPICH_{CC,CXX,FC,F77,F90})**

Set the path/name of the underlying compiler to be used.

Syntax

```
I_MPI_{CC,CXX,FC,F77,F90}=<compiler>
```

Deprecated Syntax

```
MPICH_{CC,CXX,FC,F77,F90}=<compiler>
```

Arguments

<code><compiler></code>	Specify the full path/name of compiler to be used.
-------------------------------	--

Description

Set this environment variable to select a specific compiler to be used. Specify the full path to the compiler if it is not located in the search path.

NOTE

Some compilers may require additional command line options.

NOTE

The configuration file is sourced if it exists for a specified compiler. See [Configuration Files](#) for details.

I_MPI_ROOT

Set the Intel® MPI Library installation directory path.

Syntax

```
I_MPI_ROOT=<path>
```

Arguments

<code><path></code>	Specify the installation directory of the Intel® MPI Library
---------------------------	--

Description

Set this environment variable to specify the installation directory of the Intel® MPI Library.

VT_ROOT

Set Intel® Trace Collector installation directory path.

Syntax

```
VT_ROOT=<path>
```

Arguments

<code><path></code>	Specify the installation directory of the Intel® Trace Collector
---------------------------	--

Description

Set this environment variable to specify the installation directory of the Intel® Trace Collector.

I_MPI_COMPILER_CONFIG_DIR

Set the location of the compiler configuration files.

Syntax

```
I_MPI_COMPILER_CONFIG_DIR=<path>
```

Arguments

<path>	Specify the location of the compiler configuration files. The default value is <installdir>/<arch>/etc
--------	--

Description

Set this environment variable to change the default location of the compiler configuration files.

I_MPI_LINK

Select a specific version of the Intel® MPI Library for linking.

Syntax

```
I_MPI_LINK=<arg>
```

Arguments

<arg>	Version of library
opt	The optimized, single threaded version of the Intel® MPI Library
opt_mt	The optimized, multithreaded version of the Intel MPI Library
dbg	The debugging, single threaded version of the Intel MPI Library
dbg_mt	The debugging, multithreaded version of Intel MPI Library

Description

Set this variable to always link against the specified version of the Intel® MPI Library.

I_MPI_DEBUG_INFO_STRIP

Turn on/off the debug information stripping while linking applications statically.

Syntax

```
I_MPI_DEBUG_INFO_STRIP=<arg>
```

Arguments

<arg>	Binary indicator
enable yes on 1	Turn on. This is the default value
disable no off 0	Turn off

Description

Use this option to turn on/off the debug information stripping while linking the Intel® MPI Library statically. Debug information is stripped by default.

2.2. Simplified Job Startup Command

mpirun

Syntax

```
mpirun <options>
```

where `<options>:= <mpiexec.hydra options> | [<mpdboot options>] <mpiexec options>`

Arguments

<code><mpiexec.hydra options></code>	mpiexec.hydra options as described in the mpiexec.hydra section. This is the default operation mode.
<code><mpdboot options></code>	mpdboot options as described in the mpdboot command description, except <code>-n</code>
<code><mpiexec options></code>	mpiexec options as described in the mpiexec section

Description

Use this command to launch an MPI job. The `mpirun` command uses Hydra* or MPD as the underlying process managers. Hydra is the default process manager. Set the `I_MPI_PROCESS_MANAGER` environment variable to change the default value.

The `mpirun` command detects if the MPI job is submitted from within a session allocated using a job scheduler like Torque*, PBS Pro*, LSF*, Parallelnavi* NQS*, SLURM*, Univa* Grid Engine*, or LoadLeveler*. The `mpirun` command extracts the host list from the respective environment and uses these nodes automatically according to the above scheme.

In this case, you do not need to create the `mpd.hosts` file. Allocate the session using a job scheduler installed on your system, and use the `mpirun` command inside this session to run your MPI job.

Example

```
$ mpirun -n <# of processes> ./myprog
```

This command invokes the `mpiexec.hydra` command which uses the Hydra Process Manager by default.

Hydra* Specification

The `mpirun` command silently ignores the MPD specific options for compatibility reasons if you select Hydra* as the active process manager. The following table provides the list of silently ignored and unsupported MPD* options. Avoid these unsupported options if the Hydra* process manager is used.

Ignored mpdboot Options	Ignored mpiexec Options	Unsupported mpdboot Options	Unsupported mpiexec Options
<code>--locons</code>	<code>-[g]envuser</code>	<code>--user=<user> -u <user></code>	<code>-a</code>
<code>--remcons</code>	<code>-[g]envexcl</code>	<code>--mpd=<mpdcmd> -m <mpdcmd></code>	
<code>--ordered -o</code>	<code>-m</code>	<code>--shell -s</code>	
<code>--maxbranch=<maxbranch> -b <maxbranch></code>	<code>-ifhn <interface/hostname></code>	<code>-1</code>	
<code>--parallel-startup -p</code>	<code>-ecfn <filename></code>	<code>--ncpus=<ncpus></code>	

	-tvsu		
--	-------	--	--

MPD* Specification

If you select MPD* as the process manager, the `mpirun` command automatically starts an independent ring of the `mpd` daemons, launches an MPI job, and shuts down the `mpd` ring upon job termination.

The first non-`mpdboot` option (including `-n` or `-np`) delimits the `mpdboot` and the `mpiexec` options. All options up to this point, excluding the delimiting option, are passed to the `mpdboot` command. All options from this point on, including the delimiting option, are passed to the `mpiexec` command.

All configuration files and environment variables applicable to the `mpdboot` and `mpiexec` commands also apply to the `mpirun` command.

The set of hosts is defined by the following rules, which are executed in this order:

1. All host names from the `mpdboot` host file (either `mpd.hosts` or the file specified by the `-f` option).
2. All host names returned by the `mpdtrace` command, if there is an `mpd` ring running.
3. The local host (a warning is issued in this case).

I_MPI_MPIRUN_CLEANUP

Control the environment cleanup after the `mpirun` command.

Syntax

`I_MPI_MPIRUN_CLEANUP=<value>`

Arguments

<code><value></code>	Define the option
<code>enable yes on 1</code>	Enable the environment cleanup
<code>disable no off 0</code>	Disable the environment cleanup. This is the default value

Description

Use this environment variable to define whether to clean up the environment upon the `mpirun` completion. The cleanup includes the removal of the eventual stray service process, temporary files, and so on.

I_MPI_PROCESS_MANAGER

Select a process manager to be used by the `mpirun` command.

Syntax

`I_MPI_PROCESS_MANAGER=<value>`

Arguments

<code><value></code>	String value
<code>hydra</code>	Use Hydra* process manager. This is the default value
<code>mpd</code>	Use MPD* process manager

Description

Set this environment variable to select the process manager to be used by the `mpirun` command.

NOTE

You can run each process manager directly by invoking the `mpiexec` command for MPD* and the `mpiexec.hydra` command for Hydra*.

I_MPI_YARN

Set this variable when running on a YARN*-managed cluster.

Arguments

<value>	Binary indicator
enable yes on 1	Enable YARN support
disable no off 0	Disable YARN support. This is the default value.

Description

Set this environment variable to make Hydra request resources from the YARN cluster manager prior to running an MPI job. Use this functionality only when you launch MPI on a YARN-managed cluster with Llama* installed (for example on cluster with the Cloudera* Distribution for Hadoop*).

Usage Example

Verify that YARN is configured to work properly with Llama (refer to the [Llama documentation](#) for the specific configuration details), and the Apache* Thrift* installation is available on the cluster.

1. Make sure Llama is started on the same host where YARN is running, or start it by issuing the following command as the `llama` user:

```
$ llama [--verbose &]
```

2. Make sure passwordless `ssh` is configured on the cluster.
3. Set the `I_MPI_YARN` environment variable:

```
$ export I_MPI_YARN=1
```

4. Either set `I_MPI_THRIFT_PYTHON_LIB` to point to Thrift's Python* modules, or add these modules explicitly to `PYTHONPATH`.
5. Set `I_MPI_LLAMA_HOST/I_MPI_LLAMA_PORT` to point to the Llama server host/port (by default, it is `localhost:15000`, so you can skip this step, if launching MPI from the same host where the Llama service is running).
6. Launch an MPI job as usual (do not specify `hosts` or `machinefile` explicitly - the resources will be automatically allocated by YARN):

```
$ mpirun -n 16 -ppn 2 [other IMPI options] <application>
```

NOTE

The functionality is available with the Hydra process manager only.

2.3. Scalable Process Management System (Hydra) Commands

mpiexec.hydra

The `mpiexec.hydra` is a more scalable alternative to the MPD* process manager.

Syntax

```
mpiexec.hydra <g-options> <l-options> <executable>
or
mpiexec.hydra <g-options> <l-options> <executable1> : \
<l-options> <executable2>
```

Arguments

<code><g-options></code>	Global options that apply to all MPI processes
<code><l-options></code>	Local options that apply to a single arg-set
<code><executable></code>	<code>./a.out</code> or path/name of the executable file

Description

Use the `mpiexec.hydra` utility to run MPI applications without the MPD ring.

Use the first short command-line syntax to start all MPI processes of the `<executable>` with the single set of arguments. For example, the following command executes `a.out` over the specified processes and hosts:

```
$ mpiexec.hydra -f <hostsfile> -n <# of processes> ./a.out
```

where:

- `<# of processes>` specifies the number of processes on which to run the `a.out` executable
- `<hostsfile>` specifies a list of hosts on which to run the `a.out` executable

Use the second long command-line syntax to set different argument sets for different MPI program runs. For example, the following command executes two different binaries with different argument sets:

```
$ mpiexec.hydra -f <hostsfile> -env <VAR1> <VAL1> -n 2 ./a.out : \
-env <VAR2> <VAL2> -n 2 ./b.out
```

NOTE

If there is no "." in the PATH environment variable on all nodes of the cluster, specify `<executable>` as `./a.out` instead of `a.out`.

NOTE

You need to distinguish global options from local options. In a command-line syntax, place the local options after the global options.

2.3.1. Global Options

-hostfile <hostfile> or -f <hostfile>

Use this option to specify host names on which to run the application. If a host name is repeated, this name is used only once.

See also the `_MPI_HYDRA_HOST_FILE` environment variable for more details.

NOTE

Use the `-perhost`, `-ppn`, `-grr`, and `-rr` options to change the process placement on the cluster nodes.

- Use the `-perhost`, `-ppn`, and `-grr` options to place consecutive MPI processes on every host using the round robin scheduling.
- Use the `-rr` option to place consecutive MPI processes on different hosts using the round robin scheduling.

`-machinefile <machine file> or -machine <machine file>`

Use this option to control the process placement through the `<machine file>`. To define the total number of processes to start, use the `-n` option.

When you are pinning within a machine, the option `-binding=map` is available within the machine file for each line. For example:

```
$ cat ./machinefile
node0:2 binding=map=0,3
node1:2 binding=map=[2,8]
node0:1 binding=map=8
$ mpiexec.hydra -machinefile ./machinefile -n 5 -l numactl --show
[4] policy: default
[4] preferred node: current
[4] physcpubind: 8
[4] cpubind: 0
[4] nodebind: 0
[4] membind: 0 1
[0] policy: default
[0] preferred node: current
[0] physcpubind: 0
[0] cpubind: 0
[0] nodebind: 0
[0] membind: 0 1
[1] policy: default
[1] preferred node: current
[1] physcpubind: 3
[1] cpubind: 1
[1] nodebind: 1
[1] membind: 0 1
[3] policy: default
[3] preferred node: current
[3] physcpubind: 3
[3] cpubind: 1
[3] nodebind: 1
[3] membind: 0 1
[2] policy: default
[2] preferred node: current
[2] physcpubind: 1
[2] cpubind: 1
[2] nodebind: 1
[2] membind: 0 1
```

`-genv <ENVVAR> <value>`

Use this option to set the `<ENVVAR>` environment variable to the specified `<value>` for all MPI processes.

`-genvall`

Use this option to enable propagation of all environment variables to all MPI processes.

-genvnone

Use this option to suppress propagation of any environment variables to any MPI processes.

-genvlist <list of genv var names>

Use this option to pass a list of environment variables with their current values. *<list of genv var names>* is a comma separated list of environment variables to be sent to all MPI processes.

-pmi-connect <mode>

Use this option to choose the caching mode of Process Management Interface* (PMI) message. Possible values for *<mode>* are:

- `nocache` - do not cache PMI messages.
- `cache` - cache PMI messages on the local `pmi_proxy` management processes to minimize PMI requests. Cached information is propagated to the child management processes.
- `lazy-cache` - cache mode with on-request propagation of the PMI information.
- `alltoall` - information is automatically exchanged between all `pmi_proxy` before any get request can be done.

The `lazy-cache` mode is the default mode.

See the [_MPI_HYDRA_PMI_CONNECT](#) environment variable for more details.

-perhost <# of processes >, -ppn <# of processes >, or -grr <# of processes>

Use this option to place the indicated number of consecutive MPI processes on every host in the group using round robin scheduling. See the [_MPI_PERHOST](#) environment variable for more details.

-rr

Use this option to place consecutive MPI processes on different hosts using the round robin scheduling. This option is equivalent to `-perhost 1`. See the [_MPI_PERHOST](#) environment variable for more details.

(SDK only) -trace [<profiling_library>] or -t [<profiling_library>]

Use this option to profile your MPI application using the indicated *<profiling_library>*. If you do not specify *<profiling_library>*, the default profiling library `libVT.so` is used.

Set the [_MPI_JOB_TRACE_LIBS](#) environment variable to override the default profiling library.

(SDK only) -mps

Use this option to collect statistics from your MPI application using internal Intel MPI statistics and additional collector, which can collect hardware counters, memory consumption, internal MPI imbalance and OpenMP* imbalance if the application runs with Intel OpenMP implementation. When you use this option, two text files are generated: `stats.txt` and `app_stat.txt`. The `stats.txt` file contains the Intel® MPI Library native statistics and the `app_stat.txt` file contains the statistics information for your application provided by MPI Performance Snapshot. These files can be analyzed by `mps` utility. By using the `mps` utility, you can simplify the analysis of the Intel MPI statistics.

For example, to collect the statistics, use the following command:

```
$ mpirun -mps -n 2 ./myApp
```

For more information, see the [Native Statistics Format](#).

NOTE

1. Before running an application with this option, set up the environment by using `mpsvars.sh` from the Intel® Trace Analyzer and Collector installation directory. For detailed description, see *MPI*

Performance Snapshot for Linux* OS User's Guide at

`<installdir>/itac_latest/doc/MPI_Perf_Snapshot_User_Guide.pdf` in Intel® Parallel Studio XE Cluster Edition.

2. If you have used the options `-trace` or `-check_mpi` in the command, the `-mps` option is ignored.

(SDK only) `-check_mpi [<checking_library>]`

Use this option to check your MPI application using the indicated `<checking_library>`. If you do not specify `<checking_library>`, the default checking library `libVTmc.so` is used.

Set the `_MPI_JOB_CHECK_LIBS` environment variable to override the default checking library.

(SDK only) `-trace-pt2pt`

Use this option to collect the information about point-to-point operations.

(SDK only) `-trace-collectives`

Use this option to collect the information about collective operations.

NOTE

Use the `-trace-pt2pt` and `-trace-collectives` to reduce the size of the resulting trace file or the number of message checker reports. These options work with both statically and dynamically linked applications.

`-configfile <filename>`

Use this option to specify the file `<filename>` that contains the command-line options. Blank lines and lines that start with `'#'` as the first character are ignored.

`-branch-count <num>`

Use this option to restrict the number of child management processes launched by the `mpiexec.hydra` command, or by each `pmi_proxy` management process.

See the `_MPI_HYDRA_BRANCH_COUNT` environment variable for more details.

`-pmi-aggregate` or `-pmi-noaggregate`

Use this option to switch on or off, respectively, the aggregation of the PMI requests. The default value is `-pmi-aggregate`, which means the aggregation is enabled by default.

See the `_MPI_HYDRA_PMI_AGGREGATE` environment variable for more details.

`-tv`

Use this option to run `<executable>` under the TotalView* debugger. For example:

```
$ mpiexec.hydra -tv -n <# of processes> <executable>
```

See [Environment Variables](#) for information on how to select the TotalView* executable file.

NOTE

TotalView* uses `rsh` by default. If you want to use `ssh`, set the value of the `TVDSVRLAUNCHCMD` environment variable to `ssh`.

NOTE

The TotalView* debugger can display message queue state of your MPI program. To enable this feature, do the following steps:

1. Run your `<executable>` with the `-tv` option.

```
$ mpiexec.hydra -tv -n <# of processes> <executable>
```

2. Answer Yes to the question about stopping the `mpiexec.hydra` job.

To display the internal state of the MPI library textually, select the **Tools > Message Queue** command. If you select the **Process Window Tools > Message Queue Graph** command, the TotalView* environment variable displays a window that shows a graph of the current message queue state. For more information, see the [TotalView*](#) environment variable.

-tva <pid>

Use this option to attach the TotalView* debugger to an existing Intel® MPI Library job. Use the `mpiexec.hydra` process id as `<pid>`. You can use the following command:

```
$ mpiexec.hydra -tva <pid>
```

-gdb

Use this option to run `<executable>` under the GNU* debugger. You can use the following command:

```
$ mpiexe.hydra -gdb -n <# of processes> <executable>
```

-gdba <pid>

Use this option to attach the GNU* debugger to the existing Intel® MPI Library job. You can use the following command:

```
$ mpiexec.hydra -gdba <pid>
```

-gtool

Use this option to launch tools such as Intel® VTune™ Amplifier XE, Valgrind*, and GNU* Debugger on specified ranks through the `mpiexec.hydra` command.

NOTE

You can not use the `-gdb` option and `-gtool` option simultaneously except that you have not specified any of debugger to be launched with the `-gtool` option.

Syntax

```
-gtool "<command line for a tool 1>:<ranks set 1>[=lanuch mode 1][@arch 1];  
<command line for a tool 2>:<ranks set 2>[=exclusive][@arch 2]; ... ;<command line  
for a tool n>:<ranks set n>[=exclusive][@arch n]" <executable>
```

Or

```
$ mpiexec.hydra -n <# of processes> -gtool "<command line for a tool 1>:<ranks  
set 1>[=launch mode 1][@arch 1]" -gtool "<command line for a tool 2>:<ranks set  
2>[=launch mode 2][@arch 2]" ... -gtool "<command line for a tool n>:<ranks set  
n>[=launch mode 3][@arch n]" <executable>
```

NOTE

In the syntax, the separator ";" and the `-gtool` option can be replaced with each other.

Arguments

<code><arg></code>	Parameters
--------------------------	------------

<code><rank set></code>	<p>Specify the ranks range which are involved in the tool execution. Ranks are separated with a comma or “-” symbol can be used for a set of contiguous ranks.</p> <hr/> <p>NOTE</p> <p>If you specify incorrect rank index, the corresponding warning is printed and a tool continues working for valid ranks.</p> <hr/>
<code>[=launch mode]</code>	Specify the launch mode.
<code>[@arch]</code>	<p>Specify the architecture on which a tool applies. For a given <code><rank set></code>, if you specify this argument, the tool is only applied for those ranks which have been allocated on hosts with the specific architecture. This parameter is optional. For the values of <code>[@arch]</code>, see the argument table of I_MPI_PLATFORM for the detail value descriptions.</p> <p>If you launch the debugger on Intel® Xeon Phi™ coprocessor, setting <code>[@arch]</code> is required. See examples for details.</p>

NOTE

Rank sets cannot intersect; for example, the `-gtool` option with missed parameter or the same `[@arch]` parameter. However, you can specify the same rank sets with clearly specified different `[@arch]` parameters so that these two sets are not intersect. You should specify which tools to apply within a single `mpiexec.hydra` launch. It may happen that some tools cannot work at the same time or such usages may lead to incorrect results.

The following table shows the parameter values for `[=launch mode]`

Arguments

Value	Description
<code>exclusive</code>	Specify this value to prevent launching a tool for more than one rank per host. The value is optional.
<code>attach</code>	Specify this value to attach the tool from a <code>-gtool</code> option to an executable file. If you use debuggers or other tools which can attach to a process in a debugger manner, you need to specify this value in <code>[=launch mode]</code> . The current implementation has been tested for debuggers only.
<code>node-wide</code>	Specify this value to apply the tool from a <code>-gtool</code> option to a higher level than an executable file (to <code>pmi_proxy</code> daemon). This parameter will be implemented in future release.

You can specify several values for each tool. In this case, separate the tools with a comma sign “,”. For example:

```
$ mpiexec.hydra -gtool "gdb-ia:all=attach,exclusive; /usr/bin/gdb:all=exclusive, attach@knc" -host <hostname> -n 2 <app> : -host <hostname-mic0> -n 2 <mic-app>
```

In this example, the Intel version of GNU* GDB (`gdb-ia`) and the native GNU* GDB for Intel® Xeon Phi™ coprocessor are launched only for one rank on a host as well as for one rank on a coprocessor.

Examples

The following command examples demonstrate different scenarios of using the `-gtool` option:

1. Launch Intel® VTune™ Amplifier XE and Valgrind* through the `mpiexec.hydra` command:


```
$ mpiexec.hydra -n 16 -gtool "amplxe-cl -collect advanced-hotspots -
analyze-system -r result1:5,3,7-9=exclusive@nhm;valgrind -log-file=log_%p
:0,1,10-12@wsm" a.out
```

Use this command to apply `amplxe-cl` to a rank with a minimal index allocated on the hosts in Intel® microarchitecture code name Nehalem from the given rank set. At the same time Valgrind* is applied for all ranks allocated on the hosts in Intel® microarchitecture code name Westmere from the specified rank set. Valgrind results are written to files with names `log_<process ID>`.

2. Launch GNU* Debugger (GDB*) through the `mpiexec.hydra` command:

```
$ mpiexec.hydra -n 16 -gtool "gdb:3,5,7-9=attach" a.out
```

Use this command to apply `gdb` to the given rank set.

3. Set up an environment variable for a specific rank set:

```
$ mpiexec.hydra -n 16 -gtool "env VARIABLE1=value1 VARIABLE2=value2:3,5,7-
9; env VARIABLE3=value3:0,11" a.out
```

Use this command to set `VARIABLE1` and `VARIABLE2` for 3,5,7,8,9 ranks and establish `VARIABLE3` for ranks with numbers 0 and 11.

4. Debug an application on selected ranks through the `mpiexec.hydra` command.

When configure the debugger through the `mpiexec.hydra` command, you can use the following options for corresponding debuggers. You can also launch any other debuggers by using the `-gtool` option.

- o standard GNU* Debugger (GDB*): `gdb`
- o Intel version of GNU* Debugger: `gdb-ia`
- o native GNU* Debugger for Intel® Xeon Phi™ coprocessor: `gdb`

The `-gtool` option can simultaneously support debuggers for the host with Intel® Xeon Phi™ coprocessor. In this case, you must specify `@arch` parameter for ranks allocated on Intel® Xeon Phi™ machines. Set `@arch=knc` in the command for architectures with Intel® Xeon Phi™ coprocessor.

NOTE

When debugging on hosts with Intel® Xeon Phi™ Coprocessor or on heterogeneous clusters, you need to enable the Intel® Xeon Phi™ coprocessor recognition by setting the `I_MPI_MIC` environment variable. See [I_MPI_MIC](#) on how to specify this environment variable.

If you have specified `@arch`, you can specify rank sets with intersections.

If the debugging session has been launched on a binary compatible cluster, the parameter `@arch` can be omitted. The value `generic` comprises all platforms except for the platform with Intel® Xeon Phi™ Coprocessor.

If the debugging session has been launched on a heterogeneous cluster and the debugger for Intel® Xeon Phi™ coprocessor has not been specified, `/usr/bin/gdb` is launched on the coprocessor by default.

For example:

- a. `$ mpiexec.hydra -n 16 -gtool "gdb:3,5=attach;gdb-ia:7-9=attach" a.out`

In this case the standard GNU* Debugger (gdb) is launched for ranks with numbers 3, 5, 7, 8, 9.

- b. `$ mpiexec.hydra -gtool "gdb-ia:all=attach@generic;
/tmp/gdb:all=attach@knc" -host <hostname> -n 8 <host-app> : -host
<hostname-mic0> -n 8 <mic-app>`

In this case the Intel version of GNU* GDB (gdb-ia) is launched for all hostname ranks. While native GNU* GDB for Intel® Xeon Phi™ coprocessor is applied to all ranks allocated on a coprocessor. This example demonstrates the fact that rank sets can intersect if you have specified different architectures for those rank sets.

- c. `$ mpiexec.hydra -gtool "gdb-ia:3,5,7-9" -host <hostname> -n 8 <host-app> : -host <hostname-mic0> -n 8 <mic-app>`

In this case, the Intel version of GNU* GDB (gdb-ia) is launched for all ranks from the rank set which are allocated on a hostname machine. Any rank from the rank set that is assigned to the coprocessor uses the native GNU* GDB for Intel® Xeon Phi™ coprocessor (/usr/bin/gdb) automatically.

4. Apply a tool for a certain rank through the <machine file> option.

In this example, suppose the `m_file` has the following contents:

```
hostname_1:2
hostname_2:3
hostname_3:1
```

The following command line demonstrates how to use the `-machinefile` option to apply a tool:

```
$ mpiexec.hydra -n 6 -machinefile m_file -gtool "amplxe-cl -collect  
advanced-hotspots -analyze-system -r  
result1:5,3=exclusive@nhm;valgrind:0,1@wsm" a.out
```

In this example, the use of `-machinefile` option means that ranks with indices 0 and 1 is allocated on `hostname_1` machine; rank with index 3 is allocated on `hostname_2` machine and rank number 5 - on `hostname_3` machine. After that, `amplxe-cl` is applied only ranks 3 and 5 (since these ranks belong to different machines, `exclusive` option passes both of them) in case if `hostname_2` and `hostname_3` machines have Intel® microarchitecture code name Nehalem architecture. At the same time the Valgrind* tool is applied to both ranks allocated on `hostname_1` machine in case if it has Intel® microarchitecture code name Westmere architecture.

5. Show how the ranks are distributed across the cluster nodes. Use the `z show map` command with `-gtool` option. To see the full function of `z` commands, use the `z help` command.

```
$ mpiexec.hydra -gtool "gdb-ia:0,2,3,9,10=attach;/tmp/gdb:5,6=attach@knc"  
-host <hostname_1> -n 4 <host-app> : -host <hostname_1-mic0> -n 4 <mic-app>  
: -host <hostname_2> -n 4 <host-app>  
[0,2,3,5,6,9,10] (mpigdb) z show map  
[0,2,3]: hostname_1  
[5,6]: hostname_1-mic0  
[9,10]: hostname_2
```

```
[0,2,3,5,6,9,10] (mpigdb) z help
z <positive number(s) up to 11 or all> - Sets ranks for debugging
z show map - Shows ranks distribution across the cluster nodes
z help - Shows help information
[0,2,3,5,6,9,10] (mpigdb)
```

-gtoolfile <gtool_config_file>

Use this option to launch tools such as Intel® VTune™ Amplifier XE, Valgrind*, and GNU* Debugger on specified ranks through the `mpiexec.hydra` command.

Example

If `gtool_config_file` contains the following settings:

```
env VARIABLE1=value1 VARIABLE2=value2:3,5,7-9; env VARIABLE3=value3:0,11
env VARIABLE4=value4:1,12
```

then the following command sets `VARIABLE1` and `VARIABLE2` for 3,5,7,8,9 ranks and establishes `VARIABLE3` for ranks with numbers 0 and 11, while `VARIABLE4` is stated for the first and the twelfth ranks.

```
$ mpiexec.hydra -n 16 -gtoolfile gtool_config_file a.out
```

NOTE

The options and the environment variable `-gtool`, `-gtoolfile` and `I_MPI_GTOOL` are mutually exclusive. The options `-gtool` and `-gtoolfile` are of the same priority. The first specified option in a command line is effective and the second specified option is ignored. Both `-gtool` and `-gtoolfile` options have higher priority than the `I_MPI_GTOOL` environment variable. Thus, use the `I_MPI_GTOOL` environment variable if you have not specified neither `-gtool` or `-gtoolfile` options in the `mpiexec.hydra` command.

-nolocal

Use this option to avoid running the `<executable>` on the host where the `mpiexec.hydra` is launched. You can use this option on clusters that deploy a dedicated master node for starting the MPI jobs and a set of dedicated compute nodes for running the actual MPI processes.

-hosts <nodelist>

Use this option to specify a particular `<nodelist>` on which to run the MPI processes. For example, the following command runs the executable `a.out` on hosts `host1` and `host2`:

```
$ mpiexec.hydra -n 2 -ppn 1 -hosts host1,host2 ./a.out
```

NOTE

If `<nodelist>` consists of only one node, this option is interpreted as a local option. See [Local Options](#) for details.

-iface <interface>

Use this option to choose the appropriate network interface. For example, if the IP emulation of your InfiniBand* network is configured to `ib0`, you can use the following command.

```
$ mpiexec.hydra -n 2 -iface ib0 ./a.out
```

See the `I_MPI_HYDRA_IFACE` environment variable for more details.

-demux <mode>

Use this option to set the polling mode for multiple I/O. The default is `poll`.

Arguments

<code><spec></code>	Define the polling mode for multiple I/O
<code>poll</code>	Set <code>poll</code> as the polling mode. This is the default value.
<code>select</code>	Set <code>select</code> as the polling mode.

See the [_MPI_HYDRA_DEMUX](#) environment variable for more details.

-enable-x or -disable-x

Use this option to control the Xlib* traffic forwarding. The default value is `-disable-x`, which means the Xlib traffic is not forwarded.

-l

Use this option to insert the MPI process rank at the beginning of all lines written to the standard output.

-tune [*<arg>*]

where:

`<arg>= {<dir_name>, <configuration_file>}`.

Use this option to optimize the Intel® MPI Library performance by using the data collected by the `mpitune` utility.

NOTE

Use the `mpitune` utility to collect the performance tuning data before using this option.

If `<arg>` is not specified, the best-fit tune options are selected for the given configuration. The default location of the configuration file is `<installdir>/<arch>/etc` directory.

To specify a different location for the configuration file, set `<arg>=<dir_name>`.

To specify a different configuration file, set `<arg>=<configuration_file>`.

-ilp64

Use this option to preload the ILP64 interface. See [ILP64 Support](#) for more details.

-s <spec>

Use this option to direct standard input to the specified MPI processes.

Arguments

<code><spec></code>	Define MPI process ranks
<code>all</code>	Use all processes
<code><l>, <m>, <n></code>	Specify an exact list and use processes <code><l></code> , <code><m></code> and <code><n></code> only. The default value is zero
<code><k>, <l>-<m>, <n></code>	Specify a range and use processes <code><k></code> , <code><l></code> through <code><m></code> , and <code><n></code>

-noconf

Use this option to disable processing of the `mpiexec.hydra` configuration files described in [Configuration Files](#).

-ordered-output

Use this option to avoid intermingling of data output from the MPI processes. This option affects both the standard output and the standard error streams.

NOTE

When using this option, end the last output line of each process with the end-of-line (\n) character. Otherwise the application may stop responding.

-path <directory>

Use this option to specify the path to the *<executable>* file.

-cleanup

Use this option to create a temporary file containing information about the launched processes. The file name is `mpiexec_${username}_${PPID}.log`, where PPID is a parent process PID. This file is created in the temporary directory selected by the *-tmpdir* option. This file is used by the `mpicleanup` utility. If a job terminates successfully, the `mpiexec.hydra` command automatically removes this file.

See the [I_MPI_HYDRA_CLEANUP](#) environment variable for more details.

-tmpdir

Use this option to set a directory for temporary files.

See the [I_MPI_TMPDIR](#) environment variable for more details.

-version or -V

Use this option to display the version of the Intel® MPI Library.

-info

Use this option to display build information of the Intel® MPI Library. When this option is used, the other command line arguments are ignored.

-use-app-topology <value>

Use this option to performs rank placement on the assumption of transferred data from the stats file and cluster topology. You can use the following command:

```
$ mpiexec.hydra -use-app-topology ./stats.txt <...> ./my_app
```

Arguments

<code><value></code>	The path to the native Intel MPI statistics file level 1 or higher
----------------------------	--

NOTE

The hydra PM uses API of `libmpitune.so` the same way as `mpitune_rank_placement` in the static method and uses the resulting host list for rank assignment.

See [I_MPI_HYDRA_USE_APP_TOPOLOGY](#) and [Topology Awareness Application Tuning](#) for more details.

-localhost

Use this option to explicitly specify the local host name for the launching node.

Example:

```
$ mpiexec.hydra -localhost <localhost_ip> -machinefile <file> -n 2 ./a.out
```

Bootstrap Options

-bootstrap <bootstrap server>

Use this option to select a built-in bootstrap server to use. A bootstrap server is the basic remote node access mechanism that is provided by the system. Hydra supports multiple runtime bootstrap servers such as `ssh`, `rsh`, `pdsh`, `fork`, `persist`, `slurm`, `ll`, `lsf`, `sge`, or `jmi` to launch the MPI processes. The default bootstrap server is `ssh`. By selecting `slurm`, `ll`, `lsf`, or `sge`, you use the corresponding `srun`, `llspawn.stdio`, `blaunch`, or `qsh` internal job scheduler utility to launch service processes under the respective selected job scheduler (SLURM*, LoadLeveler*, LSF*, and SGE*).

Arguments

<arg>	String parameter
<code>ssh</code>	Use secure shell. This is the default value
<code>rsh</code>	Use remote shell
<code>pdsh</code>	Use parallel distributed shell
<code>pbsdsh</code>	Use Torque* and PBS* <code>pbsdsh</code> command
<code>fork</code>	Use fork call
<code>persist</code>	Use Hydra persist server
<code>slurm</code>	Use SLURM* <code>srun</code> command
<code>ll</code>	Use LoadLeveler* <code>llspawn.stdio</code> command
<code>lsf</code>	Use LSF <code>blaunch</code> command
<code>sge</code>	Use Univa* Grid Engine* <code>qsh</code> command
<code>jmi</code>	Use Job Manager Interface (tighter integration)

To enable tighter integration with the SLURM* or PBS Pro* job manager, use the `jmi` bootstrap server. Tighter integration includes registration of the process identifiers by the respective job managers. This configuration enables better resource accounting by the respective job manager, and better node cleanup upon job termination.

NOTE

Some bootstrap servers that use parallel startup of remote processes (`slurm` and `pdsh`) might not work in heterogeneous environment, for example, when `I_MPI_MIC` is enabled.

See the [-bootstrap jmi](#) description and the `I_MPI_HYDRA_BOOTSTRAP` environment variable for details.

-bootstrap-exec <bootstrap server>

Use this option to set the executable to be used as a bootstrap server. The default bootstrap server is `ssh`. For example:

```
$ mpiexec.hydra -bootstrap-exec <bootstrap_server_executable> \
-f hosts.file -env <VAR1> <VAL1> -n 2 ./a.out
```

See the `I_MPI_HYDRA_BOOTSTRAP` environment variable for more details.

-bootstrap-exec-args <args>

Use this option to provide the additional parameters to the bootstrap server executable file.

```
$ mpiexec.hydra -bootstrap-exec-args <arguments> -n 2 ./a.out
```

See the [I_MPI_HYDRA_BOOTSTRAP_EXEC_EXTRA_ARGS](#) environment variable for more details.

-bootstrap persist

Use this option to launch MPI processes using Hydra persist server. Before running a job, start these servers on each host:

```
$ hydra_persist&
```

NOTE

Do not start the services under the root account. A server can be shutdown using the Linux shell kill command.

-bootstrap jmi

Use this option to enable tight integration with the SLURM* or PBS Pro* job schedulers. Tighter integration is implemented using a particular job scheduler application programming interface or utility. If you specify this option, the default `libjmi.so` library is loaded. You can overwrite the default library name through the `I_MPI_HYDRA_JMI_LIBRARY` environment variable.

See the [I_MPI_HYDRA_JMI_LIBRARY](#) environment variable for more details.

Binding Options**-binding**

Use this option to pin or bind MPI processes to a particular processor and avoid undesired process migration. In the following syntax, the quotes may be omitted for a one-member list. Each parameter corresponds to a single pinning property.

This option is supported on both Intel® and non-Intel microprocessors, but it may perform additional optimizations for Intel microprocessors than it performs for non-Intel microprocessors.

Syntax

```
-binding "<parameter>=<value>[;<parameter>=<value> ...]"
```

Parameters

pin	Pinning switch
enable yes on 1	Turn on the pinning property. This is the default value
disable no off 0	Turn off the pinning property

cell	Pinning resolution
unit	Basic processor unit (logical CPU)
core	Processor core in multi-core system

map	Process mapping
spread	The processes are mapped consecutively to separate processor cells. Thus, the processes do not share the common resources of the adjacent cells.
scatter	The processes are mapped to separate processor cells. Adjacent processes are mapped upon the cells that are the most remote in the multi-core topology.
bunch	The processes are mapped to separate processor cells by #processes/#sockets processes per socket. Each socket processor portion is a set of the cells that are the closest in the multi-core topology.
p_0, p_1, \dots, p_n	<p>The processes are mapped upon the separate processors according to the processor specification on the p_0, p_1, \dots, p_n list: the i^{th} process is mapped upon the processor p_i, where</p> <p>p_i takes one of the following values:</p> <ul style="list-style-type: none"> • processor number like n • range of processor numbers like $n-m$ • -1 for no pinning of the corresponding process
$[m_0, m_1, \dots, m_n]$	<p>The i^{th} process is mapped upon the processor subset defined by m_i hexadecimal mask using the following rule:</p> <p>The j^{th} processor is included into the subset m_i if the j^{th} bit of m_i equals 1.</p>

domain	Processor domain set on a node
cell	Each domain of the set is a single processor cell (unit or core).
core	Each domain of the set consists of the processor cells that share a particular core.
cache1	Each domain of the set consists of the processor cells that share a particular level 1 cache.
cache2	Each domain of the set consists of the processor cells that share a particular level 2 cache.
cache3	Each domain of the set consists of the processor cells that share a particular level 3 cache.
cache	The set elements of which are the largest domains among cache1, cache2, and cache3
socket	Each domain of the set consists of the processor cells that are located on a particular socket.
node	All processor cells on a node are arranged into a single domain.
$\langle size \rangle[:\langle layout \rangle]$	Each domain of the set consists of $\langle size \rangle$ processor cells. $\langle size \rangle$ may have the

	<p>following values:</p> <ul style="list-style-type: none"> • <code>auto</code> - domain size = #cells/#processes • <code>omp</code> - domain size = <code>OMP_NUM_THREADS</code> environment variable value • positive integer - exact value of the domain size
	<p>NOTE</p> <p>Domain size is limited by the number of processor cores on the node.</p>
	<p>Each member location inside the domain is defined by the optional <code><layout></code> parameter value:</p> <ul style="list-style-type: none"> • <code>compact</code> - as close with others as possible in the multi-core topology • <code>scatter</code> - as far away from others as possible in the multi-core topology • <code>range</code> - by BIOS numbering of the processors <p>If <code><layout></code> parameter is omitted, <code>compact</code> is assumed as the value of <code><layout></code></p>

<code>order</code>	Linear ordering of the domains
<code>compact</code>	Order the domain set so that adjacent domains are the closest in the multi-core topology
<code>scatter</code>	Order the domain set so that adjacent domains are the most remote in the multi-core topology
<code>range</code>	Order the domain set according to the BIOS processor numbering

<code>offset</code>	Domain list offset
<code><n></code>	Integer number of the starting domain among the linear ordered domains. This domain gets number zero. The numbers of other domains will be cyclically shifted.

Communication Subsystem Options

-rmk <RMK>

Use this option to select a resource management kernel to be used. Intel® MPI Library only supports `pbs`. See the `_MPI_HYDRA_RMK` environment variable for more details.

Other Options

-verbose or -v

Use this option to print debug information from `mpiexec.hydra`, such as:

- Service processes arguments
- Environment variables and arguments passed to start an application
- PMI requests/responses during a job life cycle

See the `_MPI_HYDRA_DEBUG` environment variable for more details.

-print-rank-map

Use this option to print out the MPI rank mapping.

-print-all-exitcodes

Use this option to print the exit codes of all processes.

2.3.2. Local Options

-n <# of processes> or -np <# of processes>

Use this option to set the number of MPI processes to run with the current argument set.

-env <ENVVAR> <value>

Use this option to set the `<ENVVAR>` environment variable to the specified `<value>` for all MPI processes in the current arg-set.

-envall

Use this option to propagate all environment variables in the current arg-set.

See the `_MPI_HYDRA_ENV` environment variable for more details.

-envnone

Use this option to suppress propagation of any environment variables to the MPI processes in the current arg-set.

-envlist <list of env var names>

Use this option to pass a list of environment variables with their current values. `<list of env var names>` is a comma separated list of environment variables to be sent to the MPI processes.

-host <nodename>

Use this option to specify a particular `<nodename>` on which the MPI processes are to be run. For example, the following command executes `a.out` on hosts `host1` and `host2`:

```
$ mpiexec.hydra -n 2 -host host1 ./a.out : -n 2 -host host2 ./a.out
```

-path <directory>

Use this option to specify the path to the `<executable>` file to be run in the current arg-set.

-wdir <directory>

Use this option to specify the working directory in which the `<executable>` file runs in the current arg-set.

-umask <umask>

Use this option to perform the `umask <umask>` command for the remote `<executable>` file.

-hostos <host OS>

Use this option to specify an operating system installed on a particular host. MPI processes are launched on each host in accordance with this option specified. The default value is `linux`.

Arguments

<code><arg></code>	String parameter
<code>linux</code>	The host with Linux* OS installed. This is the default value
<code>windows</code>	The host with Windows* OS installed.

NOTE

The option is used in conjunction with `-host` option. For instance, the following command runs the executable `a.exe` on `host1` and `b.out` on `host2`:

```
$ mpiexec.hydra -n 1 -host host1 -hostos windows a.exe : -n 1 -host host2 \ -
hostos linux ./a.out
```

2.3.3. Extended Device Control Options**-rdma**

Use this option to select an RDMA-capable network fabric. The application attempts to use the first available RDMA-capable network fabric from the list: `dapl, ofa`. If no such fabric is available, another fabric from the list `tcp, tmi, ofi` is used. This option is equivalent to the setting: `-genv I_MPI_FABRICS_LIST dapl, ofa, tcp, tmi, ofi -genv I_MPI_FALLBACK 1`.

-RDMA

Use this option to select an RDMA-capable network fabric. The application attempts to use the first available RDMA-capable network fabric from the list: `dapl, ofa`. The application fails if no such fabric is found. This option is equivalent to the setting: `-genv I_MPI_FABRICS_LIST dapl, ofa -genv I_MPI_FALLBACK 1`.

-dapl

Use this option to select a DAPL-capable network fabric. The application attempts to use a DAPL-capable network fabric. If no such fabric is available, another fabric from the list `tcp, tmi, ofa, ofi` is used. This option is equivalent to the setting: `-genv I_MPI_FABRICS_LIST dapl, tcp, tmi, ofa, ofi -genv I_MPI_FALLBACK 1`.

-DAPL

Use this option to select a DAPL-capable network fabric. The application fails if no such fabric is found. This option is equivalent to the setting: `-genv I_MPI_FABRICS_LIST dapl -genv I_MPI_FALLBACK 0`.

-ib

Use this option to select an OFA-capable network fabric. The application attempts to use an OFA-capable network fabric. If no such fabric is available, another fabric from the list `dapl, tcp, tmi, ofi` is used. This option is equivalent to the setting: `-genv I_MPI_FABRICS_LIST ofa, dapl, tcp, tmi, ofi -genv I_MPI_FALLBACK 1`.

-IB

Use this option to select an OFA-capable network fabric. The application fails if no such fabric is found. This option is equivalent to the setting: `-genv I_MPI_FABRICS_LIST ofa -genv I_MPI_FALLBACK 0`.

-tmi

Use this option to select a TMI-capable network fabric. The application attempts to use a TMI-capable network fabric. If no such fabric is available, another fabric from the list `dapl,tcp,ofa,ofi` is used. This option is equivalent to the setting: `-genv I_MPI_FABRICS_LIST tmi,dapl,tcp,ofa,ofi -genv I_MPI_FALLBACK 1`.

-TMI

Use this option to select a TMI-capable network fabric. The application fails if no such fabric is found. This option is equivalent to the setting: `-genv I_MPI_FABRICS_LIST tmi -genv I_MPI_FALLBACK 0`.

-mx

Use this option to select the Myrinet MX* network fabric. The application attempts to use the Myrinet MX* network fabric. If no such fabric is available, another fabric from the list `dapl,tcp,ofa,ofi` is used. This option is equivalent to the setting: `-genv I_MPI_FABRICS_LIST tmi,dapl,tcp,ofa,ofi -genv I_MPI_TMI_PROVIDER mx -genv I_MPI_DAPL_PROVIDER mx -genv I_MPI_FALLBACK 1`.

-MX

Use this option to select Myrinet MX* network fabric. The application fails if no such fabric is found. This option is equivalent to the setting: `-genv I_MPI_FABRICS_LIST tmi -genv I_MPI_TMI_PROVIDER mx -genv I_MPI_FALLBACK 0`.

-psm

Use this option to select a PSM-capable network fabric: Intel® True Scale Fabric or Intel® Omni-Path Fabric in PSM-compatibility mode. The application attempts to use a PSM-capable network fabric. If no such fabric is available, another fabric from the list `dapl,tcp,ofa,ofi` is used. This option is equivalent to the setting: `-genv I_MPI_FABRICS_LIST tmi,dapl,tcp,ofa,ofi -genv I_MPI_TMI_PROVIDER psm -genv I_MPI_FALLBACK 1`.

-PSM

Use this option to select a PSM-capable network fabric: Intel® True Scale Fabric or Intel® Omni-Path Fabric in PSM-compatibility mode. The application fails if no such fabric is found. This option is equivalent to the setting: `-genv I_MPI_FABRICS_LIST tmi -genv I_MPI_TMI_PROVIDER psm -genv I_MPI_FALLBACK 0`.

-psm2

Use this option to select the Intel® Omni-Path Fabric. The application attempts to use the Intel® Omni-Path Fabric. If no such fabric is available, another fabric from the list `dapl,tcp,ofa,ofi` is used. This option is equivalent to the setting: `-genv I_MPI_FABRICS_LIST tmi,dapl,tcp,ofa,ofi -genv I_MPI_TMI_PROVIDER psm2 -genv I_MPI_FALLBACK 1`.

-PSM2

Use this option to select the Intel® Omni-Path Fabric. The application fails if no such fabric is found. This option is equivalent to the setting: `-genv I_MPI_FABRICS_LIST tmi -genv I_MPI_TMI_PROVIDER psm2 -genv I_MPI_FALLBACK 0`.

-ofi

Use this option to select an OpenFabrics Interfaces* (OFI*) capable network fabric. The application attempts to use an OFI-capable network fabric. If no such fabric is available, another fabric from the list `tmi,dapl,tcp,ofa` is used. This option is equivalent to the setting: `-genv I_MPI_FABRICS_LIST ofi,tmi,dapl,tcp,ofa, -genv I_MPI_FALLBACK 1`.

-OFI

Use this option to select an OFI-capable network fabric. The application fails if no such fabric is found. This option is equivalent to the setting: `-genv I_MPI_FABRICS_LIST ofi -genv I_MPI_FALLBACK 0`.

2.3.4. Environment Variables

I_MPI_HYDRA_HOST_FILE

Set the host file to run the application.

Syntax

`I_MPI_HYDRA_HOST_FILE=<arg>`

Deprecated Syntax

`HYDRA_HOST_FILE=<arg>`

Arguments

<code><arg></code>	String parameter
<code><hostsfile></code>	The full or relative path to the host file

Description

Set this environment variable to specify the hosts file.

I_MPI_HYDRA_DEBUG

Print out the debug information.

Syntax

`I_MPI_HYDRA_DEBUG=<arg>`

Arguments

<code><arg></code>	Binary indicator
<code>enable yes on 1</code>	Turn on the debug output
<code>disable no off 0</code>	Turn off the debug output. This is the default value

Description

Set this environment variable to enable the debug mode.

I_MPI_HYDRA_ENV

Control the environment propagation.

Syntax

`I_MPI_HYDRA_ENV=<arg>`

Arguments

<code><arg></code>	String parameter
<code>all</code>	Pass all environment to all MPI processes

Description

Set this environment variable to control the environment propagation to the MPI processes. By default, the entire launching node environment is passed to the MPI processes. Setting this variable also overwrites environment variables set by the remote shell.

I_MPI_JOB_TIMEOUT, I_MPI_MPIEXEC_TIMEOUT (MPIEXEC_TIMEOUT)

Set the timeout period for `mpiexec.hydra`.

Syntax

```
I_MPI_JOB_TIMEOUT=<timeout>
```

```
I_MPI_MPIEXEC_TIMEOUT=<timeout>
```

Deprecated Syntax

```
MPIEXEC_TIMEOUT=<timeout>
```

Arguments

<code><timeout></code>	Define <code>mpiexec.hydra</code> timeout period in seconds
<code><n> >= 0</code>	The default timeout value is zero, which means no timeout.

Description

Set this environment variable to make `mpiexec.hydra` terminate the job in `<timeout>` seconds after its launch. The `<timeout>` value should be greater than zero. Otherwise the environment variable setting is ignored.

NOTE

Set the `I_MPI_JOB_TIMEOUT` environment variable in the shell environment before executing the `mpiexec.hydra` command. Do not use the `-genv` or `-env` options to set the `<timeout>` value. Those options are used for passing environment variables to the MPI process environment.

I_MPI_JOB_TIMEOUT_SIGNAL (MPIEXEC_TIMEOUT_SIGNAL)

Define the signal to be sent when a job is terminated because of a timeout.

Syntax

```
I_MPI_JOB_TIMEOUT_SIGNAL=<number>
```

Deprecated Syntax

```
MPIEXEC_TIMEOUT_SIGNAL=<number>
```

Arguments

<code><number></code>	Define signal number
<code><n> > 0</code>	The default value is 9 (SIGKILL)

Description

Define a signal number sent to stop the MPI job if the timeout period specified by the `I_MPI_JOB_TIMEOUT` environment variable expires. If you set a signal number unsupported by the system, the `mpiexec.hydra` operation prints a warning message and continues the task termination using the default signal number 9 (SIGKILL).

I_MPI_JOB_ABORT_SIGNAL

Define a signal to be sent to all processes when a job is terminated unexpectedly.

Syntax

`I_MPI_JOB_ABORT_SIGNAL=<number>`

Arguments

<code><number></code>	Define signal number
<code><n> > 0</code>	The default value is 9 (SIGKILL)

Description

Set this environment variable to define a signal for task termination. If you set an unsupported signal number, `mpiexec.hydra` prints a warning message and uses the default signal 9 (SIGKILL).

I_MPI_JOB_SIGNAL_PROPAGATION (MPIEXEC_SIGNAL_PROPAGATION)

Control signal propagation.

Syntax

`I_MPI_JOB_SIGNAL_PROPAGATION=<arg>`

Deprecated Syntax

`MPIEXEC_SIGNAL_PROPAGATION=<arg>`

Arguments

<code><arg></code>	Binary indicator
<code>enable yes on 1</code>	Turn on propagation
<code>disable no off 0</code>	Turn off propagation. This is the default value

Description

Set this environment variable to control propagation of the signals (SIGINT, SIGALRM, and SIGTERM). If you enable signal propagation, the received signal is sent to all processes of the MPI job. If you disable signal propagation, all processes of the MPI job are stopped with the default signal 9 (SIGKILL).

I_MPI_HYDRA_BOOTSTRAP

Set the bootstrap server.

Syntax

`I_MPI_HYDRA_BOOTSTRAP=<arg>`

Arguments

<code><arg></code>	String parameter
<code>ssh</code>	Use secure shell. This is the default value
<code>rsh</code>	Use remote shell

pdsh	Use parallel distributed shell
pbsdsh	Use Torque* and PBS* pbsdsh command
fork	Use fork call
slurm	Use SLURM* srun command
ll	Use LoadLeveler* llspawn.stdio command
lsf	Use LSF* blaunch command
sge	Use Univa* Grid Engine* qrsh command
jmi	Use Job Manager Interface (tighter integration)

Description

Set this environment variable to specify the bootstrap server.

NOTE

Set the `I_MPI_HYDRA_BOOTSTRAP` environment variable in the shell environment before executing the `mpiexec.hydra` command. Do not use the `-env` option to set the `<arg>` value. This option is used for passing environment variables to the MPI process environment.

I_MPI_HYDRA_BOOTSTRAP_EXEC

Set the executable file to be used as a bootstrap server.

Syntax

`I_MPI_HYDRA_BOOTSTRAP_EXEC=<arg>`

Arguments

<code><arg></code>	String parameter
<code><executable></code>	The name of the executable file

Description

Set this environment variable to specify the executable file to be used as a bootstrap server.

I_MPI_HYDRA_BOOTSTRAP_EXEC_EXTRA_ARGS

Set additional arguments for the bootstrap server.

Syntax

`I_MPI_HYDRA_BOOTSTRAP_EXEC_EXTRA_ARGS=<arg>`

Arguments

<code><arg></code>	String parameter
<code><args></code>	Additional bootstrap server arguments

Description

Set this environment variable to specify additional arguments for the bootstrap server.

I_MPI_HYDRA_BOOTSTRAP_AUTOFORK

Control the usage of `fork` call for local processes.

Syntax

`I_MPI_HYDRA_BOOTSTRAP_AUTOFORK = <arg>`

Arguments

<arg>	String parameter
<code>enable yes on 1</code>	Use <code>fork</code> for the local processes. This is default value for <code>ssh</code> , <code>rsh</code> , <code>ll</code> , <code>lsf</code> , and <code>pbsdsh</code> bootstrap servers
<code>disable no off 0</code>	Do not use <code>fork</code> for the local processes. This is default value for the <code>sge</code> bootstrap server

Description

Set this environment variable to control usage of `fork` call for the local processes.

NOTE

This option is not applicable to `slurm`, `pdsh`, `persist`, and `jmi` bootstrap servers.

I_MPI_HYDRA_RMK

Use the resource management kernel.

Syntax

`I_MPI_HYDRA_RMK=<arg>`

Arguments

<arg>	String parameter
<rmk>	Resource management kernel. The only supported value is <code>pbs</code>

Description

Set this environment variable to use the `pbs` resource management kernel. Intel® MPI Library only supports `pbs`.

I_MPI_HYDRA_PMI_CONNECT

Define the processing method for PMI messages.

Syntax

`I_MPI_HYDRA_PMI_CONNECT=<value>`

Arguments

<value>	The algorithm to be used
<code>nocache</code>	Do not cache PMI messages
<code>cache</code>	Cache PMI messages on the local <code>pmi_proxy</code> management processes to minimize the number of PMI requests. Cached information is automatically propagated to child management processes

lazy-cache	cache mode with on-demand propagation. This is the default value.
alltoall	Information is automatically exchanged between all <code>PMI_proxy</code> before any get request can be done

Description

Use this environment variable to select the `PMI` messages processing method.

I_MPI_PERHOST

Define the default settings for the `-perhost` option in the `mpiexec` and `mpiexec.hydra` command.

Syntax

`I_MPI_PERHOST=<value>`

Arguments

<code><value></code>	Define a value that is used for the <code>-perhost</code> option by default
<code>integer > 0</code>	Exact value for the option
<code>all</code>	All logical CPUs on the node
<code>allcores</code>	All cores (physical CPUs) on the node. This is the default value.

Description

Set this environment variable to define the default setting for the `-perhost` option. The `-perhost` option implied with the respective value if the `I_MPI_PERHOST` environment variable is defined.

I_MPI_JOB_TRACE_LIBS

Choose the libraries to preload through the `-trace` option.

Syntax

`I_MPI_JOB_TRACE_LIBS=<arg>`

Deprecated Syntax

`MPIEXEC_TRACE_LIBS=<arg>`

Arguments

<code><arg></code>	String parameter
<code><list></code>	Blank separated list of the libraries to preload. The default value is <code>vt</code>

Description

Set this environment variable to choose an alternative library for preloading through the `-trace` option.

I_MPI_JOB_CHECK_LIBS

Choose the libraries to preload through the `-check_mpi` option.

Syntax

`I_MPI_JOB_CHECK_LIBS=<arg>`

Arguments

<code><arg></code>	String parameter
<code><list></code>	Blank separated list of the libraries to preload. The default value is <code>vtmc</code>

Description

Set this environment variable to choose an alternative library for preloading through the `-check_mpi` option.

I_MPI_HYDRA_BRANCH_COUNT

Set the hierarchical branch count.

Syntax

`I_MPI_HYDRA_BRANCH_COUNT = <num>`

Arguments

<code><num></code>	Number
<code><n> >= 0</code>	<ul style="list-style-type: none"> The default value is <code>-1</code> if less than 128 nodes are used. This value also means that there is no hierarchical structure The default value is <code>32</code> if more than 127 nodes are used

Description

Set this environment variable to restrict the number of child management processes launched by the `mpiexec.hydra` operation or by each `pmi_proxy` management process.

I_MPI_HYDRA_PMI_AGGREGATE

Turn on/off aggregation of the PMI messages.

Syntax

`I_MPI_HYDRA_PMI_AGGREGATE=<arg>`

Arguments

<code><arg></code>	Binary indicator
<code>enable yes on 1</code>	Enable PMI message aggregation. This is the default value
<code>disable no off 0</code>	Disable PMI message aggregation

Description

Set this environment variable to enable/disable aggregation of PMI messages .

I_MPI_HYDRA_GDB_REMOTE_SHELL

Set the remote shell command to run GNU* debugger.

Syntax

`I_MPI_HYDRA_GDB_REMOTE_SHELL=<arg>`

Arguments

<code><arg></code>	String parameter
<code>ssh</code>	Secure Shell (SSH). This is the default value
<code>rsh</code>	Remote shell (RSH)

Description

Set this environment variable to specify the remote shell command to run the GNU* debugger on the remote machines. You can use this environment variable to specify any shell command that has the same syntax as SSH or RSH.

I_MPI_HYDRA_JMI_LIBRARY

Define the default setting of the JMI library.

Syntax

`I_MPI_HYDRA_JMI_LIBRARY=<value>`

Arguments

<code><value></code>	Define a string value, name, or path to JMI dynamic library
<code>libjmi_slurm.so.1.1 libjmi_pbs.so.1.0</code>	Set the library name or full path to library name. The default value is <code>libjmi.so</code>

Description

Set this environment variable to define the JMI library to be loaded by the Hydra* processor manager. Set the full path to the library if the path is not mentioned in the `LD_LIBRARY_PATH` environment variable. If you use the `mpirun` command, you do not need to set this environment variable. The JMI library is automatically detected and set.

I_MPI_HYDRA_IFACE

Set the network interface.

Syntax

`I_MPI_HYDRA_IFACE=<arg>`

Arguments

<code><arg></code>	String parameter
<code><network interface></code>	The network interface configured in your system

Description

Set this environment variable to specify the network interface to use. For example, use `-iface ib0`, if the IP emulation of your InfiniBand* network is configured on `ib0`.

I_MPI_HYDRA_DEMUX

Set the demultiplexer (demux) mode.

Syntax

`I_MPI_HYDRA_DEMUX=<arg>`

Arguments

<code><arg></code>	String parameter
<code>poll</code>	Set <code>poll</code> as the multiple I/O demultiplexer (<code>demux</code>) mode engine. This is the default value.
<code>select</code>	Set <code>select</code> as the multiple I/O demultiplexer (<code>demux</code>) mode engine

Description

Set this environment variable to specify the multiple I/O `demux` mode engine. The default value is `poll`.

I_MPI_HYDRA_CLEANUP

Control the creation of the default `mpicleanup` input file.

Syntax

`I_MPI_HYDRA_CLEANUP=<value>`

Arguments

<code><value></code>	Binary indicator
<code>enable yes on 1</code>	Enable the <code>mpicleanup</code> input file creation
<code>disable no off 0</code>	Disable the <code>mpicleanup</code> input file creation. This is the default value

Description

Set the `I_MPI_HYDRA_CLEANUP` environment variable to create the input file for the `mpicleanup` utility.

**I_MPI_TMPDIR
(TMPDIR)**

Set the temporary directory.

Syntax

`I_MPI_TMPDIR=<arg>`

Arguments

<code><arg></code>	String parameter
<code><path></code>	Set the temporary directory. The default value is <code>/tmp</code>

Description

Set this environment variable to specify the temporary directory to store the `mpicleanup` input file.

I_MPI_JOB_RESPECT_PROCESS_PLACEMENT

Specify whether to use the job scheduler provided process-per-node parameter.

Syntax

`I_MPI_JOB_RESPECT_PROCESS_PLACEMENT=<arg>`

Arguments

<code><value></code>	Binary indicator
----------------------------	------------------

<code>enable yes on 1</code>	Use the process placement provided by job scheduler. This is the default value
<code>disable no off 0</code>	Do not use the process placement provided by job scheduler

Description

If you set `I_MPI_JOB_RESPECT_PROCESS_PLACEMENT=enable`, then Hydra process manager uses PPN provided by job scheduler.

If you set `I_MPI_JOB_RESPECT_PROCESS_PLACEMENT=disable`, then Hydra process manager uses PPN provided in command line option or using `I_MPI_PERHOST` environment variable.

I_MPI_GTOOL

Specify the tools to be launched for selected ranks.

Syntax

```
I_MPI_GTOOL="<command line for a tool 1>:<ranks set 1>[=exclusive][@arch 1];
<command line for a tool 2>:<ranks set 2>[=exclusive][@arch 2]; ... ;<command line
for a tool n>:<ranks set n>[=exclusive][@arch n]"
```

Arguments

<code><arg></code>	Parameters
<code><command line for a tool></code>	Specify a tool along with its parameters
<code><rank set></code>	Specify the ranks range which is involved in the tool execution. Ranks are separated with a comma or “-” symbol can be used for a set of contiguous ranks. NOTE If you specify incorrect rank index, a tool prints the corresponding warning and a tool continues working for valid ranks.
<code>[=exclusive]</code>	Specify this parameter to prevent launching a tool for more than one rank per host. This parameter is optional.
<code>[@arch]</code>	Specify the architecture on which a tool will applies. For a given <code><rank set></code> , if you specify this parameter, the tool is only applied for those ranks which have been allocated on hosts with the specific architecture. This parameter is optional. For the values of <code>[@arch]</code> , see the argument table of I_MPI_PLATFORM for the detail value descriptions. If you launch the debugger on Intel® Xeon Phi™ coprocessor, setting <code>[@arch]</code> is required. See examples for details.

Description

Use this option to launch the tools such as Intel® VTune™ Amplifier XE, Valgrind*, and GNU* Debugger on specified ranks.

Examples

The following command examples demonstrate different scenarios of using the `I_MPI_GTOOL` environment variable:

1. Launch Intel® VTune™ Amplifier XE and Valgrind* by setting the `I_MPI_GTOOL` environment variable:

```
$ export I_MPI_GTOOL="amplxe-cl -collect advanced-hotspots -analyze-system
-r result1:5,3,7-9=exclusive@nhm;valgrind -log-file=log_%p :0,1,10-12@wsm"
$ mpiexec.hydra -n 16 a.out
```

Use this command to apply `amplxe-cl` to a rank with a minimal index allocated on the hosts in Intel® microarchitecture code name Nehalem from the given rank set. At the same time Valgrind* is applied for all ranks allocated on the hosts in Intel® microarchitecture code name Westmere from the specified rank set. Valgrind results are written to files with names `log_<process ID>`.

2. Launch GNU* Debugger (GDB*) by setting the `I_MPI_GTOOL` environment variable:

```
$ mpiexec.hydra -n 16 -genv I_MPI_GTOOL="gdb:3,5,7-9" a.out
```

Use this command to apply `gdb` to the given rank set.

NOTE

The options and the environment variable `-gtool`, `-gtoolfile` and `I_MPI_GTOOL` are mutually exclusive. The options `-gtool` and `-gtoolfile` are of the same priority. The first specified option in a command line is effective and the second specified option is ignored. Both `-gtool` and `-gtoolfile` options have higher priority than the `I_MPI_GTOOL` environment variable. Thus, use the `I_MPI_GTOOL` environment variable if you have not specified neither `-gtool` or `-gtoolfile` options in the `mpiexec.hydra` command.

I_MPI_HYDRA_USE_APP_TOPOLOGY

Syntax

`I_MPI_HYDRA_USE_APP_TOPOLOGY=<value>`

Arguments

<code><value></code>	The path to the native Intel MPI statistics file level 1 or higher
----------------------------	--

Description

If you define `I_MPI_HYDRA_USE_APP_TOPOLOGY`, hydra process manager (PM) performs rank placement on the assumption of transferred data from the stats file and cluster topology.

```
$ mpiexec.hydra -use-app-topology ./stats.txt <...> ./my_app
```

NOTE

The hydra PM uses API of `libmpitune.so` the same way as `mpitune_rank_placement` in the static method and uses the resulting host list for rank assignment.

See the description of [-use-app-topology](#) and [Topology Awareness Application Tuning](#) for more details.

2.3.5. Cleaning up Utility

mpicleanup

Clean up the environment after an abnormally terminated MPI run under the `mpiexec.hydra` process manager.

Syntax

```
mpicleanup [ -i <input_file> | -t -f <hostsfile> ] [ -r <rshcmd> ] \
[ -b <branch_count> ] [-p] [-s | -d] [-h] [-V]
```

or

```
mpicleanup [ --input <input_file> | --total --file <hostsfile> ] \
[ --rsh <rshcmd> ] [ --branch <branch_count> ] [ --parallel ] \
[ --silent | --verbose ] [ --help ] [ --version ]
```

Arguments

-i <input_file> --input <input_file>	Specify the input file generated by <code>mpiexec.hydra</code> . The default value is <code>mpiexec_\${username}_\${PPID}.log</code> located in the temporary directory determined by the values of the <code>I_MPI_TMPDIR</code> or <code>TMPDIR</code> environment variables, or in the <code>/tmp</code> directory.
-t --total	Use the total mode to stop all user processes on the specified machines. This option is not supported for the <code>root</code> user.
-f <hostsfile> --file <hostsfile>	Specify the file containing the list of machines to clean up.
-r <rshcmd> --rsh <rshcmd>	Specify the remote shell to use. The default shell is <code>ssh</code> .
-b <branch_count> --branch <branch_count>	Define the number of the child processes. The default value is 32.
-p --parallel	Use the parallel launch mode. This option is only applicable if all hosts are available. Otherwise a part of machines may stay in an undefined state.
-s --silent	Suppress extra output generation.
-d --verbose	Output verbose information.
-h --help	Display a help message.
-V --version	Display Intel® MPI Library version information.

Description

Use this command to clean up the environment after an abnormal MPI job termination.

For example, use the following command to stop processes mentioned in the input file generated by the prior `mpiexec.hydra` invocation:

```
$ mpicleanup
```

or

```
$ mpicleanup --input /path/to/input.file
```

Use the following command to stop all your user processes on the machines specified in the `hostsfile` file:

```
$ mpicleanup --file hostsfile --total
```


2.3.6. Checkpoint-Restart Support

The Checkpoint-Restart feature in Intel® MPI Library is designed to be application transparent. You can access to the Checkpoint-Restart functionality through the MPI process management interface. The Checkpoint-Restart options and environment variables are applicable to the Hydra process manager only. To use the Hydra process manager, set `I_MPI_PROCESS_MANAGER=hydra` if you to change the process manager from the default value.

NOTE

The Checkpoint-Restart feature requires the OFA* network module. You can choose the OFA network module, for example, with the `I_MPI_FABRICS` environment variable by setting the value to `ofa`, or the `-ib` option.

NOTE

To enable the Checkpoint-Restart feature, set the following:

- 1 for `I_MPI_OFA_DYNAMIC_QPS` environment variable
- 0 for `I_MPI_OFA_NUM_RDMA_CONNECTIONS` environment variable

NOTE

Install the Berkeley Lab Checkpoint/Restart* (BLCR) Software for the Checkpoint-Restart function.

Global Options

-ckpoint <switch>

Arguments

<switch>	Checkpoint switch
enable yes on 1	Enables the check point function for the application started
disable no off 0	Disables the check point function for the application started. This is the default value

Use this option to enable/disable checkpoints capability. When this capability is disabled, other checkpoint options are ignored.

-ckpoint-interval <sec>

Arguments

<sec>	Interval between consecutive checkpoints in seconds
-------	---

Use this option to turn on timer driven checkpoints. See also [Timer Driven Checkpoint](#). The checkpoints are taken every <sec> seconds. If this option is not specified, signal driven checkpoint function may be used. See [Explicit Signal Driven Checkpoint](#) for more details.

-ckptpoint-preserve <N>**Arguments**

<N>	Maximal number of checkpoint images kept. The default value is 1
-----	--

Use this option while running the checkpoint function to keep last <N> checkpoints to reduce checkpoint image space. By default, only the last checkpoint is kept.

-restart

Use this option to restart an application from one of the stored checkpoints. `-ckptpointlib`, `-ckptpoint-prefix` and `-ckptpoint-num` options are meaningful for restarting. The executable name may be provided to the process manager, but is ignored. Taking checkpoints is allowed for the restarted application, so `-restart` option may be accompanied with `-ckptpoint` and other applicable checkpoint options.

-ckptpoint-num <N>**Arguments**

<N>	Identifier of the checkpoint image to restart an application with. Valid values are any number equal or lower than the last checkpoint number. The default is the last checkpoint number.
-----	---

Use this option while restarting an application. The checkpoint number <N> (counting from 0) is taken as a restart point. To determine the best choice for this value, examine the checkpoint storage directory setting with the `-ckptpoint-prefix` option.

NOTE

The number of images determined by the `-ckptpoint-preserve` option is kept at maximum.

The application will abort with an error message during startup if this checkpoint does not exist. By default, the last checkpoint is selected.

Local Options**-ckptpointlib <lib>****Arguments**

<lib>	Checkpoint-Restart system library
blcr	Berkeley Lab Checkpoint/Restart* (BLCR) Library. This is the default value

Use this option to select underlying Checkpoint-Restart system library. Only the Berkeley Lab Checkpoint/Restart* (BLCR) Library is supported.

NOTE

You need to provide the same option when using the checkpoint function, or when restarting the application.

-ckptpoint-prefix <dir>**Arguments**

<dir>	The directory to store checkpoints. The default value is <code>/tmp</code>
-------	--

Use this option to specify a directory to store checkpoints. By default, `/tmp` is used. The directory `<dir>` should be writable, otherwise an error will be raised during process launch, and the application will abort with an error message.

NOTE

You need to provide the same option when using the checkpoint function, or when restarting the application.

-ckpoint-tmp-prefix <dir>

Arguments

<code><dir></code>	The directory to store temporary checkpoints. The default value is <code>/tmp</code>
--------------------------	--

Use this option to indicate the directory to store temporary checkpoints. Checkpoints are migrated from `-ckpoint-tmp-prefix` to the directory specified in `-ckpoint-prefix`. The directory `<dir>` should be writable; otherwise the application will abort during startup with an error message. Temporary storage is not used if the option is not set.

-ckpoint-logfile <file>

Use this option for monitoring checkpoint activity. The trace is dumped into `<file>`. You should be able to write in `<file>`; otherwise the application will abort during startup with an error message. This is an optional feature.

Environment Variables

I_MPI_CKPOINT

Syntax

`I_MPI_CKPOINT=<switch>`

Arguments

<code><switch></code>	Checkpoint switch
<code>enable yes on 1</code>	Enables the check point function for the application started
<code>disable no off 0</code>	Disables the check point function for the application started. This is the default value

Description

Use this variable to turn on taking checkpoints capability. This has the same effect as the `-ckpoint` option. If you have set the `-ckpoint` option, the Hydra process manager sets the `I_MPI_CKPOINT` even if you do not set this environment variable.

I_MPI_CKPOINTLIB

Syntax

`I_MPI_CKPOINTLIB=<lib>`

Arguments

<code><lib></code>	Checkpoint-Restart system library
--------------------------	-----------------------------------

<code>bldcr</code>	Berkeley Lab Checkpoint/Restart* (BLCR) Library. This is the default value
--------------------	--

Description

Use this variable to select underlying Checkpoint-Restart system library. This has the same effect as the `-ckptlib` option.

I_MPI_CKPOINT_PREFIX**Syntax**

`I_MPI_CKPOINT_PREFIX=<dir>`

Arguments

<code><dir></code>	The directory to store checkpoints. The default value is <code>/tmp</code>
--------------------------	--

Description

Use this variable to specify a directory to store checkpoints. This has the same effect as the `-ckpt-prefix` option.

I_MPI_CKPOINT_TMP_PREFIX**Syntax**

`I_MPI_CKPOINT_TMP_PREFIX=<dir>`

Arguments

<code><dir></code>	The directory to store temporary checkpoints
--------------------------	--

Description

Use this variable to indicate storage of temporary checkpoints while `-ckpt-prefix` indicates permanent storage. This has the same effect as the `-ckpt-tmp-prefix` option.

I_MPI_CKPOINT_INTERVAL**Syntax**

`I_MPI_CKPOINT_INTERVAL=<sec>`

Arguments

<code><sec></code>	Interval between consecutive checkpoints in seconds
--------------------------	---

Description

Use this variable to turn on timer driven checkpoints. This has the same effect as the `-ckpt-interval` option.

I_MPI_CKPOINT_PRESERVE**Syntax**

`I_MPI_CKPOINT_PRESERVE=<N>`

Arguments

<code><N></code>	Maximal number of checkpoint images kept. The default value is 1
------------------------	--

Description

Use this option while running the checkpoint function to keep last `<N>` checkpoints to reduce checkpoint image space. This has the same effect as the `-ckpt-preserve` option.

I_MPI_CHKPOINT_LOGFILE

Syntax

`I_MPI_CHKPOINT_LOGFILE=<file>`

Arguments

<code><file></code>	The file keeping the trace for checkpoint activity
---------------------------	--

Description

Use this option for checkpoint activity monitoring. The trace is dumped into `<file>`. This has the same effect as the `-ckpt-logfile` option.

I_MPI_CHKPOINT_NUM

Syntax

`I_MPI_CHKPOINT_NUM=<N>`

Arguments

<code><N></code>	Number of checkpoint image to restart an application with
------------------------	---

Description

Use this option while restarting application. This has the same effect as the `-ckpt-num` option.

I_MPI_RESTART

Syntax

`I_MPI_RESTART=<switch>`

Arguments

<code><switch></code>	Restart switch
<code>enable yes on 1</code>	Enables the restart of the application from one of the stored checkpoints.
<code>disable no off 0</code>	Disables the restart of the application. This is the default value.

Description

Use this variable to restart an application from one of the stored checkpoints. Using this variable has the same effect as `-restart` option.

Running MPI Applications

The checkpoint-restart feature is available with the Hydra process launcher (`mpiexec.hydra`). The launcher provides two mutually exclusive methods of taking checkpoints:

- By timers
- By explicit signal

You can provide directory paths where checkpoints can be stored temporarily and permanently.

Timer Driven Checkpoint

In the following example, a checkpoint is taken every 3600 seconds (=1hour). The checkpoints are stored in a directory called `ckptdir`. Each node generates one checkpoint which is named by the node number and number of that checkpoint.

```
user@head $ mpiexec.hydra -ckpoint on -ckpoint-prefix /home/user/ckptdir -
ckpoint-interval 3600 -ckpointlib blcr -n 32 -f hosts /home/user/myapp
```

Explicit Signal Driven Checkpoint

In the following example, an application is started and then an explicit signal (`SIGUSR1`) is passed to the application to take a checkpoint. The checkpoints are stored in a directory called `ckptdir`.

```
user@head $ mpiexec.hydra -ckpoint on -ckpoint-prefix /home/user/ckptdir -
ckpointlib blcr -n 32 -f hosts /home/user/myapp
```

...

```
user@head $ kill -s SIGUSR1 <PID of mpiexec.hydra>
```

It is necessary and sufficient for you to signal the `mpiexec.hydra` process on node head.

Using Local Storage

In the following example, there are two locations for storing checkpoints.

- Temporary location: indicated in the argument to `-ckpoint-tmp-prefix`
- Permanent location: indicated in the argument to `-ckpoint--prefix`

```
user@head $ mpiexec.hydra -ckpoint on -ckpoint-tmp-prefix
/ssd/user/ckptdir -ckpoint-prefix /home/user/ckptdir -ckpointlib blcr -n
32 -f hosts /home/user/myapp
```

Restarting MPI Applications

The following is an example of restarting an application from checkpoint number `<N>`.

```
user@head $ mpiexec.hydra -restart -ckpoint-prefix /home/user/ckptdir -
ckpointlib blcr -ckpoint-num <N> -n 32 -f hosts
```

When restarting, you need to revise the "hosts" file to eliminate any dead or unavailable nodes. Also, providing the executable name is not necessary when restarting because it is already stored in the checkpoint images.

Viewing Checkpoint Activity in Log File

The following is an example of launching an MPI job and specifying a checkpoint log file so that you can watch the checkpoint activity.

```
user@head $ mpiexec.hydra -ckpoint on -ckpoint-logfile /home/user/ckpt.log -
ckpoint-tmp-prefix /ssd/user/ckptdir -ckpoint-prefix /home/user/ckptdir -
ckpointlib blcr -n 32 -f hosts /home/user/myapp
```

The following output is a sample log:

```
[Mon Dec 19 13:31:36 2011] cst-linux Checkpoint log initialized (master mpiexec
pid 10687, 48 processes, 6 nodes)
[Mon Dec 19 13:31:36 2011] cst-linux Permanent checkpoint storage:
/mnt/lustre/user
[Mon Dec 19 13:31:36 2011] cst-linux Temporary checkpoint storage: /tmp
[Mon Dec 19 13:32:06 2011] cst-linux Started checkpoint number 0 ...
[Mon Dec 19 13:33:00 2011] cst-linux Finished checkpoint number 0.
[Mon Dec 19 13:33:00 2011] cst-linux Moving checkpoint 0 from /tmp to
/mnt/lustre/user ...
[Mon Dec 19 13:38:00 2011] cst-linux Moved checkpoint 0 from /tmp to
/mnt/lustre/user
```

Automatic Cleanup of Previous Checkpoints

Checkpoint images are large; thus, Intel® MPI Library only keeps the last useful checkpoint by default. The following is an example to keep `<N>` previous checkpoints. The flag is `-ckptpoint-preserve <N>`. The default value of `-ckptpoint-preserve` is 1 (keep only the last checkpoint).

```
user@head $ mpiexec.hydra -ckptpoint on -ckptpoint-preserve <N> -ckptpoint-tmp-prefix
/ssd/user/ckptdir -ckptpoint-prefix /home/user/ckptdir -ckptpointlib blcr -n 32 -f
hosts /home/user/myapp
```

2.4. Hetero Operating System Cluster Support

The Intel® MPI Library provides support for heterogeneous Windows-Linux environment. Hydra process manager available on Windows* OS and Linux* OS is used to provide possibility for Intel MPI Library on Linux* OS and Windows* OS to cooperate within one job. For more information about hydra process manager, see [Scalable Process Management System Commands](#).

To run Linux-Windows operating system (OS) mixed job, do the following:

- Ensure the Intel MPI Library is installed and operable on each node of your job.
- The `-demux=select` and `I_MPI_FABRICS=shm:tcp` are supported for the operating system mixed run.
- Set the `-bootstrap` option. The default value in operating system mixed run mode is `-bootstrap service`. To enable such configuration, ensure the hydra service is launched on Windows* OS on each node of an MPI job. Provide the ssh connectivity between Linux and Windows machines for the `-bootstrap ssh`.
- Use `-hostos` option to specify an operating system installed on a particular host.
- Use `I_MPI_ROOT` and `PATH` local environment variables to overwrite incorrect settings for adjacent operating system environment inheritance.

For example, the following command runs `IMB-MPI1` job under Windows-Linux heterogeneous environment:

```
$ mpirun -demux select -genv I_MPI_FABRICS shm:tcp -env I_MPI_ROOT \
<windows_installdir> -env PATH <windows_installdir>\\<arch>\\bin -hostos \
windows -host <win_host> -n 2 <windows_installdir>\\<arch>\\bin\\IMB-MPI1 : \ -
host <lin_host> -n 3 <linux_installdir>/<arch>/bin/IMB-MPI1
```

2.5. Intel® Xeon Phi™ Coprocessor Support

This topic concentrates on the Intel® MPI Library specifics related to the support of the Intel® Xeon Phi™ Coprocessor (codename: Knights Corner) based on Intel® Many Integrated Core Architecture (Intel® MIC Architecture).

2.5.1. Usage Model

To use the Intel MPI Library on Intel® Xeon Phi™ Coprocessor (codename: Knights Corner), ensure that:

- Each host and each Intel® Xeon Phi™ coprocessor must have a unique IP address and symbolic name, which is the same as to classic cluster.
- Password-less access between host and Intel® Xeon Phi™ Coprocessor by ssh is established.

If the connection fails, check the possible causes and solutions in the following table:

Possible Cause	Solution
The version of Intel® MIC Software Stack you used is out of date.	Install a newer version.
The <code>iptables</code> service is running on the host.	Stop that service.
The route is incomplete.	Add the missing routes.

Refer to the system administrator and *Intel® MIC Software Stack readme* to configure the settings for the IP connectivity.

When using Intel MPI Library on an Intel Xeon Phi coprocessor, consider Intel Xeon Phi coprocessor card to be another cluster node with a different Intel® architecture. The way that MPI features work for the Intel Xeon Phi coprocessor is similar to the way they work for an Intel® Xeon processor.

For example, MPI libraries may be available on both Intel Xeon processor and Intel Xeon Phi coprocessor through an NFS share that has the same path for Intel Xeon processor host and Intel Xeon Phi coprocessor; MPI tasks may be started from Intel Xeon processor host or Intel Xeon Phi coprocessor, and so on.

To build an application for running on the Intel Xeon Phi coprocessor and the host node, go through the following steps:

1. Establish environment settings for the compiler and for the Intel MPI Library:

```
(host)$ . <compiler_installdir>/bin/compilervars.sh intel64em64t
```

```
(host)$ . <mpi_installdir>/intel64em64t/bin/mpivars.sh
```

2. Build your application for Intel MIC Architecture, for example:

```
(host)$ mpiicc -mmic test.c -o test_hello.mic
```

3. Build your application for Intel 64 Architecture, for example:

```
(host)$ mpiicc test.c -o test_hello
```

To run an application on the Intel Xeon Phi coprocessor and the host node, go through the following steps:

1. Ensure that NFS is properly set-up between the hosts and the Intel Xeon Phi coprocessor, which is the recommended way for using Intel MPI Library on Intel MIC Architecture.

For information on how to set up NFS on the Intel Xeon Phi coprocessor, see

<http://software.intel.com/en-us/articles/intel-mpi-library-for-linux-kb/all/> or

<http://software.intel.com/mic-developer> .

2. Establish environment settings for the Intel MPI Library:

```
(host)$ . <mpi_installdir>/intel64em64t/bin/mpivars.sh
```

3. Launch the executable file from host, for example:

```
(host)$ export I_MPI_MIC=1
```

```
(host)$ mpirun -n 2 -host <host ID> ./test_hello : -n 2 -host <coprocessor ID> ./test_hello.mic
```

NOTE

See *Intel® MPI Library for Linux® OS Reference Manual* for `-configfile`, `-hostfile` and `-machinefile` options which also can be used.

To run the application on Intel Xeon Phi coprocessor only, follow the steps described above except for the step of building the application for Intel 64 Architecture. Meanwhile, ensure that the hostfile only contains the Intel Xeon Phi coprocessor name.

For more details, see <http://software.intel.com/en-us/articles/intel-mpi-library-for-linux-kb/all/> and <http://software.intel.com/mic-developer>.

2.5.2. Environment Variables

I_MPI_MIC

Syntax

`I_MPI_MIC=<value>`

Arguments

<code><value></code>	Intel® Xeon Phi™ recognition
<code>enable yes on 1</code>	Enable the Intel Xeon Phi coprocessor recognition
<code>disable no off 0</code>	Disables the Intel Xeon Phi coprocessor recognition. This is the default value

Description

Set this environment variable to control whether the Intel Xeon processor of the Intel® MPI Library tries to detect and work with the Intel® MIC Architecture components.

If the value of environment variable `I_MPI_MIC` is `enable`, the default value of environment variable `I_MPI_SSHM` is `enable`.

NOTE

This is a provisional variable and is only temporarily introduced, until the architecture detection and other related matters are clarified.

I_MPI_MIC_PREFIX

Syntax

`I_MPI_MIC_PREFIX=<value>`

Arguments

<code><value></code>	Specify a string as the prefix of an Intel Xeon Phi coprocessor file name. The default value is an empty string
----------------------------	---

Description

Set this environment variable to add a prefix to a host executable file name to get a corresponding Intel Xeon Phi coprocessor executable file name.

For example, set different locations as the value for the `I_MPI_MIC_PREFIX` environment variable to distinguish Intel MIC Architecture and Intel® 64 Architecture executable files:

```
(host)$ mpiicc test.c -o test_hello
(host)$ mpiicc -mmic test.c -o ./MIC/test_hello
(host)$ export I_MPI_MIC=1
(host)$ export I_MPI_MIC_PREFIX=./MIC/
```

```
(host)$ mpirun -n 4 -hostfile <hostfile> test_hello
```

In the example, `./test_hello` binary is launched on Intel® 64 Architecture nodes and `./MIC/test_hello` binary is launched on Intel Xeon Phi coprocessor nodes.

I_MPI_MIC_POSTFIX

Syntax

```
I_MPI_MIC_POSTFIX=<value>
```

Arguments

<code><value></code>	Specify a string as the postfix of an Intel Xeon Phi coprocessor file name. The default value is an empty string
----------------------------	--

Description

Set this environment variable to add a postfix to a host executable name to get a corresponding Intel Xeon Phi coprocessor executable file name.

For example, set different names as the value for the `I_MPI_MIC_POSTFIX` environment variable to distinguish Intel Xeon Phi coprocessor and Intel 64 Architecture executable files:

```
(host)$ mpiicc test.c -o test_hello
(host)$ mpiicc -mmic test.c -o test_hello.mic
(host)$ export I_MPI_MIC=1
(host)$ export I_MPI_MIC_POSTFIX=.mic
(host)$ mpirun -n 4 -hostfile <hostfile> test_hello
```

In the example, `test_hello` binary is launched on Intel 64 Architecture nodes and `test_hello.mic` binary on Intel Xeon Phi coprocessor nodes.

I_MPI_DAPL_PROVIDER_LIST

Syntax

```
I_MPI_DAPL_PROVIDER_LIST=<primary provider>[,<local secondary provider>
[,<remote secondary provider>]]
```

Arguments

<code><primary provider></code>	Provides the best latency and available on all network segments (cross box and within box)
<code><local secondary provider></code>	Provides the best bandwidth for local configurations. The distance is less than the value of the <code>I_MPI_DAPL_LOCALITY_THRESHOLD</code> environment variable
<code><remote secondary provider></code>	Provides best bandwidth for remote configurations. The distance greater than the value of the <code>I_MPI_DAPL_LOCALITY_THRESHOLD</code> environment variable

Description

Use this variable to define the DAPL providers to load.

With Intel® Manycore Platform Software Stack (Intel® MPSS), the arguments of `I_MPI_DAPL_PROVIDER_LIST` are set as the following values:

- `<primary provider>`- CCL-direct
- `<local secondary provider>`- IBSCIF

- `<remote secondary provider>- CCL-proxy`

Thus, the setting is `I_MPI_DAPL_PROVIDER_LIST=<CCL-direct>[,<IBSCIF>[,<CCL-proxy>]]`

The following configuration is an example with the default `dat.conf` provided with Intel MPSS:

```
I_MPI_DAPL_PROVIDER_LIST=ofa-v2-mlx4_0-1u,ofa-v2-scif0,ofa-v2-mcm-1
```

You can adjust the threshold for secondary provider through the

`I_MPI_DAPL_DIRECT_COPY_THRESHOLD` environment variable (`<secondary provider threshold>`):

```
I_MPI_DAPL_DIRECT_COPY_THRESHOLD=<primary provider direct copy
threshold>[,<secondary provider threshold>]
```

`<primary provider direct copy threshold>` has to be lower than `<secondary provider threshold>`.

If the environment variable `I_MPI_DAPL_PROVIDER_LIST` contains a list of values, then the syntax of the following environment variables may be extended by the values related to all corresponding providers.

- `I_MPI_DAPL_DIRECT_COPY_THRESHOLD`
- `I_MPI_DAPL_TRANSLATION_CACHE`
- `I_MPI_DAPL_TRANSLATION_CACHE_AVL_TREE`
- `I_MPI_DAPL_CONN_EVD_SIZE`
- `I_MPI_DAPL_RDMA_RNDV_WRITE`

If you set only single value, this value applies to all providers. In case of mismatch or incorrect values, the default value is used for all providers.

For example:

```
export I_MPI_DAPL_PROVIDER_LIST=ofa-v2-mlx4_0-1,ofa-v2-scif0
```

```
export I_MPI_DAPL_TRANSLATION_CACHE=enable,disable
```

This `I_MPI_DAPL_TRANSLATION_CACHE` setting turns on the memory registration cache for the first provider; but turns it off for the second one.

I_MPI_DAPL_LOCALITY_THRESHOLD

Define the threshold to switch from local secondary provider to remote secondary provider.

Syntax

```
I_MPI_DAPL_LOCALITY_THRESHOLD=<value>
```

Arguments

<code><value></code>	Defines the threshold to switch from local secondary provider to remote secondary provider. See I_MPI_DAPL_PROVIDER_LIST for the detail options descriptions
See the following description	The default value depends on Non-Uniform Memory Architecture (NUMA) configuration and DAPL*/Intel® Manycore Platform Software Stack (Intel® MPSS) version.

Description

This value from range [10; 255] is associated with NUMA distances inside the host and extrapolated for Intel Xeon Phi coprocessors. The default value relates to clusters and jobs, and is defined with the formulas shown below.

Starting from DAPL* 2.1.3, there is an auto adjustment logic which uses the formula: $255 - d_{\max} + d_{\min}$.

The following table shows the list for formulas that you should use for the

`I_MPI_DAPL_LOCALITY_THRESHOLD` default value in difference cases:

Formula for the Default Value	Mode
$d + (255 - 1 - d_{\max})$	Two ranks are executed on different Intel Xeon Phi coprocessors of the same host
$d + (255 - 1 - d_{\max}) / 2$	One rank is executed on Intel Xeon Phi coprocessor and another on the same host
d	Both ranks are executed inside the host
255	Two ranks are executed on different hosts

In the above formula:

- d - host system NUMA distance.
- d_{\min} - minimal host system NUMA distance.
- d_{\max} - maximal host system NUMA distance.

Before DAPL 2.1.3, `I_MPI_DAPL_LOCALITY_THRESHOLD` is equal 255 by default, what conforms to a local secondary provider will be chosen for all ranks which are executed inside the box; otherwise a remote secondary provider will be chosen.

`I_MPI_ENV_PREFIX_LIST`

Define the prefixes of environment variables for the intended platforms.

Syntax

```
I_MPI_ENV_PREFIX_LIST=[platform:prefix][,...]
```

Argument

<code>platform</code>	The intended platform (string). Options: <code>htn</code> , <code>nhm</code> , <code>wsm</code> , <code>snb</code> , <code>ivb</code> See I_MPI_PLATFORM for the detail options descriptions
<code>prefix</code>	A prefix (string) for a name of an environment variable to be used for the intended platform

Description

Set this environment variable to define the prefix of environment variables to be used for the intended platform.

If you specify a prefix in `I_MPI_ENV_PREFIX_LIST` for an environment variable, the prefixed environment variable overrides the respective non-prefixed environment variable on the intended platform.

If you do not specify `I_MPI_ENV_PREFIX_LIST`, environment variables are applied to all platforms.

NOTE

Use the lower case when you specify the platform names.

Examples

1. `I_MPI_ENV_PREFIX_LIST=platform:prefix`
`<NAME>=value` is applied to all systems.
`<prefix>_<NAME>=value` defines `<NAME>=value` for all `<platform>` systems.
2. Assume that some machines are on the Intel® microarchitecture code name Sandy Bridge based platform, and the rest machines are on other architectures based platforms. The environment variable `OMP_NUM_THREADS` value is 3 on all platforms.

To set `OMP_NUM_THREADS=5` for the ranks on the Intel® microarchitecture code name Sandy Bridge based platform, specify the prefix in `I_MPI_ENV_PREFIX_LIST` for `OMP_NUM_THREADS` with the following configurations:

```
I_MPI_ENV_PREFIX_LIST=snb:<prefix>
OMP_NUM_THREADS=3
<prefix>_OMP_NUM_THREADS=5
```

2.5.3. Compiler Commands

The following table lists available MPI compiler commands and Intel® Composer XE 2013 for Linux* OS for Intel® MIC Architecture, languages, and application binary interfaces (ABIs) that they support.

Compiler Command	Default Compiler	Supported Language(s)	Supported ABI(s)
<code>mpiicc</code>	<code>icc</code>	C	64 bit
<code>mpiicpc</code>	<code>icpc</code>	C++	64 bit
<code>mpiifort</code>	<code>ifort</code>	Fortran77/Fortran 95	64 bit

The compiler commands have the following common features:

- The compiler commands reside in the `<installdir>/intel64em64t/bin` directory.
- The environment settings should be established by sourcing the `<installdir>/intel64em64t/bin/mpivars.sh` script. If you need to establish environment settings for different library configurations, you can pass one of the following arguments to the `mpivars.sh` script to switch to corresponding configurations:
 - o `debug`
 - o `release`
 - o `debug_mt`
 - o `release_mt`

Multi-threaded optimized library is chosen by default.

- To compile a heterogeneous MPI application, compile it twice: one time for Intel® 64 Architecture and another time for Intel® MIC Architecture.
- To distinguish targeted architectures, the scripts parse the underlying compiler options. If they detect the compiler options that target Intel® MIC Architecture (such as `-mmic`) is currently used by Intel® Composer XE 2013 for Linux* OS for Intel® MIC Architecture, they create an Intel® MIC Compiler executable file. Otherwise, they create an Intel® Xeon processor executable file.

- GNU* Compiler use requires that the compiler be specified with `-cc/-cxx/-fc/-f77/-f90` options or through the environment variables described in the Reference Manual. For example:

```
(host)$ mpicc -cc=/usr/linux-k10m-4.7/bin/x86_64-k10m-linux-gcc -mmic
test.c -o test_hello.mic
```

NOTE

Use different file names and/or locations to distinguish Intel® MIC Architecture and Intel® 64 Architecture executable files.

2.6. Multipurpose Daemon Commands

mpd

Start Multipurpose daemon* (MPD).

Syntax

```
mpd [ --help ] [ -V ] [ --version ] [ --host=<host> --port=<portnum> ] \
[ --noconsole ] [ --trace ] [ --echo ] [ --daemon ] [ --bulletproof ] \
[ --i fhn <interface/hostname> ] [ --listenport <listenport> ]
```

Arguments

<code>--help</code>	Display a help message
<code>-V --version</code>	Display the Intel® MPI Library version information
<code>-h <host> -p <portnum> --host=<host> --port=<portnum></code>	Specify the host and port to be used for entering an existing ring. The <code>--host</code> and <code>--port</code> options must be specified together
<code>-n --noconsole</code>	Do not create a console at startup
<code>-t --trace</code>	Print internal MPD trace information
<code>-e --echo</code>	Print a port number at startup to which other <code>mpds</code> may connect
<code>-d --daemon</code>	Start <code>mpd</code> in daemon mode. By default, the interactive mode is enabled
<code>--bulletproof</code>	Turn MPD bulletproofing on
<code>--ifhn=<interface/hostname></code>	Specify <code><interface/hostname></code> to use for MPD communications
<code>-l <listenport> --listenport=<listenport></code>	Specify the <code>mpd</code> listening port

Description

Multipurpose daemon* (MPD) is the Intel® MPI Library process management system for starting parallel jobs. Before running a job, start `mpd` daemons on each host and connect them into a ring. Long parameter names may be abbreviated to their first letters by using only one hyphen and no equal sign. For example,

```
$ mpd -h masterhost -p 4268 -n
```

is equivalent to

```
$ mpd --host=masterhost --port=4268 --noconsole
```

If a file named `.mpd.conf` is available in the user's home directory, only the user can have read and write privileges. The file must minimally contain a line with `secretword=<secretword>`. If you want to run MPD as root, create the `mpd.conf` file in the `/etc` directory instead of `.mpd.conf` in the root's home directory to run `mpd` as root. Avoid starting the MPD ring under the root account.

NOTE

Multipurpose daemon* (MPD) has been deprecated since Intel® MPI Library 5.0 release. You can use scalable process management system (Hydra) to start parallel jobs.

mpdboot

Start `mpd` ring.

Syntax

```
mpdboot [ -h ] [ -V ] [ -n <#nodes> ] [ -f <hostsfile> ] [ -r <rshcmd> ] \
[ -u <user> ] [ -m<mpdcmd> ] [ --locons ] [ --remcons ] \
[ -s ] [ -d ] [ -v ] [ -l ] [ --ncpus=<ncpus> ] [ -o ] \
[ -b <maxbranch> ] [ -p ]
```

or

```
mpdboot [ --help ] [ --version ] [ --totalnum=<#nodes> ] \
[ --file=<hostsfile> ] [ --rsh=<rshcmd> ] [ --user=<user> ] \
[ --mpd=<mpdcmd> ] [ --locons ] [ --remcons ] [ --shell ] \
[ --debug ] [ --verbose ] [ -l ] [ --ncpus=<ncpus> ] [ --ordered ]
[ --maxbranch=<maxbranch> ] [ --parallel-startup ]
```

Arguments

<code>-h --help</code>	Display a help message
<code>-V --version</code>	Display Intel® MPI Library version information
<code>-d --debug</code>	Print debug information
<code>-v --verbose</code>	Print more information. Show the <code><rshcmd></code> attempts
<code>-n <#nodes> </code> <code>--totalnum=<#nodes></code>	Number of nodes in <code>mpd.hosts</code> on which daemons are started
<code>-r <rshcmd> </code> <code>--rsh=<rshcmd></code>	Specify remote shell to start daemons and jobs. The default value is <code>ssh</code>
<code>-f <hostsfile> </code> <code>--file=<hostsfile></code>	Path/name of the file that has the list of machine names on which the daemons are started

<code>-l</code>	Enable starting multiple <code>mpd</code> per machine
<code>-m <mpdcmd> </code> <code>--mpd=<mpdcms></code>	Specify the full path name of the <code>mpd</code> on the remote hosts
<code>-s --shell</code>	Specify the shell
<code>-u <user> --</code> <code>user=<user></code>	Specify the user
<code>--locons</code>	Do not create local MPD consoles
<code>--remcons</code>	Do not create remote MPD consoles
<code>--ncpus=<ncpus></code>	Indicate how many processors to use on the local machine (other nodes are listed in the hosts file)
<code>-o --ordered</code>	Start all the <code>mpd</code> daemons in the order as specified in the <code>mpd.hosts</code> file
<code>-b <maxbranch> </code> <code>--maxbranch=<maxbranch></code>	Use this option to indicate the maximum number of the <code>mpd</code> daemons to enter the <code>mpd</code> ring under another. This helps to control the parallelism of the <code>mpd</code> ring start. The default value is four
<code>-p --parallel-startup</code>	Use this option to allow parallel fast starting of <code>mpd</code> daemons under one local root. No daemon checking is performed. This option also supports shells which do not transfer the output from the remote commands

Description

Start the `mpd` daemons on the specified number of nodes by providing a list of node names in `<mpd.hosts>`.

The `mpd` daemons are started using the `ssh` command by default. If the `ssh` connectivity is not enabled, use the `-r rsh` option to switch over to `rsh`. Make sure that all nodes in the cluster can connect to each other through the `ssh` command without a password or, if the `-r rsh` option is used, through the `rsh` command without a password.

NOTE

The `mpdboot` command spawns an MPD daemon on the host machine, even if the machine name is not listed in the `mpd.hosts` file.

mpdexit

Shut down a single `mpd` daemon.

Syntax

```
mpdexit [ --help ] [ -V ] [--version ] <mpdid>
```

Arguments

<code>--help</code>	Display a help message
<code>-V --version</code>	Display Intel® MPI Library version information
<code><mpdid></code>	Specify the <code>mpd</code> daemon to kill

Description

Use this command to cause the single `mpd` daemon to exit. Use `<mpdid>` obtained through the `mpdtrace -l` command in the form `<hostname>_<port number>`.

mpdallexit

Shut down all `mpd` daemons on all nodes.

Syntax

```
mpdallexit [ --help ] [ -V ] [ --version ]
```

Arguments

<code>--help</code>	Display a help message
<code>-V</code> <code>--version</code>	Display Intel® MPI Library version information

Description

Use this command to shut down all MPD rings you own.

mpdcleanup

Clean up the environment after an `mpd` crash.

Syntax

```
mpdcleanup [ -h ] [ -V ] [ -f <hostsfile> ] [ -r <rshcmd> ] [ -u <user> ] \
[ -c <cleancmd> ] [ -a ]
```

or

```
mpdcleanup [ --help ] [ --version ] [ --file=<hostsfile> ] \
[ --rsh=<rshcmd> ] [ --user=<user> ] [ --clean=<cleancmd> ] \
[ --all ]
```

Arguments

<code>-h</code> <code>--help</code>	Display a help message
<code>-V</code> <code>--version</code>	Display Intel® MPI Library version information
<code>-f <hostsfile></code> <code>--file=<hostsfile></code>	Specify the file containing a list of machines to clean up
<code>-r <rshcmd></code> <code>--rsh=<rshcmd></code>	Specify the remote shell to use
<code>-u <user></code> <code>--user=<user></code>	Specify the user
<code>-c <cleancmd></code> <code>--clean=<cleancmd></code>	Specify the command to use for removing the UNIX* socket. The default command is <code>/bin/rm -f</code>
<code>-a</code> <code>--all</code>	Kill all <code>mpd</code> daemons related to the current settings of the <code>I_MPI_JOB_CONTEXT</code> environment variable on all hosts specified in <code><hostsfile></code>

Description

Use this command to clean up the environment after an `mpd` crash. It removes the UNIX* socket on local and remote machines or kills all `mpd` daemons related to the current environment controlled by the `I_MPI_JOB_CONTEXT` environment variable.

For instance, use the following command to remove the UNIX sockets on machines specified in the `hostsfile` file:

```
$ mpdcleanup --file=hostsfile
```

Use the following command to kill the `mpd` daemons on the machines specified in the `hostsfile` file:

```
$ mpdcleanup --file=hostsfile --all
```

mpdtrace

Determine whether `mpd` is running.

Syntax

```
mpdtrace [ --help ] [ -V ] [ --version ] [ -l ]
```

Arguments

<code>--help</code>	Display a help message
<code>-V</code> <code>--version</code>	Display Intel® MPI Library version information
<code>-l</code>	Show MPD identifiers instead of the hostnames

Description

Use this command to list the hostnames or identifiers of all `mpds` in the ring. The output identifiers have the form `<hostname>_<port number>`.

mpdcheck

Check for configuration problems on the host or print configuration information about this host.

Syntax

```
mpdcheck [ -v ] [ -l ] [ -h ] [ --help ] [ -V ] [ --version ]
```

```
mpdcheck -pc [ -v ] [ -l ]
```

```
mpdcheck -f <host_file> [ -ssh ] [ -v ] [ -l ]
```

```
mpdcheck -s [ -v ] [ -l ]
```

```
mpdcheck -c <server_host> <server_port> [ -v ] [ -l ]
```

Arguments

<code>-h</code> <code>--help</code>	Display a help message
<code>-V</code> <code>--version</code>	Display Intel® MPI Library version information
<code>-pc</code>	Print configuration information about a local host
<code>-f <host_file></code>	Print information about the hosts listed in <code><host_file></code>
<code>-ssh</code>	Invoke testing of <code>ssh</code> on each remote host. Use in conjunction with the <code>-f</code> option

<code>-s</code>	Run <code>mpdcheck</code> as a server on one host
<code>-c <server_host> <server_port></code>	Run <code>mpdcheck</code> as a client on the current or different host. Connect to the <code><server_host> <server_port></code>
<code>-l</code>	Print diagnostic messages in extended format
<code>-v</code>	Print the actions that <code>mpdcheck</code> is performing

Description

Use this command to check configuration problems on the cluster nodes. Any output line that starts with ******* indicates a potential problem.

If you have problems running parallel jobs through `mpd` on one or more hosts, try to run the script once on each of those hosts.

mpdringtest

Test the MPD ring.

Syntax

```
mpdringtest [ --help ] [ -V ] [ --version ] <number of loops>
```

Arguments

<code>--help</code>	Display a help message
<code>-V --version</code>	Display Intel® MPI Library version information
<code><number of loops></code>	Number of loops

Description

Use this command to test how long it takes for a message to circle the `mpdring`.

mpdlistjobs

Syntax

```
mpdlistjobs [ -h ] [ -V ] [ -u <username> ] [ -a <jobalias> ] [ -j <jobid> ]
```

or

```
mpdlistjobs [ --help ] [ --version ] [ --user=<username> ] \  
[ --alias=<jobalias> ] [ --jobid=<jobid> ]
```

Arguments

<code>-h --help</code>	Display a help message
<code>-V --version</code>	Display Intel® MPI Library version information
<code>-u <username> --user=<username></code>	List jobs of a particular user
<code>-a <jobalias> -- alias=<jobalias></code>	List information about the particular job specified by <code><jobalias></code>

<code>-j <jobid> </code> <code>--jobid=<jobid></code>	List information about the particular job specified by <code><jobid></code>
---	--

Description

Use this command to list the running processes for a set of MPI jobs. All jobs for the current machine are displayed by default.

mpdsigjob

Apply a signal to a process running an application.

Syntax

```
mpdsigjob [ --help ] [ -V ] [ --version ] <sigtype> \
[-j <jobid> | -a <jobalias> ] [-s | -g ]
```

Arguments

<code>--help</code>	Display a help message
<code>-V --version</code>	Display Intel® MPI Library version information
<code><sigtype></code>	Specify the signal type to send. Valid options are <code>-j</code> or <code>-a</code> .
<code>-a <jobalias></code>	Send a signal to the job specified by <code><jobalias></code>
<code>-j <jobid></code>	Send a signal to the job specified by <code><jobid></code>
<code>-s</code>	Deliver a signal to a single user process
<code>-g</code>	Deliver a signal to a group of processes. This is the default behavior.

Description

Use this command to deliver a specific signal to the processes of a running job. The desired signal is the first argument. Specify one of two options: `-j` or `-a`.

mpdkilljob

Terminate a job.

Syntax

```
mpdkilljob [ --help ] [ -V ] [ --version ] [ <jobnum> ] [ -a <jobalias> ]
```

Arguments

<code>--help</code>	Display a help message
<code>-V --version</code>	Display Intel® MPI Library version information
<code><jobnum></code>	Kill the job specified by <code><jobnum></code>
<code>-a <jobalias></code>	Kill the job specified by <code><jobalias></code>

Description

Use this command to kill the job specified by *<jobnum>* or by *<jobalias>*. Obtain *<jobnum>* and *<jobalias>* from the `mpdlistjobs` command. The *<jobid>* field has the following format: *<jobnum>@<mpdid>*.

mpdhelp

Print brief help concerning MPD commands.

Syntax

```
mpdhelp [ -V ] [ --version ]
```

Arguments

<code>-V --version</code>	Display Intel® MPI Library version information
-----------------------------	--

Description

Use this command to obtain a brief help message concerning MPD commands.

2.6.1. Job Startup Commands

mpiexec

Syntax

```
mpiexec <g-options> <l-options> <executable>
```

or

```
mpiexec <g-options> <l-options> <executable1> : \
<l-options> <executable2>
```

or

```
mpiexec -configfile <file>
```

Arguments

<i><g-options></i>	Global options that apply to all MPI processes
<i><l-options></i>	Local options that apply to a single arg-set
<i><executable></i>	<code>./a.out</code> or path/name of the executable file
<i><file></i>	File with command-line options

Description

Use the first command-line syntax to start all MPI processes of the *<executable>* with the single arg-set. For example, the following command executes `a.out` over the specified *<# of processes>*:

```
$ mpiexec -n <# of processes> ./a.out
```

Use the second command-line syntax to start several MPI programs or the same MPI program with different argument sets. For example, the following command runs each given executable on a different host:

```
$ mpiexec -n 2 -host host1 ./a.out : \
-n 2 -host host2 ./b.out
```

Use the third command-line syntax to read the command line from specified *<file>*. For a command with a single arg-set, the entire command should be specified on a single line in *<file>*. For a command with

multiple arg-sets, each arg-set should be specified on a single, separate line in *<file>*. Global options should always appear at the beginning of the first line in *<file>*.

MPD daemons must already be running in order for `mpirun` to succeed.

NOTE

If there is no "." in the PATH environment variable on all nodes in the cluster, specify *<executable>* as `./a.out` rather than `a.out`.

Extended Device Control Options

Use these options to select a specific fabric combination.

The exact combination of fabrics depends on the number of processes started per node.

If all processes start on one node, the Intel® MPI Library uses `shm` intra-node communication regardless of the selected option from the list in this topic.

If the number of started processes is less than or equal to the number of available nodes, the library uses the first available fabric from the list of fabrics for inter-node communication.

For other cases, the library uses `shm` for intra-node communication, and the first available fabric from the list of fabrics for inter-node communication. See [I_MPI_FABRICS](#) and [I_MPI_FABRICS_LIST](#) for more details.

The `shm` fabric is available for both Intel® and non-Intel microprocessors, but it may perform additional optimizations for Intel microprocessors than it performs for non-Intel microprocessors.

-rdma

Use this option to select an RDMA-capable network fabric for inter-node communication. The application attempts to use first available RDMA-capable network fabric from the list `dapl` or `ofa`. If no such fabric is available, other fabrics from the list `tcp` or `tmi` are used. This option is equivalent to the `-genv I_MPI_FABRICS_LIST dapl,ofa,tcp,tmi -genv I_MPI_FALLBACK 1` setting.

-RDMA

Use this option to select an RDMA-capable network fabric for inter-node communication. The application attempts to use first available RDMA-capable network fabric from the list `dapl` or `ofa`. The application fails if no such fabric is found. This option is equivalent to the `-genv I_MPI_FABRICS_LIST dapl,ofa -genv I_MPI_FALLBACK 1` setting.

-dapl

Use this option to select DAPL capable network fabric for inter-node communication. The application attempts to use DAPL capable network fabric. If no such fabric is available, another fabric from the list `tcp,tmi` or `ofa` is used. This option is equivalent to the `-genv I_MPI_FABRICS_LIST dapl,tcp,tmi,ofa -genv I_MPI_FALLBACK 1` setting.

-DAPL

Use this option to select DAPL capable network fabric for inter-node communication. The application fails if no such fabric is found. This option is equivalent to the `-genv I_MPI_FABRICS_LIST dapl -genv I_MPI_FALLBACK 0` setting.

-ib

Use this option to select OFA capable network fabric for inter-node communication. The application attempts to use OFA capable network fabric. If no such fabric is available, another fabric from the list `dapl,tcp` or `tmi` is used. This option is equivalent to the `-genv I_MPI_FABRICS_LIST ofa,dapl,tcp,tmi -genv I_MPI_FALLBACK 1` setting.

-IB

Use this option to select OFA capable network fabric for inter-node communication. The application fails if no such fabric is found. This option is equivalent to the `-genv I_MPI_FABRICS_LIST ofa -genv I_MPI_FALLBACK 0` setting.

-tmi

Use this option to select TMI capable network fabric for inter-node communication. The application attempts to use TMI capable network fabric. If no such fabric is available, another fabric from the list `dapl,tcp` or `ofa` is used. This option is equivalent to the `-genv I_MPI_FABRICS_LIST tmi,dapl,tcp,ofa -genv I_MPI_FALLBACK 1` setting.

-TMI

Use this option to select TMI capable network fabric for inter-node communication. The application will fail if no such fabric is found. This option is equivalent to the `-genv I_MPI_FABRICS_LIST tmi -genv I_MPI_FALLBACK 0` setting.

-mx

Use this option to select Myrinet MX* network fabric for inter-node communication. The application attempts to use Myrinet MX* network fabric. If no such fabric is available, another fabric from the list `dapl,tcp` or `ofa` is used. This option is equivalent to the `-genv I_MPI_FABRICS_LIST tmi,dapl,tcp,ofa -genv I_MPI_TMI_PROVIDER mx -genv I_MPI_DAPL_PROVIDER mx -genv I_MPI_FALLBACK 1` setting.

-MX

Use this option to select Myrinet MX* network fabric for inter-node communication. The application fails if no such fabric is found. This option is equivalent to the `-genv I_MPI_FABRICS_LIST tmi -genv I_MPI_TMI_PROVIDER mx -genv I_MPI_FALLBACK 0` setting.

-psm

Use this option to select Intel® True Scale Fabric for inter-node communication. The application attempts to use Intel True Scale Fabric. If no such fabric is available, another fabric from the list `dapl,tcp` or `ofa` is used. This option is equivalent to the `-genv I_MPI_FABRICS_LIST tmi,dapl,tcp,ofa -genv I_MPI_TMI_PROVIDER psm -genv I_MPI_FALLBACK 1` setting.

-PSM

Use this option to select Intel True Scale Fabric for inter-node communication. The application fails if no such fabric is found. This option is equivalent to the `-genv I_MPI_FABRICS_LIST tmi -genv I_MPI_TMI_PROVIDER psm -genv I_MPI_FALLBACK 0` setting.

Global Options

-version or -V

Use this option to display Intel® MPI Library version information.

-h or -help or --help

Use this option to display the `mpiexec` help message.

-tune [<arg >]

where:

`<arg>={<dir_name>, <configuration_file>}`.

Use this option to optimize the Intel® MPI Library performance using data collected by the `mpitune` utility. If `<arg>` is not specified, the best-fit tuning options are selected for the given configuration. The default location of the configuration file is `<installdir>/<arch>/etc` directory. You can override this default location by explicitly stating: `<arg>=<dir_name>`. The provided configuration file is used if you set `<arg>=<configuration_file>`.

See [Automatic Tuning Utility](#) for more details.

NOTE

If `<arg>` does not point to the configuration file, set the environment variable `I_MPI_FABRICS`. If `I_MPI_FABRICS` is not set, performance data will not be found and warnings will be printed.

-nolocal

Use this option to avoid running `<executable>` on the host where the `mpiexec` is launched. This option is useful for clusters that deploy a dedicated master node for starting the MPI jobs, and a set of compute nodes for running the actual MPI processes.

-perhost <# of processes>

Use this option to place the indicated number of consecutive MPI processes on every host in a group using round robin scheduling. The total number of processes to start is controlled by the `-n` option.

The `mpiexec` command controls how the ranks of the processes are allocated to the nodes in the cluster. By default, `mpiexec` uses round-robin assignment of ranks to nodes, executing consecutive MPI processes on all processor cores.

To change this default behavior, set the number of processes per host by using the `-perhost` option, and set the total number of processes by using the `-n` option. See [Local Options](#) for details. The first `<# of processes>` indicated by the `-perhost` option is executed on the first host; the next `<# of processes>` is executed on the next host, and so on.

See also the `I_MPI_PERHOST` environment variable.

-rr

Use this option to execute consecutive MPI processes on different hosts using round robin scheduling. This option is equivalent to `-perhost 1`.

-grr <# of processes>

Use this option to place the indicated number of consecutive MPI processes on every host using round robin scheduling. This option is equivalent to `-perhost <# of processes>`.

-ppn <# of processes>

Use this option to place the indicated number of consecutive MPI processes on every host using round robin scheduling. This option is equivalent to `-perhost <# of processes>`.

-machinefile <machine file>

Use this option to control the process placement through `<machine file>`. The total number of processes to start is controlled by the `-n` option.

A machine file is a list of fully qualified or short host names, one name per line. Blank lines and lines that start with `#` as the first character are ignored.

By repeating a host name, you place additional processes on this host. You can also use the following format to avoid repetition of the same host name: `<host name>:<number of processes>`. For example, the following machine file:

```
host1
host1
```



```
host2
host2
host3
```

is equivalent to:

```
host1:2
host2:2
host3
```

It is also possible to specify the network interface used for communication for each node: *<host name>:<number of processes> [ifhn=<interface_host_name>]*.

NOTE

The `-machinefile`, `-ppn`, `-rr`, and `-perhost` options are intended for process distribution. If used simultaneously, `-machinefile` takes precedence.

-configfile <filename>

Use this option to specify the file *<filename>* that contains command-line options. Blank lines and lines that start with # as the first character are ignored. For example, the configuration file contains the following commands to run the executable files `a.out` and `b.out` using the `shm:dapl` fabric over `host1` and `host2` respectively:

```
-host host1 -env I_MPI_DEBUG 2 -env I_MPI_FABRICS shm:dapl -n 2 ./a.out
-host host2 -env I_MPI_DEBUG 2 -env I_MPI_FABRICS shm:dapl -n 2 ./b.out
```

To launch a MPI application according to the parameters above, use:

```
$ mpiexec -configfile <filename>
```

NOTE

This option may only be used alone. It terminates parsing of the `mpiexec` command line.

-g<l-option>

Use this option to apply the named local option *<l-option>* globally. See [Local Options](#) for a list of all local options. During the application startup, the default value is the `-genvuser` option.

NOTE

Local options have higher priority than global options:

- The `-genv` option has the highest priority.
 - The options `-genvlist`, `-genvexcl` have lower priority than the `-genv` option.
 - The options `-genvnone`, `-genvuser`, `-genvall` have the lowest priority.
-

-genv <ENVVAR> <value>

Use this option to set the *<ENVVAR>* environment variable to the specified *<value>* for all MPI processes.

-genvuser

Use this option to propagate all user environment variables to all MPI processes, with the exception of the following system environment variables: `$HOSTNAME`, `$HOST`, `$HOSTTYPE`, `$MACHTYPE`, `$OSTYPE`. This is the default setting.

-genvall

Use this option to enable propagation of all environment variables to all MPI processes.

-genvnone

Use this option to suppress propagation of any environment variables to any MPI processes.

(SDK only) -trace [<profiling_library>] or -t [<profiling_library>]

Use this option to profile your MPI application using the indicated *<profiling_library>*. If the *<profiling_library>* is not mentioned, the default profiling library `libVT.so` is used.

Set the `I_MPI_JOB_TRACE_LIBS` environment variable to override the default profiling library.

NOTE

It is not necessary to link your application against the profiling library before execution.

(SDK only) -check_mpi [<checking_library>]

Use this option to check your MPI application using the indicated *<checking_library>*. If *<checking_library>* is not mentioned, the default checking library `libVTmc.so` is used.

Set the `I_MPI_JOB_CHECK_LIBS` environment variable to override the default checking library.

NOTE

It is not necessary to link your application against the checking library before execution.

-tv

Use this option to run *<executable>* under the TotalView* debugger. For example:

```
$ mpiexec -tv -n <# of processes> <executable>
```

See [Environment Variables](#) for information on how to select the TotalView* executable file.

NOTE

Ensure the environment variable `TVDSVRLAUNCHCMD=ssh` because the TotalView* uses `rsh` by default.

NOTE

The TotalView* debugger has a feature to displays the message queue state of your MPI program. To use the state display feature, do the following steps:

1. Run your *<executable>* with `-tv` option.

```
$ mpiexec -tv -n <# of processes> <executable>
```
2. Answer **Yes** to the question about stopping the Python* job.

To display the internal state of the MPI library textually, select **Tools > Message Queue** command. If you select the **Process Window Tools > Message Queue Graph** command, the TotalView* displays a window that shows a graph of the current message queue state. For more information, see [TotalView*](#).

-tva <jobid>

Use this option to attach the TotalView* debugger to existing *<jobid>*. For example:

```
$ mpiexec -tva <jobid>
```

-tvsu

Use this option to run *<executable>* for later attachment with the TotalView* debugger. For example:

```
$ mpiexec -tvsu -n <# of processes> <executable>
```

NOTE

To debug the running Intel® MPI job, attach the TotalView* to the Python* instance that is running the `mpiexec` script.

-gdb

Use this option to run `<executable>` under the GNU* debugger. For example:

```
$ mpiexec -gdb -n <# of processes> <executable>
```

-gdba <jobid>

Use this option to attach the GNU* debugger to the existing `<jobid>`. For example:

```
$ mpiexec -gdba <jobid>
```

-a <alias>

Use this option to assign `<alias>` to the job.

-ordered-output

Use this option to avoid intermingling of data output by the MPI processes. This option affects both the standard output and standard error streams.

NOTE

For this option to work, the last line output by each process must end with the end-of-line (`\n`) character. Otherwise the application may stop responding.

-m

Use this option to merge output lines.

-l

Use this option to insert the MPI process rank at the beginning of all lines written to the standard output.

-s <spec>

Use this option to direct standard input to the specified MPI processes.

Arguments

<code><spec></code>	Define MPI process ranks
<code>all</code>	Use all processes
<code><l>, <m>, <n></code>	Specify an exact list and use processes <code><l></code> , <code><m></code> and <code><n></code> only. The default value is zero
<code><k>, <l>-<m>, <n></code>	Specify a range and use processes <code><k></code> , <code><l></code> through <code><m></code> , and <code><n></code>

-noconf

Use this option to disable processing of the `mpiexec` configuration files described in the section [Configuration Files](#).

-ifhn <interface/hostname>

Use this option to specify the network interface for communication with the local MPD daemon; where *<interface/hostname>* is an IP address or a hostname associated with the alternative network interface.

-ecfn <filename>

Use this option to list XML exit codes to the file *<filename>*.

Local Options**-n <# of processes> or -np <# of processes>**

Use this option to set the number of MPI processes to run with the current arg-set.

-env <ENVVAR> <value>

Use this option to set the *<ENVVAR>* environment variable to specified *<value>* for all MPI processes in the current arg-set.

-envuser

Use this option to propagate all user environment variables, with the exception of the following variables: \$HOSTNAME, \$HOST, \$HOSTTYPE, \$MACHTYPE, \$OSTYPE. This is the default setting.

-envall

Use this option to propagate all environment variables in the current environment.

-envnone

Use this option to suppress propagation of any environment variables to the MPI processes in the current arg-set.

-envlist <list of env var names>

Use this option to pass a list of environment variables with their current values. *<list of env var names>* is a comma separated list of environment variables to be sent to the processes. If this option is used several times in the command line, all variables listed in the arguments are included into one list.

-envexcl <list of env var names>

Use this option to suppress propagation of the listed environment variables to the MPI processes in the current arg-set.

-host <nodename>

Use this option to specify a particular *<nodename>* on which to run MPI processes in the current argument set. For example, the following command runs the executable `a.out` on host `host1` only:

```
$ mpiexec -n 2 -host host1 ./a.out
```

-path <directory>

Use this option to specify the path to the *<executable>* to run.

-wdir <directory>

Use this option to specify the working directory in which *<executable>* is to be run in the current arg-set.

-umask <umask>

Use this option to perform the `umask <umask>` command for the remote process.

Configuration Files

The `mpiexec` configuration files specify the default options applied to all `mpiexec` commands.

If any of these files exist, their contents are prefixed to the command-line options for `mpiexec` in the following order:

System-wide `<installdir>/etc/mpiexec.conf`. The default location of the configuration file is the `<installdir>/<arch>/etc`.

User-specific `$HOME/.mpiexec.conf`

Session-specific `$PWD/mpiexec.conf`

You can override these files by defining environment variables and using command line options. You can skip these configuration files by using the `mpiexec -noconf` option.

You can create or modify these files. They contain `mpiexec` command-line options. Blank lines and lines that start with `#` are ignored. For example, to specify a default fabric, add the following line to the respective `mpiexec.conf` file:

```
-genv I_MPI_FABRICS <fabric>
```

Environment Variables**I_MPI_DEBUG**

Print out debugging information when an MPI program starts running.

Syntax

```
I_MPI_DEBUG=<level>[,<flags>]
```

Arguments

<level>	Indicate level of debug information provided
0	Output no debugging information. This is the default value
1	Output verbose error diagnostics
2	Confirm which <code>I_MPI_FABRICS</code> was used and which Intel® MPI Library configuration was used.
3	Output effective MPI rank, <code>pid</code> and node mapping table
4	Output process pinning information
5	Output Intel MPI-specific environment variables
6	Output collective operation algorithms settings
> 6	Add extra levels of debug information

<flags>	Comma-separated list of debug flags
---------	-------------------------------------

pid	Show process id for each debug message
tid	Show thread id for each debug message for multithreaded library
time	Show time for each debug message
datetime	Show time and date for each debug message
host	Show host name for each debug message
level	Show level for each debug message
scope	Show scope for each debug message
line	Show source line number for each debug message
file	Show source file name for each debug message
nofunc	Do not show routine name
norank	Do not show rank
flock	Synchronize debug output from different process or threads
nobuf	Do not use buffered I/O for debug output

Description

Set this environment variable to control the output of the debugging information.

NOTE

Set the same `<level>` value for all ranks.

You can specify the output file name for debug information by setting the `I_MPI_DEBUG_OUTPUT` environment variable.

Each printed line has the following format:

```
[<identifier>] <message>
```

where

- `<identifier>` identifies the MPI process that produced the message. The `<identifier>` is an MPI process rank if `<level>` is an unsigned number. If the '+' sign is added in front of the `<level>` number, the `<identifier>` contains a `rank#pid@hostname` tuple. Here, `rank` is the MPI process rank; `pid` is the UNIX process id; and `hostname` is the host name as defined at process launch time.
- `<message>` contains the debugging output.

For example, the following command:

```
$ mpiexec -n 1 -env I_MPI_DEBUG 2 ./a.out
```

may produce the following output:

```
[0] MPI startup(): shared memory data transfer mode
```

while the command

```
$ mpiexec -n 1 -env I_MPI_DEBUG +2 ./a.out
```

or

```
$ mpiexec -n 1 -env I_MPI_DEBUG 2,pid,host ./a.out
```

may produce the following output:

```
[0#1986@mpiclust001] MPI startup(): shared memory data transfer mode
```

NOTE

Compiling with `mpiicc -g` adds a considerable amount of printed debug information.

I_MPI_DEBUG_OUTPUT

Set output file name for debug information.

Syntax

```
I_MPI_DEBUG_OUTPUT=<arg>
```

Arguments

<arg>	String value
stdout	Output to stdout - default value
stderr	Output to stderr
<file_name>	Specify the output file name for debug information (a maximum of 256 symbols)

Description

Set this environment variable if you want to split output of debug information from the output produced by an application. If you use format like `%r`, `%p` or `%h`, rank, pid or host name is added to the file name accordingly.

I_MPI_PERHOST

Define the default settings for the `-perhost` option in the `mpiexec` command.

Syntax

```
I_MPI_PERHOST=<value>
```

Arguments

<value>	Define the default process layout
<n> > 0	<n> processes per node
all	All logical CPUs on a node
allcores	All cores (physical CPUs) on a node

Description

Set this environment variable to define the default setting for the `-perhost` option. If `-perhost` is explicitly called in the command line, the `I_MPI_PERHOST` environment variable has no effect. The `-perhost` option assumes the value of the `I_MPI_PERHOST` environment variable if this environment variable is defined.

NOTE

When `I_MPI_PERHOST` is defined together with `mpiexec -host` option, `I_MPI_PERHOST` is ignored.

I_MPI_PRINT_VERSION

Print library version information.

Syntax

`I_MPI_PRINT_VERSION=<arg>`

Arguments

<code><arg></code>	Binary indicator
<code>enable yes on 1</code>	Print library version information.
<code>disable no off 0</code>	No action. This is the default value.

Description

Set this environment variable to enable/disable printing of Intel® MPI library version information when an MPI application starts running.

**(SDK only) I_MPI_JOB_TRACE_LIBS
(MPIEXEC_TRACE_LIBS)**

Choose the libraries to preload through the `-trace` option.

Syntax

`I_MPI_JOB_TRACE_LIBS=<arg>`

Deprecated Syntax

`MPIEXEC_TRACE_LIBS=<arg>`

Arguments

<code><arg></code>	String parameter
<code><list></code>	Blank separated list of libraries to preload. The default value is <code>vt</code>

Description

Set this environment variable to choose an alternative library for preloading by the `-trace` option.

(SDK only) I_MPI_JOB_CHECK_LIBS

Choose the libraries to preload through the `-check_mpi` option.

Syntax

`I_MPI_JOB_CHECK_LIBS=<arg>`

Arguments

<code><arg></code>	String parameter
<code><list></code>	Blank separated list of libraries to preload. The default value is <code>vtmc</code>

Description

Set this environment variable to choose an alternative library for preloading by the `-check_mpi` option.

I_MPI_JOB_STARTUP_TIMEOUT

Set the `mpiexec` job startup timeout.

Syntax

```
I_MPI_JOB_STARTUP_TIMEOUT=<timeout>
```

Arguments

<code><timeout></code>	Define <code>mpiexec</code> job startup timeout period in seconds
<code><n> >= 0</code>	The default timeout value is 20 seconds

Description

Set this environment variable to make `mpiexec` wait for the job to start in `<timeout>` seconds after its launch. The `<timeout>` value should be greater than zero. Otherwise the environment variable setting is ignored and a warning message is printed. Setting this environment variable may make sense on large clusters with a lot of nodes where the job startup time may exceed the default value.

NOTE

Set the `I_MPI_JOB_STARTUP_TIMEOUT` environment variable in the shell environment before executing the `mpiexec` command. Do not use the `-genv` or `-env` options for setting the `<timeout>` value. Those options are used only for passing environment variables to the MPI process environment.

I_MPI_JOB_TIMEOUT (MPIEXEC_TIMEOUT)

Set the `mpiexec` timeout.

Syntax

```
I_MPI_JOB_TIMEOUT=<timeout>
```

Deprecated Syntax

```
MPIEXEC_TIMEOUT=<timeout>
```

Arguments

<code><timeout></code>	Define <code>mpiexec</code> timeout period in seconds
<code><n> >= 0</code>	The default timeout value is zero, meaning no timeout

Description

Set this environment variable to make `mpiexec` terminate the job in `<timeout>` seconds after its launch. The `<timeout>` value should be greater than zero. Otherwise the environment variable setting is ignored.

NOTE

Set the `I_MPI_JOB_TIMEOUT` environment variable in the shell environment before executing the `mpiexec` command. Do not use the `-genv` or `-env` options for setting the `<timeout>` value. Those options are used only for passing environment variables to the MPI process environment.

I_MPI_JOB_TIMEOUT_SIGNAL (MPIEXEC_TIMEOUT_SIGNAL)

Define a signal to be used when a job is terminated because of a timeout.

Syntax

```
I_MPI_JOB_TIMEOUT_SIGNAL=<number>
```

Deprecated Syntax

```
MPIEXEC_TIMEOUT_SIGNAL=<number>
```

Arguments

<number>	Define signal number
<n> > 0	The default value is 9 (SIGKILL)

Description

Define a signal number for task termination upon the timeout period specified by the environment variable `I_MPI_JOB_TIMEOUT`. If you set a signal number unsupported by the system, `mpiexec` prints a warning message and continues task termination using the default signal number 9 (SIGKILL).

I_MPI_JOB_ABORT_SIGNAL

Define a signal to be sent to all processes when a job is terminated unexpectedly.

Syntax

```
I_MPI_JOB_ABORT_SIGNAL=<number>
```

Arguments

<number>	Define signal number
<n> > 0	The default value is 9 (SIGKILL)

Description

Set this environment variable to define a signal for task termination. If you set an unsupported signal number, `mpiexec` prints a warning message and uses the default signal 9 (SIGKILL).

I_MPI_JOB_SIGNAL_PROPAGATION (MPIEXEC_SIGNAL_PROPAGATION)

Control signal propagation.

Syntax

```
I_MPI_JOB_SIGNAL_PROPAGATION=<arg>
```

Deprecated Syntax

```
MPIEXEC_SIGNAL_PROPAGATION=<arg>
```

Arguments

<arg>	Binary indicator
enable yes on 1	Turn on propagation.
disable no off 0	Turn off propagation. This is the default value

Description

Set this environment variable to control propagation of the signals (`SIGINT`, `SIGALRM`, and `SIGTERM`) that may be received by the MPD daemons. If signal propagation is enabled, the received signal is sent to all processes of the MPI job. If signal propagation is disabled, all processes of the MPI job are stopped with the default signal 9 (`SIGKILL`).

I_MPI_OUTPUT_CHUNK_SIZE

Set the size of the `stdout/stderr` output buffer.

Syntax

`I_MPI_OUTPUT_CHUNK_SIZE=<size>`

Arguments

<code><size></code>	Define output chunk size in kilobytes
<code><n> > 0</code>	The default chunk size value is 1 KB

Description

Set this environment variable to increase the size of the buffer used to intercept the standard output and standard error streams from the processes. If the `<size>` value is not greater than zero, the environment variable setting is ignored and a warning message is displayed.

Use this setting for applications that create a significant amount of output from different processes. With the `-ordered-output mpiexec` option, this setting helps to prevent the output from garbling.

NOTE

Set the `I_MPI_OUTPUT_CHUNK_SIZE` environment variable in the shell environment before executing the `mpiexec` command. Do not use the `-genv` or `-env` options for setting the `<size>` value. Those options are used only for passing environment variables to the MPI process environment.

I_MPI_PMI_EXTENSIONS

Turn on/off the use of the Intel® MPI Library Process Management Interface (PMI) extensions.

Syntax

`I_MPI_PMI_EXTENSIONS=<arg>`

Arguments

<code><arg></code>	Binary indicator
<code>enable yes on 1</code>	Turn on the PMI extensions
<code>disable no off 0</code>	Turn off the PMI extensions

Description

The Intel® MPI Library automatically detects if your process manager supports the PMI extensions. If supported, the extensions substantially decrease task startup time. Set `I_MPI_PMI_EXTENSIONS` to `disable` if your process manager does not support these extensions.

I_MPI_PMI_LIBRARY

Specify the name to third party implementation of the PMI library.

Syntax

```
I_MPI_PMI_LIBRARY=<name>
```

Arguments

<name>	Full name of the third party PMI library
--------	--

Description

Set `I_MPI_PMI_LIBRARY` to specify the name of third party PMI library. When you set this environment variable, provide full name of the library with full path to it.

I_MPI_JOB_FAST_STARTUP (I_MPI_PMI_FAST_STARTUP)

Turn on/off the faster Intel® MPI Library process startup algorithm.

Syntax

```
I_MPI_JOB_FAST_STARTUP=<arg>
```

Deprecated Syntax

```
I_MPI_PMI_FAST_STARTUP=<arg>
```

Arguments

<arg>	Binary indicator
enable yes on 1	Turn on the algorithm for fast startup. This is the default value
disable no off 0	Turn off the algorithm for fast startup

Description

The new algorithm significantly decreases the application startup time. Some DAPL providers may be overloaded during startup of large number of processes (greater than 512). To avoid this problem, turn off this algorithm by setting the `I_MPI_JOB_FAST_STARTUP` environment variable to `disable`.

TOTALVIEW

Select a particular TotalView* executable file to use.

Syntax

```
TOTALVIEW=<path>
```

Arguments

<path>	Path/name of the TotalView* executable file instead of the default <code>totalview</code>
--------	---

Description

Set this environment variable to select a particular TotalView* executable file.

I_MPI_PLATFORM

Select the intended optimization platform.

Syntax

```
I_MPI_PLATFORM=<platform>
```

Arguments

<code><platform></code>	Intended optimization platform (string value)
<code>auto[:min]</code>	Optimize for the oldest supported Intel® Architecture Processor across all nodes. This is the default value
<code>auto:max</code>	Optimize for the newest supported Intel® Architecture Processor across all nodes
<code>auto:most</code>	Optimize for the most numerous Intel® Architecture Processor across all nodes. In case of a tie, choose the newer platform
<code>uniform</code>	Optimize locally. The behavior is unpredictable if the resulting selection differs from node to node
<code>none</code>	Select no specific optimization
<code>htn generic</code>	Optimize for the Intel® Xeon® Processors 5400 series and other Intel® Architecture processors formerly code named Harpertown
<code>nhm</code>	Optimize for the Intel® Xeon® Processors 5500, 6500, 7500 series and other Intel® Architecture processors formerly code named Nehalem
<code>wsm</code>	Optimize for the Intel® Xeon® Processors 5600, 3600 series and other Intel® Architecture processors formerly code named Westmere
<code>snb</code>	Optimize for the Intel® Xeon® Processors E3, E5, and E7 series and other Intel® Architecture processors formerly code named Sandy Bridge
<code>ivb</code>	Optimize for the Intel® Xeon® Processors E3, E5, and E7 V2 series and other Intel® Architecture processors formerly code named Ivy Bridge
<code>knc</code>	Optimize for the Intel® Xeon® Processors (codename: Knights Corner). If Intel Xeon Phi coprocessor is present on the cluster, the value is chosen by default
<code>hsw</code>	Optimize for the Intel® Xeon® Processors E3, E5, and E7 V3 series and other Intel® Architecture processors formerly code named Haswell
<code>knl</code>	Optimize for the Intel® Xeon Phi™ Processor x200 Product Family formerly code named Knights Landing

Description

Set this variable to use the predefined platform settings. It is available for both Intel® and non-Intel microprocessors, but it may utilize additional optimizations for Intel microprocessors than it utilizes for non-Intel microprocessors.

NOTE

The values `auto:min`, `auto:max` and `auto:most` may increase the MPI job startup time.

I_MPI_PLATFORM_CHECK

Turn on/off the optimization setting similarity check.

Syntax

```
I_MPI_PLATFORM_CHECK=<arg>
```

Argument

<arg>	Binary indicator
enable yes on 1	Turns on the optimization platform similarity check. This is the default value
disable no off 0	Turns off the optimization platform similarity check

Description

Set this variable to check the optimization platform settings of all processes for similarity. If the settings are not the same on all ranks, the library terminates the program. Disabling this check may reduce the MPI job startup time.

I_MPI_THREAD_LEVEL_DEFAULT

Set this environment variable to initialize the MPI thread environment for the multi-threaded library if `MPI_Init()` call is used for initialization.

Syntax

```
I_MPI_THREAD_LEVEL_DEFAULT=<threadlevel>
```

Arguments

<threadlevel>	Define the default level of thread support
SINGLE single	Set the default level of thread support to <code>MPI_THREAD_SINGLE</code>
FUNNELED funneled	Set the default level of thread support to <code>MPI_THREAD_FUNNELED</code> . This is the default value if <code>MPI_Init()</code> call is used for initialization
SERIALIZED serialized	Set the default level of thread support to <code>MPI_THREAD_SERIALIZED</code>
MULTIPLE multiple	Set the default level of thread support to <code>MPI_THREAD_MULTIPLE</code>

Description

Set `I_MPI_THREAD_LEVEL_DEFAULT` to define the default level of thread support for the multi-threaded library if `MPI_Init()` call is used for initialization.

NOTE

The environment variable `I_MPI_THREAD_LEVEL_DEFAULT` is equivalent to the environment variable `MPICH_THREADLEVEL_DEFAULT`.

2.6.2. Configuration Files

\$HOME/.mpd.conf

This optional configuration file contains an mpd daemon password. Create it before setting up the mpd daemons. Use it to control access to the daemons by various Intel® MPI Library users.

Syntax

The file has a single line:

```
secretword=<mpd password>
```

or

```
MPD_SECRETWORD=<mpd password>
```

Description

An arbitrary `<mpd password>` string only controls access to the `mpd` daemons by various cluster users. Do not use Linux* OS login passwords here.

Place the `$HOME/.mpd.conf` file on a network-mounted file system, or replicate this file so that it is accessible as `$HOME/.mpd.conf` on all nodes of the cluster.

When `mpdboot` is executed by some non-root `<user>`, this file should have user and ownership set to `<user>` and `<<user>'s group>` accordingly. The access permissions should be set to `600` mode (only user has read and write privileges).

NOTE

`MPD_SECRETWORD` is a synonym for `secretword`.

mpd.hosts

This file has a list of node names which the `mpdboot` command uses to start `mpd` daemons.

Ensure that this file is accessible by the user who runs `mpdboot` on the node where the `mpdboot` command is actually invoked.

Syntax

The format of the `mpd.hosts` file is a list of node names, one name per line. Blank lines and the portions of any lines that follow a `#` character are ignored.

2.6.3. Environment Variables

I_MPI_JOB_CONFIG_FILE (I_MPI_MPD_CONF)

Set the path/name of the `mpd` configuration file.

Syntax

```
I_MPI_JOB_CONFIG_FILE=<path/name>
```

Deprecated Syntax

```
I_MPI_MPD_CONF=<path/name>
```

Arguments

<code><path/name></code>	Absolute path of the MPD configuration file
--------------------------------	---

Description

Set this environment variable to define the absolute path of the file that is used by the `mpdbootscript` instead of the default value `${HOME}/.mpd.conf`.

I_MPI_JOB_CONTEXT (MPD_CON_EXT)

Set a unique name for the `mpd` console file. This enables you to run several `mpd` rings under the same user account.

Syntax

```
I_MPI_JOB_CONTEXT=<tag>
```

Deprecated Syntax

```
MPD_CON_EXT=<tag>
```

Arguments

<tag>	Unique MPD identifier
-------	-----------------------

Description

Set this environment variable to different unique values to allow several MPD rings to co-exist. Each MPD ring is associated with a separate `I_MPI_JOB_CONTEXT` value. Once this environment variable is set, you can start one MPD ring and work with it without affecting other available MPD rings. Set the appropriate `I_MPI_JOB_CONTEXT` value to work with a particular MPD ring. See [Simplified Job Startup Command](#) to learn about an easier way to run several Intel® MPI Library jobs at once.

I_MPI_JOB_TAGGED_PORT_OUTPUT

Turn on/off the use of the tagged mpd port output.

Syntax

```
I_MPI_JOB_TAGGED_PORT_OUTPUT=<arg>
```

Arguments

<arg>	Binary indicator
enable yes on 1	Turn on the tagged output. This is the default value
disable no off 0	Turn off the tagged output

Description

The tagged output format works at the `mpdboot` stage and prevents a failure during startup due to unexpected output from a remote shell like `ssh.mpdboot` sets this environment variable to 1 automatically. Set `I_MPI_JOB_TAGGED_PORT_OUTPUT` to `disable` if you do not want to use the new format.

I_MPI_MPD_CHECK_PYTHON

Toggle the Python* versions check at the MPD ring startup stage.

Syntax

```
I_MPI_MPD_CHECK_PYTHON=<arg>
```

Arguments

<arg>	Binary indicator
enable yes on 1	Check for Python version compatibility
disable no off 0	Do not check the Python version compatibility. This is the default value

Description

Set this environment variable to `enable` compatibility checking of Python versions installed on the cluster nodes. This may lead to increased MPD ring startup time. The MPD behavior is undefined if incompatible Python versions are installed on the cluster.

If `I_MPI_MPD_CHECK_PYTHON` is set to `enable` and the compatibility check fails, `mpdboot` exits abnormally and print a diagnostic message. An MPD ring is not started.

I_MPI_MPD_RSH

Set the remote shell to start `mpd` daemons.

Syntax

`I_MPI_MPD_RSH = <arg>`

Arguments

<code><arg></code>	String parameter
<code><remote shell></code>	The remote shell

Description

Set this environment variable to define the default setting for the `--rsh mpdboot` option. If `--rsh` is explicitly called in the command line, the `I_MPI_MPD_RSH` environment variable has no effect. If the `--rsh` option is not explicitly defined, it assumes the value of the `I_MPI_MPD_RSH` environment variable.

I_MPI_MPD_TMPDIR

Set a temporary directory for the MPD subsystem.

Syntax

`I_MPI_MPD_TMPDIR=<arg>`

`TMPDIR=<arg>`

Arguments

<code><arg></code>	String parameter
<code><directory name></code>	A string that points to a scratch space location. The default value is <code>/tmp</code>

Description

Set one of these environment variables to specify an alternative scratch space location. The MPD subsystem creates its own files in the directory specified by these environment variables. If both environment variables point to valid directories, the value of the `TMPDIR` environment variable is ignored.

NOTE

The `mpd2.console_*` file path length is limited in some operating systems. If you get the following diagnostic message: `socket.error: AF_UNIX path too long`, you need to decrease the length of the `<directory name>` string to avoid this issue.

NOTE

If `<arg>` points to a distributed file system (PANFS, PVFS, etc.), the `mpd` demons may not start. If this happens, set the `I_MPI_MPD_TMPDIR` and `TMPDIR` to point to a standard file system, such as `ext2`, `ext3`, or `NFS`.

I_MPI_MPD_CLEAN_LOG

Control the removal of the log file upon MPD termination.

Syntax

`I_MPI_MPD_CLEAN_LOG=<value>`

Arguments

<code><value></code>	Define the value
<code>enable yes on 1</code>	Remove the log file
<code>disable no off 0</code>	Keep the log file. This is the default value

Description

Set this environment variable to define the `mpdallexit` behavior. If you enable this environment variable, the `mpdallexit` removes the log file created during its execution. If you disable this environment variable, the `mpdallexit` keeps the log file.

2.7. Processor Information Utility

cpuinfo

The `cpuinfo` utility provides processor architecture information.

Syntax

`cpuinfo [[-]<options>]]`

Arguments

<code><options></code>	Sequence of one-letter options. Each option controls a specific part of the output data
<code>g</code>	General information about single cluster node shows: <ul style="list-style-type: none"> the processor product name the number of packages/sockets on the node core and threads numbers on the node and within each package SMT mode enabling
<code>i</code>	Logical processors identification table identifies threads, cores, and packages of each logical processor accordingly. <ul style="list-style-type: none"> <i>Processor</i> - logical processor number. <i>Thread Id</i> - unique processor identifier within a core. <i>Core Id</i> - unique core identifier within a package. <i>Package Id</i> - unique package identifier within a node.
<code>d</code>	Node decomposition table shows the node contents. Each entry contains the information on packages, cores, and logical processors. <ul style="list-style-type: none"> <i>Package Id</i> - physical package identifier.

	<ul style="list-style-type: none"> • <i>Cores Id</i> - list of core identifiers that belong to this package. • <i>Processors Id</i> - list of processors that belong to this package. This list order directly corresponds to the core list. A group of processors enclosed in brackets belongs to one core.
c	<p>Cache sharing by logical processors shows information of sizes and processors groups, which share particular cache level.</p> <ul style="list-style-type: none"> • Size - cache size in bytes. • Processors - a list of processor groups enclosed in the parentheses those share this cache or no sharing otherwise.
s	<p>Microprocessor signature hexadecimal fields (Intel platform notation) show signature values:</p> <ul style="list-style-type: none"> • extended family • extended model • family • model • type • stepping
f	<p>Microprocessor feature flags indicate what features the microprocessor supports. The Intel platform notation is used.</p>
A	Equivalent to <code>gidcsf</code>
<code>gidc</code>	Default sequence
?	Utility usage info

Description

The `cpuinfo` utility prints out the processor architecture information that can be used to define suitable process pinning settings. The output consists of a number of tables. Each table corresponds to one of the single options listed in the arguments table.

NOTE

The architecture information is available on systems based on the Intel® 64 architectures.

The `cpuinfo` utility is available for both Intel microprocessors and non-Intel microprocessors, but it may provide only partial information about non-Intel microprocessors.

Examples

`cpuinfo` output for the processor of Intel® Xeon® E5-2697 v2:

```
$ cpuinfo A
```

Intel(R) processor family information utility, Version 5.1 Build 20150225 (build id: 11316)
 Copyright (C) 2005-2015 Intel Corporation. All rights reserved.

===== Processor composition =====

Processor name : Intel(R) Xeon(R) E5-2697 v2
 Packages(sockets) : 2
 Cores : 24
 Processors(CPUs) : 48
 Cores per package : 12
 Threads per core : 2

===== Processor identification =====

Processor	Thread Id.	Core Id.	Package Id.
0	0	0	0
1	0	1	0
2	0	2	0
3	0	3	0
4	0	4	0
5	0	5	0
6	0	8	0
7	0	9	0
8	0	10	0
9	0	11	0
10	0	12	0
11	0	13	0
12	0	0	1
13	0	1	1
14	0	2	1
15	0	3	1
16	0	4	1
17	0	5	1
18	0	8	1
19	0	9	1
20	0	10	1
21	0	11	1
22	0	12	1
23	0	13	1
24	1	0	0
25	1	1	0
26	1	2	0
27	1	3	0
28	1	4	0
29	1	5	0
30	1	8	0
31	1	9	0
32	1	10	0
33	1	11	0
34	1	12	0
35	1	13	0
36	1	0	1
37	1	1	1
38	1	2	1
39	1	3	1
40	1	4	1
41	1	5	1
42	1	8	1
43	1	9	1
44	1	10	1
45	1	11	1
46	1	12	1
47	1	13	1

===== Placement on packages =====

Package Id.	Core Id.	Processors
0	0,1,2,3,4,5,8,9,10,11,12,13	(0,24) (1,25) (2,26) (3,27) (4,28) (5,29) (6,30) (7,31) (8,32) (9,33) (10,34) (11,35)
1	0,1,2,3,4,5,8,9,10,11,12,13	(12,36) (13,37) (14,38) (15,39) (16,40) (17,41) (18,42) (19,43) (20,44) (21,45) (22,46) (23,47)

===== Cache sharing =====

Cache	Size	Processors
L1	32 KB	(0,24) (1,25) (2,26) (3,27) (4,28) (5,29) (6,30) (7,31) (8,32) (9,33) (10,34) (11,35) (12,36) (13,37) (14,38) (15,39) (16,40) (17,41) (18,42) (19,43) (20,44) (21,45) (22,46) (23,47)
L2	256 KB	(0,24) (1,25) (2,26) (3,27) (4,28) (5,29) (6,30) (7,31) (8,32) (9,33) (10,34) (11,35) (12,36) (13,37) (14,38) (15,39) (16,40) (17,41) (18,42) (19,43) (20,44) (21,45) (22,46) (23,47)
L3	30 MB	(0,1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,20,21,22,23,24,25,26,27,28,29,30,31,32,33,34,35) (12,13,14,15,16,17,18,19,20,21,22,23,24,25,26,27,28,29,30,31,32,33,34,35)

===== Processor Signature =====

xFamily	xModel	Type	Family	Model	Stepping
00	3	0	6	e	4

===== Processor Feature Flags =====

SSE3	PCLMULQ	DTES64	MONITOR	DS-CPL	VMX	SMX	EIST	TM2	SSSE3	CNXT-ID	FMA	CX16	xTFR
1	1	1	1	1	1	1	1	1	1	0	0	1	1

3. Tuning Reference

The Intel® MPI Library provides an automatic tuning utility to help you select optimal values for many environment variables that can be used to influence program behavior and performance at run time.

3.1. Using mpitune Utility

mpitune

Use the `mpitune` utility to find optimal settings for the Intel® MPI Library relevant to your cluster configuration or your application.

Syntax

```
mpitune [ -a "<application command line>" ] [ -of <file-name> ] \  
[ -t "<test_cmd_line>" ] [ -cm ] [ -d ] [ -D ] \  
[ -dl [d1[,d2...[,dN]]] ] [ -fl [f1[,f2...[,fN]]] ] [ -er ] \  
[ -hf <hostsfile> ] [ -h ] [ -hr {min:max|min:|:max} ] \  
[ -i <count> ] [ -mr {min:max|min:|:max} ] [ -od <outputdir> ] \  
[ -odr <outputdir> ] [ -r <rshcmd> ] [ -pr {min:max|min:|:max} ] \  
[ -sf [file-path] ] [ -ss ] [ -s ] [ -td <dir-path> ] \  
[ -tl <minutes> ] [ -mh ] [ -os <opt1,...,optN> ] \  
[ -oe <opt1,...,optN> ] [ -V ] [ -vi {percent} | -vix {X factor} ] \  
[ -zb ] [ -t ] [ -so ] [ -ar "reg-expr" ] [ -trf <appoutfile> ] \  
[ -m {base|optimized} ] [ -avd {min|max} ] [ -pm {mpd|hydra} ] \  
[ -co ] [ -sd ] [ -soc ]
```

or

```
mpitune [ --application "<app_cmd_line>" ] [ --output-file <file-name> ] \  
[ --test "<test_cmd_line>" ] [ --cluster-mode ] [ --debug ] \  
[ --distinct ] [ --device-list [d1[,d2,...[,dN]]] ] \  
[ --fabric-list [f1[,f2...[,fN]]] ] [ --existing-ring ] \  
[ --host-file <hostsfile> ] [ --help ] \  
[ --host-range {min:max|min:|:max} ] [ --iterations <count> ] \  
[ --message-range {min:max|min:|:max} ] \  
[ --output-directory <outputdir> ] \  
[ --output-directory-results <outputdir> ] [ --rsh <rshcmd> ] \  
[ --ppn-range {min:max|min:|:max} |  
--perhost-range {min:max|min:|:max} ] \  
[ --session-file [file-path] ] [ --show-session ] [ --silent ] \  
[ --temp-directory <dir-path> ] [ --time-limit <minutes> ] \  
[ --master-host ] [ --options-set <opt1,...,optN> ] \  
[ --options-exclude <opt1,...,optN> ] [ --version ] \  
[ --valuable-improvement | --valuable-improvement-x {X factor} ] \  
[ --zero-based ] [ --trace ] [ --scheduler-only ] \  
[ --application-regexp "reg-expr" ]
```

```
[ --test-regexp-file <appoutfile> ] [ --model {base|optimized} ] \
[ --application-value-direction {min|max} ] \
[ --process-manager {mpd|hydra} ] [ -co ] [ -sd ] [ -soc ]
```

Arguments

<code>-a \"<app_cmd_line>\" --application \"<app_cmd_line>\"</code>	Switch on the application-specific mode. Quote the full command line as shown, including the backslashes.
<code>-of <file-name> --output-file <file-name></code>	Specify the name of the application configuration file to be generated in the application-specific mode. By default, use the file name <code>\$PWD/app.conf</code> .
<code>-t \"<test_cmd_line>\" --test \"<test_cmd_line>\"</code>	Replace the default Intel® MPI Benchmarks by the indicated benchmarking program in the cluster-specific mode. Quote the full command line as shown including the backslashes.
<code>-cm {exclusive full} --cluster-mode {exclusive full}</code>	Set the cluster usage mode <ul style="list-style-type: none"> <code>full</code> - maximum number of tasks are executed. This is the default mode. <code>exclusive</code> - only one task is executed on the cluster at a time.
<code>-d --debug</code>	Print out the debug information.
<code>-D --distinct</code>	Tune all options separately from each other. This argument is applicable only for the cluster-specific mode.
<code>-dl [d1[,d2...[,dN]] --device-list [d1[,d2,...[,dN]]]</code>	Select the device(s) you want to tune. Any previously set fabrics are ignored.. By default, use all devices listed in the <code><installdir>/<arch>/etc/devices.xml</code> file.
<code>-fl [f1[,f2...[,fN]] --fabric-list [f1[,f2...[,fN]]]</code>	Select the fabric(s) you want to tune. Any previously set devices are ignored. By default, use all fabrics listed in the <code><installdir>/<arch>/etc/fabrics.xml</code> file.
<code>-er --existing-ring</code>	Use an existing MPD ring. By default, a new MPD ring is created. This argument is applicable only if <code>I_MPI_PROCESS_MANAGER</code> is set to <code>mpd</code> .
<code>-hf <hostsfile> --host-file <hostsfile></code>	Specify an alternative host file name. By default, use the <code>\$PWD/mpd.hosts</code> .
<code>-h --help</code>	Display the help message.
<code>-hr {min:max min: :max} --host-range {min:max min: :max}</code>	Set the range of hosts used for testing. The default minimum value is 1. The default maximum value is the number of hosts defined by the <code>mpd.hosts</code> or the existing MPD ring. The <code>min:</code> or <code>:max</code> format uses the default values as appropriate.
<code>-i <count> </code>	Define how many times to run each tuning step. Higher iteration counts

<code>--iterations <count></code>	increase the tuning time, but may also increase the accuracy of the results. The default value is 3.
<code>-mr {min:max min: :max}</code> <code>--message-range {min:max min: :max}</code>	Set the message size range. The default minimum value is 0. The default maximum value is 4194304 (4mb). By default, the values are given in bytes. They can also be given in the following format: 16kb, 8mb or 2gb. The min: or :max format uses the default values as appropriate.
<code>-od <outputdir> </code> <code>--output-directory <outputdir></code>	Specify the directory name for all output files: log-files, session-files, local host-files and report-files. By default, use the current directory. This directory should be accessible from all hosts.
<code>-odr <outputdir> </code> <code>--output-directory-results <outputdir></code>	Specify the directory name for the resulting configuration files. By default, use the current directory in the application-specific mode and the <installdir>/<arch>/etc in the cluster-specific mode. If <installdir>/<arch>/etc is unavailable, \$PWD is used as the default value in the cluster-specific mode.
<code>-r <rshcmd> --rsh <rshcmd></code>	Specify the remote shell used to start daemons (as applicable) and jobs. The default value is ssh.
<code>-pr {min:max min: :max}</code> <code>--ppn-range {min:max min: :max}</code> <code>--perhost-range {min:max min: :max}</code>	Set the maximum number of processes per host. The default minimum value is 1. The default maximum value is the number of cores of the processor. The min: or :max format uses the default values as appropriate.
<code>-sf [file-path] </code> <code>--session-file [file-path]</code>	Continue the tuning process starting from the state saved in the file-path session file.
<code>-ss </code> <code>--show-session</code>	Show information about the session file and exit. This option works only jointly with the -sf option.
<code>-s --silent</code>	Suppress all diagnostics.
<code>-td <dir-path> </code> <code>--temp-directory <dir-path></code>	Specify a directory name for the temporary data. Use \$PWD/mpitunertemp by default. This directory should be accessible from all hosts.
<code>-tl <minutes> </code> <code>--time-limit <minutes></code>	Set mpitune execution time limit in minutes. The default value is 0, which means no limitations.
<code>-mh </code> <code>--master-host</code>	Dedicate a single host to run the mpitune.
<code>-os <opt1,...,optN> </code> <code>--options-set <opt1,...,optN></code>	Use mpitune to tune the only required options you have set in the option values
<code>-oe <opt1,...,optN> </code> <code>--options-exclude</code>	Exclude the settings of the indicated Intel® MPI Library options from the tuning process.

<opt1,...,optN>	
-V --version	Print out the version information.
-vi {percent} > --valuable-improvement {percent} -vix {X factor} --valuable-improvement- x {X factor}	Control the threshold for performance improvement. The default threshold is 3%.
-zb --zero-based	Set zero as the base for all options before tuning. This argument is applicable only for the cluster-specific mode.
-t --trace	Print out error information such as error codes and tuner trace back.
-so --scheduler-only	Create the list of tasks to be executed, display the tasks, and terminate execution.
-ar "reg-expr" --application-regexp "reg-expr"	Use <code>reg-expr</code> to determine the performance expectations of the application. This option is applicable only for the application-specific mode. The <code>reg-expr</code> setting should contain only one group of numeric values which is used by <code>mpitune</code> for analysis. Use backslash for symbols when setting the value of this argument in accordance with the operating system requirements.
-trf <appoutfile> --test-regexp-file <appoutfile>	Use a test output file to check the correctness of the regular expression. This argument is applicable only for the cluster-specific mode when you use the <code>-ar</code> option.
-m {base optimized} --model {base optimized}	Specify the search model: <ul style="list-style-type: none"> • Set <code>base</code> to use the old model. • Set <code>optimized</code> to use the new faster search model. This is the default value.
-avd {min max} --application-value- direction {min max}	Specify the direction of the value optimization : <ul style="list-style-type: none"> • Set <code>min</code> to specify that lower is better. For example, use this value when optimizing the wall time. • Set <code>max</code> to specify that higher is better. For example, use this value when optimizing the solver ratio.
-pm {mpd hydra} --process-manager {mpd hydra}	Specify the process manager used to run the benchmarks. The default value is <code>hydra</code> .
-co --collectives- only	Tune collective operations only.
-sd --save-defaults	Use <code>mpitune</code> to save the default values of the Intel® MPI Library options.

<code>-soc --skip-options-check</code>	Specify whether to check the command line options.
--	--

Deprecated Options

Deprecated Option	New Option
<code>--outdir</code>	<code>-od --output-directory</code>
<code>--verbose</code>	<code>-d --debug</code>
<code>--file</code>	<code>-hf --host-file</code>
<code>--logs</code>	<code>-lf --log-file</code>
<code>--app</code>	<code>-a --application</code>

Description

Use the `mpitune` utility to create a set of Intel® MPI Library configuration files that contain optimal settings for a particular cluster or application. You can reuse these configuration files in the `mpirun` job launcher by using the `-tune` option. If configuration files from previous `mpitune` sessions exist, `mpitune` creates a copy of the existing files before starting execution.

The MPI tuner utility operates in two modes:

- Cluster-specific, evaluating a given cluster environment using either the Intel® MPI Benchmarks or a user-provided benchmarking program to find the most suitable configuration of the Intel® MPI Library. This mode is used by default.
- Application-specific, evaluating the performance of a given MPI application to find the best configuration for the Intel® MPI Library for the particular application. Application tuning is enabled by the `--application` command line option.

3.1.1. Cluster Specific Tuning

To find the optimal settings for tuning your cluster, run the `mpitune` utility once after the Intel® MPI Library installation and after every cluster configuration change (processor or memory upgrade, network re-configuration, and so on.). To get the list of settings, run the utility under the user account that was used for the Intel® MPI Library installation, or appropriately set the tuner data directory through the `--output-directory` option and the results directory through the `--output-directory-results` option.

If there are any configuration files in the `<installdir>/<arch>/etc` directory, the recorded Intel® MPI Library configuration settings are used automatically by `mpirun` with the `-tune` option.

For example:

- Collect configuration settings for the cluster hosts listed in the `./mpd.hosts` file by using the Intel® MPI Benchmarks

```
$ mpitune
```

- Use the optimal recorded values when running on the cluster

```
$ mpirun -tune -n 32 ./myprog
```

The job launcher finds a proper set of configuration options based on the following execution conditions: communication fabrics, number of hosts and processes, etc. If you have write access permission for `<installdir>/<arch>/etc`, all generated files are saved in this directory; otherwise the current working directory is used.

NOTE

When you use the `-tune` option in the cluster specific mode (such as, without the tuning configuration file name), you need to explicitly select the communication device or fabric, the number of processes per node, and the total number of processes. For example:

```
$ mpirun -tune -genv I_MPI_FABRICS shm:dapl -ppn 8 -n 32 ./myprog
```

Replacing the Default Benchmark

This tuning feature is an extension of the cluster-specific tuning mode in which you specify a benchmarking application that is used for tuning.

The Intel® MPI Benchmarks executable files, which are more optimized for Intel microprocessors than for non-Intel microprocessors, are used by default. This may result in different tuning settings on Intel microprocessors than on non-Intel microprocessors.

For example:

1. Collect the configuration settings for the cluster hosts listed in the `./mpd.hosts` file by using the desired benchmarking program

```
$ mpitune --test \"benchmark -param1 -param2\"
```

2. Use the optimal recorded values for your cluster

```
$ mpirun -tune -n 32 ./myprog
```

3.1.2. Application Specific Tuning

Run the tuning process for any MPI application by specifying its command line to the tuner. Performance is measured as inversed execution time of the given application. To reduce the overall tuning time, use the shortest representative application workload that is applicable to the configuration (fabric, rank placement, and so on.).

NOTE

In the application specific mode, you can achieve the best tuning results using a similar command line and environment.

For example:

Collect configuration settings for the given application

```
$ mpitune --application \"mpirun -n 32 ./myprog\" -of ./myprog.conf
```

Use the optimal recorded values for your application

```
$ mpirun -tune ./myprog.conf -n 32 ./myprog
```

Based on the default tuning rules, the automated tuning utility evaluates a full set of the library configuration parameters to minimize the application execution time. By default, all generated files are saved in the current working directory.

The resulting application configuration file contains the optimal Intel® MPI Library parameters for this application and configuration only. To tune the Intel® MPI Library for the same application in a different configuration (number of hosts, workload, and so on.), rerun the automated tuning utility with the desired configuration.

NOTE

By default, the automated tuning utility overwrites the existing application configuration files. If you want to keep various application and configuration files, you should use a naming convention to save the different versions and select the correct file when you need it.

Fast Application Tuning

This topic describes how to use fast application tuning to find optimal settings for the Intel® MPI Library.

Syntax

```
--fast [<value>] or -f [<value>]
```

or

```
I_MPI_TUNE_FAST=<value>
```

Arguments

<value>	Binary indicator
enable yes on 1	Enable the fast application tuning. This value is only applicable for the application tuning mode
disable no off 0	Disable the fast application tuning. This is the default value

Description

If you set `I_MPI_TUNE_FAST` to `enable`, the `mpitune` utility executes alternate fast application tuning procedure. The fast application tuning uses the same resulting configuration file used in previous tuning.

Example 1

```
$ mpitune --application \"mpirun ...\" --fast
```

Example 2

```
$ export I_MPI_TUNE_FAST=enable
```

```
$ mpitune --application \"mpirun ...\"
```

NOTE

If you use the `--help` option and the `--fast` option together, the help message of `app_tune` is printed.

Topology Awareness Application Tuning

This section describes how to use the `mpitune` utility to perform the topology awareness application tuning. If you want to perform this tuning in the dynamic method, see the description of [-use-app-topology](#) and `I_MPI_HYDRA_USE_APP_TOPOLOGY`.

I_MPI_TUNE_RANK_PLACEMENT

Syntax

```
--rank-placement [<value>] or -rp [<value>]
```

or

```
I_MPI_TUNE_RANK_PLACEMENT=<value>
```

Arguments

<value>	Binary indicator
enable yes on 1	Switch the <code>mpitune</code> utility to the topology tuning tool. This value is only applicable for the application tuning mode. For example, when you define the <code>--application</code> option or application communication graph (ACG) and optionally hardware topology graph (HTG) are passed by additional options <code>-acg</code> and <code>-htg</code>

disable no off 0	Switch off the topology tuning tool. This is the default value
------------------------------	--

Description

If you set `I_MPI_TUNE_RANK_PLACEMENT` to enable, the `mpitune` wrapper executes alternate topology tool (`mpitune_rank_placement`).

Example

```
$ mpitune --application \"mpirun ... \" --rank-placement
$ mpitune --application \"mpirun ... \" --rank-placement enable
$ mpitune --application \"mpirun ... \" -rp -acg <path to acg_file> -htg <path to htg_file>
```

The result is a host file and a configuration file with the record for automatic usage of the host file.

NOTE

If you use the `--help` option and the `--rank-placement` option together, the help message of `mpitune_rank_placement` is printed.

I_MPI_TUNE_APPLICATION_STATISTICS**Syntax**

```
--application-statistics [<value>] or -s [<value>]
```

or

```
I_MPI_TUNE_APPLICATION_STATICSTICS=<value>
```

Arguments

<value>	Path to the native Intel MPI statistics file level 1 or higher. This setting is only applicable together with <code>--rank-placement</code> option
---------	--

Description

To pass the Intel MPI statistics file to `mpitune` (`mpitune_rank_placement`), it reduces the tuning time.

Example

```
$ mpitune -rp -s <path to statistics file>
```

I_MPI_TUNE_APPLICATION_COMMUNICATION_GRAPH**Syntax**

```
--application-communication-graph [<value>] or -acg [<value>]
```

or

```
I_MPI_TUNE_APPLICATION_COMMUNICATION_GRAPH=<value>
```

Arguments

<value>	Path to the ACG file. This setting is only applicable together with <code>--rank-placement</code> option
---------	--

Description

To pass the ACG file to `mpitune` (`mpitune_rank_placement`), it reduces the tuning time.

Example

```
$ mpitune -rp -acg <path to acg_file>
```

I_MPI_TUNE_HARDWARE_TOPOLOGY_GRAPH

Syntax

```
--hardware-topology-graph [<value>] or -htg [<value>]
```

or

```
I_MPI_TUNE_HARDWARE_TOPOLOGY_GRAPH=<value>
```

Arguments

<value>	Path to the file with the description of hardware topology graph. This setting is only applicable together with <code>--rank-placement</code> option
---------	--

Description

To pass the HTG file to `mpitune` (`mpitune_rank_placement`), it reduces the tuning time.

Example

```
$ mpitune -rp -acg <path to acg_file> -htg <path to htg_file>
```

3.1.3. Tuning Utility Output

Upon completion of the tuning process, the Intel® MPI Library tuning utility records the chosen values in the configuration file in the following format:

```
-genv I_MPI_DYNAMIC_CONNECTION 1
-genv I_MPI_ADJUST_REDUCE 1:0-8
```

The Intel MPI Library tuning utility ignores any environment variables that have no effect on the application when the difference between probes is at the noise level (1%). In this case, the utility does not set the environment variable and preserves the default library heuristics.

In the case of an tuning application that has significant run-to-run performance variation, the Intel MPI Library tuning utility might select divergent values for the same environment variable under the same conditions. To improve decision accuracy, increase the number of iterations for each test run with the `--iterations` command line option. The default value for the number of iterations is 3.

3.2. Process Pinning

Use this feature to pin a particular MPI process to a corresponding CPU within a node and avoid undesired process migration. This feature is available on operating systems that provide the necessary kernel interfaces.

3.2.1. Default Settings for Process Pinning

If you do not specify values for any process pinning environment variables, the following default settings are used. For more details about these settings, see [Environment Variables](#) and [Interoperability with OpenMP API](#).

- `I_MPI_PIN=on`
- `I_MPI_PIN_MODE=pm`
- `I_MPI_PIN_RESPECT_CPUSET=on`
- `I_MPI_PIN_RESPECT_HCA=on`
- `I_MPI_PIN_CELL=unit`
- `I_MPI_PIN_DOMAIN=auto:compact`
- `I_MPI_PIN_ORDER=compact`

3.2.2. Processor Identification

The following schemes are used to identify logical processors in a system:

- System-defined logical enumeration
- Topological enumeration based on three-level hierarchical identification through triplets (package/socket, core, thread)

The number of a logical CPU is defined as the corresponding position of this CPU bit in the kernel affinity bit-mask. Use the `cpuinfo` utility, provided with your Intel MPI Library installation, or the `cat /proc/cpuinfo` command to find out the logical CPU numbers.

The three-level hierarchical identification uses triplets that provide information about processor location and their order. The triplets are hierarchically ordered (package, core, and thread).

See the example for one possible processor numbering scenario with two sockets, four cores (two cores per socket), and eight logical processors (two processors per core).

NOTE

Logical and topological enumerations are not the same.

Table 3.2-1 Logical Enumeration

0	4	1	5	2	6	3	7
---	---	---	---	---	---	---	---

Table 3.2-2 Hierarchical Levels

Socket	0	0	0	0	1	1	1	1
Core	0	0	1	1	0	0	1	1
Thread	0	1	0	1	0	1	0	1

Table 3.2-3 Topological Enumeration

0	1	2	3	4	5	6	7
---	---	---	---	---	---	---	---

Use the `cpuinfo` utility to identify the correspondence between the logical and topological enumerations. See [Processor Information Utility](#) for more details.

3.2.3. Environment Variables

I_MPI_PIN

Turn on/off process pinning.

Syntax

`I_MPI_PIN=<arg>`

Arguments

<code><arg></code>	Binary indicator
<code>enable yes on 1</code>	Enable process pinning. This is the default value
<code>disable no off 0</code>	Disable processes pinning

Description

Set this environment variable to control the process pinning feature of the Intel® MPI Library.

I_MPI_PIN_MODE

Choose the pinning method.

Syntax

`I_MPI_PIN_MODE=<pinmode>`

Arguments

<code><pinmode></code>	Choose the CPU pinning mode
<code>mpd pm</code>	Pin processes inside the process manager involved (Multipurpose Daemon*/MPD or Hydra*). This is the default value
<code>lib</code>	Pin processes inside the Intel MPI Library

Description

Set the `I_MPI_PIN_MODE` environment variable to choose the pinning method. This environment variable is valid only if the `I_MPI_PIN` environment variable is enabled.

If you set the `I_MPI_PIN_MODE` environment variable to `mpd|pm` to make the `mpd` daemon or the Hydra process launcher pin processes through system specific means, if they are available. The pinning is done before the MPI process launch. Therefore, it is possible to co-locate the process CPU and memory. This pinning method has an advantage over a system with Non-Uniform Memory Architecture (NUMA) like SGI* Altix*. Under NUMA, a processor can access its own local memory faster than non-local memory.

If you set the `I_MPI_PIN_MODE` environment variable to `lib`, the Intel® MPI Library pins the processes. This mode does not offer the capability to co-locate the CPU and memory for a process.

I_MPI_PIN_PROCESSOR_LIST

(`I_MPI_PIN_PROCS`)

Define a processor subset and the mapping rules for MPI processes within this subset.

Syntax

`I_MPI_PIN_PROCESSOR_LIST=<value>`

The environment variable value has the following syntax forms:

1. `<proclist>`

2.

`[<procset>][:[grain=<grain>][,shift=<shift>][,preoffset=<preoffset>][,postoffset=<postoffset>]]`

3. `[<procset>][:map=<map>]`

The following paragraphs provide detail descriptions for the values of these syntax forms.

Deprecated Syntax

`I_MPI_PIN_PROCS=<proclist>`

NOTE

The `postoffset` keyword has `offset` alias.

NOTE

The second form of the pinning procedure has three steps:

1. Cyclic shift of the source processor list on `preoffset*grain` value.
2. Round robin shift of the list derived on the first step on `shift*grain` value.
3. Cyclic shift of the list derived on the second step on the `postoffset*grain` value.

NOTE

The `grain`, `shift`, `preoffset`, and `postoffset` parameters have a unified definition style.

This environment variable is available for both Intel® and non-Intel microprocessors, but it may perform additional optimizations for Intel microprocessors than it performs for non-Intel microprocessors.

Syntax

```
I_MPI_PIN_PROCESSOR_LIST=<proclist>
```

Arguments

<proclist>	A comma-separated list of logical processor numbers and/or ranges of processors. The process with the i-th rank is pinned to the i-th processor in the list. The number should not exceed the amount of processors on a node.
<l>	Processor with logical number <l>.
<l>-<m>	Range of processors with logical numbers from <l> to <m>.
<k>,<l>-<m>	Processors <k>, as well as <l> through <m>.

Syntax

```
I_MPI_PIN_PROCESSOR_LIST=[<procset>][:[grain=<grain>][,shift=<shift>][,preoffset=<preoffset>][,postoffset=<postoffset>]]
```

Arguments

<procset>	Specify a processor subset based on the topological numeration. The default value is <code>allcores</code> .
<code>all</code>	All logical processors. Specify this subset to define the number of CPUs on a node.
<code>allcores</code>	All cores (physical CPUs). Specify this subset to define the number of cores on a node. This is the default value. If Intel® Hyper-Threading Technology is disabled, <code>allcores</code> equals to <code>all</code> .
<code>allsockets</code>	All packages/sockets. Specify this subset to define the number of sockets on a node.

<grain>	Specify the pinning granularity cell for a defined <procset>. The minimal <grain> value is a single element of the <procset>. The maximal <grain> value is the number of <procset> elements in a socket. The <grain> value must be a multiple of the <procset> value. Otherwise, the minimal <grain> value is assumed. The default value is the minimal <grain> value.
---------	--

<code><shift></code>	Specify the granularity of the round robin scheduling shift of the cells for the <code><procset></code> . <code><shift></code> is measured in the defined <code><grain></code> units. The <code><shift></code> value must be positive integer. Otherwise, no shift is performed. The default value is no shift, which is equal to 1 normal increment
<code><preoffset></code>	Specify the cyclic shift of the processor subset <code><procset></code> defined before the round robin shifting on the <code><preoffset></code> value. The value is measured in the defined <code><grain></code> units. The <code><preoffset></code> value must be non-negative integer. Otherwise, no shift is performed. The default value is no shift.
<code><postoffset></code>	Specify the cyclic shift of the processor subset <code><procset></code> derived after round robin shifting on the <code><postoffset></code> value. The value is measured in the defined <code><grain></code> units. The <code><postoffset></code> value must be non-negative integer. Otherwise no shift is performed. The default value is no shift.

The following table displays the values for `<grain>`, `<shift>`, `<preoffset>`, and `<postoffset>` options:

<code><n></code>	Specify an explicit value of the corresponding parameters. <code><n></code> is non-negative integer.
<code>fine</code>	Specify the minimal value of the corresponding parameter.
<code>core</code>	Specify the parameter value equal to the amount of the corresponding parameter units contained in one core.
<code>cache1</code>	Specify the parameter value equal to the amount of the corresponding parameter units that share an L1 cache.
<code>cache2</code>	Specify the parameter value equal to the amount of the corresponding parameter units that share an L2 cache.
<code>cache3</code>	Specify the parameter value equal to the amount of the corresponding parameter units that share an L3 cache.
<code>cache</code>	The largest value among <code>cache1</code> , <code>cache2</code> , and <code>cache3</code> .
<code>socket sock</code>	Specify the parameter value equal to the amount of the corresponding parameter units contained in one physical package/socket.
<code>half mid</code>	Specify the parameter value equal to <code>socket/2</code> .
<code>third</code>	Specify the parameter value equal to <code>socket/3</code> .
<code>quarter</code>	Specify the parameter value equal to <code>socket/4</code> .
<code>octavo</code>	Specify the parameter value equal to <code>socket/8</code> .

Syntax

```
I_MPI_PIN_PROCESSOR_LIST=[<procset>] [:map=<map>]
```

Arguments

<code><map></code>	The mapping pattern used for process placement.
<code>bunch</code>	The processes are mapped as close as possible on the sockets.

scatter	The processes are mapped as remotely as possible so as not to share common resources: FSB, caches, and core.
spread	The processes are mapped consecutively with the possibility not to share common resources.

Description

Set the `I_MPI_PIN_PROCESSOR_LIST` environment variable to define the processor placement. To avoid conflicts with different shell versions, the environment variable value may need to be enclosed in quotes.

NOTE

This environment variable is valid only if `I_MPI_PIN` is enabled.

The `I_MPI_PIN_PROCESSOR_LIST` environment variable has the following different syntax variants:

- **Explicit processor list.** This comma-separated list is defined in terms of logical processor numbers. The relative node rank of a process is an index to the processor list such that the *i*-th process is pinned on *i*-th list member. This permits the definition of any process placement on the CPUs.

For example, process mapping for `I_MPI_PIN_PROCESSOR_LIST=p0,p1,p2,...,pn` is as follows:

Rank on a node	0	1	2	...	n-1	N
Logical CPU	p0	p1	p2	...	pn-1	Pn

- **grain/shift/offset mapping.** This method provides cyclic shift of a defined `grain` along the processor list with steps equal to `shift*grain` and a single shift on `offset*grain` at the end. This shifting action is repeated `shift` times.

For example: `grain = 2` logical processors, `shift = 3` grains, `offset = 0`.

Legend:

gray - MPI process grains

A) red - processor grains chosen on the 1st pass

B) cyan - processor grains chosen on the 2nd pass

C) green - processor grains chosen on the final 3rd pass

D) Final map table ordered by MPI ranks

A)

0 1			2 3			...	2n-2 2n-1		
0 1	2 3	4 5	6 7	8 9	10 11	...	6n-6 6n-5	6n-4 6n-3	6n-2 6n-1

B)

0 1	2n 2n+1		2 3	2n+2 2n+3		...	2n-2 2n-1	4n-2 4n-1	
0 1	2 3	4 5	6 7	8 9	10 11	...	6n-6 6n-5	6n-4 6n-3	6n-2 6n-1

							5	3	1
--	--	--	--	--	--	--	---	---	---

C)

0 1	2n 2n+1	4n 4n+1	2 3	2n+2 2n+3	4n+2 4n+3	...	2n-2 2n-1	4n-2 4n-1	6n-2 6n-1
0 1	2 3	4 5	6 7	8 9	10 11	...	6n-6 6n-5	6n-4 6n-3	6n-2 6n-1

D)

0 1	2 3	...	2n-2 2n-1	2n 2n+1	2n+2 2n+3	...	4n-2 4n-1	4n 4n+1	4n+2 4n+3	...	6n-2 6n-1
0 1	6 7	...	6n-6 6n-5	2 3	8 9	...	6n-4 6n-3	4 5	10 11	...	6n-2 6n-1

- Predefined mapping scenario. In this case popular process pinning schemes are defined as keywords selectable at runtime. There are two such scenarios: `bunch` and `scatter`.

In the `bunch` scenario the processes are mapped proportionally to sockets as closely as possible. This mapping makes sense for partial processor loading. In this case the number of processes is less than the number of processors.

In the `scatter` scenario the processes are mapped as remotely as possible so as not to share common resources: FSB, caches, and cores.

In the example, there are two sockets, four cores per socket, one logical CPU per core, and two cores per shared cache.

Legend:

gray - MPI processes

cyan - 1st socket processors

green - 2nd socket processors

Same color defines a processor pair sharing a cache

0	1	2			3	4		
0	1	2	3		4	5	6	7

`bunch` scenario for 5 processes

0	4	2	6		1	5	3	7
0	1	2	3		4	5	6	7

`scatter` scenario for full loading

Examples

To pin the processes to CPU0 and CPU3 on each node globally, use the following command:

```
$ mpirun -genv I_MPI_PIN_PROCESSOR_LIST 0,3 \
-n <# of processes> <executable>
```

To pin the processes to different CPUs on each node individually (CPU0 and CPU3 on host1 and CPU0, CPU1 and CPU3 on host2), use the following command:

```
$ mpirun -host host1 -env I_MPI_PIN_PROCESSOR_LIST 0,3 \
-n <# of processes> <executable> : \
-host host2 -env I_MPI_PIN_PROCESSOR_LIST 1,2,3 \
-n <# of processes> <executable>
```

To print extra debug information about process pinning, use the following command:

```
$ mpirun -genv I_MPI_DEBUG 4 -m -host host1 \
-env I_MPI_PIN_PROCESSOR_LIST 0,3 -n <# of processes> <executable> :\
-host host2 -env I_MPI_PIN_PROCESSOR_LIST 1,2,3 \ -n <# of processes>
<executable>
```

NOTE

If the number of processes is greater than the number of CPUs used for pinning, the process list is wrapped around to the start of the processor list.

I_MPI_PIN_PROCESSOR_EXCLUDE_LIST

Define a subset of logical processors to be excluded for the pinning capability on the intended hosts.

Syntax

```
I_MPI_PIN_PROCESSOR_EXCLUDE_LIST=<proclist>
```

Arguments

<proclist>	A comma-separated list of logical processor numbers and/or ranges of processors.
<l>	Processor with logical number <l>.
<l>--<m>	Range of processors with logical numbers from <l> to <m>.
<k>,<l>--<m>	Processors <k>, as well as <l> through <m>.

Description

Set this environment variable to define the logical processors that Intel® MPI Library does not use for pinning capability on the intended hosts. Logical processors are numbered as in `/proc/cpuinfo`.

I_MPI_PIN_CELL

Set this environment variable to define the pinning resolution granularity. `I_MPI_PIN_CELL` specifies the minimal processor cell allocated when an MPI process is running.

Syntax

```
I_MPI_PIN_CELL=<cell>
```

Arguments

<cell>	Specify the resolution granularity
unit	Basic processor unit (logical CPU)
core	Physical processor core

Description

Set this environment variable to define the processor subset used when a process is running. You can choose from two scenarios:

- all possible CPUs in a node (`unit` value)
- all cores in a node (`core` value)

The environment variable has effect on both pinning types:

- one-to-one pinning through the `I_MPI_PIN_PROCESSOR_LIST` environment variable
- one-to-many pinning through the `I_MPI_PIN_DOMAIN` environment variable

The default value rules are:

- If you use `I_MPI_PIN_DOMAIN`, then the cell granularity is `unit`.
- If you use `I_MPI_PIN_PROCESSOR_LIST`, then the following rules apply:
 - When the number of processes is greater than the number of cores, the cell granularity is `unit`.
 - When the number of processes is equal to or less than the number of cores, the cell granularity is `core`.

NOTE

The `core` value is not affected by the enabling/disabling of Intel® Hyper-Threading Technology in a system.

I_MPI_PIN_RESPECT_CPUSET

Respect the process affinity mask.

Syntax

`I_MPI_PIN_RESPECT_CPUSET=<value>`

Arguments

<code><value></code>	Binary indicator
<code>enable yes on 1</code>	Respect the process affinity mask. This is the default value
<code>disable no off 0</code>	Do not respect the process affinity mask

Description

If you set `I_MPI_PIN_RESPECT_CPUSET=enable`, the Hydra process launcher uses its process affinity mask on each intended host to determine logical processors for applying Intel MPI Library pinning capability.

If you set `I_MPI_PIN_RESPECT_CPUSET=disable`, the Hydra process launcher does not use its process affinity mask to determine logical processors for applying Intel MPI Library pinning capability.

I_MPI_PIN_RESPECT_HCA

In the presence of Infiniband architecture* host channel adapter (IBA* HCA), adjust the pinning according to the location of IBA HCA.

Syntax

`I_MPI_PIN_RESPECT_HCA=<value>`

Arguments

<code><value></code>	Binary indicator
----------------------------	------------------

enable yes on 1	Use the location of IBA HCA if available. This is the default value
disable no off 0	Do not use the location of IBA HCA

Description

If you set `I_MPI_PIN_RESPECT_HCA=enable`, the Hydra process launcher uses the location of IBA HCA on each intended host for applying Intel MPI Library pinning capability.

If you set `I_MPI_PIN_RESPECT_HCA=disable`, the Hydra process launcher does not use the location of IBA HCA on each intended host for applying Intel MPI Library pinning capability.

3.2.4. Interoperability with OpenMP* API

I_MPI_PIN_DOMAIN

The Intel® MPI Library provides an additional environment variable to control process pinning for hybrid Intel MPI Library applications. This environment variable is used to define a number of non-overlapping subsets (domains) of logical processors on a node, and a set of rules on how MPI processes are bound to these domains by the following formula: *one MPI process per one domain*. See the picture.

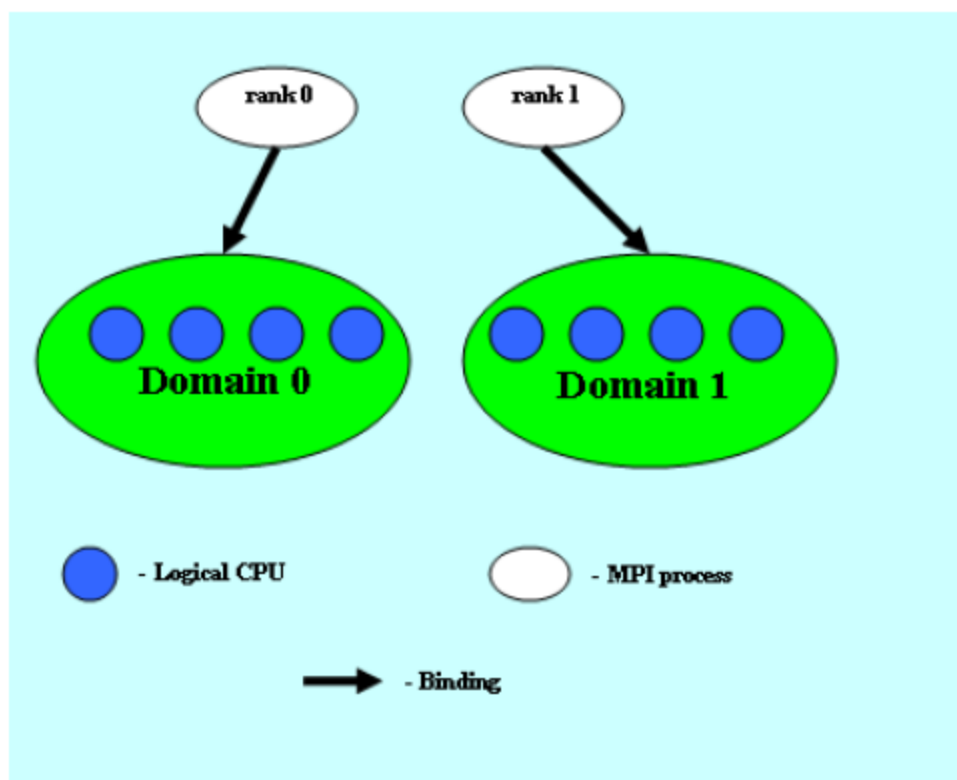


Figure 3.2-1 Domain Example

Each MPI process can create a number of children threads for running within the corresponding domain. The process threads can freely migrate from one logical processor to another within the particular domain.

If the `I_MPI_PIN_DOMAIN` environment variable is defined, then the `I_MPI_PIN_PROCESSOR_LIST` environment variable setting is ignored.

If the `I_MPI_PIN_DOMAIN` environment variable is not defined, then MPI processes are pinned according to the current value of the `I_MPI_PIN_PROCESSOR_LIST` environment variable.

The `I_MPI_PIN_DOMAIN` environment variable has the following syntax forms:

- Domain description through multi-core terms `<mc-shape>`
- Domain description through domain size and domain member layout `<size>[:<layout>]`

- Explicit domain description through bit mask *<masklist>*

The following tables describe these syntax forms.

Multi-core Shape

`I_MPI_PIN_DOMAIN=<mc-shape>`

<i><mc-shape></i>	Define domains through multi-core terms.
<code>core</code>	Each domain consists of the logical processors that share a particular core. The number of domains on a node is equal to the number of cores on the node.
<code>socket</code> <code>sock</code>	Each domain consists of the logical processors that share a particular socket. The number of domains on a node is equal to the number of sockets on the node. This is the recommended value.
<code>node</code>	All logical processors on a node are arranged into a single domain.
<code>cache1</code>	Logical processors that share a particular level 1 cache are arranged into a single domain.
<code>cache2</code>	Logical processors that share a particular level 2 cache are arranged into a single domain.
<code>cache3</code>	Logical processors that share a particular level 3 cache are arranged into a single domain.
<code>cache</code>	The largest domain among <code>cache1</code> , <code>cache2</code> , and <code>cache3</code> is selected.

Explicit Shape

`I_MPI_PIN_DOMAIN=<size>[:<layout>]`

<i><size></i>	Define a number of logical processors in each domain (domain size)
<code>omp</code>	The domain size is equal to the <code>OMP_NUM_THREADS</code> environment variable value. If the <code>OMP_NUM_THREADS</code> environment variable is not set, each node is treated as a separate domain.
<code>auto</code>	The domain size is defined by the formula <code>size=#cpu/#proc</code> , where <code>#cpu</code> is the number of logical processors on a node, and <code>#proc</code> is the number of the MPI processes started on a node
<i><n></i>	The domain size is defined by a positive decimal number <i><n></i>

<i><layout></i>	Ordering of domain members. The default value is <code>compact</code>
<code>platform</code>	Domain members are ordered according to their BIOS numbering (platform-depended numbering)
<code>compact</code>	Domain members are located as close to each other as possible in terms of common resources (cores, caches, sockets, and so on). This is the default value
<code>scatter</code>	Domain members are located as far away from each other as possible in terms of common resources (cores, caches, sockets, and so on)

Explicit Domain Mask

`I_MPI_PIN_DOMAIN=<masklist>`

<code><masklist></code>	Define domains through the comma separated list of hexadecimal numbers (domain masks)
<code>[m₁, . . . , m_n]</code>	<p>For <code><masklist></code>, each <code>m_i</code> is a hexadecimal bit mask defining an individual domain. The following rule is used: the <code>ith</code> logical processor is included into the domain if the corresponding <code>m_i</code> value is set to 1. All remaining processors are put into a separate domain. BIOS numbering is used.</p> <hr/> <p>NOTE</p> <p>To ensure that your configuration in <code><masklist></code> is parsed correctly, use square brackets to enclose the domains specified by the <code><masklist></code>. For example:</p> <pre>I_MPI_PIN_DOMAIN=[0x55, 0xaa]</pre> <hr/>

NOTE

These options are available for both Intel® and non-Intel microprocessors, but they may perform additional optimizations for Intel microprocessors than they perform for non-Intel microprocessors.

NOTE

To pin OpenMP* processes or threads inside the domain, the corresponding OpenMP feature (for example, the `KMP_AFFINITY` environment variable for Intel® Composer XE) should be used.

NOTE

The following configurations are effectively the same as if pinning is not applied:

- If you set `I_MPI_PIN_DOMAIN=auto` and a single process is running on a node (for example, due to `I_MPI_PERHOST=1`)
- `I_MPI_PIN_DOMAIN=node`

If you do not want the process to be migrated between sockets on a multi-socket platform, specify the domain size as `I_MPI_PIN_DOMAIN=socket` or smaller.

You can also use `I_MPI_PIN_PROCESSOR_LIST`, which produces a single-cpu process affinity mask for each rank (the affinity mask is supposed to be automatically adjusted in presence of IBA* HCA).

See the following model of an SMP node in the examples:

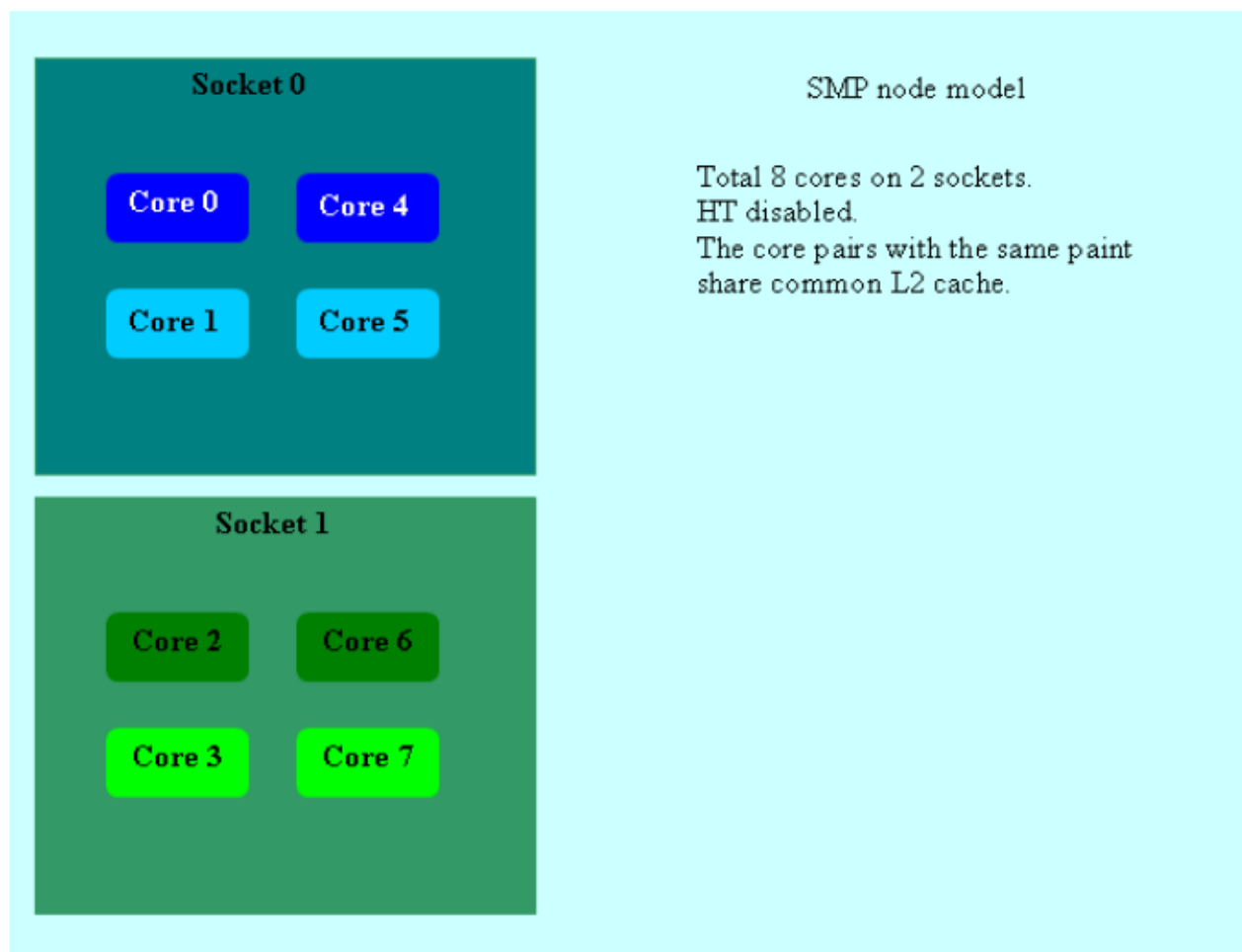


Figure 3.2-2 Model of a Node

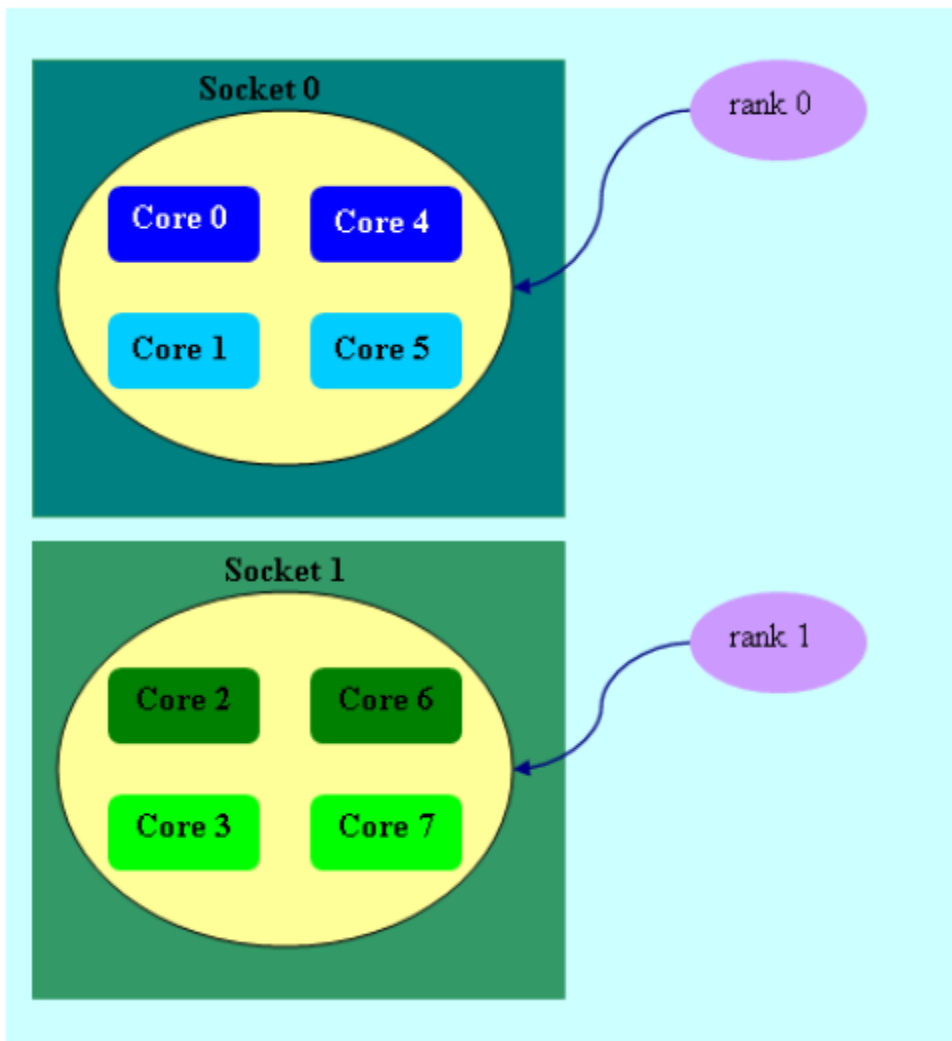


Figure 3.2-3 `mpirun -n 2 -env I_MPI_PIN_DOMAIN socket ./a.out`

In Figure 3.2-3, two domains are defined according to the number of sockets. Process rank 0 can migrate on all cores on the 0-th socket. Process rank 1 can migrate on all cores on the first socket.

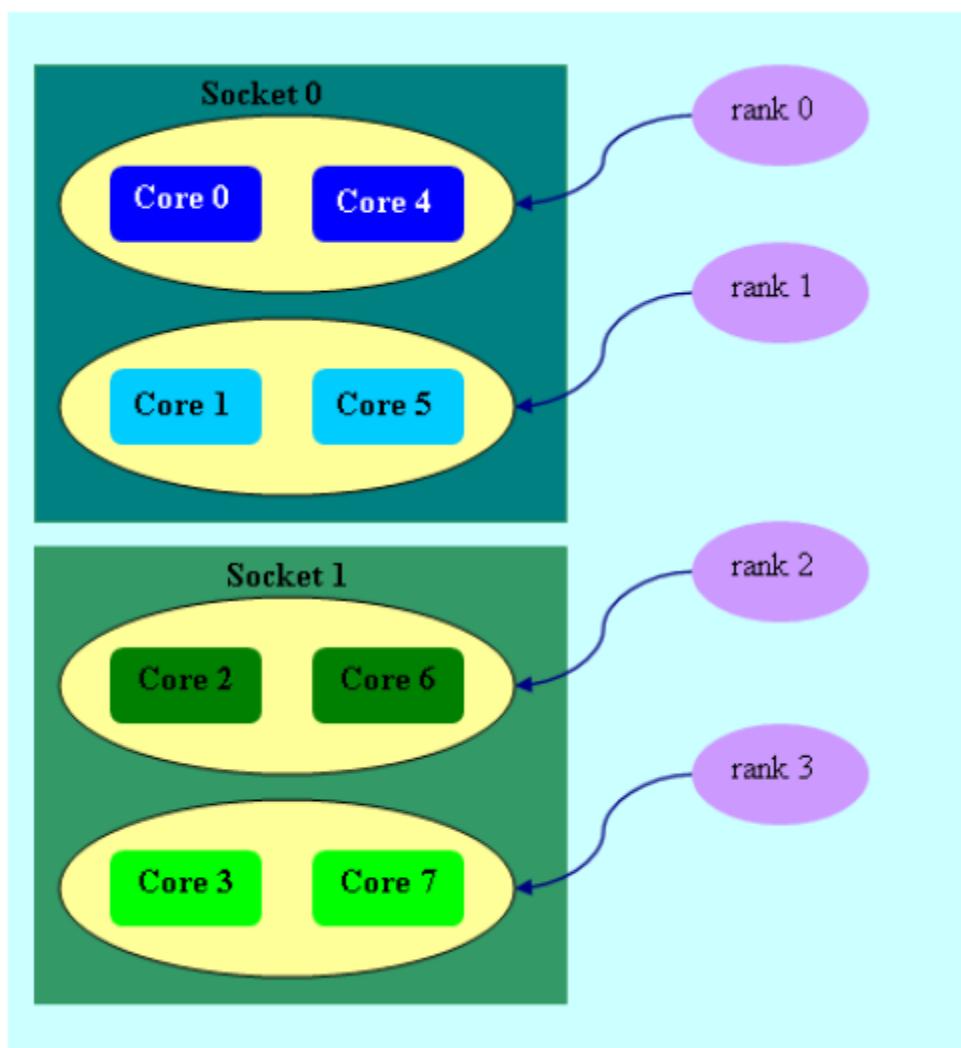


Figure 3.2-4 `mpirun -n 4 -env I_MPI_PIN_DOMAIN cache2 ./a.out`

In Figure 3.2-4, four domains are defined according to the amount of common L2 caches. Process rank 0 runs on cores {0,4} that share an L2 cache. Process rank 1 runs on cores {1,5} that share an L2 cache as well, and so on.

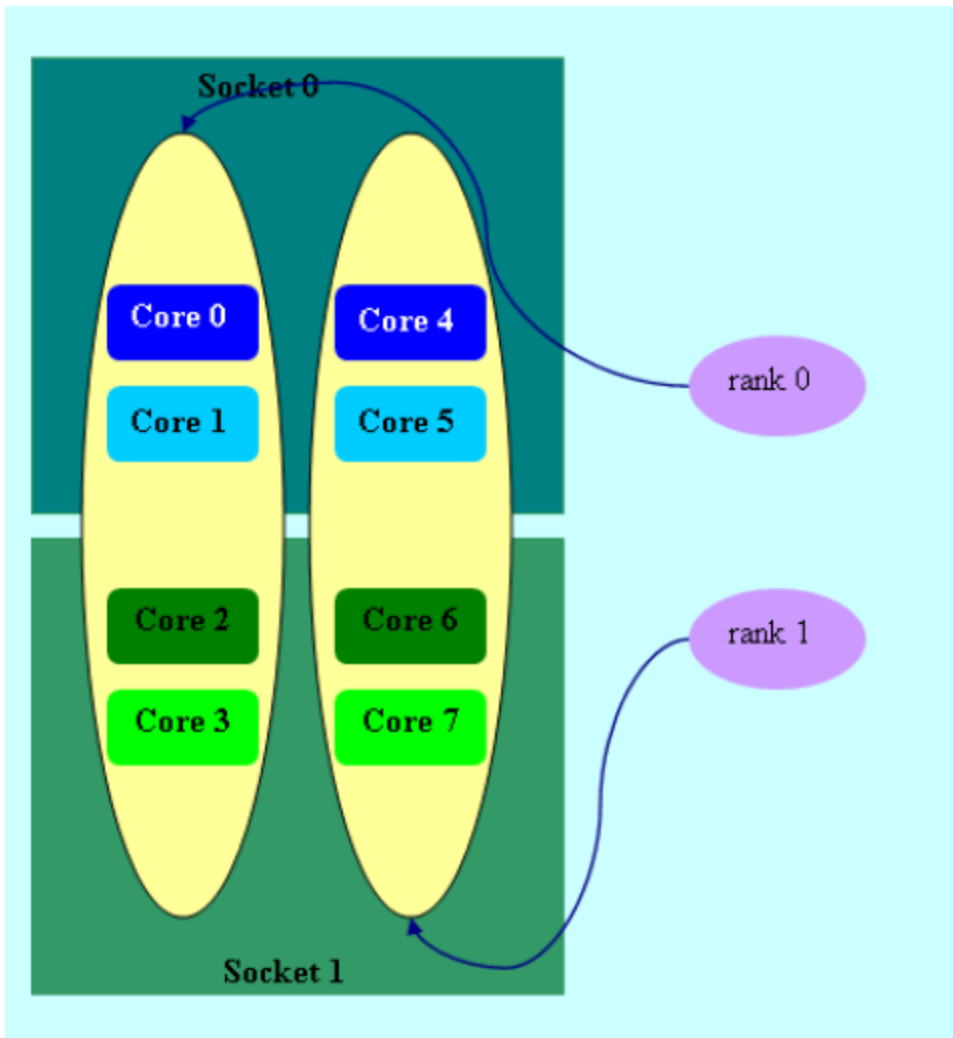


Figure 3.2-5 `mpirun -n 2 -env I_MPI_PIN_DOMAIN 4:platform ./a.out`

In Figure 3.2-5, two domains with size=4 are defined. The first domain contains cores {0,1,2,3}, and the second domain contains cores {4,5,6,7}. Domain members (cores) have consecutive numbering as defined by the `platform` option.

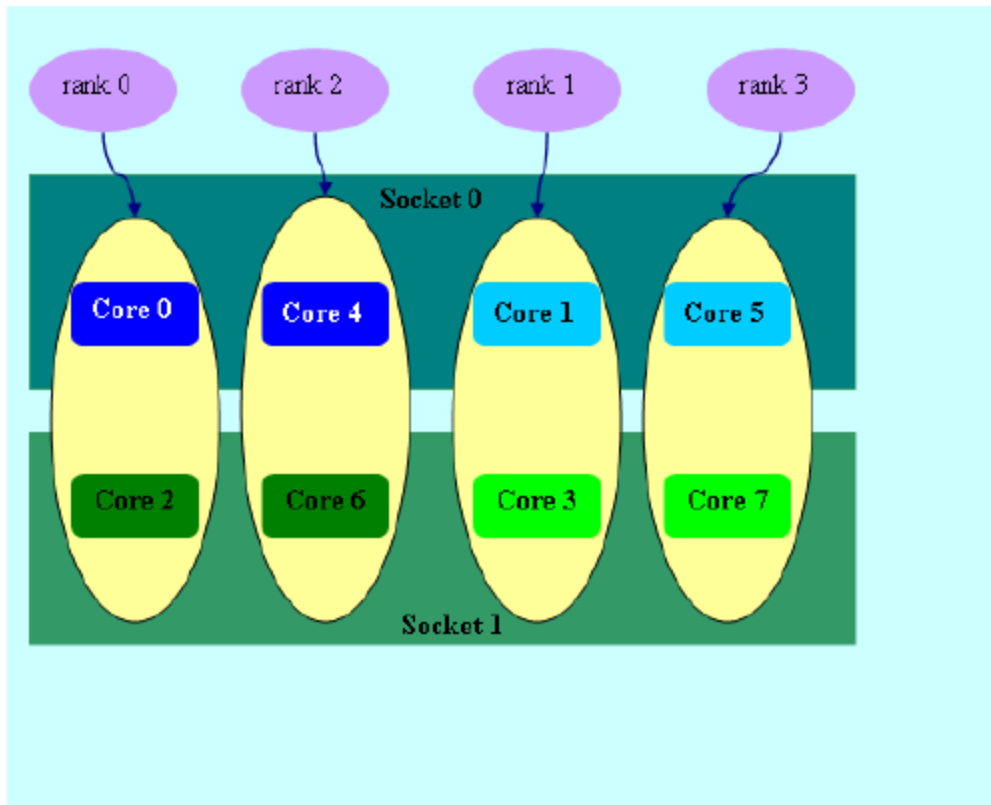


Figure 3.2-6 `mpirun -n 4 -env I_MPI_PIN_DOMAIN auto:scatter ./a.out`

In [Picture 3.2-6](#), domain size=2 (defined by the number of CPUs=8 / number of processes=4), `scatter` layout. Four domains {0,2}, {1,3}, {4,6}, {5,7} are defined. Domain members do not share any common resources.

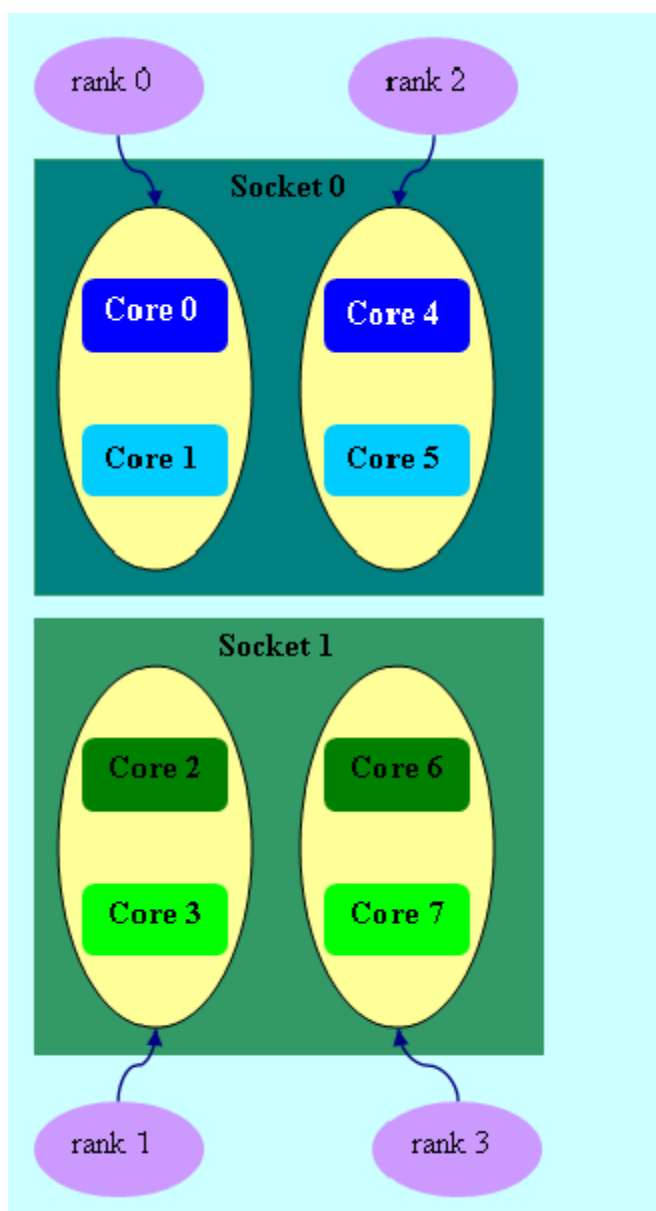


Figure 3.2-7 `setenv OMP_NUM_THREADS=2`

```
mpirun -n 4 -env I_MPI_PIN_DOMAIN omp:platform ./a.out
```

In [Figure 3.2-7](#), domain size=2 (defined by `OMP_NUM_THREADS=2`), `platform` layout. Four domains {0,1}, {2,3}, {4,5}, {6,7} are defined. Domain members (cores) have consecutive numbering.

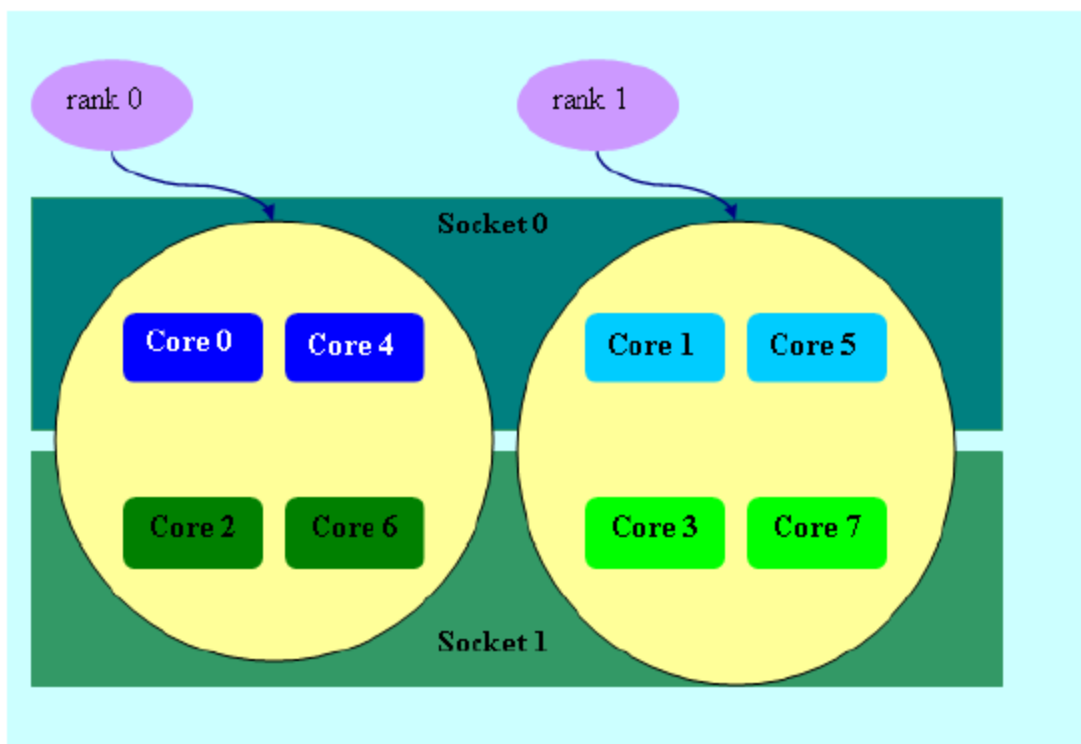


Figure 3.2-8 `mpirun -n 2 -env I_MPI_PIN_DOMAIN [0x55,0xaa] ./a.out`

In Figure 3.2-8 (the example for `I_MPI_PIN_DOMAIN=<masklist>`), the first domain is defined by the 0x55 mask. It contains all cores with even numbers {0,2,4,6}. The second domain is defined by the 0xAA mask. It contains all cores with odd numbers {1,3,5,7}.

I_MPI_PIN_ORDER

Set this environment variable to define the mapping order for MPI processes to domains as specified by the `I_MPI_PIN_DOMAIN` environment variable.

Syntax

`I_MPI_PIN_ORDER=<order>`

Arguments

<code><order></code>	Specify the ranking order
<code>range</code>	The domains are ordered according to the processor's BIOS numbering. This is a platform-dependent numbering
<code>scatter</code>	The domains are ordered so that adjacent domains have minimal sharing of common resources
<code>compact</code>	The domains are ordered so that adjacent domains share common resources as much as possible. This is the default value
<code>spread</code>	The domains are ordered consecutively with the possibility not to share common resources
<code>bunch</code>	The processes are mapped proportionally to sockets and the domains are ordered as close as possible on the sockets

Description

The optimal setting for this environment variable is application-specific. If adjacent MPI processes prefer to share common resources, such as cores, caches, sockets, FSB, use the `compact` or `bunch` values. Otherwise, use the `scatter` or `spread` values. Use the `range` value as needed. For detail information and examples about these values, see the Arguments table and the Example section of `I_MPI_PIN_ORDER` in this topic.

The options `scatter`, `compact`, `spread` and `bunch` are available for both Intel® and non-Intel microprocessors, but they may perform additional optimizations for Intel microprocessors than they perform for non-Intel microprocessors.

Example

For the following configuration:

- Two socket nodes with four cores and a shared L2 cache for corresponding core pairs.
- 4 MPI processes you want to run on the node using the following settings:
- For compact order:
`I_MPI_PIN_DOMAIN=2`
`I_MPI_PIN_ORDER=compact`

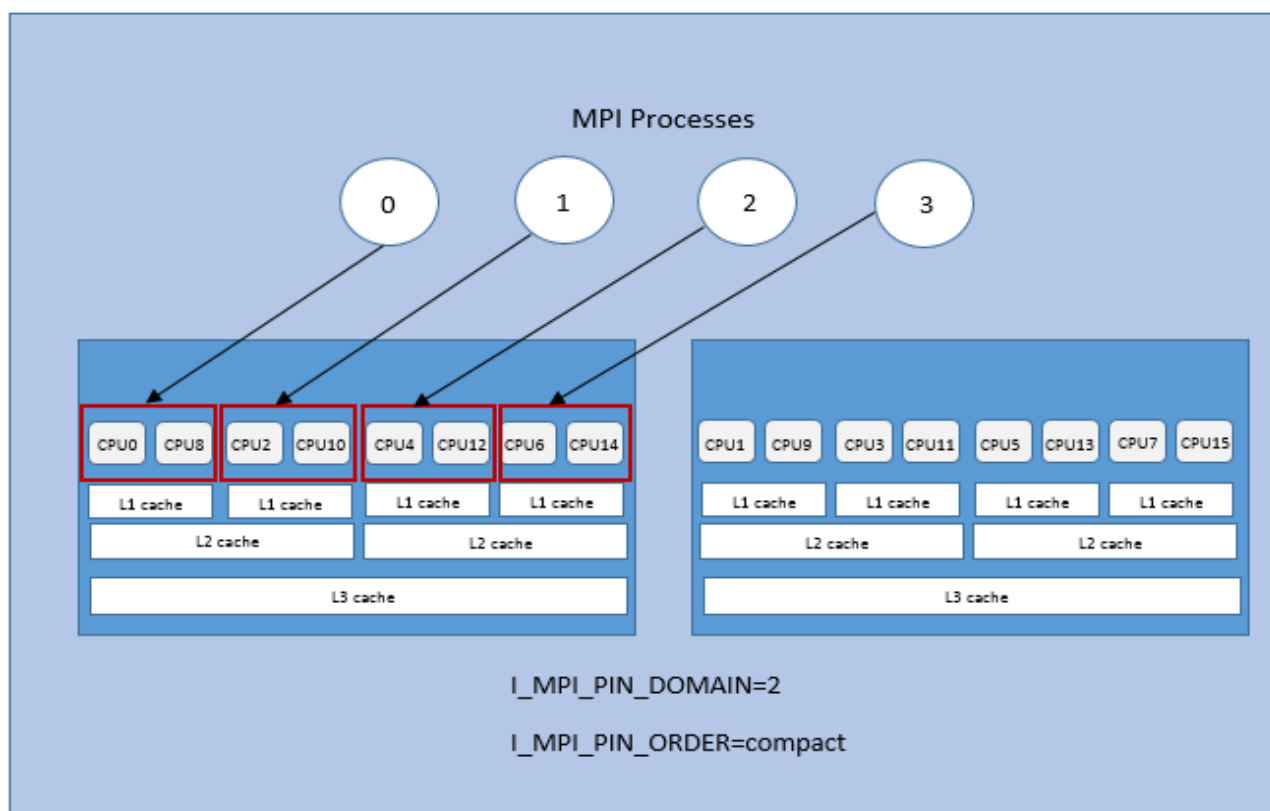


Figure 3.2-9 Compact Order Example

- For scatter order:
`I_MPI_PIN_DOMAIN=2`
`I_MPI_PIN_ORDER=scatter`

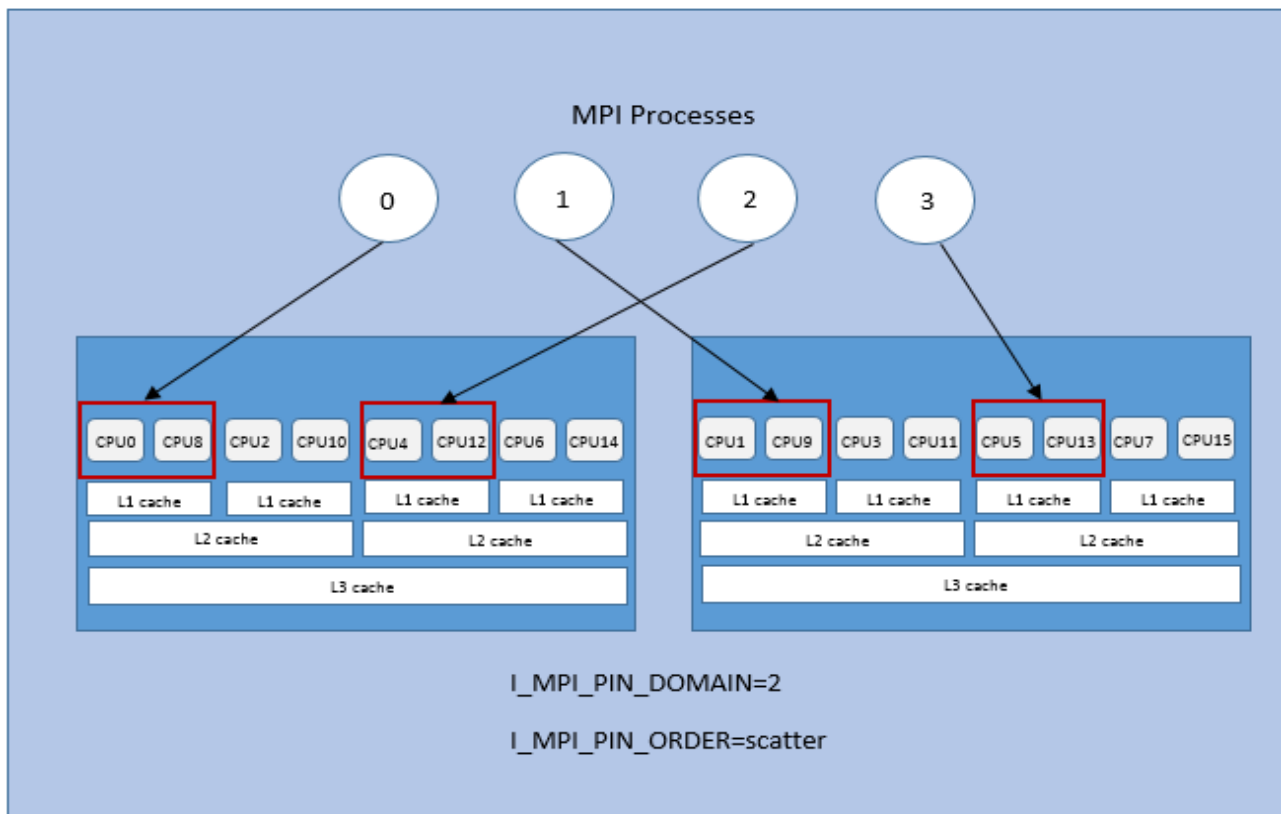


Figure 3.2-10 Scatter Order Example

- For spread order:
I_MPI_PIN_DOMAIN=2
I_MPI_PIN_ORDER=spread

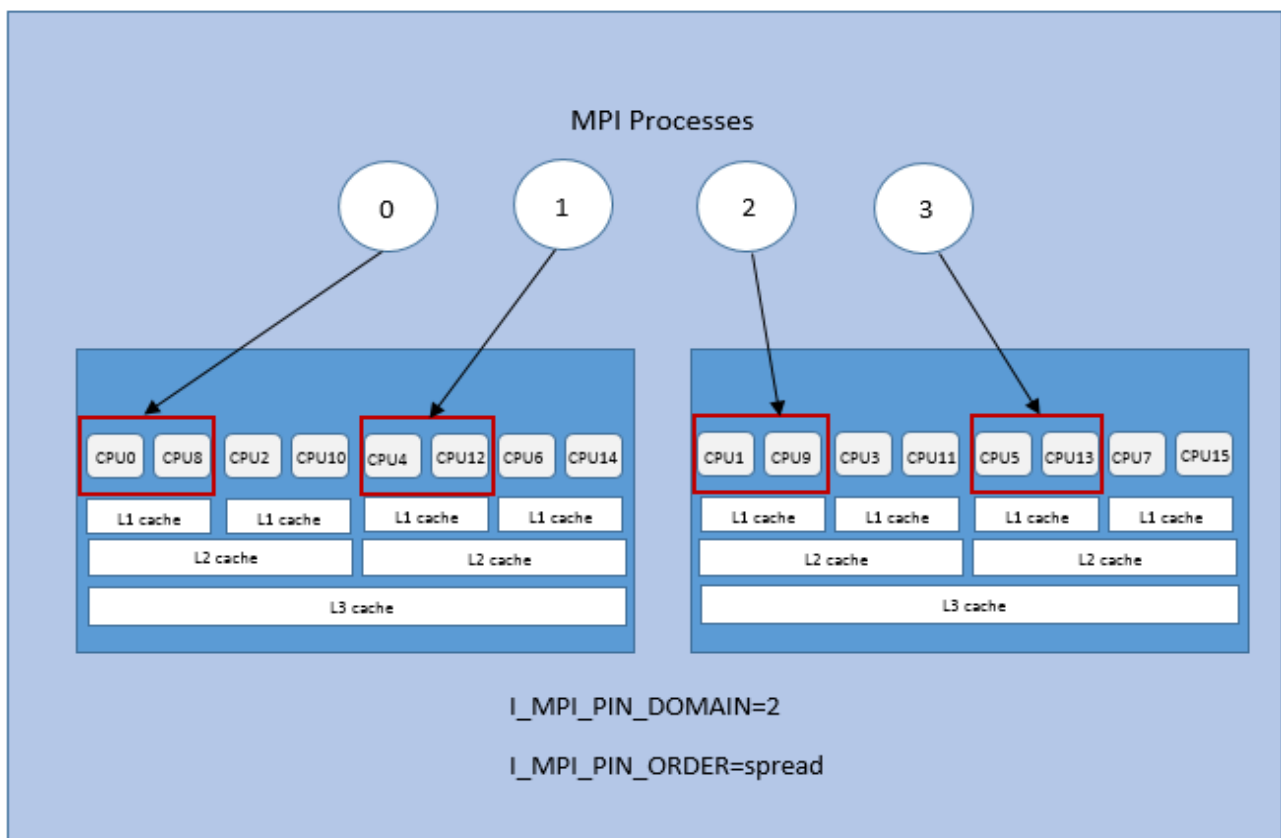
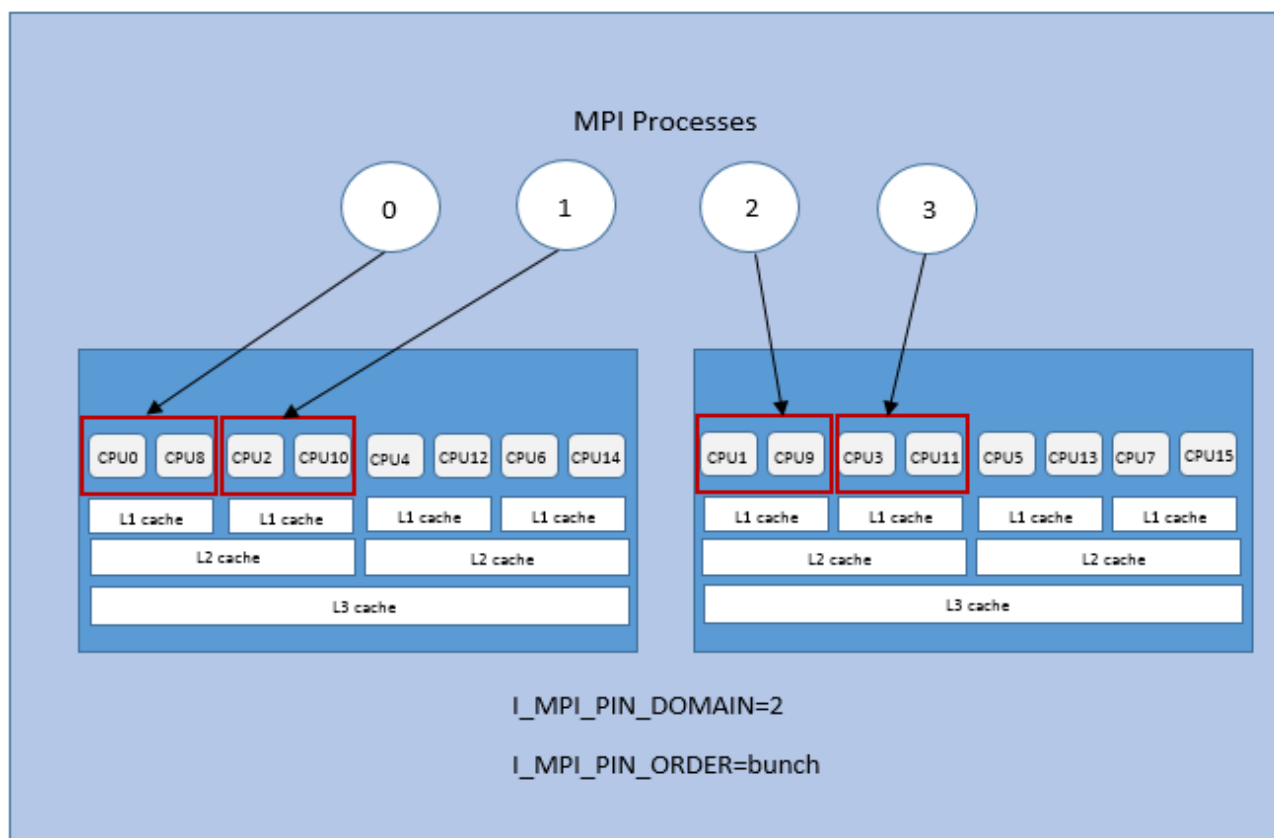


Figure 3.2-11 Spread Order Example

- For bunch order:
`I_MPI_PIN_DOMAIN=2`
`I_MPI_PIN_ORDER=bunch`

**Figure 3.2-12 Bunch Order Example**

3.3. Fabrics Control

This topic provides you with the information on how to use environment variables to control the following fabrics:

- Communication fabrics
- Shared memory fabrics
- DAPL-capable network fabrics
- DAPL UD-capable network fabrics
- TCP-capable network fabrics
- TMI-capable network fabrics
- OFA*-capable network fabrics
- OFI*-capable network fabrics

3.3.1. Communication Fabrics Control

I_MPI_FABRICS (I_MPI_DEVICE)

Select the particular network fabrics to be used.

Syntax

```
I_MPI_FABRICS=<fabric>|<intra-node fabric>:<inter-nodes fabric>
```

Where <fabric>:= {shm, dapl, tcp, tmi, ofa, ofi}

<intra-node fabric>:= {shm, dapl, tcp, tmi, ofa, ofi}

<inter-nodes fabric>:= {dapl, tcp, tmi, ofa, ofi}

Deprecated Syntax

```
I_MPI_DEVICE=<device>[:<provider>]
```

Arguments

<fabric>	Define a network fabric
shm	Shared-memory
dapl	DAPL-capable network fabrics, such as InfiniBand*, iWarp*, Dolphin*, and XPMEM* (through DAPL*)
tcp	TCP/IP-capable network fabrics, such as Ethernet and InfiniBand* (through IPoIB*)
tmi	TMI-capable network fabrics including Intel® True Scale Fabric, Myrinet*, (through Tag Matching Interface)
ofa	OFA-capable network fabric including InfiniBand* (through OFED* verbs)
ofi	OFI (OpenFabrics Interfaces*)-capable network fabric including Intel® True Scale Fabric, and TCP (through OFI* API)

Correspondence with I_MPI_DEVICE

<device>	<fabric>
sock	tcp
shm	shm
ssm	shm:tcp
rdma	dapl
rdssm	shm:dapl
<provider>	Optional DAPL* provider name (only for the rdma and the rdssm devices) I_MPI_DAPL_PROVIDER=<provider> or I_MPI_DAPL_UD_PROVIDER=<provider>

Use the <provider> specification only for the {rdma, rdssm} devices.

For example, to select the OFED* InfiniBand* device, use the following command:

```
$ mpiexec -n <# of processes> \
-env I_MPI_DEVICE rdssm:OpenIB-cma <executable>
```

For these devices, if *<provider>* is not specified, the first DAPL* provider in the `/etc/dat.conf` file is used.

Description

Set this environment variable to select a specific fabric combination. If the requested fabric(s) is not available, Intel® MPI Library can fall back to other fabric(s). See [I_MPI_FALLBACK](#) for details. If the `I_MPI_FABRICS` environment variable is not defined, Intel® MPI Library selects the most appropriate fabric combination automatically.

The exact combination of fabrics depends on the number of processes started per node.

- If all processes start on one node, the library uses `shm` intra-node communication.
- If the number of started processes is less than or equal to the number of available nodes, the library uses the first available fabric from the fabrics list for inter-node communication.
- For other cases, the library uses `shm` for intra-node communication, and the first available fabric from the fabrics list for inter-node communication. See [I_MPI_FABRICS_LIST](#) for details.

The `shm` fabric is available for both Intel® and non-Intel microprocessors, but it may perform additional optimizations for Intel microprocessors than it performs for non-Intel microprocessors.

NOTE

The combination of selected fabrics ensures that the job runs, but this combination may not provide the highest possible performance for the given cluster configuration.

For example, to select shared-memory as the chosen fabric, use the following command:

```
$ mpirun -n <# of processes> -env I_MPI_FABRICS shm <executable>
```

To select shared-memory and DAPL-capable network fabric as the chosen fabric combination, use the following command:

```
$ mpirun -n <# of processes> -env I_MPI_FABRICS shm:dapl <executable>
```

To enable Intel® MPI Library to select most appropriate fabric combination automatically, use the following command:

```
$ mpirun -n <# of procs> -perhost <# of procs per host> <executable>
```

Set the level of debug information to 2 or higher to check which fabrics have been initialized. See [I_MPI_DEBUG](#) for details. For example:

```
[0] MPI startup(): shm and dapl data transfer modes
```

or

```
[0] MPI startup(): tcp data transfer mode
```

I_MPI_FABRICS_LIST

Define a fabrics list.

Syntax

```
I_MPI_FABRICS_LIST=<fabrics list>
```

Where *<fabrics list>* := *<fabric>*, ..., *<fabric>*

<fabric> := {dapl, tcp, tmi, ofa, ofi}

Arguments

<code><fabrics list></code>	Specify a list of fabrics
<code>dapl, ofa, tcp, tmi, ofi</code>	This is the default value
<code>dapl, tcp, ofa, tmi, ofi</code>	If you specify <code>I_MPI_WAIT_MODE=enable</code> , this is the default value
<code>tmi, dapl, tcp, ofa, ofi</code>	This is the default value for nodes that have Intel® True Scale Fabric available and do not have any other type of interconnect cards. In case host has several types of HCAs, this does not apply

Description

Set this environment variable to define a list of fabrics. The library uses the fabrics list to choose the most appropriate fabrics combination automatically. For more information on fabric combination, see [I_MPI_FABRICS](#).

For example, if `I_MPI_FABRICS_LIST=dapl, tcp`, and `I_MPI_FABRICS` is not defined, and the initialization of DAPL-capable network fabrics fails, the library falls back to TCP-capable network fabric. For more information on fallback, see [I_MPI_FALLBACK](#).

[I_MPI_FALLBACK](#) ([I_MPI_FALLBACK_DEVICE](#))

Set this environment variable to enable fallback to the first available fabric.

Syntax

`I_MPI_FALLBACK=<arg>`

Deprecated Syntax

`I_MPI_FALLBACK_DEVICE=<arg>`

Arguments

<code><arg></code>	Binary indicator
<code>enable yes on 1</code>	Fall back to the first available fabric. This is the default value if you do not set the <code>I_MPI_FABRICS(I_MPI_DEVICE)</code> environment variable.
<code>disable no off 0</code>	Terminate the job if MPI cannot initialize the one of the fabrics selected by the <code>I_MPI_FABRICS</code> environment variable. This is the default value if you do not set the <code>I_MPI_FABRICS(I_MPI_DEVICE)</code> environment variable.

Description

Set this environment variable to control fallback to the first available fabric.

If you set `I_MPI_FALLBACK` to `enable` and an attempt to initialize a specified fabric fails, the library uses the first available fabric from the list of fabrics. See [I_MPI_FABRICS_LIST](#) for details.

If you set `I_MPI_FALLBACK` to `disable` and an attempt to initialize a specified fabric fails, the library terminates the MPI job.

NOTE

If you set `I_MPI_FABRICS` and `I_MPI_FALLBACK=enable`, the library falls back to fabrics with higher numbers in the fabrics list. For example, if `I_MPI_FABRICS=dapl`, `I_MPI_FABRICS_LIST=ofa, tmi, dapl, tcp`, `I_MPI_FALLBACK=enable` and the initialization of DAPL-capable network fabrics fails, the library falls back to TCP-capable network fabric.

I_MPI_LARGE_SCALE_THRESHOLD

Change the threshold for enabling scalable optimizations.

Syntax

`I_MPI_LARGE_SCALE_THRESHOLD=<arg>`

Arguments

<code><nprocs></code>	Define the scale threshold
<code>> 0</code>	The default value is 4096

Description

This variable defines the number of processes when the DAPL UD IB extension is turned on automatically.

I_MPI_EAGER_THRESHOLD

Change the eager/rendezvous message size threshold for all devices.

Syntax

`I_MPI_EAGER_THRESHOLD=<nbytes>`

Arguments

<code><nbytes></code>	Set the eager/rendezvous message size threshold
<code>> 0</code>	The default <code><nbytes></code> value is equal to 262144 bytes

Description

Set this environment variable to control the protocol used for point-to-point communication:

- Messages shorter than or equal in size to `<nbytes>` are sent using the eager protocol.
- Messages larger than `<nbytes>` are sent using the rendezvous protocol. The rendezvous protocol uses memory more efficiently.

I_MPI_INTRANODE_EAGER_THRESHOLD

Change the eager/rendezvous message size threshold for intra-node communication mode.

Syntax

`I_MPI_INTRANODE_EAGER_THRESHOLD=<nbytes>`

Arguments

<code><nbytes></code>	Set the eager/rendezvous message size threshold for intra-node communication
<code>> 0</code>	The default <code><nbytes></code> value is equal to 262144 bytes for all fabrics except shm. For shm, cutover point is equal to the value of <code>I_MPI_SHM_CELL_SIZE</code> environment variable

Description

Set this environment variable to change the protocol used for communication within the node:

- Messages shorter than or equal in size to `<nbytes>` are sent using the eager protocol.
- Messages larger than `<nbytes>` are sent using the rendezvous protocol. The rendezvous protocol uses the memory more efficiently.

If you do not set `I_MPI_INTRANODE_EAGER_THRESHOLD`, the value of `I_MPI_EAGER_THRESHOLD` is used.

`I_MPI_SPIN_COUNT`

Control the spin count value.

Syntax

`I_MPI_SPIN_COUNT=<scout>`

Arguments

<code><scout></code>	Define the loop spin count when polling fabric(s)
<code>> 0</code>	The default <code><scout></code> value is equal to 1 when more than one process runs per processor/core. Otherwise the value equals 250. The maximum value is equal to 2147483647

Description

Set the spin count limit. The loop for polling the fabric(s) spins `<scout>` times before the library releases the processes if no incoming messages are received for processing. Within every spin loop, the `shm` fabric (if enabled) is polled an extra `I_MPI_SHM_SPIN_COUNT` times. Smaller values for `<scout>` cause the Intel® MPI Library to release the processor more frequently.

Use the `I_MPI_SPIN_COUNT` environment variable for tuning application performance. The best value for `<scout>` can be chosen on an experimental basis. It depends on the particular computational environment and the application.

`I_MPI_SCALABLE_OPTIMIZATION`

Turn on/off scalable optimization of the network fabric communication.

Syntax

`I_MPI_SCALABLE_OPTIMIZATION=<arg>`

Arguments

<code><arg></code>	Binary indicator
<code>enable yes on 1</code>	Turn on scalable optimization of the network fabric communication. This is the default for 16 or more processes
<code>disable no off 0</code>	Turn off scalable optimization of the network fabric communication. This is the default value for less than 16 processes

Description

Set this environment variable to enable scalable optimization of the network fabric communication. In most cases, using optimization decreases latency and increases bandwidth for a large number of processes.

`I_MPI_WAIT_MODE`

Turn on/off wait mode.

Syntax

`I_MPI_WAIT_MODE=<arg>`

Arguments

<code><arg></code>	Binary indicator
<code>enable yes on 1</code>	Turn on the wait mode
<code>disable no off 0</code>	Turn off the wait mode. This is the default

Description

Set this environment variable to control the wait mode. If you enable this mode, the processes wait for receiving messages without polling the fabric(s). This mode can save CPU time for other tasks.

Use the Native POSIX Thread Library* with the wait mode for `shm` communications.

NOTE

To check which version of the thread library is installed, use the following command:

```
$ getconf GNU_LIBPTHREAD_VERSION
```

I_MPI_DYNAMIC_CONNECTION (I_MPI_USE_DYNAMIC_CONNECTIONS)

Control the dynamic connection establishment.

Syntax

`I_MPI_DYNAMIC_CONNECTION=<arg>`

Deprecated Syntax

`I_MPI_USE_DYNAMIC_CONNECTIONS=<arg>`

Arguments

<code><arg></code>	Binary indicator
<code>enable yes on 1</code>	Turn on the dynamic connection establishment. This is the default for 64 or more processes
<code>disable no off 0</code>	Turn off the dynamic connection establishment. This is the default for less than 64 processes

Description

Set this environment variable to control dynamic connection establishment.

- If this mode is enabled, all connections are established at the time of the first communication between each pair of processes.
- If this mode is disabled, all connections are established upfront.

The default value depends on the number of processes in the MPI job. The dynamic connection establishment is off if the total number of processes is less than 64.

3.3.2. Shared Memory Control

I_MPI_SHM_CACHE_BYPASS (I_MPI_CACHE_BYPASS)

Control the message transfer algorithm for the shared memory.

Syntax

```
I_MPI_SHM_CACHE_BYPASS=<arg>
```

Deprecated Syntax

```
I_MPI_CACHE_BYPASS=<arg>
```

Arguments

<arg>	Binary indicator
enable yes on 1	Enable message transfer bypass cache. This is the default value
disable no off 0	Disable message transfer bypass cache

Description

Set this environment variable to enable/disable message transfer bypass cache for the shared memory. When you enable this feature, the MPI sends the messages greater than or equal in size to the value specified by the `I_MPI_SHM_CACHE_BYPASS_THRESHOLD` environment variable through the bypass cache. This feature is enabled by default.

`I_MPI_SHM_CACHE_BYPASS_THRESHOLDS` (`I_MPI_CACHE_BYPASS_THRESHOLDS`)

Set the message copying algorithm threshold.

Syntax

```
I_MPI_SHM_CACHE_BYPASS_THRESHOLDS=<nb_send>,<nb_recv>[,<nb_send_pk>,<nb_recv_pk>]
```

Deprecated Syntax

```
I_MPI_CACHE_BYPASS_THRESHOLDS=<nb_send>,<nb_recv>[,<nb_send_pk>,<nb_recv_pk>]
```

Arguments

<nb_send>	Set the threshold for sent messages in the following situations: <ul style="list-style-type: none"> Processes are pinned on cores that are not located in the same physical processor package Processes are not pinned
<nb_recv>	Set the threshold for received messages in the following situations: <ul style="list-style-type: none"> Processes are pinned on cores that are not located in the same physical processor package Processes are not pinned
<nb_send_pk>	Set the threshold for sent messages when processes are pinned on cores located in the same physical processor package
<nb_recv_pk>	Set the threshold for received messages when processes are pinned on cores located in the same physical processor package

Description

Set this environment variable to control the thresholds for the message copying algorithm. Intel® MPI Library uses different message copying implementations which are optimized to operate with different memory hierarchy levels. Intel® MPI Library copies messages greater than or equal in size to the defined threshold value using copying algorithm optimized for far memory access. The value of -1 disables using of those algorithms. The default values depend on the architecture and may vary among the Intel® MPI Library versions. This environment variable is valid only when `I_MPI_SHM_CACHE_BYPASS` is enabled.

This environment variable is available for both Intel and non-Intel microprocessors, but it may perform additional optimizations for Intel microprocessors than it performs for non-Intel microprocessors.

I_MPI_SHM_FBOX

Control the usage of the shared memory fast-boxes.

Syntax

`I_MPI_SHM_FBOX=<arg>`

Arguments

<code><arg></code>	Binary indicator
<code>enable yes on 1</code>	Turn on fast box usage. This is the default value.
<code>disable no off 0</code>	Turn off fast box usage.

Description

Set this environment variable to control the usage of fast-boxes. Each pair of MPI processes on the same computing node has two shared memory fast-boxes, for sending and receiving eager messages.

Turn off the usage of fast-boxes to avoid the overhead of message synchronization when the application uses mass transfer of short non-blocking messages.

I_MPI_SHM_FBOX_SIZE

Set the size of the shared memory fast-boxes.

Syntax

`I_MPI_SHM_FBOX_SIZE=<nbytes>`

Arguments

<code><nbytes></code>	The size of shared memory fast-boxes in bytes
<code>> 0</code>	The default <code><nbytes></code> value depends on the specific platform you use. The value range is from 8K to 64K typically.

Description

Set this environment variable to define the size of shared memory fast-boxes.

I_MPI_SHM_CELL_NUM

Change the number of cells in the shared memory receiving queue.

Syntax

`I_MPI_SHM_CELL_NUM=<num>`

Arguments

<code><num></code>	The number of shared memory cells
<code>> 0</code>	The default value is 128

Description

Set this environment variable to define the number of cells in the shared memory receive queue. Each MPI process has own shared memory receive queue, where other processes put eager messages. The queue is used when shared memory fast-boxes are blocked by another MPI request.

I_MPI_SHM_CELL_SIZE

Change the size of a shared memory cell.

Syntax

`I_MPI_SHM_CELL_SIZE=<nbytes>`

Arguments

<code><nbytes></code>	The size of a shared memory cell in bytes
<code>> 0</code>	The default <code><nbytes></code> value depends on the specific platform you use. The value range is from 8K to 64K typically.

Description

Set this environment variable to define the size of shared memory cells.

If you set this environment variable, `I_MPI_INTRANODE_EAGER_THRESHOLD` is also changed and becomes equal to the given value.

I_MPI_SHM_LMT

Control the usage of large message transfer (LMT) mechanism for the shared memory.

Syntax

`I_MPI_SHM_LMT=<arg>`

Arguments

<code><arg></code>	Binary indicator
<code>shm</code>	Turn on the shared memory copy LMT mechanism.
<code>direct</code>	Turn on the direct copy LMT mechanism. This is the default value
<code>disable no off 0</code>	Turn off LMT mechanism

Description

Set this environment variable to control the usage of the large message transfer (LMT) mechanism. To transfer rendezvous messages, you can use the LMT mechanism by employing either of the following implementations:

- Use intermediate shared memory queues to send messages.
- Use direct copy mechanism that transfers messages without intermediate buffer if the Linux* kernel is higher than version 3.2 which supports the cross memory attach (CMA) feature. If you set

the `I_MPI_SHM_LMT` environment variable to `direct`, but the operating system does not support CMA, then the `shm` LTM mechanism runs.

`I_MPI_SHM_LMT_BUFFER_NUM` (`I_MPI_SHM_NUM_BUFFERS`)

Change the number of shared memory buffers for the large message transfer (LMT) mechanism.

Syntax

```
I_MPI_SHM_LMT_BUFFER_NUM=<num>
```

Deprecated Syntax

```
I_MPI_SHM_NUM_BUFFERS=<num>
```

Arguments

<code><num></code>	The number of shared memory buffers for each process pair
<code>> 0</code>	The default value is 8

Description

Set this environment variable to define the number of shared memory buffers between each process pair.

`I_MPI_SHM_LMT_BUFFER_SIZE` (`I_MPI_SHM_BUFFER_SIZE`)

Change the size of shared memory buffers for the LMT mechanism.

Syntax

```
I_MPI_SHM_LMT_BUFFER_SIZE=<nbytes>
```

Deprecated Syntax

```
I_MPI_SHM_BUFFER_SIZE=<nbytes>
```

Arguments

<code><nbytes></code>	The size of shared memory buffers in bytes
<code>> 0</code>	The default <code><nbytes></code> value is equal to 32768 bytes

Description

Set this environment variable to define the size of shared memory buffers for each pair of processes.

`I_MPI_SSHM`

Control the usage of the scalable shared memory mechanism.

Syntax

```
I_MPI_SSHM =<arg>
```

Arguments

<code><arg></code>	Binary indicator
<code>enable yes on 1</code>	Turn on the usage of this mechanism

disable no off 0	Turn off the usage of this mechanism. This is the default value
------------------------	---

Description

Set this environment variable to control the usage of an alternative shared memory mechanism. This mechanism replaces the shared memory fast-boxes, receive queues and LMT mechanism.

If you set this environment variable, the `I_MPI_INTRANODE_EAGER_THRESHOLD` environment variable is changed and becomes equal to 262,144 bytes.

I_MPI_SSHM_BUFFER_NUM

Change the number of shared memory buffers for the alternative shared memory mechanism.

Syntax

`I_MPI_SSHM_BUFFER_NUM=<num>`

Arguments

<num>	The number of shared memory buffers for each process pair
> 0	The default value is 4

Description

Set this environment variable to define the number of shared memory buffers between each process pair.

I_MPI_SSHM_BUFFER_SIZE

Change the size of shared memory buffers for the alternative shared memory mechanism.

Syntax

`I_MPI_SSHM_BUFFER_SIZE=<nbytes>`

Arguments

<nbytes>	The size of shared memory buffers in bytes
> 0	The default <nbytes> value depends on the specific platform you use. The value range is from 8K to 64K typically.

Description

Set this environment variable to define the size of shared memory buffers for each pair of processes.

I_MPI_SSHM_DYNAMIC_CONNECTION

Control the dynamic connection establishment for the alternative shared memory mechanism.

Syntax

`I_MPI_SSHM_DYNAMIC_CONNECTION=<arg>`

Arguments

<arg>	Binary indicator
enable yes on 1	Turn on the dynamic connection establishment
disable no off 0	Turn off the dynamic connection establishment. This is the default value

Description

Set this environment variable to control the dynamic connection establishment.

- If this mode is enabled, all connections are established at the time of the first communication between each pair of processes.
- If this mode is disabled, all connections are established upfront.

I_MPI_SHM_BYPASS

(I_MPI_INTRANODE_SHMEM_BYPASS, I_MPI_USE_DAPL_INTRANODE)

Turn on/off the intra-node communication mode through network fabric along with `shm`.

Syntax

`I_MPI_SHM_BYPASS=<arg>`

Deprecated Syntaxes

`I_MPI_INTRANODE_SHMEM_BYPASS=<arg>`

`I_MPI_USE_DAPL_INTRANODE=<arg>`

Arguments

<code><arg></code>	Binary indicator
<code>enable yes on 1</code>	Turn on the intra-node communication through network fabric
<code>disable no off 0</code>	Turn off the intra-node communication through network fabric. This is the default

Description

Set this environment variable to specify the communication mode within the node. If the intra-node communication mode through network fabric is enabled, data transfer algorithms are selected according to the following scheme:

- Messages shorter than or equal in size to the threshold value of the `I_MPI_INTRANODE_EAGER_THRESHOLD` environment variable are transferred using shared memory.
- Messages larger than the threshold value of the `I_MPI_INTRANODE_EAGER_THRESHOLD` environment variable are transferred through the network fabric layer.

NOTE

This environment variable is applicable only when you turn on shared memory and a network fabric either by default or by setting the `I_MPI_FABRICS` environment variable to `shm:<fabric>` or an equivalent `I_MPI_DEVICE` setting. This mode is available only for `dapl` and `tcp` fabrics.

I_MPI_SHM_SPIN_COUNT

Control the spin count value for the shared memory fabric.

Syntax

`I_MPI_SHM_SPIN_COUNT=<shm_scount>`

Arguments

<code><scout></code>	Define the spin count of the loop when polling the <code>shm</code> fabric
<code>> 0</code>	When internode communication uses the <code>tcp</code> fabric, the default <code><shm_scout></code> value is equal to 100 spins When internode communication uses the <code>ofa</code> , <code>tmi</code> , <code>ofi</code> or <code>dapl</code> fabric, the default <code><shm_scout></code> value is equal to 10 spins. The maximum value is equal to 2147483647

Description

Set the spin count limit of the shared memory fabric to increase the frequency of polling. This configuration allows polling of the `shm` fabric `<shm_scout>` times before the control is passed to the overall network fabric polling mechanism.

To tune application performance, use the `I_MPI_SHM_SPIN_COUNT` environment variable. The best value for `<shm_scout>` can be chosen on an experimental basis. It depends largely on the application and the particular computation environment. An increase in the `<shm_scout>` value benefits multi-core platforms when the application uses topological algorithms for message passing.

3.3.3. DAPL-capable Network Fabrics Control

I_MPI_DAPL_PROVIDER

Define the DAPL provider to load.

Syntax

`I_MPI_DAPL_PROVIDER=<name>`

Arguments

<code><name></code>	Define the name of DAPL provider to load
---------------------------	--

Description

Set this environment variable to define the name of DAPL provider to load. This name is also defined in the `dat.conf` configuration file.

I_MPI_DAT_LIBRARY

Select the DAT library to be used for DAPL* provider.

Syntax

`I_MPI_DAT_LIBRARY=<library>`

Arguments

<code><library></code>	Specify the DAT library for DAPL provider to be used. Default values are <code>libdat.so</code> or <code>libdat.so.1</code> for DAPL* 1.2 providers and <code>libdat2.so</code> or <code>libdat2.so.2</code> for DAPL* 2.0 providers
------------------------------	--

Description

Set this environment variable to select a specific DAT library to be used for DAPL provider. If the library is not located in the dynamic loader search path, specify the full path to the DAT library. This environment variable affects only on DAPL and DAPL UD capable fabrics.

I_MPI_DAPL_TRANSLATION_CACHE (I_MPI_RDMA_TRANSLATION_CACHE)

Turn on/off the memory registration cache in the DAPL path.

Syntax

I_MPI_DAPL_TRANSLATION_CACHE=<arg>

Deprecated Syntax

I_MPI_RDMA_TRANSLATION_CACHE=<arg>

Arguments

<arg>	Binary indicator
enable yes on 1	Turn on the memory registration cache. This is the default
disable no off 0	Turn off the memory registration cache

Description

Set this environment variable to turn on/off the memory registration cache in the DAPL path.

The cache substantially increases performance, but may lead to correctness issues in certain situations. See product *Release Notes* for further details.

I_MPI_DAPL_TRANSLATION_CACHE_AVL_TREE

Enable/disable the AVL tree* based implementation of the RDMA translation cache in the DAPL path.

Syntax

I_MPI_DAPL_TRANSLATION_CACHE_AVL_TREE=<arg>

Arguments

<arg>	Binary indicator
enable yes on 1	Turn on the AVL tree based RDMA translation cache
disable no off 0	Turn off the AVL tree based RDMA translation cache. This is the default value

Description

Set this environment variable to enable the AVL tree based implementation of RDMA translation cache in the DAPL path. When the search in RDMA translation cache handles over 10,000 elements, the AVL tree based RDMA translation cache is faster than the default implementation.

I_MPI_DAPL_DIRECT_COPY_THRESHOLD (I_MPI_RDMA_EAGER_THRESHOLD, RDMA_IBA_EAGER_THRESHOLD)

Change the threshold of the DAPL direct-copy protocol.

Syntax

I_MPI_DAPL_DIRECT_COPY_THRESHOLD=<nbytes>

Deprecated Syntaxes

I_MPI_RDMA_EAGER_THRESHOLD=<nbytes>

RDMA_IBA_EAGER_THRESHOLD=<nbytes>

Arguments

<code><nbytes></code>	Define the DAPL direct-copy protocol threshold
<code>> 0</code>	The default <code><nbytes></code> value depends on the platform

Description

Set this environment variable to control the DAPL direct-copy protocol threshold. Data transfer algorithms for the DAPL-capable network fabrics are selected based on the following scheme:

- Messages shorter than or equal to `<nbytes>` are sent using the eager protocol through the internal pre-registered buffers. This approach is faster for short messages.
- Messages larger than `<nbytes>` are sent using the direct-copy protocol. It does not use any buffering but involves registration of memory on sender and receiver sides. This approach is faster for large messages.

This environment variable is available for both Intel® and non-Intel microprocessors, but it may perform additional optimizations for Intel microprocessors than it performs for non-Intel microprocessors.

NOTE

The equivalent of this variable for Intel® Xeon Phi™ Coprocessor is

`I_MIC_MPI_DAPL_DIRECT_COPY_THRESHOLD`

I_MPI_DAPL_EAGER_MESSAGE_AGGREGATION

Control the use of concatenation for adjourned MPI send requests. Adjourned MPI send requests are those that cannot be sent immediately.

Syntax

`I_MPI_DAPL_EAGER_MESSAGE_AGGREGATION=<arg>`

Arguments

<code><arg></code>	Binary indicator
<code>enable yes on 1</code>	Enable the concatenation for adjourned MPI send requests
<code>disable no off 0</code>	Disable the concatenation for adjourned MPI send requests. This is the default value

Set this environment variable to control the use of concatenation for adjourned MPI send requests intended for the same MPI rank. In some cases, this mode can improve the performance of applications, especially when `MPI_Isend()` is used with short message sizes and the same destination rank, such as:

```
for( i = 0; i< NMSG; i++)
{ret = MPI_Isend( sbuf[i], MSG_SIZE, datatype, dest , tag, \
comm, &req_send[i]);
}
```

I_MPI_DAPL_DYNAMIC_CONNECTION_MODE (I_MPI_DYNAMIC_CONNECTION_MODE, I_MPI_DYNAMIC_CONNECTIONS_MODE)

Choose the algorithm for establishing the DAPL* connections.

Syntax

```
I_MPI_DAPL_DYNAMIC_CONNECTION_MODE=<arg>
```

Deprecated Syntax

```
I_MPI_DYNAMIC_CONNECTION_MODE=<arg>
```

```
I_MPI_DYNAMIC_CONNECTIONS_MODE=<arg>
```

Arguments

<arg>	Mode selector
reject	Deny one of the two simultaneous connection requests. This is the default
disconnect	Deny one of the two simultaneous connection requests after both connections have been established

Description

Set this environment variable to choose the algorithm for handling dynamically established connections for DAPL-capable fabrics according to the following scheme:

- In the `reject` mode, if two processes initiate the connection simultaneously, one of the requests is rejected.
- In the `disconnect` mode, both connections are established, but then one is disconnected. The `disconnect` mode is provided to avoid a bug in certain DAPL* providers.

I_MPI_DAPL_SCALABLE_PROGRESS (I_MPI_RDMA_SCALABLE_PROGRESS)

Turn on/off scalable algorithm for DAPL read progress.

Syntax

```
I_MPI_DAPL_SCALABLE_PROGRESS=<arg>
```

Deprecated Syntax

```
I_MPI_RDMA_SCALABLE_PROGRESS=<arg>
```

Arguments

<arg>	Binary indicator
enable yes on 1	Turn on scalable algorithm. When the number of processes is larger than 128, this is the default value
disable no off 0	Turn off scalable algorithm. When the number of processes is less than or equal to 128, this is the default value

Description

Set this environment variable to enable scalable algorithm for the DAPL read progress. In some cases, this provides advantages for systems with many processes.

I_MPI_DAPL_BUFFER_NUM (I_MPI_RDMA_BUFFER_NUM, NUM_RDMA_BUFFER)

Change the number of internal pre-registered buffers for each process pair in the DAPL path.

Syntax

```
I_MPI_DAPL_BUFFER_NUM=<nbuf>
```

Deprecated Syntaxes

```
I_MPI_RDMA_BUFFER_NUM=<nbuf>
```

```
NUM_RDMA_BUFFER=<nbuf>
```

Arguments

<nbuf>	Define the number of buffers for each pair in a process group
> 0	The default value depends on the platform

Description

Set this environment variable to change the number of the internal pre-registered buffers for each process pair in the DAPL path.

NOTE

The more pre-registered buffers are available, the more memory is used for every established connection.

I_MPI_DAPL_BUFFER_SIZE**(I_MPI_RDMA_BUFFER_SIZE, I_MPI_RDMA_VBUF_TOTAL_SIZE)**

Change the size of internal pre-registered buffers for each process pair in the DAPL path.

Syntax

```
I_MPI_DAPL_BUFFER_SIZE=<nbytes>
```

Deprecated Syntaxes

```
I_MPI_RDMA_BUFFER_SIZE=<nbytes>
```

```
I_MPI_RDMA_VBUF_TOTAL_SIZE=<nbytes>
```

Arguments

<nbytes>	Define the size of pre-registered buffers
> 0	The default value depends on the platform

Description

Set this environment variable to define the size of the internal pre-registered buffer for each process pair in the DAPL path. The actual size is calculated by adjusting the <nbytes> to align the buffer to an optimal value.

I_MPI_DAPL_RNDV_BUFFER_ALIGNMENT**(I_MPI_RDMA_RNDV_BUFFER_ALIGNMENT, I_MPI_RDMA_RNDV_BUF_ALIGN)**

Define the alignment of the sending buffer for the DAPL direct-copy transfers.

Syntax

```
I_MPI_DAPL_RNDV_BUFFER_ALIGNMENT=<arg>
```

Deprecated Syntaxes

```
I_MPI_RDMA_RNDV_BUFFER_ALIGNMENT=<arg>
```

```
I_MPI_RDMA_RNDV_BUF_ALIGN=<arg>
```

Arguments

<code><arg></code>	Define the alignment for the sending buffer
<code>> 0</code> and a power of 2	The default value is 64

Set this environment variable to define the alignment of the sending buffer for DAPL direct-copy transfers. When a buffer specified in a DAPL operation is aligned to an optimal value, the data transfer bandwidth may be increased.

I_MPI_DAPL_RDMA_RNDV_WRITE (**I_MPI_RDMA_RNDV_WRITE, I_MPI_USE_RENDEZVOUS_RDMA_WRITE**)

Turn on/off the RDMA Write-based rendezvous direct-copy protocol in the DAPL path.

Syntax

`I_MPI_DAPL_RDMA_RNDV_WRITE=<arg>`

Deprecated Syntaxes

`I_MPI_RDMA_RNDV_WRITE=<arg>`

`I_MPI_USE_RENDEZVOUS_RDMA_WRITE=<arg>`

Arguments

<code><arg></code>	Binary indicator
<code>enable yes on 1</code>	Turn on the RDMA Write rendezvous direct-copy protocol
<code>disable no off 0</code>	Turn off the RDMA Write rendezvous direct-copy protocol

Description

Set this environment variable to select the RDMA Write-based rendezvous direct-copy protocol in the DAPL path. Certain DAPL* providers have a slow RDMA Read implementation on certain platforms. Switching on the rendezvous direct-copy protocol based on the RDMA Write operation can increase performance in these cases. The default value depends on the DAPL provider attributes.

I_MPI_DAPL_CHECK_MAX_RDMA_SIZE (**I_MPI_RDMA_CHECK_MAX_RDMA_SIZE**)

Check the value of the DAPL attribute `max_rdma_size`.

Syntax

`I_MPI_DAPL_CHECK_MAX_RDMA_SIZE=<arg>`

Deprecated Syntax

`I_MPI_RDMA_CHECK_MAX_RDMA_SIZE=<arg>`

Arguments

<code><arg></code>	Binary indicator
<code>enable yes on 1</code>	Check the value of the DAPL* attribute <code>max_rdma_size</code>
<code>disable no off 0</code>	Do not check the value of the DAPL* attribute <code>max_rdma_size</code> . This is the default value

Description

Set this environment variable to control message fragmentation according to the following scheme:

- If this mode is enabled, the Intel® MPI Library fragmentizes the messages bigger than the value of the DAPL attribute `max_rdma_size`
- If this mode is disabled, the Intel® MPI Library does not take into account the value of the DAPL attribute `max_rdma_size` for message fragmentation

I_MPI_DAPL_MAX_MSG_SIZE (I_MPI_RDMA_MAX_MSG_SIZE)

Control message fragmentation threshold.

Syntax

`I_MPI_DAPL_MAX_MSG_SIZE=<nbytes>`

Deprecated Syntax

`I_MPI_RDMA_MAX_MSG_SIZE=<nbytes>`

Arguments

<code><nbytes></code>	Define the maximum message size that can be sent through DAPL without fragmentation
<code>> 0</code>	If the <code>I_MPI_DAPL_CHECK_MAX_RDMA_SIZE</code> environment variable is enabled, the default <code><nbytes></code> value is equal to the <code>max_rdma_size</code> DAPL attribute value. Otherwise the default value is <code>MAX_INT</code>

Description

Set this environment variable to control message fragmentation size according to the following scheme:

- If the `I_MPI_DAPL_CHECK_MAX_RDMA_SIZE` environment variable is set to `disable`, the Intel® MPI Library fragmentizes the messages whose sizes are greater than `<nbytes>`.
- If the `I_MPI_DAPL_CHECK_MAX_RDMA_SIZE` environment variable is set to `enable`, the Intel® MPI Library fragmentizes the messages whose sizes are greater than the minimum of `<nbytes>` and the `max_rdma_size` DAPL* attribute value.

I_MPI_DAPL_CONN_EVD_SIZE (I_MPI_RDMA_CONN_EVD_SIZE, I_MPI_CONN_EVD_QLEN)

Define the event queue size of the DAPL event dispatcher for connections.

Syntax

`I_MPI_DAPL_CONN_EVD_SIZE=<size>`

Deprecated Syntaxes

`I_MPI_RDMA_CONN_EVD_SIZE=<size>`

`I_MPI_CONN_EVD_QLEN=<size>`

Arguments

<code><size></code>	Define the length of the event queue
<code>> 0</code>	The default value is <code>2*number of processes + 32</code> in the MPI job

Description

Set this environment variable to define the event queue size of the DAPL event dispatcher that handles connection related events. If this environment variable is set, the minimum value between *<size>* and the value obtained from the provider is used as the size of the event queue. The provider is required to supply a queue size that equal or larger than the calculated value.

I_MPI_DAPL_SR_THRESHOLD

Change the threshold of switching send/recv to *rdma* path for DAPL wait mode.

Syntax

I_MPI_DAPL_SR_THRESHOLD=*<arg>*

Arguments

<i><nbytes></i>	Define the message size threshold of switching send/recv to <i>rdma</i>
<i>>= 0</i>	The default <i><nbytes></i> value is 256 bytes

Description

Set this environment variable to control the protocol used for point-to-point communication in DAPL wait mode:

- Messages shorter than or equal in size to *<nbytes>* are sent using DAPL send/recv data transfer operations.
- Messages greater in size than *<nbytes>* are sent using DAPL RDMA WRITE or RDMA WRITE immediate data transfer operations.

I_MPI_DAPL_SR_BUF_NUM

Change the number of internal pre-registered buffers for each process pair used in DAPL wait mode for send/recv path.

Syntax

I_MPI_DAPL_SR_BUF_NUM=*<nbuf>*

Arguments

<i><nbuf></i>	Define the number of send/recv buffers for each pair in a process group
<i>> 0</i>	The default value is 32

Description

Set this environment variable to change the number of the internal send/recv pre-registered buffers for each process pair.

I_MPI_DAPL_RDMA_WRITE_IMM (I_MPI_RDMA_WRITE_IMM)

Enable/disable RDMA Write with immediate data InfiniBand (IB) extension in DAPL wait mode.

Syntax

I_MPI_DAPL_RDMA_WRITE_IMM=*<arg>*

Deprecated syntax

I_MPI_RDMA_WRITE_IMM=*<arg>*

Arguments

<code><arg></code>	Binary indicator
<code>enable yes on 1</code>	Turn on RDMA Write with immediate data IB extension
<code>disable no off 0</code>	Turn off RDMA Write with immediate data IB extension

Description

Set this environment variable to utilize RDMA Write with immediate data IB extension. The algorithm is enabled if this environment variable is set and a certain DAPL provider attribute indicates that RDMA Write with immediate data IB extension is supported.

I_MPI_DAPL_DESIRED_STATIC_CONNECTIONS_NUM

Define the number of processes that establish DAPL static connections at the same time.

Syntax

`I_MPI_DAPL_DESIRED_STATIC_CONNECTIONS_NUM=<num_procesess>`

Arguments

<code><num_procesess></code>	Define the number of processes that establish DAPL static connections at the same time
<code>> 0</code>	The default <code><num_procesess></code> value is equal to 256

Description

Set this environment variable to control the algorithm of DAPL static connection establishment.

If the number of processes in the MPI job is less than or equal to `<num_procesess>`, all MPI processes establish the static connections simultaneously. Otherwise, the processes are distributed into several groups. The number of processes in each group is calculated to be close to `<num_procesess>`. Then static connections are established in several iterations, including intergroup connection setup.

I_MPI_CHECK_DAPL_PROVIDER_COMPATIBILITY

Enable/disable the check that the same DAPL provider is selected by all ranks.

Syntax

`I_MPI_CHECK_DAPL_PROVIDER_COMPATIBILITY=<arg>`

Arguments

<code><arg></code>	Binary indicator
<code>enable yes on 1</code>	Turn on the check that the DAPL provider is the same on all ranks. This is default value
<code>disable no off 0</code>	Turn off the check that the DAPL provider is the same on all ranks

Description

Set this variable to make a check if the DAPL provider is selected by all MPI ranks. If this check is enabled, Intel® MPI Library checks the name of DAPL provider and the version of DAPL. If these parameters are not

the same on all ranks, Intel MPI Library does not select the RDMA path and may fall to sockets. Turning off the check reduces the execution time of `MPI_Init()`. It may be significant for MPI jobs with a large number of processes.

3.3.4. DAPL UD-capable Network Fabrics Control

I_MPI_DAPL_UD

Enable/disable using DAPL UD path.

Syntax

`I_MPI_DAPL_UD=<arg>`

Arguments

<code><arg></code>	Binary indicator
<code>enable yes on 1</code>	Turn on using DAPL UD IB extension
<code>disable no off 0</code>	Turn off using DAPL UD IB extension. This is the default value

Description

Set this environment variable to enable DAPL UD path for transferring data. The algorithm is enabled if you set this environment variable and a certain DAPL provider attribute indicates that UD IB extension is supported.

I_MPI_DAPL_UD_PROVIDER

Define the DAPL provider to work with IB UD transport.

Syntax

`I_MPI_DAPL_UD_PROVIDER=<name>`

Arguments

<code><name></code>	Define the name of DAPL provider to load
---------------------------	--

Description

Set this environment variable to define the name of DAPL provider to load. This name is also defined in the `dat.conf` configuration file. Make sure that specified DAPL provider supports UD IB extension.

I_MPI_DAPL_UD_DIRECT_COPY_THRESHOLD

Change the message size threshold of the DAPL UD direct-copy protocol.

Syntax

`I_MPI_DAPL_UD_DIRECT_COPY_THRESHOLD=<nbytes>`

Arguments

<code><nbytes></code>	Define the DAPL UD direct-copy protocol threshold
<code>> 0</code>	The default <code><nbytes></code> value is equal to 16456 bytes

Description

Set this environment variable to control the DAPL UD direct-copy protocol threshold. Data transfer algorithms for the DAPL-capable network fabrics are selected based on the following scheme:

- Messages shorter than or equal to `<nbytes>` are sent using the eager protocol through the internal pre-registered buffers. This approach is faster for short messages.
- Messages larger than `<nbytes>` are sent using the direct-copy protocol. It does not use any buffering but involves registration of memory on sender and receiver sides. This approach is faster for large messages.

This environment variable is available for both Intel® and non-Intel microprocessors, but it may perform additional optimizations for Intel microprocessors than it performs for non-Intel microprocessors.

I_MPI_DAPL_UD_RECV_BUFFER_NUM

Change the number of the internal pre-registered UD buffers for receiving messages.

Syntax

`I_MPI_DAPL_UD_RECV_BUFFER_NUM=<nbuf>`

Arguments

<code><nbuf></code>	Define the number of buffers for receiving messages
<code>> 0</code>	The default value is $16 + n \cdot 4$ where n is a total number of process in MPI job

Description

Set this environment variable to change the number of the internal pre-registered buffers for receiving messages. These buffers are common for all connections or process pairs.

NOTE

The pre-registered buffers use up memory, however they help avoid the loss of packets.

I_MPI_DAPL_UD_SEND_BUFFER_NUM

Change the number of internal pre-registered UD buffers for sending messages.

Syntax

`I_MPI_DAPL_UD_SEND_BUFFER_NUM=<nbuf>`

Arguments

<code><nbuf></code>	Define the number of buffers for sending messages
<code>> 0</code>	The default value is $16 + n \cdot 4$ where n is a total number of process in MPI job

Description

Set this environment variable to change the number of the internal pre-registered buffers for sending messages. These buffers are common for all connections or process pairs.

NOTE

The pre-registered buffers use up memory, however they help avoid the loss of packets.

I_MPI_DAPL_UD_ACK_SEND_POOL_SIZE

Change the number of ACK UD buffers for sending messages.

Syntax

`I_MPI_DAPL_UD_ACK_SEND_POOL_SIZE=<nbuf>`

Arguments

<code><nbuf></code>	Define the number of ACK UD buffers for sending messages
<code>> 0</code>	The default value is 256

Description

Set this environment variable to change the number of the internal pre-registered ACK buffers for sending service messages. These buffers are common for all connections or process pairs.

I_MPI_DAPL_UD_ACK_RECV_POOL_SIZE

Change the number of ACK UD buffers for receiving messages.

Syntax

`I_MPI_DAPL_UD_ACK_RECV_POOL_SIZE=<nbuf>`

Arguments

<code><nbuf></code>	Define the number of ACK UD buffers for receiving messages
<code>> 0</code>	The default value is $512+n*4$, where n is total number of process in MPI job

Description

Set this environment variable to change the number of the internal pre-registered ACK buffers for receiving service messages. These buffers are common for all connections or process pairs.

I_MPI_DAPL_UD_TRANSLATION_CACHE

Turn on/off the memory registration cache in the DAPL UD path.

Syntax

`I_MPI_DAPL_UD_TRANSLATION_CACHE=<arg>`

Arguments

<code><arg></code>	Binary indicator
<code>enable yes on 1</code>	Turn on the memory registration cache. This is the default
<code>disable no off 0</code>	Turn off the memory registration cache

Description

Set this environment variable to turn off the memory registration cache in the DAPL UD path. Using the cache substantially improves performance. See product *Release Notes* for further details.

I_MPI_DAPL_UD_TRANSLATION_CACHE_AVL_TREE

Enable/disable the AVL* tree based implementation of RDMA translation cache in the DAPL UD path.

Syntax

`I_MPI_DAPL_UD_TRANSLATION_CACHE_AVL_TREE=<arg>`

Arguments

<code><arg></code>	Binary indicator
--------------------------	------------------

enable yes on 1	Turn on the AVL tree based RDMA translation cache
disable no off 0	Turn off the AVL tree based RDMA translation cache. This is the default value

Description

Set this environment variable to enable the AVL tree based implementation of RDMA translation cache in the DAPL UD path. When the search in RDMA translation cache handles over 10,000 elements, the AVL tree based RDMA translation cache is faster than the default implementation.

I_MPI_DAPL_UD_REQ_EVD_SIZE

Define the event queue size of the DAPL UD event dispatcher for sending data transfer operations.

Syntax

I_MPI_DAPL_UD_REQ_EVD_SIZE=<size>

Arguments

<size>	Define the length of the event queue
> 0	The default value is 2,000

Description

Set this environment variable to define the event queue size of the DAPL event dispatcher that handles completions of sending DAPL UD data transfer operations (DTO). If this environment variable is set, the minimum value between <size> and the value obtained from the provider is used as the size of the event queue. The provider is required to supply a queue size that is at least equal to the calculated value, but it can also provide a larger queue size.

I_MPI_DAPL_UD_CONN_EVD_SIZE

Define the event queue size of the DAPL UD event dispatcher for connections.

Syntax

I_MPI_DAPL_UD_CONN_EVD_SIZE=<size>

Arguments

<size>	Define the length of the event queue
> 0	The default value is 2*number of processes + 32

Description

Set this environment variable to define the event queue size of the DAPL event dispatcher that handles connection related events. If this environment variable is set, the minimum value between <size> and the value obtained from the provider is used as the size of the event queue. The provider is required to supply a queue size that is at least equal to the calculated value, but it can also provide a larger queue size.

I_MPI_DAPL_UD_RECV_EVD_SIZE

Define the event queue size of the DAPL UD event dispatcher for receiving data transfer operations.

Syntax

I_MPI_DAPL_UD_RECV_EVD_SIZE=<size>

Arguments

<code><size></code>	Define the length of the event queue
<code>> 0</code>	The default value depends on the number UD and ACK buffers

Description

Set this environment variable to define the event queue size of the DAPL event dispatcher that handles completions of receiving DAPL UD data transfer operations (DTO). If this environment variable is set, the minimum value between `<size>` and the value obtained from the provider is used as the size of the event queue. The provider is required to supply a queue size that is at least equal to the calculated value, but it can also provide a larger queue size.

I_MPI_DAPL_UD_RNDV_MAX_BLOCK_LEN

Define maximum size of block that is passed at one iteration of DAPL UD direct-copy protocol.

Syntax

`I_MPI_DAPL_UD_RNDV_MAX_BLOCK_LEN=<nbytes>`

Arguments

<code><arg></code>	Define maximum size of block that is passed at one iteration of DAPL UD direct-copy protocol
<code>> 0</code>	The default value is 1,048,576

Set this environment variable to define the maximum size of memory block that is passed at one iteration of DAPL UD direct-copy protocol. If the size of message in direct-copy protocol is greater than given value, the message will be divided in several blocks and passed in several operations.

I_MPI_DAPL_UD_RNDV_BUFFER_ALIGNMENT

Define the alignment of the sending buffer for the DAPL UD direct-copy transfers.

Syntax

`I_MPI_DAPL_UD_RNDV_BUFFER_ALIGNMENT=<arg>`

Arguments

<code><arg></code>	Define the alignment of the sending buffer
<code>> 0</code> and a power of 2	The default value is 16

Set this environment variable to define the alignment of the sending buffer for DAPL direct-copy transfers. When a buffer specified in a DAPL operation is aligned to an optimal value, this may increase data transfer bandwidth.

I_MPI_DAPL_UD_RNDV_COPY_ALIGNMENT_THRESHOLD

Define threshold where alignment is applied to send buffer for the DAPL UD direct-copy transfers.

Syntax

`I_MPI_DAPL_UD_RNDV_COPY_ALIGNMENT_THRESHOLD=<nbytes>`

Arguments

<code><nbytes></code>	Define send buffer alignment threshold
-----------------------------	--

> 0 and a power of 2	The default value is 32,768
----------------------	-----------------------------

Set this environment variable to define the threshold where the alignment of the sending buffer is applied in DAPL direct-copy transfers. When a buffer specified in a DAPL operation is aligned to an optimal value, this may increase data transfer bandwidth.

I_MPI_DAPL_UD_RNDV_DYNAMIC_CONNECTION

Control the algorithm of dynamic connection establishment for DAPL UD endpoints used in the direct copy protocol.

Syntax

I_MPI_DAPL_UD_RNDV_DYNAMIC_CONNECTION=<arg>

Arguments

<arg>	Binary indicator
enable yes on 1	Turns on the dynamic connection mode. This is the default value
disable no off 0	Turns off the dynamic connections mode

Set this variable to control the dynamic connection establishment of DAPL UD endpoints used in the direct copy protocol.

If you disable the dynamic connection mode, all possible connections are established during the MPI startup phase.

If you enable the mode, the connection is established when an application calls the MPI function to pass the data from one process to another and invokes the communication between the two processes.

NOTE

For the RNDV dynamic connection mode, the size of the messages passed in the data is larger than the value you set in the I_MPI_DAPL_UD_DIRECT_COPY_THRESHOLD environment variable.

I_MPI_DAPL_UD_EAGER_DYNAMIC_CONNECTION

Control the algorithm of the dynamic connection establishment for DAPL UD endpoints used in eager protocol.

Syntax

I_MPI_DAPL_UD_EAGER_DYNAMIC_CONNECTION=<arg>

Arguments

<arg>	Binary indicator
enable yes on 1	Turn on the dynamic connection mode. If the number of processes is over 64, this is the default value
disable no off 0	Turn off the dynamic connections mode

Set this variable to control the dynamic connection establishment of DAPL UD endpoints involved in eager protocol. Eager protocol is used to transfer messages through internal pre-registered buffers.

If you disable this mode, all possible connections are established during MPI startup phase.

If you enable this mode, the connection is established when an application calls the MPI function to pass the data from one process to another and invokes the communication between the two processes.

NOTE

For the eager dynamic connection mode, the size of the messages passed in the data is shorter than or equal to the value you set in the `I_MPI_DAPL_UD_DIRECT_COPY_THRESHOLD` environment variable.

I_MPI_DAPL_UD_DESIRED_STATIC_CONNECTIONS_NUM

Define the number of processes that establish DAPL static connections at the same time.

Syntax

```
I_MPI_DAPL_UD_DESIRED_STATIC_CONNECTIONS_NUM=<num_procesess>
```

Arguments

<code><num_procesess></code>	Define the number of processes that establish DAPL UD static connections at the same time
<code>> 0</code>	The default value is equal to 200

Description

Set this environment variable to control the algorithm of DAPL UD static connections establishment.

If the number of processes in an MPI job is less than or equal to `<num_procesess>`, all MPI processes establish the static connections simultaneously. Otherwise, the processes are distributed into several groups. The number of processes in each group is calculated to be close to `<num_procesess>`. Then static connections are established in several iterations, including intergroup connection setup.

I_MPI_DAPL_UD_RDMA_MIXED

Control the use of the DAPL UD/RDMA mixed communication.

Syntax

```
I_MPI_DAPL_UD_RDMA_MIXED =<arg>
```

Arguments

<code><arg></code>	Binary indicator
<code>enable yes on 1</code>	Turn on the use of DAPL UD/RDMA mixed communication
<code>disable no off 0</code>	Turn off the use of DAPL UD/RDMA mixed communication. This is the default value

Description

Set this environment variable to enable the DAPL UD/RDMA mixed mode for transferring data. In the DAPL UD/RDMA mixed mode, small messages are passed through the UD transport and large messages are passed through the RDMA transport. If you set the `I_MPI_DAPL_UD_RDMA_MIXED` environment variable and a certain DAPL provider attribute indicates that UD IB extension is supported, the DAPL UD/RDMA mixed mode is enabled.

The following set of `I_MPI_DAPL_UD*` environment variables also controls the DAPL UD/RDMA mixed mode:

- `I_MPI_DAPL_UD_PROVIDER`
- `I_MPI_DAPL_UD_EAGER_DYNAMIC_CONNECTION`
- `I_MPI_DAPL_UD_RNDV_DYNAMIC_CONNECTION`

- I_MPI_DAPL_UD_DIRECT_COPY_THRESHOLD
- I_MPI_DAPL_UD_RECV_BUFFER_NUM
- I_MPI_DAPL_UD_SEND_BUFFER_NUM
- I_MPI_DAPL_UD_NUMBER_CREDIT_UPDATE
- I_MPI_DAPL_UD_ACK_SEND_POOL_SIZE
- I_MPI_DAPL_UD_ACK_RECV_POOL_SIZE
- I_MPI_DAPL_UD_RESENT_TIMEOUT
- I_MPI_DAPL_UD_MAX_MSG_SIZE
- I_MPI_DAPL_UD_SEND_BUFFER_SIZE
- I_MPI_DAPL_UD_REQ_EVD_SIZE
- I_MPI_DAPL_UD_REQUEST_QUEUE_SIZE
- I_MPI_DAPL_UD_MULTIPLE_EAGER_SEND
- I_MPI_DAPL_UD_NA_SBUF_LIMIT
- I_MPI_DAPL_UD_RECV_EVD_SIZE
- I_MPI_DAPL_UD_CONNECTION_TIMEOUT
- I_MPI_DAPL_UD_PORT
- I_MPI_DAPL_UD_CREATE_CONN_QUAL,
- I_MPI_DAPL_UD_FINALIZE_RETRY_COUNT
- I_MPI_DAPL_UD_FINALIZE_TIMEOUT
- I_MPI_DAPL_UD_TRANSLATION_CACHE
- I_MPI_DAPL_UD_TRANSLATION_CACHE_AVL_TREE
- I_MPI_DAPL_UD_TRANSLATION_CACHE_MAX_ENTRY_NUM
- I_MPI_DAPL_UD_TRANSLATION_CACHE_MAX_MEMORY_SIZE
- I_MPI_DAPL_UD_PKT_LOSS_OPTIMIZATION
- I_MPI_DAPL_UD_DFACTOR
- I_MPI_DAPL_UD_DESIRED_STATIC_CONNECTIONS_NUM
- I_MPI_DAPL_UD_CONN_EVD_SIZE
- I_MPI_DAPL_UD_RNDV_BUFFER_ALIGNMENT
- I_MPI_DAPL_UD_RNDV_COPY_ALIGNMENT_THRESHOLD

The following set of environment variables is specific for DAPL UD/RDMA mixed mode:

- I_MPI_DAPL_UD_MAX_RDMA_SIZE
- I_MPI_DAPL_UD_MAX_RDMA_DTOS

I_MPI_DAPL_UD_MAX_RDMA_SIZE

Control the maximum message size that can be sent though the RDMA for DAPL UD/RDMA mixed mode.

Syntax

```
I_MPI_DAPL_UD_MAX_RDMA_SIZE = <nbytes>
```

Arguments

<nbytes>	Define the maximum message size that can be sent through RDMA without fragmentation
> 0	The default <nbytes> value is 4 MB

Description

Set this environment variable to define the maximum message size that can be sent through RDMA protocol for the DAPL UD/RDMA mixed mode. If the message size is greater than this value, this message is divided into several fragments and is sent by several RDMA operations.

I_MPI_DAPL_UD_MAX_RDMA_DTOS

Control the maximum number of uncompleted RDMA operations per connection for the DAPL UD/RDMA mixed mode.

Syntax

```
I_MPI_DAPL_UD_MAX_RDMA_DTOS=<arg>
```

Arguments

<arg>	Define the maximum number of RDMA operations per connection
> 0	The default <arg> value is 8

Description

Set this environment variable to define the maximum number of RDMA operations per connection for the DAPL UD/RDMA mixed mode.

3.3.5. TCP-capable Network Fabrics Control**I_MPI_TCP_NETMASK
(I_MPI_NETMASK)**

Choose the network interface for MPI communication over TCP-capable network fabrics.

Syntax

```
I_MPI_TCP_NETMASK=<arg>
```

Arguments

<arg>	Define the network interface (string parameter)
<interface_mnemonic>	Mnemonic of the network interface: <code>ib</code> or <code>eth</code>
<code>ib</code>	Use IPoIB* network interface
<code>eth</code>	Use Ethernet network interface. This is the default value
<interface_name>	Name of the network interface Usually the UNIX* driver name followed by the unit number

<code><network_address></code>	Network address. Trailing zero bits imply a netmask
<code><network_address/ <netmask></code>	Network address. The <code><netmask></code> value specifies the netmask length
<code><list of interfaces></code>	A colon separated list of network addresses and interface names

Description

Set this environment variable to choose the network interface for MPI communication over TCP-capable network fabrics. If you specify a list of interfaces, the first available interface on the node is used for communication.

Examples

- Use the following setting to select the IP over InfiniBand* (IPoIB) fabric:
`I_MPI_TCP_NETMASK=ib`
- Use the following setting to select the specified network interface for socket communications:
`I_MPI_TCP_NETMASK=ib0`
- Use the following setting to select the specified network for socket communications. This setting implies the 255.255.0.0 netmask:
`I_MPI_TCP_NETMASK=192.169.0.0`
- Use the following setting to select the specified network for socket communications with netmask set explicitly:
`I_MPI_TCP_NETMASK=192.169.0.0/24`
- Use the following setting to select the specified network interfaces for socket communications:
`I_MPI_TCP_NETMASK=192.169.0.5/24:ib0:192.169.0.0`

I_MPI_TCP_BUFFER_SIZE

Change the size of the TCP socket buffers.

Syntax

`I_MPI_TCP_BUFFER_SIZE=<nbytes>`

Arguments

<code><nbytes></code>	Define the size of the TCP socket buffers
<code>> 0</code>	The default <code><nbytes></code> value is equal to default value of the TCP socket buffer size on your Linux system.

Description

Set this environment variable to manually define the size of the TCP socket buffers. The TCP socket buffer size is restricted by the existing TCP settings on your Linux system.

Use the `I_MPI_TCP_BUFFER_SIZE` environment variable for tuning your application performance for a given number of processes.

NOTE

TCP socket buffers of a large size can require more memory for an application with large number of processes. Alternatively, TCP socket buffers of a small size can considerably decrease the bandwidth of each socket connection especially for 10 Gigabit Ethernet and IPoIB (see [L_MPI_TCP_NETMASK](#) for details).

I_MPI_TCP_POLLING_MODE

Set this environment variable to define a polling mode.

Syntax

`I_MPI_TCP_POLLING_MODE=<mode>`

Arguments

<code><mode></code>	Specify the polling mode
<code>poll</code>	The polling mode based on the <code>poll()</code> function. This is <i>the</i> default value
<code>epoll[:edge]</code>	The polling mode based on the <code>epoll()</code> function as an edge-triggered interface
<code>epoll:level</code>	The polling mode based on the <code>epoll()</code> function as a level-triggered interface

Set this environment variable to select the polling mode for the `tcp` fabric.

Use the `I_MPI_TCP_POLLING_MODE` environment variable for tuning application performance. You can choose the best polling mode on an experimental basis. The best mode depends on the specific application and on the number of processes. The `epoll` polling mode is a preferable mode in the following situations:

- for large number of processes
- for APP client-server type
- for `MPI_ANY_SOURCE` tag matching

3.3.6. TMI-capable Network Fabrics Control**I_MPI_TMI_LIBRARY**

Select the TMI library to be used.

Syntax

`I_MPI_TMI_LIBRARY=<library>`

Arguments

<code><library></code>	Specify a TMI library to be used instead of the default <code>libtmi.so</code>
------------------------------	--

Description

Set this environment variable to select a specific TMI library. Specify the full path to the TMI library if the library does not locate in the dynamic loader search path.

I_MPI_TMI_PROVIDER

Define the name of the TMI provider to load.

Syntax

`I_MPI_TMI_PROVIDER=<name>`

Arguments

<code><name></code>	The name of the TMI provider to load
---------------------------	--------------------------------------

Description

Set this environment variable to define the name of the TMI provider to load. The name must also be defined in the `tmi.conf` configuration file.

I_MPI_TMI_NBITS_RANK

Defines the number of the bits that can be reserved for the storage of MPI rank values at the TMI level.

Syntax

`I_MPI_TMI_NBITS_RANK=<num_bits>`

Arguments

<code><num_bits></code>	The number of the bits reserved for the MPI rank storage
<code><=32 and > 0</code>	The default value is 24

Description

The value of `I_MPI_TMI_NBITS_RANK` specifies how many MPI ranks can be referenced and distinguished at TMI level. Thus, if you specify the default value for this environment variable, `I_MPI_TMI_NBITS_RANK=24`, the number of ranks allowed for running a job is $2^{24}=16\text{M}$ ranks.

NOTE

The value of `I_MPI_TMI_NBITS_RANK` is related to the `MPI_TAG_UB`. The larger value you specify for `I_MPI_TMI_NBITS_RANK`, the less tag value `MPI_TAG_UB` is supported. The less value you specify for `I_MPI_TMI_NBITS_RANK`, the larger tag value `MPI_TAG_UB` is supported. The correct MPI application should always query `MPI_TAG_UB` for the largest supported tag value.

I_MPI_TMI_DSEND

Control the capability of the direct send in the TMI netmod.

Syntax

`I_MPI_TMI_DSEND=<arg>`

Arguments

<code><arg></code>	Binary indicator
<code>enable yes on 1</code>	Enable the direct send. This is default value
<code>disable no off 0</code>	Disable the direct send

Description

Use the direct send capability to block `MPI_Send` calls only. Before using the direct send capability, ensure that you use it for single-threaded MPI applications and check if you have selected TMI as the network fabrics (setting `I_MPI_FABRICS=tmi`).

NOTE

The direct send capability is only supported in the TMI version 1.1 or higher. If you use a lower TMI version, the specified value of `I_MPI_TMI_DSEND` is ignored.

I_MPI_TMI_DRECV

Control the capability of the direct receive in the TMI fabric.

Syntax

`I_MPI_TMI_DRECV=<arg>`

Arguments

<code><arg></code>	Binary indicator
<code>enable yes on 1</code>	Enable the direct receive. This is default value
<code>disable no off 0</code>	Disable the direct receive

Description

Use the direct receive capability to block `MPI_Recv` calls only. Before using the direct receive capability, ensure that you use it for single-threaded MPI applications and check if you have selected TMI as the network fabrics (setting `I_MPI_FABRICS=tmi`).

3.3.7. OFA-capable Network Fabrics Control

I_MPI_OFA_NUM_ADAPTERS

Set the number of connection adapters.

Syntax

`I_MPI_OFA_NUM_ADAPTERS=<arg>`

Arguments

<code><arg></code>	Define the maximum number of connection adapters used
<code>>0</code>	Use the specified number of adapters. The default value is 1

Description

Set the number of the adapters that are used. If the number is greater than the available number of adapters, all the available adaptors are used.

I_MPI_OFA_ADAPTER_NAME

Set the name of adapter that is used.

Syntax

`I_MPI_OFA_ADAPTER_NAME=<arg>`

Arguments

<code><arg></code>	Define the name of adapter
<code>Name</code>	Use the specified adapter. By default, any adapter can be used

Description

Set the name of adaptor to be used. If the adapter with specified name does not exist, the library returns an error. This has effect only if `I_MPI_OFA_NUM_ADAPTERS=1`.

I_MPI_OFA_NUM_PORTS

Set the number of used ports on each adapter.

Syntax

```
I_MPI_OFA_NUM_PORTS=<arg>
```

Arguments

<arg>	Define the number of ports that are used on each adapter
> 0	Use the specified number of ports. The default value is 1

Description

Set the number of used ports on each adaptor. If the number is greater than the available number of ports, all the available ports are used.

I_MPI_OFA_NUM_RDMA_CONNECTIONS

Set the maximum number of connections that can use the `rdma` exchange protocol.

Syntax

```
I_MPI_OFA_NUM_RDMA_CONNECTIONS=<num_conn>
```

Arguments

<num_conn>	Define the maximum number of connections that can use the <code>rdma</code> exchange protocol
>= 0	Create the specified number of connections that use the <code>rdma</code> exchange protocol. All other processes use the send/ receive exchange protocol
-1	Create $\log_2(\text{number of processes})$ <code>rdma</code> connections
>= number of processes	Create <code>rdma</code> connections for all processes. This is the default value

Description

There are two exchange protocols between two processes: send/receive and `rdma`. This environment variable specifies the maximum amount of connections that use `rdma` protocol.

RDMA protocol is faster but requires more resources. For a large application, you can limit the number of connections that use the `rdma` protocol so that only processes that actively exchange data use the `rdma` protocol.

I_MPI_OFA_SWITCHING_TO_RDMA

Set the number of messages that a process should receive before switching this connection to RDMA exchange protocol.

Syntax

```
I_MPI_OFA_SWITCHING_TO_RDMA=<number>
```

Arguments

<number>	Define the number of messages that the process receives before switching to use the <code>rdma</code> protocol
>= 0	If this process receives <number> of messages, start using the <code>rdma</code> protocol

Description

Count the number of messages received from the specific process. The connection achieved the specified number tries to switch to `rdma` protocol for exchanging with that process. The connection will not switch to `rdma` protocol if the maximum number of connections that use the `rdma` exchange protocol defined in `I_MPI_OFA_NUM_RDMA_CONNECTIONS` has been reached.

I_MPI_OFA_RAIL_SCHEDULER

Set the method of choosing rails for short messages.

Syntax

`I_MPI_OFA_RAIL_SCHEDULER=<arg>`

Arguments

<code><arg></code>	Mode selector
<code>ROUND_ROBIN</code>	Next time use next rail
<code>FIRST_RAIL</code>	Always use the first rail for short messages
<code>PROCESS_BIND</code>	Always use the rail specific for process

Description

Set the method of choosing rails for short messages. The algorithms are selected according to the following scheme:

- In the `ROUND_ROBIN` mode, the first message is sent using the first rail; the next message is sent using the second rail, and so on.
- In the `FIRST_RAIL` mode, the first rail is always used for short messages.
- In the `PROCESS_BIND` mode, the process with the smallest rank uses the first rail, and the next uses the second rail.

I_MPI_OFA_TRANSLATION_CACHE

Turn on/off the memory registration cache.

Syntax

`I_MPI_OFA_TRANSLATION_CACHE=<arg>`

Arguments

<code><arg></code>	Binary indicator
<code>enable yes on 1</code>	Turn on the memory registration cache. This is the default
<code>disable no off 0</code>	Turn off the memory registration cache

Description

Set this environment variable to turn on/off the memory registration cache.

The cache substantially increases performance, but may lead to correctness issues in certain situations. See product *Release Notes* for further details.

I_MPI_OFA_TRANSLATION_CACHE_AVL_TREE

Enable/disable the AVL tree* based implementation of the RDMA translation cache.

Syntax

I_MPI_OFA_TRANSLATION_CACHE_AVL_TREE=<arg>

Arguments

<arg>	Binary indicator
enable yes on 1	Turn on the AVL tree based RDMA translation cache
disable no off 0	Turn off the AVL tree based RDMA translation cache. This is the default value

Description

Set this environment variable to enable the AVL tree based implementation of RDMA translation cache in the OFA path. When the search in RDMA translation cache handles over 10,000 elements, the AVL tree based RDMA translation cache is faster than the default implementation.

I_MPI_OFA_DYNAMIC_QPS

Control the library to create queue pairs (QPs) dynamically.

Syntax

I_MPI_OFA_DYNAMIC_QPS=<arg>

Arguments

<arg>	Binary indicator
enable yes on 1	Create QPs dynamically. This is the default value if the number of processes is larger than or equal 2,000
disable no off 0	Create all QPs during the initial stage. This is the default value if the number of processes is less than 2,000

Description

Set this environment variable to turn on dynamic creation of QPs.

I_MPI_OFA_PACKET_SIZE

Set the size of the packet used for sending.

Syntax

I_MPI_OFA_PACKET_SIZE=<arg>

Arguments

<arg>	Define the size of packet in bytes
> 0	Use the specified packet size. The default value is 8192

Description

Set the packet size in bytes. If the number is negative, the size is set to 8.

I_MPI_OFA_LIBRARY

Set the name of the used OFA library.

Syntax

```
I_MPI_OFA_LIBRARY=<arg>
```

Arguments

<arg>	Define the name of the OFA library
Name	Use the specified library. By default, the name is <code>libibverbs.so</code>

Description

Set the name of the InfiniBand* (IB*) library. If the library with the specified name does not exist, an error is returned.

I_MPI_OFA_NONSWITCH_CONF

Define the nonstandard template for port connections.

Syntax

```
I_MPI_OFA_NONSWITCH_CONF=<arg>
```

Arguments

<arg>	Define the template for port connections
Name	Use the specified template

Description

The nodes in clusters are normally connected so that port_i of a node can access port_i of all other nodes. Use this environment variable if ports are not connected in this way. The following example is the template format:

```
host1@port11#port12#...#host2@port21#port22....
```

Port_i defines the port used to send from host_i to host_j

For example:

```
node1@1#1#2#node2@2#1#1#node3@1#2#1#
```

This sample specifies the following configuration:

- Port1 of node1 connected to port2 of node2
- Port2 of node1 connected to port1 of node3
- Port1 of node2 connected to port2 of node3
- Port2 of node2 connected to port1 of node2
- Port1 of node3 connected to port2 of node1
- Port2 of node3 connected to port1 of node2

Port1 is always used to communicate with itself (loopback).

Failover Support in the OFA* Device

The Intel® MPI Library recognizes the following events to detect hardware issues:

- `IBV_EVENT_QP_FATAL` Error occurred on a QP and it transitioned to error state

- `IBV_EVENT_QP_REQ_ERR` Invalid request local work queue error
- `IBV_EVENT_QP_ACCESS_ERR` Local access violation error
- `IBV_EVENT_PATH_MIG_ERR` A connection failed to migrate to the alternate path
- `IBV_EVENT_CQ_ERR` CQ is in error (CQ overrun)
- `IBV_EVENT_SRQ_ERR` Error occurred on an SRQ
- `IBV_EVENT_PORT_ERR` Link became unavailable on a port
- `IBV_EVENT_DEVICE_FATAL` CA is in FATAL state

Intel® MPI Library stops using a port or the whole adapter for communications if one of these issues is detected. The communications continue through an available port or adapter, if the application is running in multi-rail mode. The application is aborted if no healthy ports/adapters are available.

Intel® MPI Library also recognizes the following event

- `IBV_EVENT_PORT_ACTIVE` Link became active on a port

The event indicates that the port can be used again and is enabled for communications.

3.3.8. OFI*-capable Network Fabrics Control

`I_MPI_OFI_LIBRARY`

Select the OpenFabrics Interfaces* (OFI*) library to be used.

Syntax

`I_MPI_OFI_LIBRARY=<library>`

Arguments

<code><library></code>	Specify an OFI library to be used instead of the default <code>libfabric.so</code>
------------------------------	--

Description

Set this environment variable to select a specific OFI library. Specify the full path to the OFI library if the library is not located in the dynamic loader search path.

`I_MPI_OFI_PROVIDER`

Define the name of the OFI provider to load.

Syntax

`I_MPI_OFI_PROVIDER=<name>`

Arguments

<code><name></code>	The name of the OFI provider to load
---------------------------	--------------------------------------

Description

Set this environment variable to define the name of the OFI provider to load. If you do not specify this variable, the OFI library chooses the provider automatically. You can check all available providers by using the `I_MPI_OFI_PROVIDER_DUMP` environment variable.

`I_MPI_OFI_PROVIDER_DUMP`

Control the capability of printing information about all OFI providers and their attributes from an OFI library.

Syntax

```
I_MPI_OFI_PROVIDER_DUMP=<arg>
```

Arguments

<arg>	Binary indicator
enable yes on 1	Print the list of all OFI providers and their attributes from an OFI library
disable no off 0	No action. This is default value

Description

Set this environment variable to control the capability of printing information about all OFI providers and their attributes from an OFI library.

3.4. Collective Operation Control

Each collective operation in the Intel® MPI Library supports a number of communication algorithms. In addition to highly optimized default settings, the library provides a way to control the algorithm selection explicitly: `I_MPI_ADJUST` environment variable family, which is described in the following section.

The environment variable is available for both Intel® and non-Intel microprocessors, but it may perform additional optimizations for Intel microprocessors than it performs for non-Intel microprocessors.

Each collective operation in the Intel® MPI Library supports a number of communication algorithms. In addition to highly optimized default settings, the library provides two ways to control the algorithm selection explicitly: the novel `I_MPI_ADJUST` environment variable family and the deprecated `I_MPI_MSG` environment variable family. They are described in the following sections.

These environment variables are available for both Intel® and non-Intel microprocessors, but they may perform additional optimizations for Intel microprocessors than they perform for non-Intel microprocessors.

3.4.1. I_MPI_ADJUST Family

`I_MPI_ADJUST_<opname>`

Control collective operation algorithm selection.

Syntax

```
I_MPI_ADJUST_<opname>="<algid>[:<conditions>][;<algid>:<conditions>[...]]"
```

Arguments

<algid>	Algorithm identifier
>= 0	The default value of zero selects the reasonable settings

<conditions>	A comma separated list of conditions. An empty list selects all message sizes and process combinations
<l>	Messages of size <l>
<l>-<m>	Messages of size from <l> to <m>, inclusive

<code><l>@<p></code>	Messages of size <code><l></code> and number of processes <code><p></code>
<code><l>-<m>@<p>-<q></code>	Messages of size from <code><l></code> to <code><m></code> and number of processes from <code><p></code> to <code><q></code> , inclusive

Description

Set this environment variable to select the desired algorithm(s) for the collective operation `<opname>` under particular conditions. Each collective operation has its own environment variable and algorithms.

Table 3.4-1 Environment Variables, Collective Operations, and Algorithms

Environment Variable	Collective Operation	Algorithms
<code>I_MPI_ADJUST_ALLGATHER</code>	<code>MPI_Allgather</code>	<ol style="list-style-type: none"> 1. Recursive doubling 2. Bruck's 3. Ring 4. Topology aware Gather + Bcast 5. Knomial
<code>I_MPI_ADJUST_ALLGATHERV</code>	<code>MPI_Allgatherv</code>	<ol style="list-style-type: none"> 1. Recursive doubling 2. Bruck's 3. Ring 4. Topology aware Gather + Bcast
<code>I_MPI_ADJUST_ALLREDUCE</code>	<code>MPI_Allreduce</code>	<ol style="list-style-type: none"> 1. Recursive doubling 2. Rabenseifner's 3. Reduce + Bcast 4. Topology aware Reduce + Bcast 5. Binomial gather + scatter 6. Topology aware binomial gather + scatter 7. Shumilin's ring 8. Ring 9. Knomial
<code>I_MPI_ADJUST_ALLTOALL</code>	<code>MPI_Alltoall</code>	<ol style="list-style-type: none"> 1. Bruck's 2. Isend/Irecv + waitall 3. Pair wise exchange 4. Plum's

I_MPI_ADJUST_ALLTOALLV	MPI_Alltoallv	<ol style="list-style-type: none"> 1. Isend/Irecv + waitall 2. Plum's
I_MPI_ADJUST_ALLTOALLW	MPI_Alltoallw	Isend/Irecv + waitall
I_MPI_ADJUST_BARRIER	MPI_Barrier	<ol style="list-style-type: none"> 1. Dissemination 2. Recursive doubling 3. Topology aware dissemination 4. Topology aware recursive doubling 5. Binominal gather + scatter 6. Topology aware binominal gather + scatter
I_MPI_ADJUST_BCAST	MPI_Bcast	<ol style="list-style-type: none"> 1. Binomial 2. Recursive doubling 3. Ring 4. Topology aware binomial 5. Topology aware recursive doubling 6. Topology aware ring 7. Shumilin's 8. Knomial
I_MPI_ADJUST_EXSCAN	MPI_Exscan	<ol style="list-style-type: none"> 1. Partial results gathering 2. Partial results gathering regarding layout of processes
I_MPI_ADJUST_GATHER	MPI_Gather	<ol style="list-style-type: none"> 1. Binomial 2. Topology aware binomial 3. Shumilin's
I_MPI_ADJUST_GATHERV	MPI_Gatherv	<ol style="list-style-type: none"> 1. Linear 2. Topology aware linear 3. Knomial

I_MPI_ADJUST_REDUCE_SCATTER	MPI_Reduce_scatter	<ol style="list-style-type: none"> 1. Recursive halving 2. Pair wise exchange 3. Recursive doubling 4. Reduce + Scatterv 5. Topology aware Reduce + Scatterv
I_MPI_ADJUST_REDUCE	MPI_Reduce	<ol style="list-style-type: none"> 1. Shumilin's 2. Binomial 3. Topology aware Shumilin's 4. Topology aware binomial 5. Rabenseifner's 6. Topology aware Rabenseifner's 7. Knomial
I_MPI_ADJUST_SCAN	MPI_Scan	<ol style="list-style-type: none"> 1. Partial results gathering 2. Topology aware partial results gathering
I_MPI_ADJUST_SCATTER	MPI_Scatter	<ol style="list-style-type: none"> 1. Binomial 2. Topology aware binomial 3. Shumilin's
I_MPI_ADJUST_SCATTERV	MPI_Scatterv	<ol style="list-style-type: none"> 1. Linear 2. Topology aware linear
I_MPI_ADJUST_IALLGATHER	MPI_Iallgather	<ol style="list-style-type: none"> 1. Recursive doubling 2. Bruck's 3. Ring
I_MPI_ADJUST_IALLGATHERV	MPI_Iallgatherv	<ol style="list-style-type: none"> 1. Recursive doubling 2. Bruck's 3. Ring
I_MPI_ADJUST_IALLREDUCE	MPI_Iallreduce	<ol style="list-style-type: none"> 1. Recursive doubling 2. Rabenseifner's

		<ol style="list-style-type: none"> 3. Reduce + Bcast 4. Ring (patarasuk) 5. Knomial 6. Binomial
I_MPI_ADJUST_IALLTOALL	MPI_Ialltoall	<ol style="list-style-type: none"> 1. Bruck's 2. Isend/Irecv + Waitall 3. Pairwise exchange
I_MPI_ADJUST_IALLTOALLV	MPI_Ialltoallv	Isend/Irecv + Waitall
I_MPI_ADJUST_IALLTOALLW	MPI_Ialltoallw	Isend/Irecv + Waitall
I_MPI_ADJUST_IBARRIER	MPI_Ibarrier	Dissemination
I_MPI_ADJUST_IBCAST	MPI_Ibcast	<ol style="list-style-type: none"> 1. Binomial 2. Recursive doubling 3. Ring 4. Knomial
I_MPI_ADJUST_IEXSCAN	MPI_Iexscan	Recursive doubling
I_MPI_ADJUST_IGATHER	MPI_Igather	<ol style="list-style-type: none"> 1. Binomial 2. Knomial
I_MPI_ADJUST_IGATHERV	MPI_Igatherv	Linear
I_MPI_ADJUST_IREDUCE_SCATTER	MPI_Ireduce_scatter	<ol style="list-style-type: none"> 1. Recursive halving 2. Pairwise 3. Recursive doubling
I_MPI_ADJUST_IREDUCE	MPI_Ireduce	<ol style="list-style-type: none"> 1. Rabenseifner's 2. Binomial 3. Knomial
I_MPI_ADJUST_ISCAN	MPI_Iscan	Recursive Doubling
I_MPI_ADJUST_ISCATTER	MPI_Iscatter	<ol style="list-style-type: none"> 1. Binomial 2. Knomial

I_MPI_ADJUST_ISCATTERV	MPI_Iscatterv	Linear
------------------------	---------------	--------

The message size calculation rules for the collective operations are described in the table. In the following table, "n/a" means that the corresponding interval $\langle l \rangle - \langle m \rangle$ should be omitted.

Table 3.4-2 Message Collective Functions

Collective Function	Message Size Formula
MPI_Allgather	recv_count*recv_type_size
MPI_Allgatherv	total_recv_count*recv_type_size
MPI_Allreduce	count*type_size
MPI_Alltoall	send_count*send_type_size
MPI_Alltoallv	n/a
MPI_Alltoallw	n/a
MPI_Barrier	n/a
MPI_Bcast	count*type_size
MPI_Exscan	count*type_size
MPI_Gather	recv_count*recv_type_size if MPI_IN_PLACE is used, otherwise send_count*send_type_size
MPI_Gatherv	n/a
MPI_Reduce_scatter	total_recv_count*type_size
MPI_Reduce	count*type_size
MPI_Scan	count*type_size
MPI_Scatter	send_count*send_type_size if MPI_IN_PLACE is used, otherwise recv_count*recv_type_size
MPI_Scatterv	n/a

Examples

Use the following settings to select the second algorithm for MPI_Reduce operation:

```
I_MPI_ADJUST_REDUCE=2
```

Use the following settings to define the algorithms for MPI_Reduce_scatter operation:

```
I_MPI_ADJUST_REDUCE_SCATTER="4:0-100,5001-10000;1:101-3200,2:3201-5000;3"
```

In this case, algorithm 4 is used for the message sizes between 0 and 100 bytes and from 5001 and 10000 bytes, algorithm 1 is used for the message sizes between 101 and 3200 bytes, algorithm 2 is used for the message sizes between 3201 and 5000 bytes, and algorithm 3 is used for all other messages.

I_MPI_ADJUST_REDUCE_SEGMENT

Syntax

```
I_MPI_ADJUST_REDUCE_SEGMENT=<block_size>|<algid>:<block_size>[,<algid>:<block_size>[...]]
```

Arguments

<algid>	Algorithm identifier
1	Shumilin's algorithm
3	Topology aware Shumilin's algorithm
<block_size>	Size of a message segment in bytes
> 0	The default value is 14000

Description

Set an internal block size to control `MPI_Reduce` message segmentation for the specified algorithm. If the <algid> value is not set, the <block_size> value is applied for all the algorithms, where it is relevant.

NOTE

This environment variable is relevant for Shumilin's and topology aware Shumilin's algorithms only (algorithm N1 and algorithm N3 correspondingly).

I_MPI_ADJUST_BCAST_SEGMENT

Syntax

```
I_MPI_ADJUST_BCAST_SEGMENT=<block_size>|<algid>:<block_size>[,<algid>:<block_size>[...]]
```

Arguments

<algid>	Algorithm identifier
1	Binomial
4	Topology aware binomial
7	Shumilin's
8	Knomial
<block_size>	Size of a message segment in bytes
> 0	The default value is 12288

Description

Set an internal block size to control `MPI_Bcast` message segmentation for the specified algorithm. If the <algid> value is not set, the <block_size> value is applied for all the algorithms, where it is relevant.

NOTE

This environment variable is relevant only for Binomial, Topology-aware binomial, Shumilin's and Knomial algorithms.

I_MPI_ADJUST_ALLGATHER_KN_RADIX**Syntax**

```
I_MPI_ADJUST_ALLGATHER_KN_RADIX=<radix>
```

Arguments

<radix>	An integer that specifies a radix used by the Knomial MPI_Allgather algorithm to build a knomial communication tree
> 1	The default value is 2

Description

Set this environment variable together with I_MPI_ADJUST_ALLGATHER=5 to select the knomial tree radix for the corresponding MPI_Allgather algorithm.

I_MPI_ADJUST_BCAST_KN_RADIX**Syntax**

```
I_MPI_ADJUST_BCAST_KN_RADIX=<radix>
```

Arguments

<radix>	An integer that specifies a radix used by the Knomial MPI_Bcast algorithm to build a knomial communication tree
> 1	The default value is 4

Description

Set this environment variable together with I_MPI_ADJUST_BCAST=8 to select the knomial tree radix for the corresponding MPI_Bcast algorithm.

I_MPI_ADJUST_ALLREDUCE_KN_RADIX**Syntax**

```
I_MPI_ADJUST_ALLREDUCE_KN_RADIX=<radix>
```

Arguments

<radix>	An integer that specifies a radix used by the Knomial MPI_Allreduce algorithm to build a knomial communication tree
> 1	The default value is 4

Description

Set this environment variable together with I_MPI_ADJUST_ALLREDUCE=9 to select the knomial tree radix for the corresponding MPI_Allreduce algorithm.

I_MPI_ADJUST_REDUCE_KN_RADIX

Syntax

I_MPI_ADJUST_REDUCE_KN_RADIX=<radix>

Arguments

<radix>	An integer that specifies a radix used by the Knomial MPI_Reduce algorithm to build a knomial communication tree
> 1	The default value is 4

Description

Set this environment variable together with I_MPI_ADJUST_REDUCE=7 to select the knomial tree radix for the corresponding MPI_Reduce algorithm.

I_MPI_ADJUST_GATHERV_KN_RADIX

Syntax

I_MPI_ADJUST_GATHERV_KN_RADIX=<radix>

Arguments

<radix>	An integer that specifies a radix used by the Knomial MPI_Gatherv algorithm to build a knomial communication tree
> 1	The default value is 2

Description

Set this environment variable together with I_MPI_ADJUST_GATHERV=3 to select the knomial tree radix for the corresponding MPI_Gatherv algorithm.

I_MPI_ADJUST_IALLREDUCE_KN_RADIX

Syntax

I_MPI_ADJUST_IALLREDUCE_KN_RADIX=<radix>

Arguments

<radix>	An integer that specifies a radix used by the Knomial MPI_Iallreduce algorithm to build a knomial communication tree
> 1	The default value is 4

Description

Set this environment variable together with I_MPI_ADJUST_IALLREDUCE=5 to select the knomial tree radix for the corresponding MPI_Iallreduce algorithm.

I_MPI_ADJUST_IBCAST_KN_RADIX

Syntax

I_MPI_ADJUST_IBCAST_KN_RADIX=<radix>

Arguments

<radix>	An integer that specifies a radix used by the Knomial MPI_Ibcast algorithm to build a
---------	---

	knomial communication tree
> 1	The default value is 4

Description

Set this environment variable together with `I_MPI_ADJUST_IBCAST=4` to select the knomial tree radix for the corresponding `MPI_Ibcast` algorithm.

I_MPI_ADJUST_IREDUCE_KN_RADIX

Syntax

`I_MPI_ADJUST_IREDUCE_KN_RADIX=<radix>`

Arguments

<radix>	An integer that specifies a radix used by the Knomial <code>MPI_Ireduce</code> algorithm to build a knomial communication tree
> 1	The default value is 4

Description

Set this environment variable together with `I_MPI_ADJUST_IREDUCE=3` to select the knomial tree radix for the corresponding `MPI_Ireduce` algorithm.

I_MPI_ADJUST_IGATHER_KN_RADIX

Syntax

`I_MPI_ADJUST_IGATHER_KN_RADIX=<radix>`

Arguments

<radix>	An integer that specifies a radix used by the Knomial <code>MPI_Igather</code> algorithm to build a knomial communication tree
> 1	The default value is 4

Description

Set this environment variable together with `I_MPI_ADJUST_IGATHER=2` to select the knomial tree radix for the corresponding `MPI_Igather` algorithm.

I_MPI_ADJUST_ISCATTER_KN_RADIX

Syntax

`I_MPI_ADJUST_ISCATTER_KN_RADIX=<radix>`

Arguments

<radix>	An integer that specifies a radix used by the Knomial <code>MPI_Iscatter</code> algorithm to build a knomial communication tree
> 1	The default value is 4

Description

Set this environment variable together with `I_MPI_ADJUST_ISCATTER=2` to select the knomial tree radix for the corresponding `MPI_Iscatter` algorithm.

3.4.2. I_MPI_MSG Family

These environment variables are deprecated and supported mostly for backward compatibility. Use the `I_MPI_ADJUST` environment variable family whenever possible.

I_MPI_FAST_COLLECTIVES

Control the default library behavior during selection of the most appropriate collective algorithm.

Syntax

`I_MPI_FAST_COLLECTIVES=<arg>`

Arguments

<code><arg></code>	Binary indicator
<code>enable yes on 1</code>	Fast collective algorithms are used. This is the default value
<code>disable no off 0</code>	Slower and safer collective algorithms are used

Description

The Intel® MPI Library uses advanced collective algorithms designed for better application performance by default. The implementation makes the following assumptions:

- It is safe to utilize the flexibility of the MPI standard regarding the order of execution of the collective operations to take advantage of the process layout and other opportunities.
- There is enough memory available for allocating additional internal buffers.

Set the `I_MPI_FAST_COLLECTIVES` environment variable to `disable` if you need to obtain results that do not depend on the physical process layout or other factors.

NOTE

Some optimizations controlled by this environment variable are of an experimental nature. In case of failure, turn off the collective optimizations and repeat the run.

I_MPI_BCAST_NUM_PROCS

Control `MPI_Bcast` algorithm thresholds.

Syntax

`I_MPI_BCAST_NUM_PROCS=<nproc>`

Arguments

<code><nproc></code>	Define the number of processes threshold for choosing the <code>MPI_Bcast</code> algorithm
<code>> 0</code>	The default value is 8

I_MPI_BCAST_MSG

Control `MPI_Bcast` algorithm thresholds.

Syntax

`I_MPI_BCAST_MSG=<nbytes1,nbytes2>`

Arguments

<code><nbytes1,nbytes2></code>	Define the message size threshold range (in bytes) for choosing the MPI_Bcast algorithm
<code>> 0</code> <code>nbytes2 >= nbytes1</code>	The default pair of values is 12288,524288

Description

Set these environment variables to control the selection of the three possible MPI_Bcast algorithms according to the following scheme (See [Table 3.4-1](#) for algorithm descriptions):

The first algorithm is selected if the message size is less than `<nbytes1>`, or the number of processes in the operation is less than `<nproc>`.

The second algorithm is selected if the message size is greater than or equal to `<nbytes1>` and less than `<nbytes2>`, and the number of processes in the operation is a power of two.

If none of the above conditions is satisfied, the third algorithm is selected.

I_MPI_ALLTOALL_NUM_PROCS

Control MPI_Alltoall algorithm thresholds.

Syntax

`I_MPI_ALLTOALL_NUM_PROCS=<nproc>`

Arguments

<code><nproc></code>	Define the number of processes threshold for choosing the MPI_Alltoall algorithm
<code>> 0</code>	The default value is 8

I_MPI_ALLTOALL_MSG

Control MPI_Alltoall algorithm thresholds.

Syntax

`I_MPI_ALLTOALL_MSG=<nbytes1,nbytes2>`

Arguments

<code><nbytes1,nbytes2></code>	Defines the message size threshold range (in bytes) for choosing the MPI_Alltoall algorithm
<code>> 0</code> <code>nbytes2 >= nbytes1</code>	The default pair of values is 256,32768

Description

Set these environment variables to control the selection of the three possible MPI_Alltoall algorithms according to the following scheme (See [Table 3.4-1](#) for algorithm descriptions):

The first algorithm is selected if the message size is greater than or equal to `<nbytes1>`, and the number of processes in the operation is not less than `<nproc>`.

The second algorithm is selected if the message size is greater than `<nbytes1>` and less than or equal to `<nbytes2>`, or if the message size is less than `<nbytes2>` and the number of processes in the operation is less than `<nproc>`.

If none of the above conditions is satisfied, the third algorithm is selected.

I_MPI_ALLGATHER_MSG

Control MPI_Allgather algorithm thresholds.

Syntax

I_MPI_ALLGATHER_MSG=<nbytes1,nbytes2>

Arguments

<nbytes1,nbytes2>	Define the message size threshold range (in bytes) for choosing the MPI_Allgather algorithm
> 0 nbytes2 >= nbytes1	The default pair of values is 81920,524288

Description

Set this environment variable to control the selection of the three possible MPI_Allgather algorithms according to the following scheme (See [Table 3.4-1](#) for algorithm descriptions):

The first algorithm is selected if the message size is less than <nbytes2> and the number of processes in the operation is a power of two.

The second algorithm is selected if the message size is less than <nbytes1> and number of processes in the operation is not a power of two.

If none of the above conditions is satisfied, the third algorithm is selected.

I_MPI_ALLREDUCE_MSG

Control MPI_Allreduce algorithm thresholds.

Syntax

I_MPI_ALLREDUCE_MSG=<nbytes>

Arguments

<nbytes>	Define the message size threshold (in bytes) for choosing the MPI_Allreduce algorithm
> 0	The default value is 2048

Description

Set this environment variable to control the selection of the two possible MPI_Allreduce algorithms according to the following scheme (See [Table 3.4-1](#) for algorithm descriptions):

The first algorithm is selected if the message size is less than or equal <nbytes>, or the reduction operation is user-defined, or the count argument is less than the nearest power of two less than or equal to the number of processes.

If the above condition is not satisfied, the second algorithm is selected.

I_MPI_REDS CAT_MSG

Control the MPI_Reduce_scatter algorithm thresholds.

Syntax

I_MPI_REDS CAT_MSG=<nbytes1,nbytes2>

Arguments

<nbytes>	Define the message size threshold range (in bytes) for choosing the MPI_Reduce_scatter algorithm
----------	--

> 0	The default pair of values is 512, 524288
-----	---

Description

Set this environment variable to control the selection of the three possible `MPI_Reduce_scatter` algorithms according to the following scheme (See [Table 3.4-1](#) for algorithm descriptions):

The first algorithm is selected if the reduction operation is commutative and the message size is less than `<nbytes2>`.

The second algorithm is selected if the reduction operation is commutative and the message size is greater than or equal to `<nbytes2>`, or if the reduction operation is not commutative and the message size is greater than or equal to `<nbytes1>`.

If none of the above conditions is satisfied, the third algorithm is selected.

I_MPI_SCATTER_MSG

Control `MPI_Scatter` algorithm thresholds.

Syntax

`I_MPI_SCATTER_MSG=<nbytes>`

Arguments

<code><nbytes></code>	Define the buffer size threshold range (in bytes) for choosing the <code>MPI_Scatter</code> algorithm
> 0	The default value is 2048

Description

Set this environment variable to control the selection of the two possible `MPI_Scatter` algorithms according to the following scheme (See [Table 3.4-1](#) for algorithm descriptions):

The first algorithm is selected on the intercommunicators if the message size is greater than `<nbytes>`.

If the above condition is not satisfied, the second algorithm is selected.

I_MPI_GATHER_MSG

Control `MPI_Gather` algorithm thresholds.

Syntax

`I_MPI_GATHER_MSG=<nbytes>`

Arguments

<code><nbytes></code>	Define the buffer size threshold range (in bytes) for choosing the <code>MPI_Gather</code> algorithm
> 0	The default value is 2048

Description

Set this environment variable to control the selection of the two possible `MPI_Gather` algorithms according to the following scheme (See [Table 3.4-1](#) for algorithm descriptions):

The first algorithm is selected on the intercommunicators if the message size is greater than `<nbytes>`.

If the above condition is not satisfied, the second algorithm is selected.

3.5. Miscellaneous

This topic provides the following information:

- Timer Control
- Compatibility Control
- Dynamic Process Support
- Fault Tolerance
- Statistics Gathering Mode
- ILP64 Support
- Unified Memory Management
- File System Support
- Multi-threaded memcpy Support

3.5.1. Timer Control

I_MPI_TIMER_KIND

Select the timer used by the `MPI_Wtime` and `MPI_Wtick` calls.

Syntax

`I_MPI_TIMER_KIND=<timername>`

Arguments

<code><timername></code>	Define the timer type
<code>gettimeofday</code>	If this setting is chosen, the <code>MPI_Wtime</code> and <code>MPI_Wtick</code> functions will work through the function <code>gettimeofday(2)</code> . This is the default value
<code>rdtsc</code>	If this setting is chosen, the <code>MPI_Wtime</code> and <code>MPI_Wtick</code> functions will use the high resolution RDTSC timer

Description

Set this environment variable to select either the ordinary or RDTSC timer.

The resolution of the default `gettimeofday(2)` timer may be insufficient on certain platforms.

3.5.2. Compatibility Control

I_MPI_COMPATIBILITY

Select the runtime compatibility mode.

Syntax

`I_MPI_COMPATIBILITY=<value>`

Arguments

<code><value></code>	Define compatibility mode
<code>not defined</code>	The MPI-3.0 standard compatibility. This is the default mode
<code>3</code>	The Intel® MPI Library 3.x compatible mode

4	The Intel® MPI Library 4.0.x compatible mode
---	--

Description

Set this environment variable to choose the Intel® MPI Library runtime compatible mode. By default, the library complies with the MPI-3.0 standard. If your application depends on the MPI-2.1 behavior, set the value of the environment variable `I_MPI_COMPATIBILITY` to 4. If your application depends on the pre-MPI-2.1 behavior, set the value of the environment variable `I_MPI_COMPATIBILITY` to 3.

3.5.3. Dynamic Process Support

The Intel® MPI Library provides support for the MPI-2 process model that allows creation and cooperative termination of processes after an MPI application has started. It provides the following:

- a mechanism to establish communication between the newly created processes and the existing MPI application
- a process attachment mechanism to establish communication between two existing MPI applications even when one of them does not spawn the other

The default placement of the spawned processes uses round robin scheduling. The first spawned process is placed after the last process of the parent group. A specific network fabric combination is selected using the usual fabrics selection algorithm (see [I_MPI_FABRICS](#) and [I_MPI_FABRICS_LIST](#) for details).

For example, to run a dynamic application, use the following commands:

```
$ mpirun -n 1 -gwdir <path_to_executable> -genv I_MPI_FABRICS shm:tcp
<spawn_app>
```

In the example, the `spawn_app` spawns 4 dynamic processes. If the `mpd.hosts` contains the following information:

```
host1
host2
host3
host4
```

The original spawning process is placed on `host1`, while the dynamic processes are distributed as follows: 1 - on `host2`, 2 - on `host3`, 3 - on `host4`, and 4 - again on `host1`.

To run a client-server application, use the following commands on the intended server host:

```
$ mpirun -n 1 -genv I_MPI_FABRICS shm:dapl <server_app> > <port_name>
```

and use the following commands on the intended client hosts:

```
$ mpirun -n 1 -genv I_MPI_FABRICS shm:dapl <client_app> < <port_name>
```

To run a simple `MPI_COMM_JOIN` based application, use the following commands on the intended server host:

```
$ mpirun -n 1 -genv I_MPI_FABRICS shm:ofa <join_server_app> < <port_number>
$ mpirun -n 1 -genv I_MPI_FABRICS shm:ofa <join_client_app> < <port_number>
```

3.5.4. Fault Tolerance

Intel® MPI Library provides extra functionality to enable fault tolerance support in the MPI applications. The MPI standard does not define behavior of MPI implementation if one or several processes of MPI application are abnormally aborted. By default, Intel® MPI Library aborts the whole application if any process stops.

Set the environment variable `I_MPI_FAULT_CONTINUE` to `on` to change this behavior. For example,

```
$ mpirun -env I_MPI_FAULT_CONTINUE on -n 2 ./test
```

An application can continue working in the case of MPI processes an issue if the issue meets the following requirements:

- An application sets error handler `MPI_ERRORS_RETURN` to communicator `MPI_COMM_WORLD` (all new communicators inherit error handler from it)
- An application uses master-slave model. In this case, the application is stopped only if the master is finished or does not respond
- An application uses only point-to-point communication between a master and a number of slaves. It does not use inter slave communication or MPI collective operations.
- Handle a certain MPI error code on a point-to-point operation with a particular failed slave rank for application to avoid further communication with this rank. The slave rank can be blocking/non-blocking send, receive, probe and test,
- Any communication operation can be used on subset communicator system. If an error appears in a collective operation, any communication inside this communicator will be prohibited.
- Master failure means the job stops.
- Fault Tolerance functionality is not available for spawned processes.

Environment Variables

I_MPI_FAULT_CONTINUE

Turn on/off support for fault tolerant applications.

Syntax

`I_MPI_FAULT_CONTINUE=<arg>`

Arguments

<code><arg></code>	Binary indicator
<code>enable yes on 1</code>	Turn on support for fault tolerant applications
<code>disable no off 0</code>	Turn off support for fault tolerant applications. This is default value

Description

Set this environment variable to provide support for fault tolerant applications.

Usage Model

An application sets `MPI_ERRORS_RETURN` error handler and checks the return code after each communication call. If a communication call does not return `MPI_SUCCESS`, the destination process should be marked unreachable and exclude communication with it. For example:

```
if(live_ranks[rank]) {
    mpi_err = MPI_Send(buf, count, dtype, rank, tag, MPI_COMM_WORLD);
    if(mpi_err != MPI_SUCCESS) {
        live_ranks[rank] = 0;
    }
}
```

In the case of non-blocking communications, errors can appear during wait/test operations.

3.5.5. Statistics Gathering Mode

This topic describes the Intel® MPI Library statistics gathering modes and how to use such gathering facility through environment variables. The Intel® MPI Library supports the following statistics formats:

- Native statistics format
- IPM statistics format

You can see the information about native statistic format in the topic, [Native Statistic Format](#) and the information about IPM statics format in the topic, [IPM Statistics Format](#). There is also possibility to collect both types of statistics. See [Native and IPM Statistics](#) for more details.

Native Statistics Format

The Intel® MPI Library has a built-in statistics gathering facility that collects essential performance data without disturbing the application execution. The collected information is sent to a text file. This section describes the environment variables used to control the built-in statistics gathering facility, and provides example output files.

Besides using the environment variables, you can also collect the native statistics using the MPI Performance Snapshot through the `-mps` option. For example:

```
$ mpirun -mps -n 2 ./myApp
```

See the description of `-mps` for more details.

I_MPI_STATS

Control statistics collection. Expand values of `I_MPI_STATS` environment variable additionally to existing values.

Syntax

```
I_MPI_STATS=[native:][n-] m
```

Arguments

<code>n, m</code>	Possible stats levels of the output information
1	Output the amount of data sent by each process
2	Output the number of calls and amount of transferred data
3	Output statistics combined according to the actual arguments
4	Output statistics defined by a buckets list
10	Output collective operation statistics for all communication contexts
20	Output additional time information for all MPI functions

Description

Set this environment variable to control the amount of statistics information collected and the output to the log file. No statistics are output by default.

NOTE

`n, m` are positive integer numbers. They define the range of output information. The statistics from level `n` to level `m` inclusive are output. If an `n` value is not provided, the default value is 1.

I_MPI_STATS_SCOPE

Select the subsystem(s) to collect statistics for.

Syntax

```
I_MPI_STATS_SCOPE="<subsystem>[:<ops>] [;<subsystem>[:<ops>] [...]]"
```

Arguments

<i><subsystem></i>	Define the target subsystem(s)
all	Collect statistics data for all operations. This is the default value
coll	Collect statistics data for all collective operations
p2p	Collect statistics data for all point-to-point operations

<i><ops></i>	Define the target operations as a comma separated list
Allgather	MPI_Allgather
Iallgather	MPI_Iallgather
Allgatherv	MPI_Allgatherv
Iallgatherv	MPI_Iallgatherv
Allreduce	MPI_Allreduce
Iallreduce	MPI_Iallreduce
Alltoall	MPI_Alltoall
Ialltoall	MPI_Ialltoall
Alltoallv	MPI_Alltoallv
Ialltoallv	MPI_Ialltoallv
Alltoallw	MPI_Alltoallw
Ialltoallw	MPI_Ialltoallw
Barrier	MPI_Barrier
Ibarrier	MPI_Ibarrier
Bcast	MPI_Bcast
Ibcast	MPI_Ibcast
Exscan	MPI_Exscan
Iexscan	MPI_Iexscan
Gather	MPI_Gather

Igather	MPI_Igather
Gatherv	MPI_Gatherv
Igatherv	MPI_Igatherv
Reduce_scatter	MPI_Reduce_scatter
Ireduce_scatter	MPI_Ireduce_scatter
Reduce	MPI_Reduce
Ireduce	MPI_Ireduce
Scan	MPI_Scan
Isan	MPI_Isan
Scatter	MPI_Scatter
Iscatter	MPI_Iscatter
Scatterv	MPI_Scatterv
Iscatterv	MPI_Iscatterv
Send	Standard transfers (MPI_Send, MPI_Isend, MPI_Send_init)
Sendrecv	Send-receive transfers (MPI_Sendrecv, MPI_Sendrecv_replace)
Bsend	Buffered transfers (MPI_Bsend, MPI_Ibsend, MPI_Bsend_init)
Csend	Point-to-point operations inside the collectives. This internal operation serves all collectives
Csendrecv	Point-to-point send-receive operations inside the collectives. This internal operation serves all collectives
Rsend	Ready transfers (MPI_Rsend, MPI_Irsend, MPI_Rsend_init)
Ssend	Synchronous transfers (MPI_Ssend, MPI_Issend, MPI_Ssend_init)

Description

Set this environment variable to select the target subsystem in which to collect statistics. All collective and point-to-point operations, including the point-to-point operations performed inside the collectives, are covered by default.

Examples

The default settings are equivalent to:

```
I_MPI_STATS_SCOPE="coll;p2p"
```

Use the following settings to collect statistics for the MPI_Bcast, MPI_Reduce, and all point-to-point operations:

```
I_MPI_STATS_SCOPE="p2p;coll:bcast,reduce"
```

Use the following settings to collect statistics for the point-to-point operations inside the collectives:

```
I_MPI_STATS_SCOPE=p2p:csend
```

I_MPI_STATS_BUCKETS

Identify a list of ranges for message sizes and communicator sizes that are used for collecting statistics.

Syntax

```
I_MPI_STATS_BUCKETS=<msg>[@<proc>] [,<msg>[@<proc>]]...
```

Arguments

<msg>	Specify range of message sizes in bytes
<l>	Single value of message size
<l>-<m>	Range from <l> to <m>

<proc>	Specify range of processes (ranks) for collective operations
<p>	Single value of communicator size
<p>-<q>	Range from <p> to <q>

Description

Set the I_MPI_STATS_BUCKETS environment variable to define a set of ranges for message sizes and communicator sizes.

Level 4 of the statistics provides profile information for these ranges.

If I_MPI_STATS_BUCKETS environment variable is not used, then level 4 statistics is not gathered.

If a range is not specified, the maximum possible range is assumed.

Examples

To specify short messages (from 0 to 1000 bytes) and long messages (from 50000 to 100000 bytes), use the following setting:

```
-env I_MPI_STATS_BUCKETS 0-1000,50000-100000
```

To specify messages that have 16 bytes in size and circulate within four process communicators, use the following setting:

```
-env I_MPI_STATS_BUCKETS "16@4"
```

NOTE

When the @ symbol is present, the environment variable value must be enclosed in quotes.

I_MPI_STATS_FILE

Define the statistics output file name.

Syntax

```
I_MPI_STATS_FILE=<name>
```

Arguments

<name>	Define the statistics output file name
--------	--

Description

Set this environment variable to define the statistics output file. By default, the `stats.txt` file is created in the current directory.

If this variable is not set and the statistics output file already exists, an index is appended to its name. For example, if `stats.txt` exists, the created statistics output file is named as `stats(2).txt`; if `stats(2).txt` exists, the created file is named as `stats(3).txt`, and so on.

The statistics data is blocked and ordered according to the process ranks in the `MPI_COMM_WORLD` communicator. The timing data is presented in microseconds. For example, with the following settings:

```
I_MPI_STATS=4
```

```
I_MPI_STATS_SCOPE="p2p;coll:allreduce"
```

The statistics output for a simple program that performs only one `MPI_Allreduce` operation may look as follows:

```
Intel(R) MPI Library Version 5.1
____ MPI Communication Statistics ____
Stats level: 4
P2P scope:< FULL >
Collectives scope:< Allreduce >
~~~~ Process 0 of 2 on node svlmpihead01 lifetime = 414.13
Data Transfers
Src Dst Amount(MB) Transfers
-----
000 --> 000 0.000000e+00 0
000 --> 001 7.629395e-06 2
=====
Totals 7.629395e-06 2
Communication Activity
Operation Volume(MB) Calls
-----
P2P
Csend 7.629395e-06 2
Csendrecv 0.000000e+00 0
Send 0.000000e+00 0
Sendrecv 0.000000e+00 0
Bsend 0.000000e+00 0
Rsend 0.000000e+00 0
Ssend 0.000000e+00 0
Collectives
Allreduce 7.629395e-06 2
=====
Communication Activity by actual args
P2P
Operation Dst Message size Calls
-----
Csend
1 1 4 2
Collectives
Operation Context Algo Comm size Message size Calls Cost(%)
-----
-----
Allreduce
1 0 1 2 4 2 44.96
=====
~~~~ Process 1 of 2 on node svlmpihead01 lifetime = 306.13
Data Transfers
Src Dst Amount(MB) Transfers
-----
001 --> 000 7.629395e-06 2
001 --> 001 0.000000e+00 0
=====
```

```

Totals 7.629395e-06 2
Communication Activity
Operation Volume(MB) Calls
-----
P2P
Csend 7.629395e-06 2
Csendrecv 0.000000e+00 0
Send 0.000000e+00 0
Sendrecv 0.000000e+00 0
Bsend 0.000000e+00 0
Rsend 0.000000e+00 0
Ssend 0.000000e+00 0
Collectives
Allreduce 7.629395e-06 2
=====
Communication Activity by actual args
P2P
Operation Dst Message size Calls
-----
Csend
1 0 4 2
Collectives
Operation Context Comm size Message size Calls Cost(%)
-----
Allreduce
1 0 2 4 2 37.93
=====
_____ End of stats.txt file _____

```

In the example above:

- All times are measured in microseconds.
- The message sizes are counted in bytes. **MB** means megabyte equal to 2^{20} or 1 048 576 bytes.
- The process life time is calculated as a stretch of time between `MPI_Init` and `MPI_Finalize`.
- The **Algo** field indicates the number of algorithm used by this operation with listed arguments.
- The **Cost** field represents a particular collective operation execution time as a percentage of the process life time.

Region Control

The Intel® MPI Library also supports an optional region feature. The region is an IPM statistics format feature. See [IPM Statistics Format](#) for more details about IPM. This feature requires the source code modification. The `MPI_Pcontrol` function can be used.

Region is a named part of the source code marked by the start/end points through the standard `MPI_Pcontrol` function calls. The `MPI_Pcontrol` function isn't used for the following special permanent regions:

- Main region contains statistics information about all MPI calls from `MPI_Init` to `MPI_Finalize`. The main region gets the "*" name for IPM statistics output. The default output file for this region is `stats.txt` for native statistics format.
- Complementary region contains statistics information not included into any named region. The region gets the "ipm_noregion" name in output for IPM statistics format. The default output file for this region is `stats_noregion.txt` for native statistics format.

If named regions are not used, the main regions and the complementary regions are identical and the complementary region is ignored.

Each region contains its own independent statistics information about MPI functions called inside the region.

The Intel® MPI Library supports the following types of regions:

- Discontiguous (several open and close).
- Intersected.
- Covering a subset of MPI processes (part of the `MPI_COMM_WORLD` environment variable).

A region is opened by the `MPI_Pcontrol(1, <name>)` call and closed by the `MPI_Pcontrol(-1, <name>)` call where `<name>` is a zero terminated string with the region name. The `<name>` is used in output for IPM statistics format. The default output file for the region is `stats_<name>.txt` for native statistics format.

All open regions are closed automatically inside the `MPI_Finalize` environment variable.

IPM Statistics Format

The Intel® MPI Library supports integrated performance monitoring (IPM) summary format as part of the built-in statistics gathering mechanism described above. You do not need to modify the source code or re-link your application to collect this information.

The `I_MPI_STATS_BUCKETS` environment variable is not applicable to the IPM format. The `I_MPI_STATS_ACCURACY` environment variable is available to control extra functionality.

I_MPI_STATS

Control the statistics data output format.

Syntax

`I_MPI_STATS=<level>`

Argument

<code><level></code>	Level of statistics data
<code>ipm</code>	Summary data throughout all regions
<code>ipm:terse</code>	Basic summary data

Description

Set this environment variable to `ipm` to get the statistics output that contains region summary. Set this environment variable to `ipm:terse` argument to get the brief statistics output.

I_MPI_STATS_FILE

Define the output file name.

Syntax

`I_MPI_STATS_FILE=<name>`

Argument

<code><name></code>	File name for statistics data gathering
---------------------------	---

Description

Set this environment variable to change the statistics output file name from the default name of `stats.ipm`.

If this variable is not set and the statistics output file already exists, an index is appended to its name. For example, if `stats.ipm` exists, the created statistics output file is named as `stats(2).ipm`; if `stats(2).ipm` exists, the created file is named as `stats(3).ipm`, and so on.

I_MPI_STATS_SCOPE

Define a semicolon separated list of subsets of MPI functions for statistics gathering.

Syntax

```
I_MPI_STATS_SCOPE="<subset>[;<subset>[...]]"
```

Argument

<code><subset></code>	Target subset
<code>all2all</code>	Collect statistics data for all-to-all functions types
<code>all2one</code>	Collect statistics data for all-to-one functions types
<code>attr</code>	Collect statistics data for attribute control functions
<code>comm</code>	Collect statistics data for communicator control functions
<code>err</code>	Collect statistics data for error handling functions
<code>group</code>	Collect statistics data for group support functions
<code>init</code>	Collect statistics data for initialize/finalize functions
<code>io</code>	Collect statistics data for input/output support function
<code>one2all</code>	Collect statistics data for one-to-all functions types
<code>recv</code>	Collect statistics data for receive functions
<code>req</code>	Collect statistics data for request support functions
<code>rma</code>	Collect statistics data for one sided communication functions
<code>scan</code>	Collect statistics data for scan collective functions
<code>send</code>	Collect statistics data for send functions
<code>sendrecv</code>	Collect statistics data for send/receive functions
<code>serv</code>	Collect statistics data for additional service functions
<code>spawn</code>	Collect statistics data for dynamic process functions
<code>status</code>	Collect statistics data for status control function
<code>sync</code>	Collect statistics data for barrier synchronization

time	Collect statistics data for timing support functions
topo	Collect statistics data for topology support functions
type	Collect statistics data for data type support functions

Description

Use this environment variable to define a subset or subsets of MPI functions for statistics gathering specified by the following table. A union of all subsets is used by default.

Table 3.5-1 Stats Subsets of MPI Functions

all2all MPI_Allgather MPI_Allgatherv MPI_Allreduce MPI_Alltoll MPI_Alltoallv MPI_Alltoallw MPI_Reduce_scatter MPI_iallgather MPI_iallgatherv MPI_iallreduce MPI_ialltoll MPI_ialltoallv MPI_ialltoallw MPI_Ireduce_scatter MPI_Ireduce_scatter_block all2one MPI_Gather MPI_Gatherv MPI_Reduce MPI_Igather MPI_Igatherv MPI_Ireduce attr MPI_Comm_create_keyval MPI_Comm_delete_attr MPI_Comm_free_keyval MPI_Comm_get_attr MPI_Comm_set_attr MPI_Comm_get_name MPI_Comm_set_name MPI_Type_create_keyval	recv MPI_Recv MPI_Irecv MPI_Recv_init MPI_Probe MPI_Iprobe req MPI_Start MPI_Startall MPI_Wait MPI_Waitall MPI_Waitany MPI_Waitsome MPI_Test MPI_Testall MPI_Testany MPI_Testsome MPI_Cancel MPI_Grequest_start MPI_Grequest_complete MPI_Request_get_status MPI_Request_free rma MPI_Accumulate MPI_Get MPI_Put MPI_Win_complete MPI_Win_create MPI_Win_fence MPI_Win_free MPI_Win_get_group MPI_Win_lock
--	--

MPI_Type_delete_attr	MPI_Win_post
MPI_Type_free_keyval	MPI_Win_start
MPI_Type_get_attr	MPI_Win_test
MPI_Type_get_name	MPI_Win_unlock
MPI_Type_set_attr	MPI_Win_wait
MPI_Type_set_name	MPI_Win_allocate
MPI_Win_create_keyval	MPI_Win_allocate_shared
MPI_Win_delete_attr	MPI_Win_create_dynamic
MPI_Win_free_keyval	MPI_Win_shared_query
MPI_Win_get_attr	MPI_Win_attach
MPI_Win_get_name	MPI_Win_detach
MPI_Win_set_attr	MPI_Win_set_info
MPI_Win_set_name	MPI_Win_get_info
MPI_Get_processor_name	MPI_Win_get_accumulate
comm	MPI_Win_fetch_and_op
MPI_Comm_compare	MPI_Win_compare_and_swap
MPI_Comm_create	MPI_Rput
MPI_Comm_dup	MPI_Rget
MPI_Comm_free	MPI_Raccumulate
MPI_Comm_get_name	MPI_Rget_accumulate
MPI_Comm_group	MPI_Win_lock_all
MPI_Comm_rank	MPI_Win_unlock_all
MPI_Comm_remote_group	MPI_Win_flush
MPI_Comm_remote_size	MPI_Win_flush_all
MPI_Comm_set_name	MPI_Win_flush_local
MPI_Comm_size	MPI_Win_flush_local_all
MPI_Comm_split	MPI_Win_sync
MPI_Comm_test_inter	scan
MPI_Intercomm_create	MPI_Exscan
MPI_Intercomm_merge	MPI_Scan
err	MPI_Iexscan
MPI_Add_error_class	MPI_Iscan
MPI_Add_error_code	send
MPI_Add_error_string	MPI_Send
MPI_Comm_call_errhandler	MPI_Bsend
MPI_Comm_create_errhandler	MPI_Rsend
MPI_Comm_get_errhandler	MPI_Ssend
MPI_Comm_set_errhandler	MPI_Isend
MPI_Errhandler_free	MPI_Ibsend
MPI_Error_class	MPI_Irsend
MPI_Error_string	MPI_Issend

MPI_File_call_errhandler	MPI_Send_init
MPI_File_create_errhandler	MPI_Bsend_init
MPI_File_get_errhandler	MPI_Rsend_init
MPI_File_set_errhandler	MPI_Ssend_init
MPI_Win_call_errhandler	sendrecv
MPI_Win_create_errhandler	MPI_Sendrecv
MPI_Win_get_errhandler	MPI_Sendrecv_replace
MPI_Win_set_errhandler	serv
group	MPI_Alloc_mem
MPI_Group_compare	MPI_Free_mem
MPI_Group_difference	MPI_Buffer_attach
MPI_Group_excl	MPI_Buffer_detach
MPI_Group_free	MPI_Op_create
MPI_Group_incl	MPI_Op_free
MPI_Group_intersection	spawn
MPI_Group_range_excl	MPI_Close_port
MPI_Group_range_incl	MPI_Comm_accept
MPI_Group_rank	MPI_Comm_connect
MPI_Group_size	MPI_Comm_disconnect
MPI_Group_translate_ranks	MPI_Comm_get_parent
MPI_Group_union	MPI_Comm_join
init	MPI_Comm_spawn
MPI_Init	MPI_Comm_spawn_multiple
MPI_Init_thread	MPI_Lookup_name
MPI_Finalize	MPI_Open_port
io	MPI_Publish_name
MPI_File_close	MPI_Unpublish_name
MPI_File_delete	status
MPI_File_get_amode	MPI_Get_count
MPI_File_get_atomicity	MPI_Status_set_elements
MPI_File_get_byte_offset	MPI_Status_set_cancelled
MPI_File_get_group	MPI_Test_cancelled
MPI_File_get_info	sync
MPI_File_get_position	MPI_Barrier
MPI_File_get_position_shared	MPI_Ibarrier
MPI_File_get_size	time
MPI_File_get_type_extent	MPI_Wtick
MPI_File_get_view	MPI_Wtime
MPI_File_iread_at	topo
MPI_File_iread	MPI_Cart_coords
MPI_File_iread_shared	MPI_Cart_create

MPI_File_iread_at	MPI_Cart_get
MPI_File_iread	MPI_Cart_map
MPI_File_iread_shared	MPI_Cart_rank
MPI_File_open	MPI_Cart_shift
MPI_File_preallocate	MPI_Cart_sub
MPI_File_read_all_begin	MPI_Cartdim_get
MPI_File_read_all_end	MPI_Dims_create
MPI_File_read_all	MPI_Graph_create
MPI_File_read_at_all_begin	MPI_Graph_get
MPI_File_read_at_all_end	MPI_Graph_map
MPI_File_read_at_all	MPI_Graph_neighbors
MPI_File_read_at	MPI_Graphdims_get
MPI_File_read	MPI_Graph_neighbors_count
MPI_File_read_ordered_begin	MPI_Topo_test
MPI_File_read_ordered_end	type
MPI_File_read_ordered	MPI_Get_address
MPI_File_read_shared	MPI_Get_elements
MPI_File_seek	MPI_Pack
MPI_File_seek_shared	MPI_Pack_external
MPI_File_set_atomicity	MPI_Pack_external_size
MPI_File_set_info	MPI_Pack_size
MPI_File_set_size	MPI_Type_commit
MPI_File_set_view	MPI_Type_contiguous
MPI_File_sync	MPI_Type_create_darray
MPI_File_write_all_begin	MPI_Type_create_hindexed
MPI_File_write_all_end	MPI_Type_create_hvector
MPI_File_write_all	MPI_Type_create_indexed_block
MPI_File_write_at_all_begin	MPI_Type_create_resized
MPI_File_write_at_all_end	MPI_Type_create_struct
MPI_File_write_at_all	MPI_Type_create_subarray
MPI_File_write_at	MPI_Type_dup
MPI_File_write	MPI_Type_free
MPI_File_write_ordered_begin	MPI_Type_get_contents
MPI_File_write_ordered_end	MPI_Type_get_envelope
MPI_File_write_ordered	MPI_Type_get_extent
MPI_File_write_shared	MPI_Type_get_true_extent
MPI_Register_datarep	MPI_Type_indexed
one2all	MPI_Type_size
MPI_Bcast	MPI_Type_vector
MPI_Scatter	MPI_Unpack_external
MPI_Scatterv	MPI_Unpack

MPI_Ibcast	
MPI_Isscatter	
MPI_Isscatterv	

I_MPI_STATS_ACCURACY

Use the I_MPI_STATS_ACCURACY environment variable to decrease statistics output.

Syntax

I_MPI_STATS_ACCURACY=<percentage>

Argument

<percentage>	Float threshold value
--------------	-----------------------

Description

Set this environment variable to collect data only on those MPI functions that take a larger portion of the elapsed time as a percentage of the total time spent inside all MPI calls.

Example

The following example represents a simple application code and IPM summary statistics format:

```
int main (int argc, char *argv[])
{
  int i, rank, size, nsend, nrecv;
  MPI_Init (&argc, &argv);
  MPI_Comm_rank (MPI_COMM_WORLD, &rank);
  nsend = rank;
  MPI_Wtime();
  for (i = 0; i < 200; i++)
  {
    MPI_Barrier(MPI_COMM_WORLD);
  }
  /* open "reduce" region for all processes */
  MPI_Pcontrol(1, "reduce");
  for (i = 0; i < 1000; i++)
    MPI_Reduce(&nsend, &nrecv, 1, MPI_INT, MPI_MAX, 0, MPI_COMM_WORLD);
  /* close "reduce" region */
  MPI_Pcontrol(-1, "reduce");
  if (rank == 0)
  {
    /* "send" region for 0-th process only */
    MPI_Pcontrol(1, "send");
    MPI_Send(&nsend, 1, MPI_INT, 1, 1, MPI_COMM_WORLD);
    MPI_Pcontrol(-1, "send");
  }
  if (rank == 1)
  {
    MPI_Recv(&nrecv, 1, MPI_INT, 0, 1, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
  }
  /* reopen "reduce" region */
  MPI_Pcontrol(1, "reduce");
  for (i = 0; i < 1000; i++)
    MPI_Reduce(&nsend, &nrecv, 1, MPI_INT, MPI_MAX, 0, MPI_COMM_WORLD);
  MPI_Wtime();
  MPI_Finalize ();
  return 0;
}
```

Command:

```
mpiexec -n 4 -env I_MPI_STATS ipm:terse ./a.out
```

Stats output:

```
#####
#
# command : ./a.out (completed)
# host : svlmpihead01/x86_64_Linux mpi_tasks : 4 on 1 nodes
# start : 05/25/11/05:44:13 wallclock : 0.092012 sec
# stop : 05/25/11/05:44:13 %comm : 98.94
# gbytes : 0.00000e+00 total gflop/sec : NA
#
#####
```

Command:

```
mpiexec -n 4 -env I_MPI_STATS ipm ./a.out
```

Stats output:

```
#####
#
# command : ./a.out (completed)
# host : svlmpihead01/x86_64_Linux mpi_tasks : 4 on 1 nodes
# start : 05/25/11/05:44:13 wallclock : 0.092012 sec
# stop : 05/25/11/05:44:13 %comm : 98.94
# gbytes : 0.00000e+00 total gflop/sec : NA
#
#####
# region : * [ntasks] = 4
#
# [total] <avg> min max
# entries 4 1 1 1
# wallclock 0.332877 0.0832192 0.0732641 0.0920119
# user 0.047992 0.011998 0.006999 0.019996
# system 0.013997 0.00349925 0.002999 0.004
# mpi 0.329348 0.082337 0.0723064 0.0912335
# %comm 98.9398 98.6928 99.154
# gflop/sec NA NA NA NA
# gbytes 0 0 0 0
#
#
# [time] [calls] <%mpi> <%wall>
# MPI_Init 0.236192 4 71.71 70.95
# MPI_Reduce 0.0608737 8000 18.48 18.29
# MPI_Barrier 0.027415 800 8.32 8.24
# MPI_Recv 0.00483489 1 1.47 1.45
# MPI_Send 1.50204e-05 1 0.00 0.00
# MPI_Wtime 1.21593e-05 8 0.00 0.00
# MPI_Finalize 3.33786e-06 4 0.00 0.00
# MPI_Comm_rank 1.90735e-06 4 0.00 0.00
# MPI_TOTAL 0.329348 8822 100.00 98.94
#####
# region : reduce [ntasks] = 4
#
# [total] <avg> min max
# entries 8 2 2 2
# wallclock 0.0638561 0.015964 0.00714302 0.0238571
# user 0.034994 0.0087485 0.003999 0.015997
# system 0.003999 0.00099975 0 0.002999
# mpi 0.0608799 0.01522 0.00633883 0.0231845
# %comm 95.3392 88.7417 97.1808
# gflop/sec NA NA NA NA
# gbytes 0 0 0 0
#
#
```



```

# [time] [calls] <%mpi> <%wall>
# MPI_Reduce 0.0608737 8000 99.99 95.33
# MPI_Finalize 3.33786e-06 4 0.01 0.01
# MPI_Wtime 2.86102e-06 4 0.00 0.00
# MPI_TOTAL 0.0608799 8008 100.00 95.34
#####
# region : send [ntasks] = 4
#
# [total] <avg> min max
# entries 1 0 0 1
# wallclock 2.89876e-05 7.24691e-06 1e-06 2.59876e-05
# user 0 0 0 0
# system 0 0 0 0
# mpi 1.50204e-05 3.75509e-06 0 1.50204e-05
# %comm 51.8165 0 57.7982
# gflop/sec NA NA NA
# gbytes 0 0 0 0
#
#
# [time] [calls] <%mpi> <%wall>
# MPI_Send 1.50204e-05 1 100.00 51.82
#####
# region : ipm_noregion [ntasks] = 4
#
# [total] <avg> min max
# entries 13 3 3 4
# wallclock 0.26898 0.0672451 0.0661182 0.068152
# user 0.012998 0.0032495 0.001 0.004999
# system 0.009998 0.0024995 0 0.004
# mpi 0.268453 0.0671132 0.0659676 0.068049
# %comm 99.8039 99.7721 99.8489
# gflop/sec NA NA NA
# gbytes 0 0 0 0
#
#
# [time] [calls] <%mpi> <%wall>
# MPI_Init 0.236192 4 87.98 87.81
# MPI_Barrier 0.027415 800 10.21 10.19
# MPI_Recv 0.00483489 1 1.80 1.80
# MPI_Wtime 9.29832e-06 4 0.00 0.00
# MPI_Comm_rank 1.90735e-06 4 0.00 0.00
# MPI_TOTAL 0.268453 813 100.00 99.80
#####

```

Native and IPM Statistics

The statistics in each supported format can be collected separately. To collect statistics in all formats with the maximal level of details, use the `I_MPI_STATS` environment variable.

I_MPI_STATS

Syntax

```
I_MPI_STATS=all
```

NOTE

The `I_MPI_STATS_SCOPE` environment variable is not applicable when both types of statistics are collected.

To control the amount of statistics information, use the ordinary `I_MPI_STATS` values, separated by comma.

Syntax

`I_MPI_STATS=[native:] [n-]m, ipm[:terse]`

NOTE

Currently the alias `all` corresponds to `I_MPI_STATS=native:20, ipm` and can be changed.

3.5.6. ILP64 Support

The term *ILP64* means that integer, long, and pointer data entities all occupy 8 bytes. This differs from the more conventional LP64 model in which only long and pointer data entities occupy 8 bytes while integer entities occupy 4 bytes. More information on the historical background and the programming model philosophy can be found, for example, in http://www.unix.org/version2/whatsnew/lp64_wp.html

Using ILP64

Use the following options to enable the ILP64 interface

- Use the Fortran compiler driver option `-i8` for separate compilation and the `-ilp64` option for separate linkage. For example,

```
$ mpiifort -i8 -c test.f
```

```
$ mpiifort -ilp64 -o test test.o
```
- For simple programs, use the Fortran compiler driver option `-i8` for compilation and linkage. Specifying `-i8` will automatically assume the ILP64 library. For example,

```
$ mpiifort -i8 test.f
```
- Use the `mpirun -ilp64` option to preload the ILP64 interface. For example,

```
$ mpirun -ilp64 -n 2 ./myprog
```

Known Issues and Limitations

- Data type counts and other arguments with values larger than $2^{31}-1$ are not supported.
- Special MPI types `MPI_FLOAT_INT`, `MPI_DOUBLE_INT`, `MPI_LONG_INT`, `MPI_SHORT_INT`, `MPI_2INT`, `MPI_LONG_DOUBLE_INT`, `MPI_2INTEGER` are not changed and still use a 4-byte integer field.
- Predefined communicator attributes `MPI_APPNUM`, `MPI_HOST`, `MPI_IO`, `MPI_LASTUSED`, `MPI_TAG_UB`, `MPI_UNIVERSE_SIZE`, and `MPI_WTIME_IS_GLOBAL` are returned by the functions `MPI_GET_ATTR` and `MPI_COMM_GET_ATTR` as 4-byte integers. The same holds for the predefined attributes that may be attached to the window and file objects.
- Do not use the `-i8` option to compile MPI callback functions, such as error handling functions, user-defined reduction operations.
- Do not use the `-i8` option with the deprecated functions that store or retrieve the 4-byte integer attribute (for example, `MPI_ATTR_GET`, `MPI_ATTR_PUT`, etc.). Use their recommended alternatives instead (`MPI_COMM_GET_ATTR`, `MPI_COMM_SET_ATTR`, etc.).
- If you want to use the Intel® Trace Collector with the Intel MPI ILP64 executable files, you must use a special Intel Trace Collector library. If necessary, the Intel MPI `mpiifort` compiler driver will select the correct Intel Trace Collector library automatically.

- There is currently no support for C and C++ applications.

3.5.7. Unified Memory Management

The Intel® MPI Library provides a way to replace the memory management subsystem by a user-defined package. You may optionally set the following function pointers:

- `i_malloc`
- `i_calloc`
- `i_realloc`
- `i_free`

These pointers also affect the C++ new and delete operators.

The respective standard C library functions are used by default.

The following contrived source code snippet illustrates the usage of the unified memory subsystem:

```
#include <i_malloc.h>
#include <my_malloc.h>
int main( int argc, int argv )
{
    // override normal pointers
    i_malloc = my_malloc;
    i_calloc = my_calloc;
    i_realloc = my_realloc;
    i_free = my_free;
#ifdef WIN32
    // also override pointers used by DLLs
    i_malloc_dll = my_malloc;
    i_calloc_dll = my_calloc;
    i_realloc_dll = my_realloc;
    i_free_dll = my_free;
#endif
    // now start using Intel(R) libraries
}
```

3.5.8. File System Support

The Intel® MPI Library provides loadable shared modules to provide native support for the following file systems:

- Panasas® ActiveScale® File System (PanFS)
- Parallel Virtual File System®, Version 2 (Pvfs2)
- Lustre® File System
- IBM® General Parallel File System® (GPFS®)

Set the `I_MPI_EXTRA_FILESYSTEM` environment variable to `on` to enable parallel file system support. Set the `I_MPI_EXTRA_FILESYSTEM_LIST` environment variable to request native support for the specific file system. For example, to request native support for Panasas® ActiveScale® File System, do the following:

```
$ mpiexec -env I_MPI_EXTRA_FILESYSTEM on \
-env I_MPI_EXTRA_FILESYSTEM_LIST panfs -n 2 ./test
```

Environment Variables

`I_MPI_EXTRA_FILESYSTEM`

Turn on/off native parallel file systems support.

Syntax

```
I_MPI_EXTRA_FILESYSTEM=<arg>
```

Arguments

<arg>	Binary indicator
enable yes on 1	Turn on native support for the parallel file systems
disable no off 0	Turn off native support for the parallel file systems. This is the default value

Description

Set this environment variable to enable parallel file system support. The

`I_MPI_EXTRA_FILESYSTEM_LIST` environment variable must be set to request native support for the specific file system.

I_MPI_EXTRA_FILESYSTEM_LIST

Select specific file systems support.

Syntax

```
I_MPI_EXTRA_FILESYSTEM_LIST=<fs>[, <fs>, ... , <fs>]
```

Arguments

<fs>	Define a target file system
panfs	Panasas* ActiveScale* File System
pvfs2	Parallel Virtual File System, Version 2
lustre	Lustre* File System
gpfs	IBM* General Parallel File System* (GPFS*)

Description

Set this environment variable to request support for the specific parallel file system. This environment variable is handled only if the `I_MPI_EXTRA_FYLESYSTEM` is enabled. The Intel® MPI Library will try to load shared modules to support the file systems specified by `I_MPI_EXTRA_FILESYSTEM_LIST`.

3.5.9. Multi-threaded memcpy Support

This topic provides information on how to use a multi-threaded version of *memcpy* implemented in the Intel® MPI Library for Intel® Xeon Phi™ Coprocessors. You can use this experimental feature to reach higher memory bandwidth between the ranks communicated through shared memory for some applications.

I_MPI_MT_MEMCPY

Controls usage of the multi-threaded *memcpy*.

Syntax

```
I_MPI_MT_MEMCPY=<value>
```

Arguments

<value>	Controls the usage of the multi-threaded <i>memcpy</i>
enable yes on 1	Enable the multi-threaded <i>memcpy</i> in the single threaded version of the Intel® MPI Library (MPI_THREAD_SINGLE). This configuration is ignored for the thread safe version of Intel® MPI Library
disable no off 0	Disable the usage of the multi-threaded <i>memcpy</i> . This is the default value

Description

Set this environment variable to control whether to use multi-threaded version of *memcpy* for intra-node communication.

I_MPI_MT_MEMCPY_NUM_THREADS

Change the number of threads involved in performing multi-threaded *memcpy*.

Syntax

I_MPI_MT_MEMCPY_NUM_THREADS=<num>

Arguments

<num>	The number of threads involved in performing multi-threaded <i>memcpy</i>
>0	The default value is the lesser of 8 and the number of physical cores within the MPI process pinning domain

Description

Use this environment variable to set the number of threads which perform *memcpy* operations per each MPI rank. The value 1 is equivalent to the setting I_MPI_MT_MEMCPY=disable.

I_MPI_MT_MEMCPY_THRESHOLD

Change the threshold for using multi-threaded *memcpy*.

Syntax

I_MPI_MT_MEMCPY_THRESHOLD=<nbytes>

Arguments

<nbytes>	Define the multi-threaded <i>memcpy</i> threshold in bytes
>0	The default value is 32768

Description

Set this environment variable to control the threshold for using multi-threaded *memcpy*. If the threshold is larger than the shared memory buffer size (for example, see [I_MPI_SHM_LMT_BUFFER_SIZE](#) or [I_MPI_SSHM_BUFFER_SIZE](#)), multi-threaded *memcpy* will never be used. The usage of multi-threaded *memcpy* is selected according to the following scheme:

- Buffers shorter than or equal to <nbytes> are sent using the serial version of *memcpy*. This approach is faster for short and medium buffers.
- Buffers larger than <nbytes> are sent using the multi-threaded *memcpy*. This approach is faster for large buffers.

I_MPI_MT_MEMCPY_SPIN_COUNT

Control the spin count value.

Syntax

I_MPI_MT_MEMCPY_SPIN_COUNT=<scount>

Arguments

<scount>	Define the loop spin count when a thread waits for data to copy before sleeping
>0	The default value is equal to 100000. The maximum value is equal to 2147483647

Description

Set the spin count limit for the loop for waiting for data to be copied by the thread. When the limit is exceeded and there is no data to copy, the thread goes to sleep.

Use the I_MPI_MT_MEMCPY_SPIN_COUNT environment variable for tuning application performance. The best value for <scount> can be chosen on an experimental basis. It depends on the particular computational environment and application.

4. Glossary

cell	A pinging resolution in descriptions for pinning property.
hyper-threading technology	A feature within the IA-64 and Intel® 64 family of processors, where each processor core provides the functionality of more than one logical processor.
logical processor	The basic modularity of processor hardware resource that allows a software executive (OS) to dispatch task or execute a thread context. Each logical processor can execute only one thread context at a time.
multi-core processor	A physical processor that contains more than one processor core.
multi-processor platform	A computer system made of two or more physical packages.
processor core	The circuitry that provides dedicated functionalities to decode, execute instructions, and transfer data between certain sub-systems in a physical package. A processor core may contain one or more logical processors.
physical package	The physical package of a microprocessor capable of executing one or more threads of software at the same time. Each physical package plugs into a physical socket. Each physical package may contain one or more processor cores.
processor topology	Hierarchical relationships of "shared vs. dedicated" hardware resources within a computing platform using physical package capable of one or more forms of hardware multi-threading.

5. Index

\$

\$HOME/.mpd.conf 87

[

-[g]envexcl 16

-[g]envuser 16

{

-{cc cxx fc f77 f90}=<compiler> 11

1

-1 16

A

-a 16

-a <alias> 75

B

-binding 31

-bootstrap <bootstrap server> 30

-bootstrap jmi 31

-bootstrap-exec-args 31

-branch-count <num> 22

C

-check_mpi 9

-check_mpi [<checking_library>] 22, 74

-ckpoint 49

-ckpoint-interval 50

-ckpointlib 50

-ckpoint-logfile 51

-ckpoint-num 50

-ckpoint-prefix 51

-ckpoint-preserve 50

-ckpoint-tmp-prefix 51

-cleanup 29

-compchk 11

-config=<name> 9

-configfile <filename> 22, 73

cpuinfo 90

D

-dapl 70

demux 45

-demux <mode> 28

-disable-x 28

-dynamic_log 10

E

-ecfn <filename> 16, 76

-echo 11, 62

-enable-x 28

-env <ENVVAR> <value> 76

-envall 76

-envexcl <list of env var names> 76

-envlist <list of env var names> 21, 76

-envnone 76

-envuser 76

F

-f <hostsfile> 20

-fast 10

G

-g 10

-g<l-option> 73

-gcc-version=<nnn> 11

-gdb 75

-gdba <jobid> 75

-genv 73

-genv <ENVVAR> <value> 73

-genvall 73, 74

-genvexcl 73

-genvlist 73

-genvnone 73, 74

-genvuser 73, 74

-grr <# of processes> 21, 72

-gtool 23

-gtoolfile 27

H

-h 63, 65, 66, 67, 71, 95

-help 71

--help 62, 63, 64, 65, 66, 67, 68

--help 71

--help 95

-host <nodename> 34, 76

-hostfile<hostfile> 20

-hostos 34, 55

-hosts <nodelist> 27

Hydra 16, 17, 30, 44

I	I_MPI_DAPL_PROVIDER_LIST 58
I_MPI_HYDRA_JMI_LIBRARY 44	I_MPI_DAPL_RDMA_RNDV_WRITE 141
I_MPI_{CC CXX FC F77 F90} 14	I_MPI_DAPL_RNDV_BUFFER_ALIGNMENT 140
I_MPI_{CC CXX FC F77 F90}_PROFILE 12	I_MPI_DAPL_SCALABLE_PROGRESS 139
I_MPI_ADJUST 163	I_MPI_DAPL_SR_BUF_NUM 143
I_MPI_ADJUST_<opname> 163	I_MPI_DAPL_SR_THRESHOLD 143
I_MPI_ADJUST_ALLGATHER_KN_RADIX 163, 170	I_MPI_DAPL_TRANSLATION_CACHE 136
I_MPI_ADJUST_ALLREDUCE_KN_RADIX 163, 170	I_MPI_DAPL_TRANSLATION_CACHE_AVL_TREE 137
I_MPI_ADJUST_BCAST_KN_RADIX 163, 170	I_MPI_DAPL_UD 145, 148
I_MPI_ADJUST_BCAST_SEGMENT 163	I_MPI_DAPL_UD_ACK_RECV_POOL_SIZE 147, 152
I_MPI_ADJUST_GATHERV_KN_RADIX 163, 171	I_MPI_DAPL_UD_ACK_SEND_POOL_SIZE 146, 152
I_MPI_ADJUST_IALLREDUCE_KN_RADIX 163	I_MPI_DAPL_UD_CONN_EVD_SIZE 148, 152
I_MPI_ADJUST_IBCAST_KN_RADIX 163	I_MPI_DAPL_UD_CONNECTION_TIMEOUT 152
I_MPI_ADJUST_IGATHER_KN_RADIX 163	I_MPI_DAPL_UD_CREATE_CONN_QUAL 152
I_MPI_ADJUST_IREDUCE_KN_RADIX 163	I_MPI_DAPL_UD_DESIRED_STATIC_CONNECTIONS_NUM 151, 152
I_MPI_ADJUST_ISCATTER_KN_RADIX 163	I_MPI_DAPL_UD_DFACTOR 152
I_MPI_ADJUST_REDUCE_KN_RADIX 163, 171	I_MPI_DAPL_UD_DIRECT_COPY_THRESHOLD 145, 150, 151, 152
I_MPI_ADJUST_REDUCE_SEGMENT 169	I_MPI_DAPL_UD_EAGER_DYNAMIC_CONNECTION 150, 151
I_MPI_CHECK_COMPILER 13	I_MPI_DAPL_UD_FINALIZE_RETRY_COUNT 152
I_MPI_CHECK_DAPL_PROVIDER_COMPATIBILITY 136	I_MPI_DAPL_UD_FINALIZE_TIMEOUT 152
I_MPI_CHECK_PROFILE 10, 13	I_MPI_DAPL_UD_MAX_MSG_SIZE 152
I_MPI_CKPOINT 51	I_MPI_DAPL_UD_MAX_RDMA_DTOS 152, 153
I_MPI_CKPOINT_INTERVAL 51	I_MPI_DAPL_UD_MAX_RDMA_SIZE 151, 152
I_MPI_CKPOINT_LOGFILE 51	I_MPI_DAPL_UD_MULTIPLE_EAGER_SEND 152
I_MPI_CKPOINT_NUM 51	I_MPI_DAPL_UD_NA_SBUF_LIMIT 152
I_MPI_CKPOINT_PREFIX 51	I_MPI_DAPL_UD_NUMBER_CREDIT_UPDATE 152
I_MPI_CKPOINT_PRESERVE 51	I_MPI_DAPL_UD_PKT_LOSS_OPTIMIZATION 152
I_MPI_CKPOINT_TMP_PREFIX 51	I_MPI_DAPL_UD_PORT 152
I_MPI_CKPOINTLIB 51	I_MPI_DAPL_UD_PROVIDER 145, 151
I_MPI_COMPATIBILITY 177	I_MPI_DAPL_UD_RDMA_MIXED 151
I_MPI_COMPILER_CONFIG_DIR 14	I_MPI_DAPL_UD_RECV_BUFFER_NUM 146, 152
I_MPI_DAPL_BUFFER_NUM 139	I_MPI_DAPL_UD_RECV_EVD_SIZE 148, 152
I_MPI_DAPL_BUFFER_SIZE 140, 154	I_MPI_DAPL_UD_REQ_EVD_SIZE 148, 152
I_MPI_DAPL_CHECK_MAX_RDMA_SIZE 141	I_MPI_DAPL_UD_REQUEST_QUEUE_SIZE 152
I_MPI_DAPL_CONN_EVD_SIZE 142	I_MPI_DAPL_UD_RESENT_TIMEOUT 152
I_MPI_DAPL_DESIRED_STATIC_CONNECTIONS_NUM 144	I_MPI_DAPL_UD_RNDV_BUFFER_ALIGNMENT 149, 152
I_MPI_DAPL_DIRECT_COPY_THRESHOLD 137	
I_MPI_DAPL_EAGER_MESSAGE_AGGREGATION 138	
I_MPI_DAPL_LOCALITY_THRESHOLD 59	
I_MPI_DAPL_MAX_MSG_SIZE 142	

- I_MPI_DAPL_UD_RNDV_COPY_ALIGNMENT_THRESHOLD 149, 152
- I_MPI_DAPL_UD_RNDV_DYNAMIC_CONNECTION 150, 151
- I_MPI_DAPL_UD_RNDV_MAX_BLOCK_LEN 149
- I_MPI_DAPL_UD_SEND_BUFFER_NUM 146, 152
- I_MPI_DAPL_UD_SEND_BUFFER_SIZE 152
- I_MPI_DAPL_UD_TRANSLATION_CACHE 147, 152
- I_MPI_DAPL_UD_TRANSLATION_CACHE_AVL_TREE 147, 152
- I_MPI_DAPL_UD_TRANSLATION_CACHE_MAX_ENTRY_NUM 152
- I_MPI_DAPL_UD_TRANSLATION_CACHE_MAX_MEMORY_SIZE 152
- I_MPI_DAT_LIBRARY 136
- I_MPI_DEBUG 10, 77
- I_MPI_DEBUG_INFO_STRIP 15
- I_MPI_DEBUG_OUTPUT 78, 79
- I_MPI_DYNAMIC_CONNECTION 129
- I_MPI_DYNAMIC_CONNECTION_MODE 138
- I_MPI_EAGER_THRESHOLD 127
- I_MPI_ENV_PREFIX_LIST 60
- I_MPI_EXTRA_FILESYSTEM 196
- I_MPI_EXTRA_FILESYSTEM_LIST 196, 197
- I_MPI_FABRICS 123, 126, 135
- I_MPI_FABRICS_LIST 125, 178
- I_MPI_FALLBACK 126
- I_MPI_FAULT_CONTINUE 178, 179
- I_MPI_GTOOL 46
- I_MPI_HYDRA_BOOTSTRAP 31, 39
- I_MPI_HYDRA_BOOTSTRAP_AUTOFORK 41
- I_MPI_HYDRA_BOOTSTRAP_EXEC 40
- I_MPI_HYDRA_BOOTSTRAP_EXEC_EXTRA_ARGS 40
- I_MPI_HYDRA_BRANCH_COUNT 43
- I_MPI_HYDRA_CLEANUP 45
- I_MPI_HYDRA_DEBUG 37
- I_MPI_HYDRA_DEMUX 45
- I_MPI_HYDRA_ENV 38
- I_MPI_HYDRA_GDB_REMOTE_SHELL 44
- I_MPI_HYDRA_HOST_FILE 37
- I_MPI_HYDRA_IFACE 44
- I_MPI_HYDRA_JMI_LIBRARY 31, 44
- I_MPI_HYDRA_PMI_AGGREGATE 22, 43
- I_MPI_HYDRA_PMI_CONNECT 41
- I_MPI_HYDRA_RMK 34, 41
- I_MPI_HYDRA_USE_APP_TOPOLOGY 47
- I_MPI_INTRANODE_EAGER_THRESHOLD 127, 132, 134, 135
- I_MPI_JOB_ABORT_SIGNAL 82
- I_MPI_JOB_CHECK_LIBS 43, 74, 80
- I_MPI_JOB_CONFIG_FILE 87
- I_MPI_JOB_CONTEXT 66, 88
- I_MPI_JOB_FAST_STARTUP 84
- I_MPI_JOB_RESPECT_PROCESS_PLACEMENT 46
- I_MPI_JOB_SIGNAL_PROPAGATION 39
- I_MPI_JOB_STARTUP_TIMEOUT 81
- I_MPI_JOB_TAGGED_PORT_OUTPUT 88
- I_MPI_JOB_TIMEOUT 38, 81
- I_MPI_JOB_TIMEOUT_SIGNAL 38, 39, 82
- I_MPI_JOB_TRACE_LIBS 42, 74, 80
- I_MPI_LARGE_SCALE_THRESHOLD 126
- I_MPI_LINK 15
- I_MPI_MIC 57
- I_MPI_MIC_POSTFIX 58
- I_MPI_MIC_PREFIX 57
- I_MPI_MPD_CHECK_PYTHON 88, 89
- I_MPI_MPD_CLEAN_LOG 90
- I_MPI_MPD_RSH 89
- I_MPI_MPD_TMPDIR 89
- I_MPI_MPIRUN_CLEANUP 17
- I_MPI_MT_MEMCPY 197
- I_MPI_MT_MEMCPY_NUM_THREADS 198
- I_MPI_MT_MEMCPY_SPIN_COUNT 198
- I_MPI_MT_MEMCPY_THRESHOLD 198
- I_MPI_OFA_ADAPTER_NAME 157
- I_MPI_OFA_DYNAMIC_QPS 160
- I_MPI_OFA_LIBRARY 161
- I_MPI_OFA_NONSWITCH_CONF 161
- I_MPI_OFA_NUM_ADAPTERS 157
- I_MPI_OFA_NUM_PORTS 158
- I_MPI_OFA_NUM_RDMA_CONNECTIONS 158, 159
- I_MPI_OFA_PACKET_SIZE 160
- I_MPI_OFA_RAIL_SCHEDULER 159
- I_MPI_OFA_SWITCHING_TO_RDMA 158

I_MPI_OFA_TRANSLATION_CACHE 159
I_MPI_OFA_TRANSLATION_CACHE_AVL_TREE 160
I_MPI_OFI_LIBRARY 162
I_MPI_OFI_PROVIDER 162
I_MPI_OFI_PROVIDER_DUMP 162
I_MPI_OUTPUT_CHUNK_SIZE 83
I_MPI_PERHOST 42, 79, 80
I_MPI_PIN 103, 107
I_MPI_PIN_CELL 109
I_MPI_PIN_DOMAIN 110, 111, 112, 121
I_MPI_PIN_MODE 104
I_MPI_PIN_ORDER 120
I_MPI_PIN_PROCESSOR_EXCLUDE_LIST 103
I_MPI_PIN_PROCESSOR_LIST 104, 111
I_MPI_PIN_PROCS 104
I_MPI_PIN_RESPECT_CPUSET 103
I_MPI_PIN_RESPECT_HCA 103
I_MPI_PMI_EXTENSIONS 83
I_MPI_PMI_LIBRARY 77
I_MPI_PRINT_VERSION 80
I_MPI_PROCESS_MANAGER 17, 95
I_MPI_RESTART 53
I_MPI_ROOT 14
I_MPI_SCALABLE_OPTIMIZATION 128
I_MPI_SHM_BYPASS 135
I_MPI_SHM_CACHE_BYPASS 129
I_MPI_SHM_CACHE_BYPASS_THRESHOLD 130
I_MPI_SHM_CACHE_BYPASS_THRESHOLDS 130
I_MPI_SHM_CELL_NUM 131
I_MPI_SHM_CELL_SIZE 127, 132
I_MPI_SHM_FBOX_SIZE 131
I_MPI_SHM_LMT 132
I_MPI_SHM_LMT_BUFFER_NUM 133
I_MPI_SHM_LMT_BUFFER_SIZE 133
I_MPI SOCK_SCALABLE_OPTIMIZATION 128
I_MPI_SPIN_COUNT 127, 135
I_MPI_SSHM 133
I_MPI_SSHM_BUFFER_NUM 134
I_MPI_SSHM_BUFFER_SIZE 134
I_MPI_SSHM_DYNAMIC_CONNECTION 134
I_MPI_STATS 184

I_MPI_STATS_ACCURACY 186
I_MPI_STATS_BUCKETS 183
I_MPI_STATS_FILE 183
I_MPI_STATS_SCOPE 180
I_MPI_TCP_BUFFER_SIZE 154
I_MPI_TCP_NETMASK 153
I_MPI_TCP_POLLING_MODE 155
I_MPI_THREAD_LEVEL_DEFAULT 86
I_MPI_TIMER_KIND 177
I_MPI_TMI_DSEND 156
I_MPI_TMI_LIBRARY 155
I_MPI_TMI_NBITS_RANK 156
I_MPI_TMI_PROVIDER 155
I_MPI_TMPDIR 45
I_MPI_TRACE_PROFILE 9, 13
I_MPI_TUNE_APPLICATION_COMMUNICATION_GRAPH 101
I_MPI_TUNE_APPLICATION_STATISTICS 101
I_MPI_TUNE_FAST 100
I_MPI_TUNE_HARDWARE_TOPOLOGY_GRAPH 102
I_MPI_TUNE_RANK_PLACEMENT 100
I_MPI_WAIT_MODE 125, 128
I_MPI_YARN 18
-ib 36, 70, 71
-iface <interface> 28
-ifhn <interface/hostname> 16, 62, 76
-ilp64 10
IPM 186
L
-l 66, 67, 75
large message transfer (LMT) 132, 133
LD_LIBRARY_PATH 44
libjmi.so 31
-link_mpi 10
LMT 132, 133
-localhost 30
--locons 16
M
-m 16, 63, 64, 75
-machine <machine file> 20
-machinefile <machine file> 20, 72
max_rdma_size 141

--maxbranch=<maxbranch> |-b <maxbranch> 16
 mpd 17, 62, 63, 65, 67, 87, 88, 89, 103
 mpd.hosts 63, 87
 MPD_SECRETWORD 87
 --mpd=<mpdcmd>|-m <mpdcmd> 16
 mpdallexit 65, 90
 mpdboot 16, 17, 63, 87
 mpdcheck 66
 mpdcleanup 65
 mpdexit 64
 mpdhelp 69
 mpdkilljob 68
 mpdlistjobs 67
 mpdringtest 67
 mpdsigjob 68
 mpdtrace 17, 66
 MPI_Allgather 168, 181
 MPI_Allgatherv 168, 181
 MPI_Allreduce 168, 181, 184
 MPI_Alltoall 168, 181
 MPI_Alltoallv 168, 181
 MPI_Alltoallw 168, 181
 MPI_ANY_SOURCE 155
 MPI_Barrier 168, 181
 MPI_Bcast 168, 181
 MPI_COMM_JOIN 178
 MPI_COMM_WORLD 179
 MPI_ERRORS_RETURN 179
 MPI_Exscan 168, 181
 MPI_Finalize 185
 MPI_Gather 168, 181
 MPI_Gatherv 168, 182
 MPI_Init 185
 MPI_Reduce 168, 182
 MPI_Reduce_scatter 168
 MPI_Scan 168, 182
 MPI_Scatter 168, 182
 MPI_Scatterv 168, 182
 mpicleanup 45, 48, 49
 mpiexec 16, 17, 42, 69, 70, 72, 73, 74, 75, 77, 81, 82, 83, 98
 mpiexec.hydra 16, 18, 19, 22, 23, 27, 29, 34, 38, 39, 43, 48
 mpiicc -g 79
 mpirun 16, 17, 44
 mpitune 28, 72, 94
 -mx 36, 71
 N
 -n 17, 20
 -n <# of processes> 34, 76
 --ncpus=<ncpus> 16
 -noconf 76
 -nolocal 27, 72
 -nostrip 9
 -np 17
 -np <# of processes> 34, 76
 O
 -O 10
 -ofi 37
 --ordered |-o 16
 -ordered-output 29, 75
 P
 --parallel-startup |-p 16
 PATH 9
 -path <directory> 29, 34, 77
 -perhost 20
 -perhost <# of processes> 21
 -perhost <# of processes> 72
 pmi_proxy 22
 -pmi-aggregate 22
 -pmi-connect <mode> 21
 -pmi-noaggregate 22
 -ppn <# of processes> 21
 -ppn <# of processes> 72
 -print-all-exitcodes 34
 -print-rank-map 34
 -profile=<profile_name> 9, 13
 -psm 36, 71
 R
 -rdma 35, 70
 --remcons 16
 -restart 50
 -rmk <RMK> 33
 -rr 21, 72

S

-s <spec> 28, 75

secretword 87

-shell|-s 16

-show 11

-static 9

-static_mpi 9

T

-t 9

-t [<profiling_library>] 21, 74

-tmi 36, 71

-tmpdir 29

TMPDIR 45, 89

TotalView 23, 74, 84

-trace 9, 13, 43

-trace [<profiling_library>] 21, 74

-tune 28, 98

-tune [<arg >] 72

-tv 22, 74

-tva <jobid> 75

-tva <pid> 23

TVDSVRLAUNCHCMD 74

-tvsu 17, 75

U

-umask <umask> 35, 77

-use-app-topology 29

--user=<user> | -u <user> 16

V

-v 11, 66, 67, 71

-verbose 34, 63, 98

-version 71

-version or -V 29

VT_ROOT 9, 10, 14

W

-wdir <directory> 35, 77