

Intel[®] MPI Library for Windows* OS

Developer Reference

Contents

| | |
|--|-----------|
| Legal Information | 3 |
| 1. Introduction | 4 |
| 1.1. Introducing Intel® MPI Library..... | 4 |
| 1.2. What's New..... | 4 |
| 1.3. Notational Conventions..... | 5 |
| 1.4. Related Information | 5 |
| 2. Command Reference | 6 |
| 2.1. Compilation Commands..... | 6 |
| 2.1.1. Compilation Command Options | 7 |
| 2.2. mpiexec..... | 8 |
| 2.2.1. Global Options..... | 9 |
| 2.2.2. Local Options | 18 |
| 2.3. cpuinfo | 19 |
| 2.4. impi_info..... | 21 |
| 2.5. mpitune..... | 22 |
| 2.5.1. mpitune Configuration Options | 22 |
| 3. Environment Variable Reference | 28 |
| 3.1. Compilation Environment Variables | 28 |
| 3.2. Hydra Environment Variables | 30 |
| 3.3. I_MPI_ADJUST Family Environment Variables | 35 |
| 3.4. Process Pinning..... | 54 |
| 3.4.1. Processor Identification..... | 54 |
| 3.4.2. Default Settings..... | 55 |
| 3.4.3. Environment Variables for Process Pinning | 55 |
| 3.4.4. Interoperability with OpenMP* API | 61 |
| 3.5. Environment Variables for Fabrics Control..... | 69 |
| 3.5.1. Communication Fabrics Control..... | 69 |
| 3.5.2. OFI*-capable Network Fabrics Control..... | 70 |
| 3.6. Other Environment Variables..... | 71 |
| 4. Miscellaneous | 76 |
| 4.1. User Authorization..... | 76 |

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 known as errata which may cause deviations from published specifications. Current characterized errata are available on request.

Intel technologies features and benefits depend on system configuration and may require enabled hardware, software or service activation. Learn more at Intel.com, or from the OEM or retailer.

Copies of documents which have an order number and are referenced in this document may be obtained by calling 1-800-548-4725 or by visiting www.intel.com/design/literature.htm.

Intel, the Intel logo, Xeon, and Xeon Phi are trademarks of Intel Corporation in the U.S. and/or other countries.

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

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

Copyright 2003-2018 Intel Corporation.

This software and the related documents are Intel copyrighted materials, and your use of them is governed by the express license under which they were provided to you (**License**). Unless the License provides otherwise, you may not use, modify, copy, publish, distribute, disclose or transmit this software or the related documents without Intel's prior written permission.

This software and the related documents are provided as is, with no express or implied warranties, other than those that are expressly stated in the License.

1. Introduction

This *Developer Reference* provides you with the complete reference for the Intel® MPI Library. It is intended to help an experienced user fully utilize the Intel MPI Library functionality. You can freely redistribute this document in any desired form.

1.1. Introducing Intel® MPI Library

Intel® MPI Library is a multi-fabric message passing library that implements the Message Passing Interface, v3.1 (MPI-3.1) specification. It provides a standard library across Intel® platforms that enable adoption of MPI-3.1 functions as their needs dictate.

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.

You can get the latest information of Intel® MPI Library at <https://software.intel.com/intel-mpi-library>.

1.2. What's New

This document reflects the updates for Intel® MPI Library 2019 release for Windows* OS:

The following latest changes in this document were made:

Intel MPI Library 2019

- Document overhaul to align with supported functionality.
- Removed the `I_MPI_HARD_FINALIZE`, `I_MPI_MIC`, `I_MPI_ENV_PREFIX_LIST`, `I_MPI_TUNE*`, `I_MPI_ENV_PREFIX_LIST`, `I_MPI_JOB_FAST_STARTUP`, `I_MPI_FALLBACK`, `I_MPI_DAPL*`, `I_MPI_LARGE_SCALE_THRESHOLD`, `I_MPI_OFA*`, `I_MPI_TCP*`, `I_MPI_TMI*` environment variables.
- Removed the `-hostos` option from [Local Options](#).
- Added the `I_MPI_OFI_LIBRARY_INTERNAL` environment variable to [OFI-capable Network Fabrics Control](#).
- Added an option for setting `MPI_UNIVERSE_SIZE` to [Global Options](#).
- Added new collective operations to [I_MPI_ADJUST Family Environment Variables](#).
- Added new variables `I_MPI_SHM_CELL_EXT_SIZE` and `I_MPI_SHM_CELL_EXT_NUM_TOTAL` to Shared Memory Control.
- Added [impi_info](#) utility.
- Updated [mpitune](#) utility.

Intel MPI Library 2018 Update 3

- Added new algorithms for `I_MPI_ADJUST_ALLREDUCE` to [I_MPI_ADJUST Family](#).

Intel MPI Library 2018 Update 2

- Improved shm performance with collective operations (`I_MPI_THREAD_YIELD`)
- Bug fixes

Intel MPI Library 2018 Update 1

- Minor changes.

Intel MPI Library 2018

- Removed support of the Intel® Xeon Phi™ coprocessors (formerly code named Knights Corner).

- Changes in environment variables:
 - `I_MPI_DAPL_TRANSLATION_CACHE` is now disabled by default

Intel MPI Library 2017 Update 2

- Added the environment variable `I_MPI_HARD_FINALIZE` in Other Environment Variables.

Intel MPI Library 2017 Update 1

- Topology-aware collective communication algorithms support (`I_MPI_ADJUST` Family).
- Added a new algorithm for `I_MPI_ADJUST_GATHER` and related environment variable `I_MPI_ADJUST_GATHER_SEGMENT` (`I_MPI_ADJUST` Family).
- Added the environment variable `I_MPI_PORT_RANGE` in Hydra Environment Variables.

Intel MPI Library 2017

- Document layout changes.

1.3. 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) |

1.4. 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

2.1. Compilation Commands

The following table lists the available Intel® MPI Library compiler commands with their underlying compilers and programming languages.

Intel® MPI Library Compiler Wrappers

| Compiler Command | Underlying Compiler | Supported Language(s) |
|---|---------------------|-----------------------|
| Common Compilers | | |
| mpicc.bat | cl.exe | C |
| mpicxx.bat | cl.exe | C++ |
| mpifc.bat | ifort.exe | Fortran 77/Fortran 95 |
| Microsoft* Visual C++* Compilers | | |
| mpicl.bat | cl.exe | C/C++ |
| Intel® Fortran, C++ Compilers | | |
| mpiicc.bat | icl.exe | C |
| mpiicpc.bat | icl.exe | C++ |
| mpiifort.bat | ifort.exe | Fortran 77/Fortran 95 |

NOTES:

- Compiler commands are available only in the Intel® MPI Library Software Development Kit (SDK).
- For the supported versions of the listed compilers, refer to the *Release Notes*.
- Compiler wrapper scripts are located in the `<installdir>\intel64\bin` directory.
- The environment settings can be established by running the `<installdir>\intel64\bin\mpivars.bat` script. If you need to use a specific library configuration, you can pass one of the following arguments to the `mpivars.bat` script to switch to the corresponding configuration: `release` or `debug`. The ordinary multi-threaded optimized library is chosen by default. Alternatively, you can use the `I_MPI_LIBRARY_KIND` environment variable to specify a configuration and source the script without arguments.
- Ensure that the corresponding underlying compiler is already in your `PATH`. If you use the Intel® Compilers, run the `compilervars.bat` script from the installation directory to set up the compiler environment.
- To display mini-help of a compiler command, execute it without any parameters.

2.1.1. Compilation Command Options

-profile=<profile_name>

Use this option to specify an MPI profiling library. <profile_name> is the name of the configuration file (profile) that loads the corresponding profiling library. The profiles are taken from <installdir>\<arch>\etc.

You can create your own profile as <installdir>\<arch>\etc\<profile_name>.conf. You can define the following environment variables in a configuration file:

- PROFILE_PRELIB – libraries (and paths) to load before the Intel® MPI Library
- PROFILE_POSTLIB – libraries to load after the Intel® MPI Library
- PROFILE_INCPATHS – C preprocessor arguments for any include files

For example, create a file <installdir>\<arch>\etc\myprof.conf with the following lines:

```
SET PROFILE_PRELIB=<path_to_myprof>\lib\myprof.lib
SET PROFILE_INCPATHS=-I"<paths_to_myprof>\include"
```

Use the -profile=myprof option for the relevant compiler wrapper to select this new profile.

-t or -trace

Use the -t or -trace option to link the resulting executable file against the Intel® Trace Collector library.

To use this option, include the installation path of the Intel® Trace Collector in the VT_ROOT environment variable. Source the itacvars.bat script provided in the Intel® Trace Analyzer and Collector installation folder.

-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 Intel® Fortran Compiler, you still have to use the ilp64 option for linkage.

-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, such as -zi.

/Zi, /Z7 or /ZI

Use these options to compile a program in debug mode and link the resulting executable against the debugging version of the Intel® MPI Library. See I_MPI_DEBUG for information on how to use additional debugging features with the /zi, /z7, /ZI or debug builds.

NOTE

The /ZI option is only valid for C/C++ compiler.

-O

Use this option to enable compiler optimization.

Setting this option triggers a call to the `libirc` library. Many of those library routines are more highly optimized for Intel microprocessors than for non-Intel microprocessors.

-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:

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}=<compiler>

Use this option to select the underlying compiler.

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

```
> mpiicc -cc=icl.exe -c test.c
```

For this to work, `icl.exe` should be in your `PATH`. Alternatively, you can specify the full path to the compiler.

NOTE

This option works only with the `mpiicc.bat` and the `mpifc.bat` commands.

-v

Use this option to print the compiler wrapper script version.

2.2. mpiexec

Launches an MPI job using the Hydra process manager.

Syntax

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

or

```
mpiexec <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 argument set |

| | |
|---------------------------------|--|
| <code><executable></code> | <code><name>.exe</code> or <code>path\name</code> of the executable file |
|---------------------------------|--|

Description

Use the `mpiexec` utility to run MPI applications using the Hydra process manager.

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 `test.exe` over the specified processes and hosts:

```
> mpiexec -f <hostfile> -n <# of processes> test.exe
```

where:

- `<# of processes>` specifies the number of processes on which to run the `test.exe` executable
- `<hostfile>` specifies a list of hosts on which to run the `test.exe` 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 -f <hostfile> -env <VAR1> <VAL1> -n 2 prog1.exe : ^
-env <VAR2> <VAL2> -n 2 prog2.exe
```

NOTE

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

2.2.1. Global Options

This section describes the global options of the Intel® MPI Library's Hydra process manager. Global options are applied to all arguments sets in the launch command. Argument sets are separated by a colon ':':

-usize <usize>

Use this option to set `MPI_UNIVERSE_SIZE`, which is available as an attribute of the `MPI_COMM_WORLD`.

| | |
|----------------------------|--|
| <code><size></code> | Define the universe size |
| SYSTEM | Set the size equal to the number of cores passed to <code>mpiexec</code> through the hostfile or the resource manager. |
| INFINITE | Do not limit the size. This is the default value. |
| <code><value></code> | Set the size to a numeric value ≥ 0 . |

-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 `I_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 process placement through a machine file. To define the total number of processes to start, use the `-n` option. To pin processes within a machine, use the option `binding=map` in the machine file. For example:

```
> type machinefile
node0:2 binding=map=0,3
node1:2 binding=map=[2,8]
node0:1 binding=map=8
```

For details, see the `-binding` option description.

-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.

-genvexcl <list of env var names>

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

-genvlist <list>

Use this option to pass a list of environment variables with their current values. `<list>` 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:

| | |
|---------------------------|---|
| <code><mode></code> | The caching mode 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. |
| <code>lazy-cache</code> | <code>cache</code> mode with on-request propagation of the PMI information. |
| <code>alltoall</code> | Information is automatically exchanged between all <code>pmi_proxy</code> before any get request can be done. This |

| | |
|--|----------------------|
| | is the default mode. |
|--|----------------------|

See the `I_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 specified number of consecutive MPI processes on every host in the group using round robin scheduling. See the `I_MPI_PERHOST` environment variable for more details.

NOTE

When running under a job scheduler, these options are ignored by default. To be able to control process placement with these options, disable the `I_MPI_JOB_RESPECT_PROCESS_PLACEMENT` variable.

-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 `I_MPI_PERHOST` environment variable for more details.

-trace-pt2pt

Use this option to collect the information about point-to-point operations using Intel® Trace Analyzer and Collector. The option requires that your application be linked against the Intel® Trace Collector profiling library.

-trace-collectives

Use this option to collect the information about collective operations using Intel® Trace Analyzer and Collector. The option requires that your application be linked against the Intel® Trace Collector profiling library.

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 Hydra process manager, or by each `pmi_proxy` management process.

See the `I_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 `I_MPI_HYDRA_PMI_AGGREGATE` environment variable for more details.

-nolocal

Use this option to avoid running the `<executable>` on the host where `mpiexec` 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 the MPI processes should be run. For example, the following command runs the executable `a.out` on the hosts `host1` and `host2`:

```
> mpiexec -n 2 -ppn 1 -hosts host1,host2 test.exe
```

NOTE

If `<nodelist>` contains 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 -n 2 -iface ib0 test.exe
```

See the `I_MPI_HYDRA_IFACE` environment variable for more details.

-l, -prepend-rank

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.hydra` 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.

-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.

-delegate

Use this option to enable the domain-based authorization with the delegation ability. See [User Authorization](#) for details.

-impersonate

Use this option to enable the limited domain-based authorization. You will not be able to open files on remote machines or access mapped network drives. See [User Authorization](#) for details.

-localhost

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

-localroot

Use this option to launch the root process directly from `mpiexec` if the host is local. You can use this option to launch GUI applications. The interactive process should be launched before any other process in a job. For example:

```
> mpiexec -n 1 -host <host2> -localroot interactive.exe : -n 1 -host <host1>  
background.exe
```

-localonly

Use this option to run an application on the local node only. If you use this option only for the local node, the Hydra service is not required.

-register

Use this option to encrypt the user name and password to the registry.

-remove

Use this option to delete the encrypted credentials from the registry.

-validate

Validate the encrypted credentials for the current host.

-whoami

Use this option to print the current user name.

-map <drive:|\\host\share>

Use this option to create network mapped drive on nodes before starting executable. Network drive will be automatically removed after the job completion.

-mapall

Use this option to request creation of all user created network mapped drives on nodes before starting executable. Network drives will be automatically removed after the job completion.

-logon

Use this option to force the prompt for user credentials.

-noprompt

Use this option to suppress the prompt for user credentials.

-port/-p

Use this option to specify the port that the service is listening on.

-verbose or -v

Use this option to print debug information from `mpirexec`, such as:

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

See the `I_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.

| | |
|--------------------------|---|
| <code><arg></code> | 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>slurm</code> | Use SLURM* <code>srun</code> command. |

| | |
|-----|--|
| ll | Use LoadLeveler* <code>llspawn.stdio</code> command. |
| lsf | Use LSF <code>blaunch</code> command. |
| sge | Use Univa* Grid Engine* <code>qrsh</code> command. |

-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.

NOTE

This option is related to the family of `I_MPI_PIN` environment variables, which have higher priority than the `-binding` option. Hence, if any of these variables are set, the option is ignored.

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

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

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

| Parameter | |
|-----------|--|
| map | Process mapping |
| Values | |
| 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 | 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_i takes one of the following values: <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 |
| $[m_0, m_1, \dots, m_n]$ | The i^{th} process is mapped upon the processor subset defined by m_i hexadecimal mask using the following rule: The j^{th} processor is included into the subset m_i if the j^{th} bit of m_i equals 1. |

| Parameter | |
|-----------|---|
| domain | Processor domain set on a node |
| Values | |
| 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. |
| <code><size>[:<layout>]</code> | <p>Each domain of the set consists of <code><size></code> processor cells. <code><size></code> may have the following values:</p> <ul style="list-style-type: none"> <code>auto</code> - domain size = #cells/#processes <code>omp</code> - domain size = OMP_NUM_THREADS environment variable value positive integer - exact value of the domain size <hr/> <p>NOTE</p> <p>Domain size is limited by the number of processor cores on the node.</p> <hr/> <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> |

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

| Parameter | |
|-----------|--|
| offset | Domain list offset |
| Values | |
| <n> | 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. |

2.2.2. Local Options

This section describes the local options of the Intel® MPI Library's Hydra process manager. Local options are applied only to the argument set they are specified in. Argument sets are separated by a colon ': '.

-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 argument set.

-envall

Use this option to propagate all environment variables in the current argument set. See the [I_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 argument set.

-envexcl <list of env var names>

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

-envlist <list>

Use this option to pass a list of environment variables with their current values. <list> 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 `test.exe` on hosts `host1` and `host2`:

```
> mpiexec -n 2 -host host1 test.exe : -n 2 -host host2 test.exe
```

-path <directory>

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

-wdir <directory>

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

NOTE

The option may work incorrectly if the file path contains Unicode characters.

-umask <umask>

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

2.3. cpuinfo

Provides information on processors used in the system.

Syntax

```
cpuinfo [[-]<options>]
```

Arguments

| | |
|-----------|--|
| <options> | Sequence of one-letter options. Each option controls a specific part of the output data. |
| g | <p>General information about single cluster node shows:</p> <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 |
| i | <p>Logical processors identification table identifies threads, cores, and packages of each logical processor accordingly.</p> <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. |
| d | <p>Node decomposition table shows the node contents. Each entry contains the information on packages, cores, and logical processors.</p> <ul style="list-style-type: none"> • <i>Package Id</i> - physical package identifier. • <i>Cores Id</i> - list of core identifiers that belong to this |

| | |
|------|---|
| | <p>package.</p> <ul style="list-style-type: none"> Processors <i>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> |
| gidc | 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 architecture.

The `cpuinfo` utility is available for both Intel microprocessors and non-Intel microprocessors, but it may provide only partial information about non-Intel microprocessors.

An example of the `cpuinfo` output:

```
> cpuinfo -gdcs
```

```

===== Processor composition =====
Processor name      : Intel(R) Xeon(R)  X5570
Packages(sockets)  : 2
Cores              : 8
Processors(CPUs)   : 8
Cores per package  : 4
Threads per core   : 1
===== Processor identification =====
Processor          Thread Id.      Core Id.      Package Id.
0                 0                0             0
1                 0                0             1
2                 0                1             0
3                 0                1             1
4                 0                2             0
5                 0                2             1
6                 0                3             0
7                 0                3             1
===== Placement on packages =====
Package Id.        Core Id.      Processors
0                  0,1,2,3      0,2,4,6
1                  0,1,2,3      1,3,5,7
===== Cache sharing =====
Cache   Size      Processors
L1      32 KB      no sharing
L2      256 KB     no sharing
L3      8 MB       (0,2,4,6) (1,3,5,7)
===== Processor Signature =====

```

| xFamily | xModel | Type | Family | Model | Stepping |
|---------|--------|------|--------|-------|----------|
| 00 | 1 | 0 | 6 | a | 5 |

2.4. impi_info

Provides information on available Intel® MPI Library environment variables.

Syntax

```
impi_info <options>
```

Arguments

| | |
|----------------|--|
| <options> | List of options. |
| -a -all | Show all IMPI variables. |
| -h -help | Show a help message. |
| -v -variable | Show all available variables or description of the specified variable. |
| -c -category | Show all available categories or variables of the specified category. |

Description

The `impi_info` utility provides information on environment variables available in the Intel MPI Library. For each variable, it prints out the name, the default value, and the value data type. By default, a reduced list of variables is displayed. Use the `-all` option to display all available variables with their descriptions.

The example of the `impi_info` output:

```
> impi_info
| NAME | DEFAULT VALUE | DATA TYPE |
=====
| I_MPI_THREAD_SPLIT | 0 | MPI_INT |
| I_MPI_THREAD_RUNTIME | none | MPI_CHAR |
| I_MPI_THREAD_MAX | -1 | MPI_INT |
| I_MPI_THREAD_ID_KEY | thread_id | MPI_CHAR |
```

2.5. mpitune

Tunes the Intel® MPI Library parameters for the given MPI application.

Syntax

```
mpitune <options>
```

Arguments

| | |
|----------------------|-------------------------------|
| <options> | List of options. |
| -c --config <file> | Specify a configuration file. |
| -h --help | Display the help message. |
| -v --version | Display the product version. |

Description

The `mpitune` utility allows you to automatically adjust Intel MPI Library parameters, such as collective operation algorithms, to your cluster configuration or application.

The tuner iteratively launches a benchmarking application with different configurations to measure performance and stores the results of each launch. Based on these results, the tuner generates optimal values for the parameters that are being tuned.

All tuner parameters should be specified in a configuration file, which is passed to the tuner with the `--config` option. A typical configuration file consists of the main section, specifying generic options, and search space sections for specific library parameters (for example, for specific collective operations). To comment a line, use the hash symbol `#`. A configuration file example is available at `<installdir>/etc/tune_cfg/example`.

2.5.1. mpitune Configuration Options

Application Options

-app

Sets a template for the command line to be launched to gather tuning results. The command line can contain variables declared as `@<var_name>@`. The variables are defined further on using other options.

For example:

```
-app: mpirun -np @np@ -ppn @ppn@ IMB-MPI1 -msglog 0:@logmax@ -npmin @np@ @func@
```

NOTE

The application must produce output (in `stdout` or file or any other destination) that can be parsed by the tuner to pick the value to be tuned and other variables. See the `-app-regex` and `-app-regex-legend` options below for details.

-app-regex

Sets a regular expression to be evaluated to extract the required values from the application output. Use regular expression groups to assign the values to variables. Variables and groups associations are set using the `-app-regex-legend` option.

For example, to extract the `#bytes` and `t_max[usec]` values from this output:

```
#bytes #repetitions t_min[usec] t_max[usec] t_avg[usec]
0      1000        0.06      0.06      0.06
1      1000        0.10      0.10      0.10
```

use the following configuration:

```
-app-regex: (\d+)\s+\d+\s+[\d.+-]+\s+([\d.+-]+)
```

-app-regex-legend

Specifies a list of variables extracted from the regular expression. Variables correspond to the regular expression groups. The tuner uses the last variable as the performance indicator of the launch. Use the `-tree-opt` to set the optimization direction of the indicator.

For example:

```
-app-regex-legend: size,time
```

-iter

Sets the number of iterations for each launch with a given set of parameters. Higher numbers of iterations increase accuracy of results.

For example:

```
-iter: 3
```

Search Space Options

Use these options to define a search space, which is a set of combinations of Intel® MPI Library parameters that the target application uses for launches. The library parameters are generally configured using run-time options or environment variables.

NOTE

A search space line can be very long, so line breaking is available for all the search space options. Use a backslash to break a line (see examples below).

-search

Defines the search space by defining variables declared with the `-app` option and by setting environment variables for the application launch.

For example:

```
-search: func=BCAST, \
        np=4,ppn={1,4},{,I_MPI_ADJUST_BCAST=[1,3]},logmax=5
```

The `-app` variables are defined as `<var1>=<value1>[,<var2>=<value2>][, ...]`. The following syntax is available for setting values:

| Syntax | Description | Examples |
|---|--|----------------------------------|
| <code><value></code> | Single value. Can be a number or a string. | 4 |
| <code>{<value1>[,<value2>][, ...]}</code> | List of independent values. | {2,4} |
| <code>[<start>,<end>[,<step>]]</code> | Linear range of values with the default step of 1. | [1,8,2] – expands to {1,2,4,6,8} |
| <code>(<start>,<end>[,<step>])</code> | Exponential range with the default step of 2. | (1,16) – expands to {1,2,4,8,16} |

To set environment variables for the command launch, use the following syntax:

| Syntax | Description | Examples |
|---|--|---|
| <code><variable>=<value></code> | Single variable definition. Any type of the syntax above can be used for the value: single values, lists or ranges. | I_MPI_ADJUST_BCAST=3 I_MPI_ADJUST_BCAST=[1,3] |
| <code>{,<variable>=<value>}</code> | A special case of the syntax above. When set this way, the variable default value is first used in an application launch. | {,I_MPI_ADJUST_BCAST=[1,3]} |
| <code><prefix>{<value1>[,<value2>][, ...]}</code> | Multi-value variable definition. Prefix is a common part for all the values, commonly the variable name. A value can be a singular value or a combination of values in the format: <code><prefix>(<value1>[,<value2>][, ...])</code> . Prefix is optional and a value in the combination is a string, which can utilize the list and range syntax above. | I_MPI_ADJUST_ALLREDUCE={=1, =2,(=9,_KN_RADIX=(2,8))} See below for a more complete example. |

The following example shows a more complex option definition:

```
I_MPI_ADJUST_BCAST{=1,=2,(=9,_KN_RADIX=(2,8)),(={10,11},_SHM_KN_RADIX=[2,8,2])}
```

This directive consecutively runs the target application with the following environment variables set:

```
I_MPI_ADJUST_BCAST=1
I_MPI_ADJUST_BCAST=2
I_MPI_ADJUST_BCAST=9,I_MPI_ADJUST_BCAST_KN_RADIX=2
I_MPI_ADJUST_BCAST=9,I_MPI_ADJUST_BCAST_KN_RADIX=4
I_MPI_ADJUST_BCAST=9,I_MPI_ADJUST_BCAST_KN_RADIX=8
I_MPI_ADJUST_BCAST=10,I_MPI_ADJUST_BCAST_SHM_KN_RADIX=2
I_MPI_ADJUST_BCAST=10,I_MPI_ADJUST_BCAST_SHM_KN_RADIX=4
I_MPI_ADJUST_BCAST=10,I_MPI_ADJUST_BCAST_SHM_KN_RADIX=6
I_MPI_ADJUST_BCAST=10,I_MPI_ADJUST_BCAST_SHM_KN_RADIX=8
```

```
I_MPI_ADJUST_BCAST=11, I_MPI_ADJUST_BCAST_SHM_KN_RADIX=2
I_MPI_ADJUST_BCAST=11, I_MPI_ADJUST_BCAST_SHM_KN_RADIX=4
I_MPI_ADJUST_BCAST=11, I_MPI_ADJUST_BCAST_SHM_KN_RADIX=6
I_MPI_ADJUST_BCAST=11, I_MPI_ADJUST_BCAST_SHM_KN_RADIX=8
```

-search-excl

Excludes certain combinations from the search space. The syntax is identical to that of the `-search` option. For example:

```
-search-excl: I_MPI_ADJUST_BCAST={1,2}
```

or

```
-search-excl: func=BCAST, np=4, ppn=1, I_MPI_ADJUST_BCAST=1
```

-search-only

Defines a subset of the search space to search in. Only this subset is used for application launches. The syntax is identical to the `-search` option.

This option is useful for the second and subsequent tuning sessions on a subset of parameters from the original session, without creating a separate configuration file.

Output Options

Use these options to customize the output. The tuner can produce output of two types:

- `table` – useful for verifying the tuning results, contains values from all the application launches
- `tree` – an internal output format, contains the optimal values

-table

Defines the layout for the resulting output table. The option value is a list of variables declared with the `-app` option, which are joined in colon-separated groups. Each group denotes a specific part of the table.

For example:

```
-table: func:ppn,np:size:*:time
```

The last group variables (`time`) are rendered in table cells. The second last group variables are used for building table columns (`*`, denotes all the variables not present the other variable groups). The third last group variables are used for building table rows (`size`). All other variable groups are used to make up the table label. Groups containing several variables are complex groups and produce output based on all the value combinations.

For example, the option definition above can produce the following output:

```
Label: "func=BCAST, ppn=2, np=2"
```

Legend:

```
set 0: ""
```

```
set 1: "I_MPI_ADJUST_BCAST=1"
```

```
set 2: "I_MPI_ADJUST_BCAST=2"
```

```
set 3: "I_MPI_ADJUST_BCAST=3"
```

Table:

| | set 0 | set 1 | set 2 | set 3 |
|----------|-------------|-------------|-------------|-------------|
| "size=0" | "time=0.10" | "time=0.08" | "time=0.11" | "time=0.10" |
| | "time=0.12" | "time=0.09" | "time=0.12" | "time=0.11" |
| | | "time=0.10" | | |

```

"size=4" | "time=1.12" | "time=1.11" | "time=1.94" | "time=1.72"
         | "time=1.35" | "time=1.18" | "time=1.97" | "time=1.81"
         | "time=1.38" | "time=1.23" | "time=2.11" | "time=1.89"
-----|-----|-----|-----|-----
"size=8" | "time=1.21" | "time=1.10" | "time=1.92" | "time=1.72"
         | "time=1.36" | "time=1.16" | "time=2.01" | "time=1.75"
         | "time=1.37" | "time=1.17" | "time=2.24" | "time=1.87"
-----|-----|-----|-----|-----
...

```

Cells include only unique values from all the launches for the given parameter combination. The number of launches is set with the `-iter` option.

-table-ignore

Specifies the variables to ignore from the `-table` option definition.

-tree

Defines the layout for the resulting tree of optimal values of the parameter that is tuned (for example, collective operation algorithms). The tree is rendered as a JSON structure. The option value is a list of variables declared with the `-app` option, which are joined in colon-separated groups. Each group denotes a specific part of the tree. Groups containing several variables are complex groups and produce output based on all the value combinations.

Example:

```
-tree: func:ppn,np:size*:time
```

The first two groups (`func` and `ppn,np`) make up the first two levels of the tree. The last group variables (`time`) are used as the optimization criteria and are not rendered. The second last group contains variables to be optimized (`*`, denotes all the variables not present the other variable groups). The third last group variables are used to split the search space into intervals based on the optimal values of parameters from the next group (for example, `I_MPI_ADJUST_<operation>` algorithm numbers).

For example, the option definition above can produce the following output:

```

{
  "func=BCAST":
  {
    "ppn=1,np=4":
    {
      "size=0":
        {"I_MPI_ADJUST_BCAST": "3"},
      "size=64":
        {"I_MPI_ADJUST_BCAST": "1"},

      "size=512":
        {"I_MPI_ADJUST_BCAST": "2"},

      ...
    }
  }
}

```

This tree representation is an intermediate format of tuning results and is ultimately converted to a string that the library can understand. The conversion script is specified with `-tree-postprocess` option.

-tree-ignore

Specifies the variables to ignore from the `-tree` option definition.

-tree-intervals

Specifies the maximum number of intervals where the optimal parameter value is applied. If not specified, any number of intervals is allowed.

-tree-tolerance

Specifies the tolerance level. Non-zero tolerance (for example, 0.03 for 3%) joins resulting intervals with the performance indicator value varying by the specified tolerance.

-tree-postprocess

Specifies an executable to convert the resulting JSON tree to a custom format.

-tree-opt

Specifies the optimization direction. The available values are `max` (default) and `min`.

-tree-file

Specifies a log file where the tuning results are saved.

3. Environment Variable Reference

3.1. Compilation Environment Variables

I_MPI_{CC,CXX,FC,F77,F90}_PROFILE

Specify the default profiling library.

Syntax

```
I_MPI_CC_PROFILE=<profile_name>  
I_MPI_CXX_PROFILE=<profile_name>  
I_MPI_FC_PROFILE=<profile_name>  
I_MPI_F77_PROFILE=<profile_name>  
I_MPI_F90_PROFILE=<profile_name>
```

Arguments

| | |
|----------------|--------------------------------------|
| <profile_name> | 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 for `mpiicc` or another Intel® MPI Library compiler wrapper.

I_MPI_{CC,CXX,FC,F77,F90}

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

Syntax

```
I_MPI_CC=<compiler>  
I_MPI_CXX=<compiler>  
I_MPI_FC=<compiler>  
I_MPI_F77=<compiler>  
I_MPI_F90=<compiler>
```

Arguments

| | |
|------------|--|
| <compiler> | 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.

I_MPI_ROOT

Set the Intel® MPI Library installation directory path.

Syntax

I_MPI_ROOT=<path>

Arguments

| | |
|--------|--|
| <path> | 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

| | |
|--------|--|
| <path> | 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

| | |
|--------------------------|---|
| <code><arg></code> | Version of library |
| <code>opt</code> | Multi-threaded optimized library. This is the default value |
| <code>dbg</code> | Multi-threaded debug library |

Description

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

3.2. Hydra Environment Variables

I_MPI_HYDRA_HOST_FILE

Set the host file to run the application.

Syntax

```
I_MPI_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

Set the timeout period for `mpiexec`.

Syntax

`I_MPI_JOB_TIMEOUT=<timeout>`

`I_MPI_MPIEXEC_TIMEOUT=<timeout>`

Arguments

| | |
|------------------------------|---|
| <code><timeout></code> | Define <code>mpiexec</code> timeout period in seconds |
| <code><n> ≥ 0</code> | The value of the timeout period. The default timeout value is zero, which means 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 this environment variable in the shell environment before executing the `mpiexec` command. Setting the variable through the `-genv` and `-env` options has no effect.

I_MPI_HYDRA_BOOTSTRAP

Set the bootstrap server.

Syntax

`I_MPI_HYDRA_BOOTSTRAP=<arg>`

Arguments

| | |
|--------------------------|-------------------------|
| <code><arg></code> | String parameter |
| <code>service</code> | Use hydra service agent |

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` 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.

NOTE**I_MPI_HYDRA_PMI_CONNECT**

Define the processing method for PMI messages.

Syntax

```
I_MPI_HYDRA_PMI_CONNECT=<value>
```

Arguments

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

Description

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

I_MPI_PERHOST

Define the default behavior for the `-perhost` option of the `mpiexec` command.

Syntax

```
I_MPI_PERHOST=<value>
```

Arguments

| | |
|-----------------------------|---|
| <code><value></code> | Define a value used for <code>-perhost</code> 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 behavior for the `-perhost` option. Unless specified explicitly, the `-perhost` option is implied with the value set in `I_MPI_PERHOST`.

NOTE

When running under a job scheduler, this environment variable is ignored by default. To be able to control process placement with `I_MPI_PERHOST`, disable the `I_MPI_JOB_RESPECT_PROCESS_PLACEMENT` variable.

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` 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

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

Description

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

I_MPI_HYDRA_IFACE

Set the network interface.

Syntax

```
I_MPI_HYDRA_IFACE=<arg>
```

Arguments

| | |
|---------------------|---|
| <arg> | String parameter |
| <network interface> | 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_TMPDIR

Specify a temporary directory.

Syntax

```
I_MPI_TMPDIR=<arg>
```

Arguments

| | |
|--------|--|
| <arg> | String parameter |
| <path> | Temporary directory. The default value is /tmp |

Description

Set this environment variable to specify a directory for temporary files.

I_MPI_JOB_RESPECT_PROCESS_PLACEMENT

Specify whether to use the process-per-node placement provided by the job scheduler, or set explicitly.

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 the variable is set, the Hydra process manager uses the process placement provided by job scheduler (default). In this case the `-ppn` option and its equivalents are ignored. If you disable the variable, the Hydra process manager uses the process placement set with `-ppn` or its equivalents.

I_MPI_PORT_RANGE

Specify a range of allowed port numbers.

Syntax

`I_MPI_PORT_RANGE=<range>`

Arguments

| | |
|--------------------------------------|--------------------|
| <code><range></code> | String parameter |
| <code><min>:<max></code> | Allowed port range |

Description

Set this environment variable to specify a range of the allowed port numbers for the Intel® MPI Library.

3.3. I_MPI_ADJUST Family Environment Variables**I_MPI_ADJUST_<opname>**

Control collective operation algorithm selection.

Syntax

`I_MPI_ADJUST_<opname>="<algid>[:<conditions>] [<algid>:<conditions>[...]]"`

Arguments

| | |
|----------------------------|--|
| <code><algid></code> | Algorithm identifier |
| <code>>= 0</code> | The default value of zero selects the optimized default settings |

| | |
|---------------------------------|---|
| <code><conditions></code> | 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 |
| <l>@<p> | Messages of size <l> and number of processes <p> |
| <l>-<m>@<p>-<q> | Messages of size from <l> to <m> and number of processes from <p> to <q>, 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.

Environment Variables, Collective Operations, and Algorithms

| Collective Operation | Environment Variable | Algorithms |
|----------------------|-------------------------|--|
| MPI_Allgather | I_MPI_ADJUST_ALLGATHER | <ol style="list-style-type: none"> 1. Recursive doubling 2. Bruck's 3. Ring 4. Topology aware Gather + Bcast 5. Knomial |
| MPI_Allgatherv | I_MPI_ADJUST_ALLGATHERV | <ol style="list-style-type: none"> 1. Recursive doubling 2. Bruck's 3. Ring 4. Topology aware Gather + Bcast |
| MPI_Allreduce | I_MPI_ADJUST_ALLREDUCE | <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 |

| | | |
|---------------|------------------------|--|
| | | <ol style="list-style-type: none"> 7. Shumilin's ring 8. Ring 9. Knomial 10. Topology aware SHM-based flat 11. Topology aware SHM-based Knomial 12. Topology aware SHM-based Knary |
| MPI_Alltoall | I_MPI_ADJUST_ALLTOALL | <ol style="list-style-type: none"> 1. Bruck's 2. Isend/Irecv + waitall 3. Pair wise exchange 4. Plum's |
| MPI_Alltoallv | I_MPI_ADJUST_ALLTOALLV | <ol style="list-style-type: none"> 1. Isend/Irecv + waitall 2. Plum's |
| MPI_Alltoallw | I_MPI_ADJUST_ALLTOALLW | Isend/Irecv + waitall |
| MPI_Barrier | I_MPI_ADJUST_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 7. Topology aware SHM-based flat 8. Topology aware SHM-based Knomial |

| | | |
|-------------|----------------------|--|
| | | <ol style="list-style-type: none"> 9. Topology aware SHM-based Knary |
| MPI_Bcast | I_MPI_ADJUST_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 9. Topology aware SHM-based flat 10. Topology aware SHM-based Knomial 11. Topology aware SHM-based Knary |
| MPI_Exscan | I_MPI_ADJUST_EXSCAN | <ol style="list-style-type: none"> 1. Partial results gathering 2. Partial results gathering regarding layout of processes |
| MPI_Gather | I_MPI_ADJUST_GATHER | <ol style="list-style-type: none"> 1. Binomial 2. Topology aware binomial 3. Shumilin's 4. Binomial with segmentation |
| MPI_Gatherv | I_MPI_ADJUST_GATHERV | <ol style="list-style-type: none"> 1. Linear 2. Topology aware linear |

| | | |
|--------------------|-----------------------------|---|
| | | <ol style="list-style-type: none"> 3. Knomial |
| MPI_Reduce_scatter | I_MPI_ADJUST_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 |
| MPI_Reduce | I_MPI_ADJUST_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 8. Topology aware SHM-based flat 9. Topology aware SHM-based Knomial 10. Topology aware SHM-based Knary 11. Topology aware SHM-based binomial |
| MPI_Scan | I_MPI_ADJUST_SCAN | <ol style="list-style-type: none"> 1. Partial results gathering 2. Topology aware partial results gathering |
| MPI_Scatter | I_MPI_ADJUST_SCATTER | <ol style="list-style-type: none"> 1. Binomial 2. Topology aware |

Environment Variable Reference

| | | |
|-----------------|--------------------------------------|--|
| | | <ul style="list-style-type: none"> binomial 3. Shumilin's |
| MPI_Scatterv | I_MPI_ADJUST_SCATTERV | <ul style="list-style-type: none"> 1. Linear 2. Topology aware linear |
| MPI_Iallgather | I_MPI_ADJUST_IALLGATHER | <ul style="list-style-type: none"> 1. Recursive doubling 2. Bruck's 3. Ring |
| | I_MPI_ADJUST_IALLGATHER_COMPOSITION | <ul style="list-style-type: none"> 1. Composition alpha 2. Composition beta |
| | I_MPI_ADJUST_IALLGATHER_NETWORK | <ul style="list-style-type: none"> 1. Recursive doubling 2. Bruck's 3. Ring |
| | I_MPI_ADJUST_IALLGATHER_NODE | <ul style="list-style-type: none"> 1. Recursive doubling 2. Bruck's 3. Ring |
| MPI_Iallgatherv | I_MPI_ADJUST_IALLGATHERV | <ul style="list-style-type: none"> 1. Recursive doubling 2. Bruck's 3. Ring |
| | I_MPI_ADJUST_IALLGATHERV_COMPOSITION | <ul style="list-style-type: none"> 1. Composition alpha 2. Composition beta |
| | I_MPI_ADJUST_IALLGATHERV_NETWORK | <ul style="list-style-type: none"> 1. Recursive doubling 2. Bruck's 3. Ring |
| | I_MPI_ADJUST_IALLGATHERV_NODE | <ul style="list-style-type: none"> 1. Recursive doubling |

| | | |
|----------------|-------------------------------------|--|
| | | <ol style="list-style-type: none"> 2. Bruck's 3. Ring |
| MPI_Iallreduce | I_MPI_ADJUST_IALLREDUCE | <ol style="list-style-type: none"> 1. Recursive doubling 2. Rabenseifner's 3. Reduce + Bcast 4. Ring 5. Knomial 6. Binomial 7. Reduce scatter allgather 8. SMP 9. Nreduce |
| | I_MPI_ADJUST_IALLREDUCE_COMPOSITION | <ol style="list-style-type: none"> 1. Composition alpha 2. Composition beta |
| | I_MPI_ADJUST_IALLREDUCE_NETWORK | <ol style="list-style-type: none"> 1. Naïve (Reduce + Bcast) 2. Recursive doubling 3. Reduce scatter allgather 4. Nreduce 5. Rabenseifner's 6. Ring 7. SMP 8. Knomial |
| | I_MPI_ADJUST_IALLREDUCE_NODE | <ol style="list-style-type: none"> 1. Naïve 2. Recursive doubling 3. Reduce scatter allgather 4. Nreduce 5. Rabenseifner's 6. Ring 7. SMP |

Environment Variable Reference

| | | |
|----------------|-------------------------------------|--|
| | | 8. Knomial |
| MPI_Ialltoall | I_MPI_ADJUST_IALLTOALL | <ol style="list-style-type: none"> 1. Bruck's 2. Isend/Irecv + Waitall 3. Pairwise exchange |
| | I_MPI_ADJUST_IALLTOALL_COMPOSITION | <ol style="list-style-type: none"> 1. Composition alpha 2. Composition beta |
| | I_MPI_ADJUST_IALLTOALL_NETWORK | <ol style="list-style-type: none"> 1. Pairwise exchange 2. Bruck's 3. Inplace 4. Isend/Irecv + Waitall |
| | I_MPI_ADJUST_IALLTOALL_NODE | <ol style="list-style-type: none"> 1. Pairwise exchange 2. Bruck's 3. Inplace 4. Isend/Irecv + Waitall |
| MPI_Ialltoallv | I_MPI_ADJUST_IALLTOALLV | Isend/Irecv + Waitall |
| | I_MPI_ADJUST_IALLTOALLV_COMPOSITION | <ol style="list-style-type: none"> 1. Composition alpha 2. Composition beta |
| | I_MPI_ADJUST_IALLTOALLV_NETWORK | <ol style="list-style-type: none"> 1. Isend/Irecv + Waitall 2. Inplace |
| | I_MPI_ADJUST_IALLTOALLV_NODE | <ol style="list-style-type: none"> 1. Isend/Irecv + Waitall 2. Inplace |
| MPI_Ialltoallw | I_MPI_ADJUST_IALLTOALLW | Isend/Irecv + Waitall |
| | I_MPI_ADJUST_IALLTOALLW_COMPOSITION | <ol style="list-style-type: none"> 1. Composition |

| | | |
|--------------|-----------------------------------|--|
| | | alpha 2. Composition beta |
| | I_MPI_ADJUST_IALLTOALLW_NETWORK | 1. Isend/Irecv + Waitall 2. Inplace |
| | I_MPI_ADJUST_IALLTOALLW_NODE | 1. Isend/Irecv + Waitall 2. Inplace |
| MPI_Ibarrier | I_MPI_ADJUST_IBARRIER | Dissemination |
| | I_MPI_ADJUST_IBARRIER_COMPOSITION | 1. Composition alpha 2. Composition beta |
| | I_MPI_ADJUST_IBARRIER_NETWORK | Dissemination |
| | I_MPI_ADJUST_IBARRIER_NODE | Dissemination |
| MPI_Ibcast | I_MPI_ADJUST_IBCAST | 1. Binomial 2. Recursive doubling 3. Ring 4. Knomial 5. SMP 6. Tree knomial 7. Tree kary |
| | I_MPI_ADJUST_IBCAST_COMPOSITION | 1. Composition alpha 2. Composition beta |
| | I_MPI_ADJUST_IBCAST_NETWORK | 1. Binomial 2. Scatter recursive doubling allgather 3. Ring allgather 4. SMP 5. Tree knomial 6. Tree kary |

Environment Variable Reference

| | | |
|--------------|-----------------------------------|--|
| | | 7. Knomial |
| | I_MPI_ADJUST_IBCAST_NODE | <ol style="list-style-type: none"> 1. Binomial 2. Scatter recursive doubling allgather 3. Ring allgather 4. SMP 5. Tree knomial 6. Tree kary 7. Knomial |
| MPI_Iexscan | I_MPI_ADJUST_IEXSCAN | Recursive doubling |
| | I_MPI_ADJUST_IEXSCAN_COMPOSITION | <ol style="list-style-type: none"> 1. Composition alpha 2. Composition beta |
| | I_MPI_ADJUST_IEXSCAN_NETWORK | Recursive doubling |
| | I_MPI_ADJUST_IEXSCAN_NODE | Recursive doubling |
| MPI_Igather | I_MPI_ADJUST_IGATHER | <ol style="list-style-type: none"> 1. Binomial 2. Knomial |
| | I_MPI_ADJUST_IGATHER_COMPOSITION | <ol style="list-style-type: none"> 1. Composition alpha 2. Composition beta |
| | I_MPI_ADJUST_IGATHER_NETWORK | <ol style="list-style-type: none"> 1. Binomial 2. Knomial |
| | I_MPI_ADJUST_IGATHER_NODE | <ol style="list-style-type: none"> 1. Binomial 2. Knomial |
| MPI_Igatherv | I_MPI_ADJUST_IGATHERV | <ol style="list-style-type: none"> 1. Linear 2. Linear ssend |
| | I_MPI_ADJUST_IGATHERV_COMPOSITION | <ol style="list-style-type: none"> 1. Composition alpha 2. Composition beta |

| | | |
|---------------------------|--|---|
| | I_MPI_ADJUST_IGATHERV_NETWORK | <ol style="list-style-type: none"> 1. Linear 2. Linear ssend |
| | I_MPI_ADJUST_IGATHERV_NODE | <ol style="list-style-type: none"> 1. Linear 2. Linear ssend |
| MPI_Ireduce_scatter | I_MPI_ADJUST_IREDUCE_SCATTER | <ol style="list-style-type: none"> 1. Recursive halving 2. Pairwise 3. Recursive doubling |
| | I_MPI_ADJUST_IREDUCE_SCATTER_COMPOSITION | <ol style="list-style-type: none"> 1. Composition alpha 2. Composition beta |
| | I_MPI_ADJUST_IREDUCE_SCATTER_NETWORK | <ol style="list-style-type: none"> 1. Noncommutative 2. Pairwise 3. Recursive doubling 4. Recursive halving |
| | I_MPI_ADJUST_IREDUCE_SCATTER_NODE | <ol style="list-style-type: none"> 1. Noncommutative 2. Pairwise 3. Recursive doubling 4. Recursive halving |
| MPI_Ireduce_scatter_block | I_MPI_ADJUST_IREDUCE_SCATTER_BLOCK | <ol style="list-style-type: none"> 1. Recursive halving 2. Pairwise 3. Recursive doubling |
| | I_MPI_ADJUST_IREDUCE_SCATTER_BLOCK_COMPOSITION | <ol style="list-style-type: none"> 1. Composition alpha 2. Composition beta |
| | I_MPI_ADJUST_IREDUCE_SCATTER_BLOCK_NETWORK | <ol style="list-style-type: none"> 1. Noncommutative 2. Pairwise 3. Recursive doubling 4. Recursive halving |

Environment Variable Reference

| | | |
|-------------|---|---|
| | I_MPI_ADJUST_IREDUCE_SCATTER_BLOCK_NODE | <ol style="list-style-type: none"> 1. Noncommutative 2. Pairwise 3. Recursive doubling 4. Recursive halving |
| MPI_Ireduce | I_MPI_ADJUST_IREDUCE | <ol style="list-style-type: none"> 1. Rabenseifner's 2. Binomial 3. Knomial 4. SMP |
| | I_MPI_ADJUST_IREDUCE_COMPOSITION | <ol style="list-style-type: none"> 1. Composition alpha 2. Composition beta |
| | I_MPI_ADJUST_IREDUCE_NETWORK | <ol style="list-style-type: none"> 1. Binomial 2. Rabenseifner's 3. SMP 4. Knomial |
| | I_MPI_ADJUST_IREDUCE_NODE | <ol style="list-style-type: none"> 1. Binomial 2. Rabenseifner's 3. SMP 4. Knomial |
| MPI_Iscan | I_MPI_ADJUST_ISCAN | <ol style="list-style-type: none"> 1. Recursive Doubling 2. SMP |
| | I_MPI_ADJUST_ISCAN_COMPOSITION | <ol style="list-style-type: none"> 1. Composition alpha 2. Composition beta |
| | I_MPI_ADJUST_ISCAN_NETWORK | <ol style="list-style-type: none"> 1. Recursive Doubling 2. SMP |
| | I_MPI_ADJUST_ISCAN_NODE | <ol style="list-style-type: none"> 1. Recursive Doubling 2. SMP |

| | | |
|----------------|------------------------------------|---|
| MPI_Isscatter | I_MPI_ADJUST_ISCATTER | 1. Binomial 2. Knomial |
| | I_MPI_ADJUST_ISCATTER_COMPOSITION | 1. Composition alpha 2. Composition beta |
| | I_MPI_ADJUST_ISCATTER_NETWORK | 1. Binomial 2. Knomial |
| | I_MPI_ADJUST_ISCATTER_NODE | 1. Binomial 2. Knomial |
| MPI_Isscatterv | I_MPI_ADJUST_ISCATTERV | Linear |
| | I_MPI_ADJUST_ISCATTERV_COMPOSITION | 1. Composition alpha 2. Composition beta |
| | I_MPI_ADJUST_ISCATTERV_NETWORK | Linear |
| | I_MPI_ADJUST_ISCATTERV_NODE | 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.

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 <code>MPI_Allgather</code> 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

| | |
|----------------------------|--|
| <code><radix></code> | An integer that specifies a radix used by the Knomial <code>MPI_Bcast</code> algorithm to build a knomial communication tree |
| <code>> 1</code> | 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

| | |
|----------------------------|--|
| <code><radix></code> | An integer that specifies a radix used by the Knomial <code>MPI_Allreduce</code> algorithm to build a knomial communication tree |
| <code>> 1</code> | 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

| | |
|----------------------------|---|
| <code><radix></code> | An integer that specifies a radix used by the Knomial <code>MPI_Reduce</code> algorithm to build a knomial communication tree |
| <code>> 1</code> | 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

| | |
|----------------------------|--|
| <code><radix></code> | An integer that specifies a radix used by the Knomial <code>MPI_Ireduce</code> algorithm to build a knomial communication tree |
| <code>> 1</code> | 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

| | |
|----------------------------|--|
| <code><radix></code> | An integer that specifies a radix used by the Knomial <code>MPI_Igather</code> algorithm to build a knomial communication tree |
| <code>> 1</code> | 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

| | |
|----------------------------|---|
| <code><radix></code> | An integer that specifies a radix used by the Knomial <code>MPI_Iscatter</code> algorithm to build a knomial communication tree |
| <code>> 1</code> | 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.

`I_MPI_ADJUST_<COLLECTIVE>_SHM_KN_RADIX`**Syntax**

```
I_MPI_ADJUST_<COLLECTIVE>_SHM_KN_RADIX=<radix>
```

Arguments

| | |
|---------|---|
| <radix> | An integer that specifies a radix used by the Knomial or Knary SHM-based algorithm to build a knomial or knary communication tree |
| > 0 | <ul style="list-style-type: none"> • If you specify the environment variables <code>I_MPI_ADJUST_BCAST_SHM_KN_RADIX</code> and <code>I_MPI_ADJUST_BARRIER_SHM_KN_RADIX</code>, the default value is 3 • If you specify the environment variables <code>I_MPI_ADJUST_REDUCE_SHM_KN_RADIX</code> and <code>I_MPI_ADJUST_ALLREDUCE_SHM_KN_RADIX</code>, the default value is 4 |

Description

This environment variable includes the following variables:

- `I_MPI_ADJUST_BCAST_SHM_KN_RADIX`
- `I_MPI_ADJUST_BARRIER_SHM_KN_RADIX`
- `I_MPI_ADJUST_REDUCE_SHM_KN_RADIX`
- `I_MPI_ADJUST_ALLREDUCE_SHM_KN_RADIX`

Set this environment variable to select the knomial or knary tree radix for the corresponding tree SHM-based algorithms. When you build a knomial communication tree, the specified value is used as the power for 2 to generate resulting radix ($2^{\langle\text{radix}\rangle}$). When you build a knary communication tree, the specified value is used for the radix.

I_MPI_COLL_INTRANODE

Syntax

`I_MPI_COLL_INTRANODE=<mode>`

Arguments

| | |
|--------|--|
| <mode> | Intranode collectives type |
| pt2pt | Use only point-to-point communication-based collectives |
| shm | Enables shared memory collectives. This is the default value |

Description

Set this environment variable to switch intranode communication type for collective operations. If there is large set of communicators, you can switch off the SHM-collectives to avoid memory overconsumption.

I_MPI_COLL_INTRANODE_SHM_THRESHOLD

Syntax

`I_MPI_COLL_INTRANODE_SHM_THRESHOLD=<nbytes>`

Arguments

| | |
|----------|--|
| <nbytes> | Define the maximal data block size processed by shared memory collectives. |
| > 0 | Use the specified size. The default value is 16384 bytes. |

Description

Set this environment variable to define the size of shared memory area available for each rank for data placement. Messages greater than this value will *not* be processed by SHM-based collective operation, but will be processed by point-to-point based collective operation. The value must be a multiple of 4096.

I_MPI_ADJUST_GATHER_SEGMENT

Syntax

I_MPI_ADJUST_GATHER_SEGMENT=<block_size>

Arguments

| | |
|--------------|---|
| <block_size> | Size of a message segment in bytes. |
| > 0 | Use the specified size. The default value is 16384 bytes. |

Description

Set an internal block size to control the `MPI_Gather` message segmentation for the binomial algorithm with segmentation.

3.4. 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.

3.4.1. 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 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 where there are 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.

Logical Enumeration

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 0 | 4 | 1 | 5 | 2 | 6 | 3 | 7 |
|---|---|---|---|---|---|---|---|

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 |

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.

3.4.2. Default Settings

If you do not specify values for any process pinning environment variables, the default settings below are used. For 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`

3.4.3. Environment Variables for Process Pinning

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_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.

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

| | |
|----------------------------------|---|
| <code><proclist></code> | A comma-separated list of logical processor numbers and/or ranges of processors. The process with the <i>i</i> -th rank is pinned to the <i>i</i> -th processor in the list. The number should not exceed the amount of processors on a node. |
| <code><l></code> | Processor with logical number <code><l></code> . |
| <code><l>-<m></code> | Range of processors with logical numbers from <code><l></code> to <code><m></code> . |

| | |
|---|--|
| <code><k>, <l>-<m></code> | Processors <code><k></code> , as well as <code><l></code> through <code><m></code> . |
|---|--|

Syntax

```
I_MPI_PIN_PROCESSOR_LIST=[<procset>] [: [grain=<grain>] [, shift=<shift>] [, preoffset=<preoffset>] [, postoffset=<postoffset>]
```

Arguments

| | |
|------------------------------|--|
| <code><procset></code> | 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. |

| | |
|---------------------------------|---|
| <code><grain></code> | Specify the pinning granularity cell for a defined <code><procset></code> . The minimal <code><grain></code> value is a single element of the <code><procset></code> . The maximal <code><grain></code> value is the number of <code><procset></code> elements in a socket. The <code><grain></code> value must be a multiple of the <code><procset></code> value. Otherwise, the minimal <code><grain></code> value is assumed. The default value is the minimal <code><grain></code> 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 <i><grain></i> units. The <i><postoffset></i> 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:

| | |
|------------------|--|
| <i><n></i> | Specify an explicit value of the corresponding parameters. <i><n></i> is non-negative integer. |
| fine | Specify the minimal value of the corresponding parameter. |
| core | Specify the parameter value equal to the amount of the corresponding parameter units contained in one core. |
| cache1 | Specify the parameter value equal to the amount of the corresponding parameter units that share an L1 cache. |
| cache2 | Specify the parameter value equal to the amount of the corresponding parameter units that share an L2 cache. |
| cache3 | Specify the parameter value equal to the amount of the corresponding parameter units that share an L3 cache. |
| cache | The largest value among <i>cache1</i> , <i>cache2</i> , and <i>cache3</i> . |
| socket sock | Specify the parameter value equal to the amount of the corresponding parameter units contained in one physical package/socket. |
| half mid | Specify the parameter value equal to <i>socket/2</i> . |
| third | Specify the parameter value equal to <i>socket/3</i> . |
| quarter | Specify the parameter value equal to <i>socket/4</i> . |
| octavo | Specify the parameter value equal to <i>socket/8</i> . |

Syntax

```
I_MPI_PIN_PROCESSOR_LIST=[<procset>] [:map=<map>]
```

Arguments

| | |
|--------------------|--|
| <i><map></i> | The mapping pattern used for process placement. |
| bunch | 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 |
|-----|-----|-----|-----|-----|-------|-----|-----------|-----------|-----------|

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:

```
> mpiexec -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:

```
> mpiexec -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 the process pinning, use the following command:

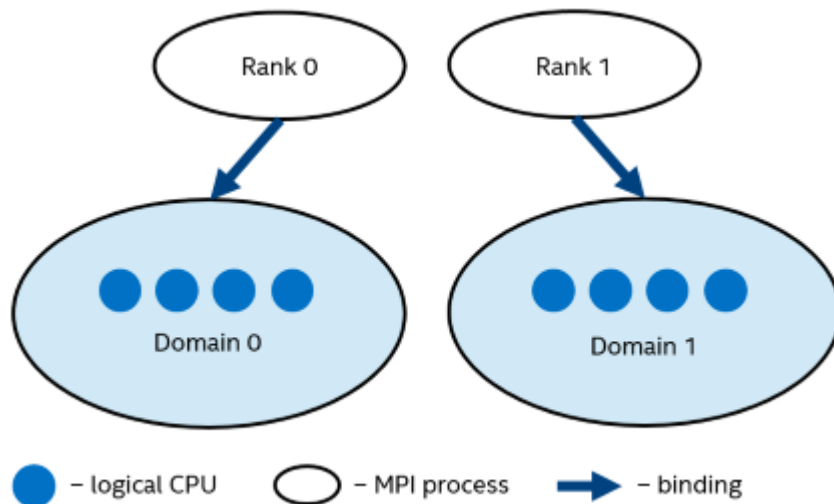
```
> mpiexec -genv I_MPI_DEBUG=4 -m -host host1 -env
I_MPI_PIN_PROCESSOR_LIST=0,3 -n <# of processes> <executable> :^
-host host2 -env I_MPI_PIN_PROCESSOR_LIST=1,2,3 -n <# of
processes> <executable>
```

3.4.4. Interoperability with OpenMP* API

I_MPI_PIN_DOMAIN

Intel® MPI Library provides an additional environment variable to control process pinning for hybrid MPI/OpenMP* 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 below.

Figure 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>`

| | |
|-------------------------------|--|
| <code><mc-shape></code> | 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 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>numa</code> | Each domain consists of the logical processors that share a particular NUMA node. The number of domains on a machine is equal to the number of NUMA nodes on the machine. |
| <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. |

NOTE

If `Cluster on Die` is disabled on a machine, the number of NUMA nodes equals to the number of sockets. In this case, pinning for `I_MPI_PIN_DOMAIN = numa` is equivalent to pinning for `I_MPI_PIN_DOMAIN = socket`.

Explicit Shape

`I_MPI_PIN_DOMAIN=<size>[:<layout>]`

| | |
|---------------------------|--|
| <code><size></code> | 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. |
| auto | The domain size is defined by the formula $size = \#cpu / \#proc$, where $\#cpu$ is the number of logical processors on a node, and $\#proc$ is the number of the MPI processes started on a node |
| <n> | The domain size is defined by a positive decimal number <n> |

| | |
|----------|---|
| <layout> | Ordering of domain members. The default value is compact |
| platform | Domain members are ordered according to their BIOS numbering (platform-depended numbering) |
| compact | 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 |
| scatter | 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>

| | |
|-----------------------|---|
| <masklist> | Define domains through the comma separated list of hexadecimal numbers (domain masks) |
| [m_1, \dots, m_n] | <p>For <masklist>, each m_i is a hexadecimal bit mask defining an individual domain. The following rule is used: the i^{th} logical processor is included into the domain if the corresponding m_i value is set to 1. All remaining processors are put into a separate domain. BIOS numbering is used.</p> <p>NOTE</p> <p>To ensure that your configuration in <masklist> is parsed correctly, use square brackets to enclose the domains specified by the <masklist>. For example: I_MPI_PIN_DOMAIN=[55, aa]</p> |

NOTE

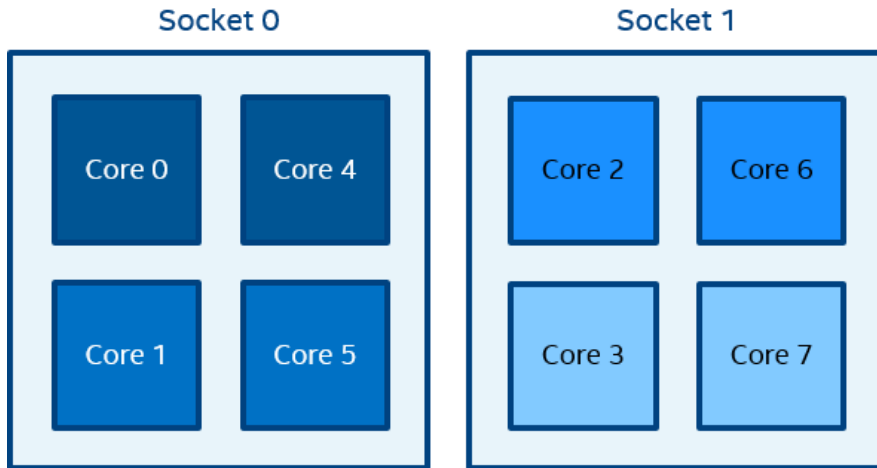
These options are available for both Intel® and non-Intel microprocessors, but they may perform additional optimizations for Intel microprocessors than they perform for non-Intel microprocessors.

NOTE

To pin OpenMP* processes or threads inside the domain, the corresponding OpenMP feature (for example, the `KMP_AFFINITY` environment variable for Intel® compilers) should be used.

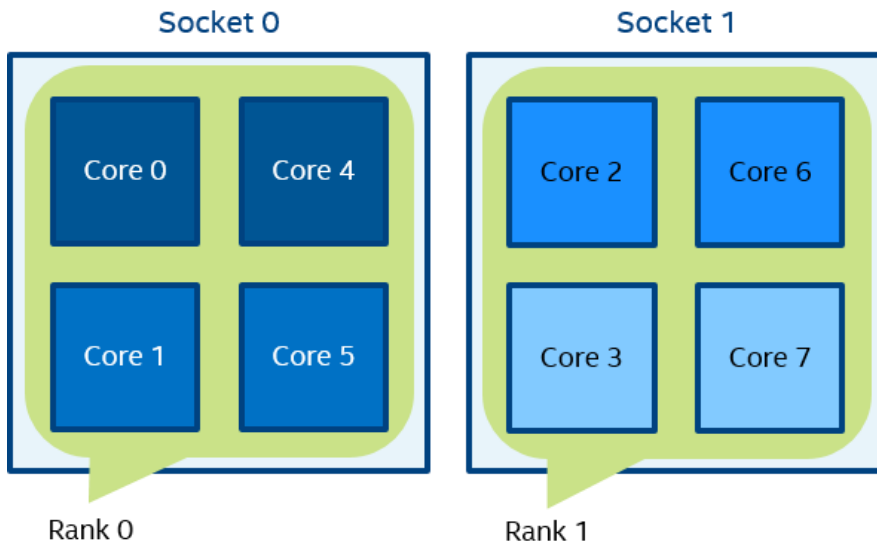
See the following model of a symmetric multiprocessing (SMP) node in the examples:

Figure 2 Model of a Node



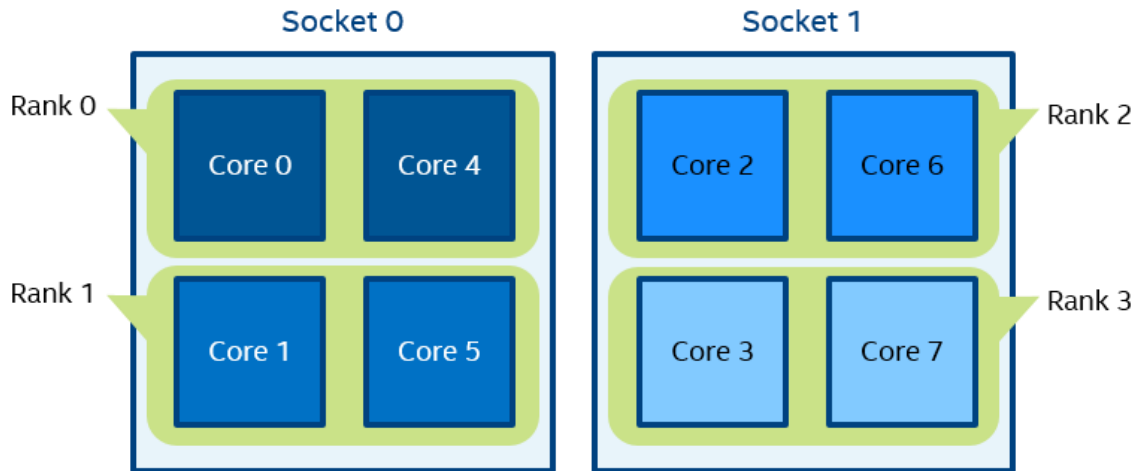
The figure above represents the SMP node model with a total of 8 cores on 2 sockets. Intel® Hyper-Threading Technology is disabled. Core pairs of the same color share the L2 cache.

Figure 3 `mpiexec -n 2 -env I_MPI_PIN_DOMAIN socket test.exe`



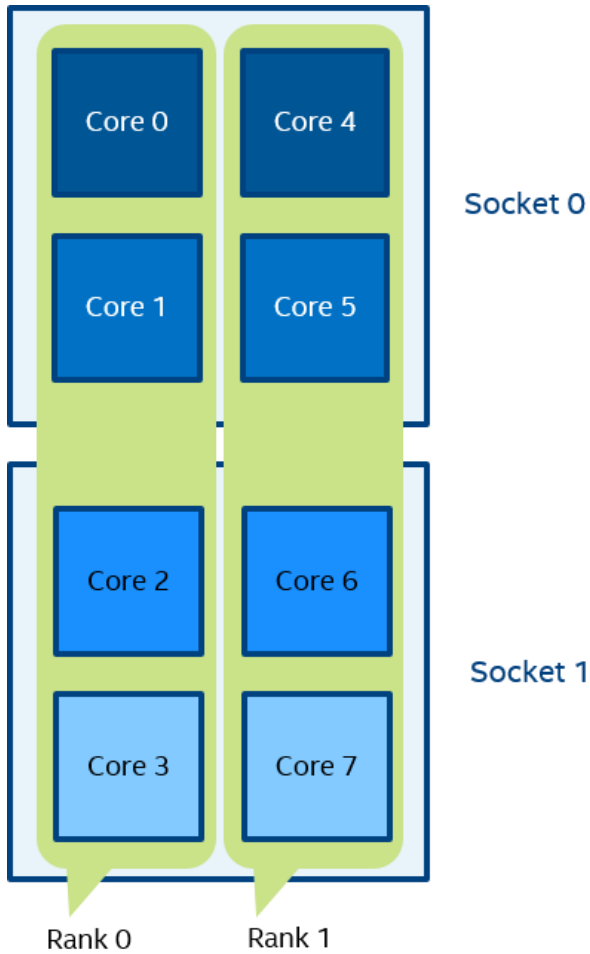
In Figure 3, two domains are defined according to the number of sockets. Process rank 0 can migrate on all cores on the 0-th socket. Process rank 1 can migrate on all cores on the first socket.

Figure 4 `mpiexec -n 4 -env I_MPI_PIN_DOMAIN cache2 test.exe`



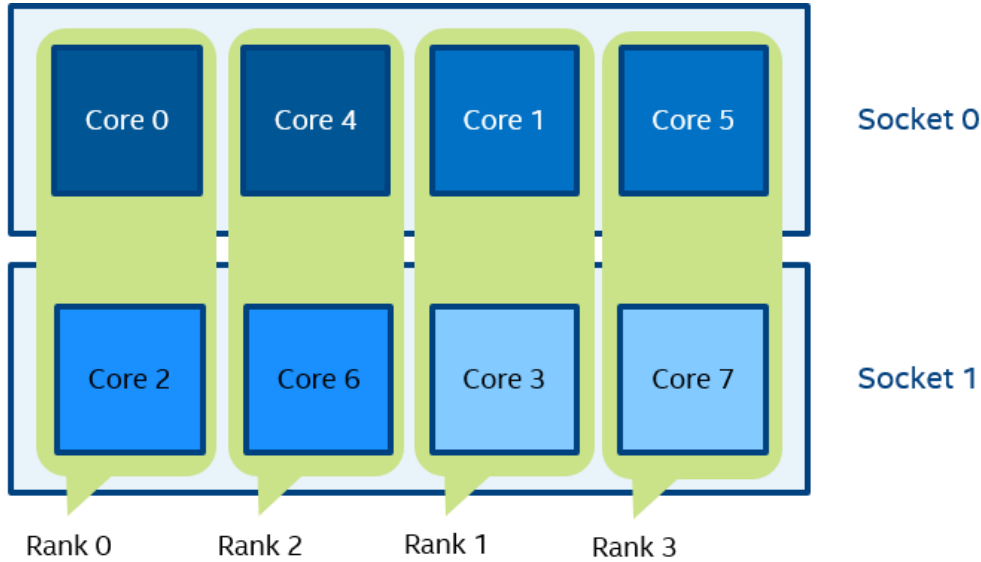
In Figure 4, four domains are defined according to the amount of common L2 caches. Process rank 0 runs on cores {0,4} that share an L2 cache. Process rank 1 runs on cores {1,5} that share an L2 cache as well, and so on.

Figure 5 `mpiexec -n 2 -env I_MPI_PIN_DOMAIN 4:platform test.exe`



In Figure 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.

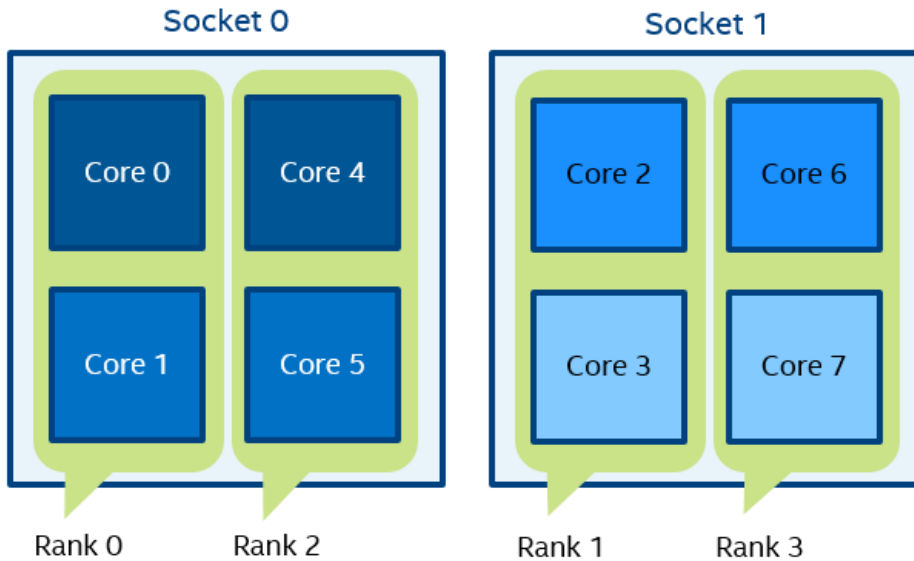
Figure 6 `mpirun -n 4 -env I_MPI_PIN_DOMAIN auto:scatter test.exe`



In Figure 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.

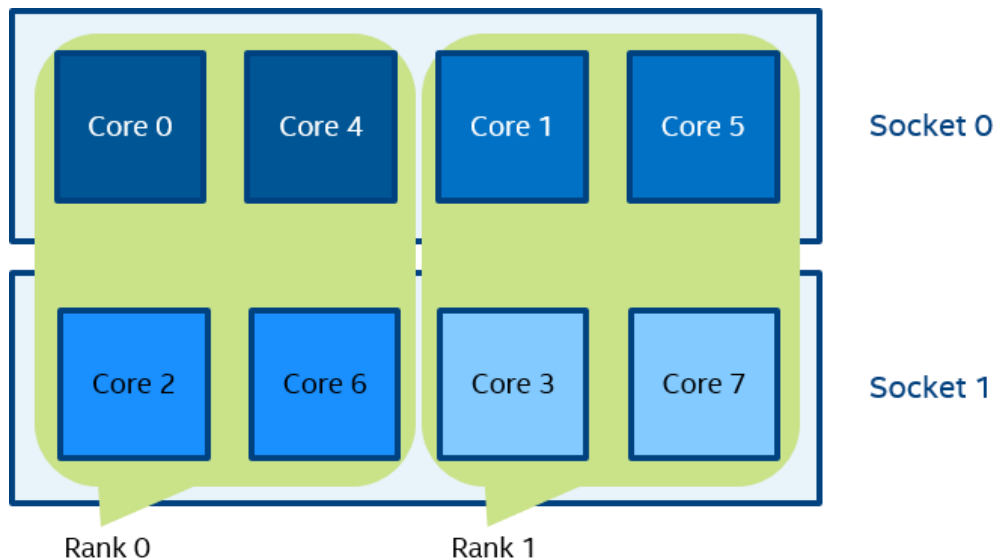
Figure 7 `set OMP_NUM_THREADS=2`

```
mpirun -n 4 -env I_MPI_PIN_DOMAIN omp:platform test.exe
```



In Figure 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 8 `mpiexec -n 2 -env I_MPI_PIN_DOMAIN [55,aa] test.exe`



In Figure 8 (the example for `I_MPI_PIN_DOMAIN=<masklist>`), the first domain is defined by the 55 mask. It contains all cores with even numbers {0,2,4,6}. The second domain is defined by the AA 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.

Examples

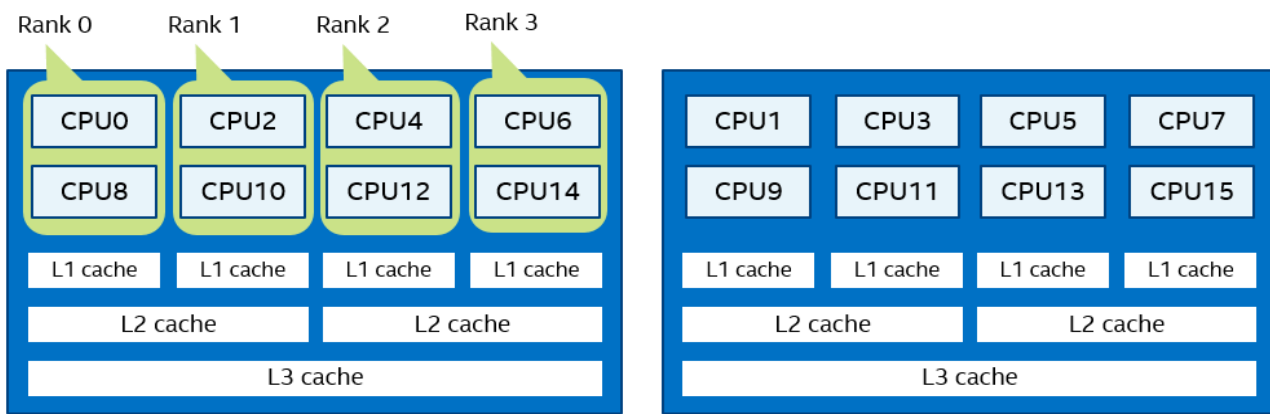
For the following configuration:

- 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 settings below.

Compact order:

```
I_MPI_PIN_DOMAIN=2
I_MPI_PIN_ORDER=compact
```

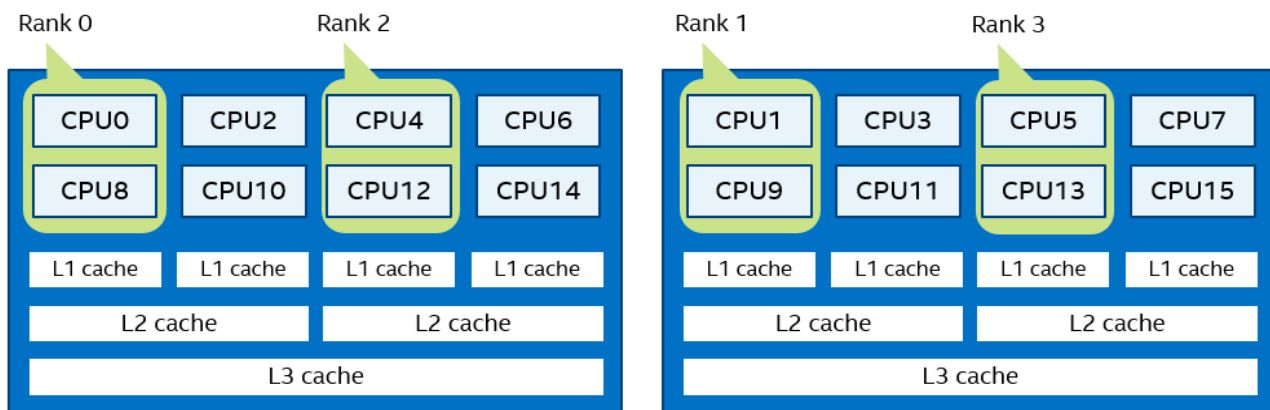
Figure 9 Compact Order Example



Scatter order:

```
I_MPI_PIN_DOMAIN=2
I_MPI_PIN_ORDER=scatter
```

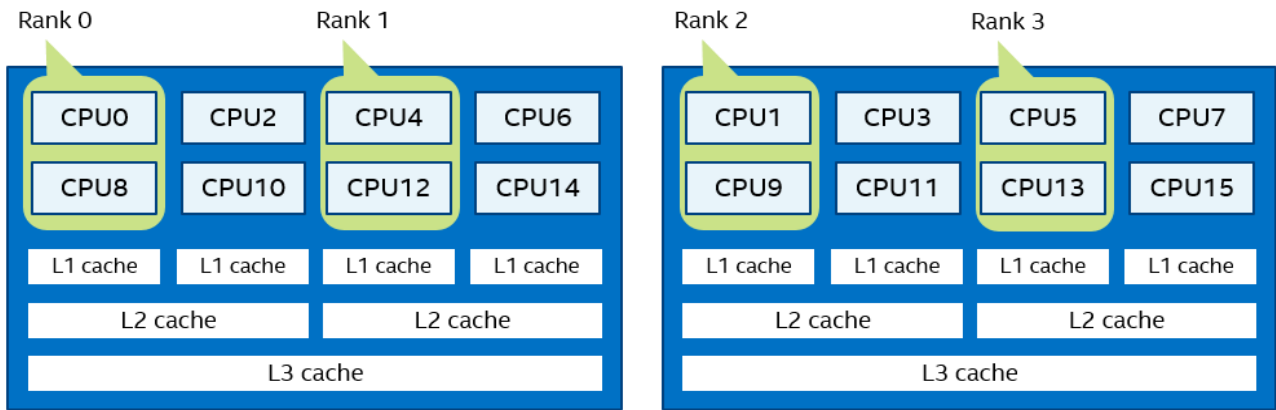
Figure 10 Scatter Order Example



Spread order:

```
I_MPI_PIN_DOMAIN=2
I_MPI_PIN_ORDER=spread
```

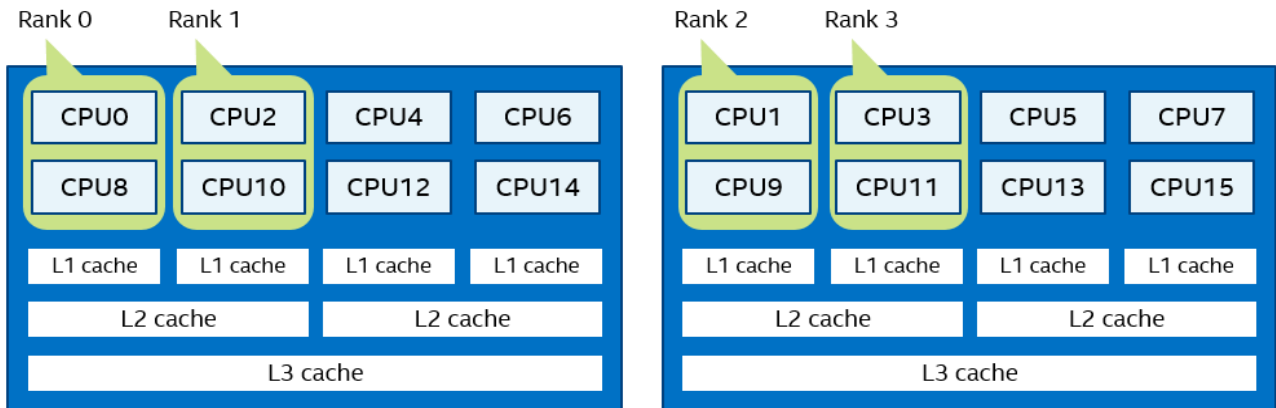
Figure 11 Spread Order Example



Bunch order:

```
I_MPI_PIN_DOMAIN=2
I_MPI_PIN_ORDER=bunch
```

Figure 12 Bunch Order Example



3.5. Environment Variables for Fabrics Control

3.5.1. Communication Fabrics Control

I_MPI_FABRICS

Select the particular fabrics to be used.

Syntax

```
I_MPI_FABRICS=ofi | shm:ofi
```

Arguments

| | |
|----------|--|
| <fabric> | Define a network fabric. |
| ofi | OpenFabrics Interfaces* (OFI)-capable network fabrics, |

| | |
|--|---|
| | such as Intel® True Scale Fabric, Intel® Omni-Path Architecture, InfiniBand*, and Ethernet (through OFI API). |
|--|---|

Description

Set this environment variable to select a specific fabric combination.

NOTE

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

3.5.2. OFI*-capable Network Fabrics Control

I_MPI_OFI_DRECV

Control the capability of the direct receive in the OFI fabric.

Syntax

`I_MPI_OFI_DRECV=<arg>`

Arguments

| | |
|-------------------------------------|--|
| <code><arg></code> | Binary indicator |
| <code>enable yes on 1</code> | Enable direct receive. This is the default value |
| <code>disable no off 0</code> | Disable 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 OFI as the network fabric by setting `I_MPI_FABRICS=ofi`.

I_MPI_OFI_LIBRARY_INTERNAL

Control the usage of `libfabric*` shipped with the Intel® MPI Library.

Syntax

`I_MPI_OFI_LIBRARY_INTERNAL=<arg>`

Arguments

| | |
|-------------------------------------|--|
| <code><arg></code> | Binary indicator |
| <code>enable yes on 1</code> | Use <code>libfabric</code> from the Intel MPI Library |
| <code>disable no off 0</code> | Do not use <code>libfabric</code> from the Intel MPI Library |

Description

Set this environment variable to disable or enable usage of `libfabric` from the Intel MPI Library. The variable must be set before sourcing the `mpivars.bat` script.

Example

```
> set I_MPI_OFI_LIBRARY_INTERNAL=1
> call <installdir>\intel64\bin\mpivars.bat
```

Setting this variable is equivalent to passing the `-ofi_internal` option to the `mpivars.bat` script.

For more information, refer to the Intel® MPI Library Developer Guide, section *Running Applications > Libfabric* Support*.

3.6. Other 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 I_MPI_FABRICS was used and which Intel® MPI Library configuration was used. |
| 3 | Output effective MPI rank, pid and node mapping table. |
| 4 | Output process pinning information. |
| 5 | Output environment variables specific to Intel MPI Library. |
| 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 print debugging information about the application.

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>` is the MPI process rank, by default. If you add the '+' sign in front of the `<level>` number, the `<identifier>` assumes the following format: `rank#pid@hostname`. Here, `rank` is the MPI process rank, `pid` is the process ID, and `hostname` is the host name. If you add the '-' sign, `<identifier>` is not printed at all.
- `<message>` contains the debugging output.

The following examples demonstrate possible command lines with the corresponding output:

```
> mpiexec -n 1 -env I_MPI_DEBUG=2 test.exe
...
[0] MPI startup(): shared memory data transfer mode
```

The following commands are equal and produce the same output:

```
> mpiexec -n 1 -env I_MPI_DEBUG=+2 test.exe
> mpiexec -n 1 -env I_MPI_DEBUG=2,pid,host test.exe
...
[0#1986@mpiclust001] MPI startup(): shared memory data transfer mode
```

NOTE

Compiling with the `/zi`, `/ZI` or `/Z7` option 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 <code>stdout</code> . This is the default value. |
| stderr | Output to <code>stderr</code> . |
| <file_name> | Specify the output file name for debug information (the maximum file name length is 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, process ID or host name is added to the file name accordingly.

I_MPI_NETMASK

Choose the network interface for MPI communication over sockets.

Syntax

I_MPI_NETMASK=<arg>

Arguments

| | |
|---------------------------|--|
| <arg> | Define the network interface (string parameter) |
| <interface_mnemonic> | Mnemonic of the network interface: <code>ib</code> or <code>eth</code> |
| ib | Select IPoIB* |
| eth | Select Ethernet. This is the default value |
| <network_address> | Network address. The trailing zero bits imply netmask |
| <network_address/netmask> | Network address. The <netmask> value specifies the netmask length |
| <list of interfaces> | A colon separated list of network addresses or interface mnemonics |

Description

Set this environment variable to choose the network interface for MPI communication over sockets in the `sock` and `ssm` communication modes. If you specify a list of interfaces, the first available interface on the node will be used for communication.

Examples

1. Use the following setting to select the IP over InfiniBand* (IPoIB) fabric:

```
I_MPI_NETMASK=ib
```

```
I_MPI_NETMASK=eth
```

2. Use the following setting to select a particular network for socket communications. This setting implies the 255.255.0.0 netmask:

```
I_MPI_NETMASK=192.169.0.0
```

3. Use the following setting to select a particular network for socket communications with netmask set explicitly:

```
I_MPI_NETMASK=192.169.0.0/24
```

4. Use the following setting to select the specified network interfaces for socket communications:

```
I_MPI_NETMASK=192.169.0.5/24:ib0:192.169.0.0
```

NOTE

If the library cannot find any suitable interface by the given value of `I_MPI_NETMASK`, the value will be used as a substring to search in the network adapter's description field. And if the substring is found in the description, this network interface will be used for socket communications. For example, if `I_MPI_NETMASK=myri` and the description field contains something like Myri-10G adapter, this interface will be chosen.

I_MPI_REMOVED_VAR_WARNING

Print out a warning if a removed environment variable is set.

Syntax

```
I_MPI_REMOVED_VAR_WARNING=<arg>
```

Arguments

| | |
|------------------------|--|
| <arg> | Binary indicator |
| enable yes on 1 | Print out the warning. This is the default value |
| disable no off 0 | Do not print the warning |

Description

Use this environment variable to print out a warning if a removed environment variable is set. Warnings are printed regardless of whether `I_MPI_DEBUG` is set.

I_MPI_LIBRARY_KIND

Specify the Intel® MPI Library configuration.

Syntax

```
I_MPI_LIBRARY_KIND=<value>
```

Arguments

| | |
|---------|------------------|
| <value> | Binary indicator |
|---------|------------------|

| | |
|---------|---|
| release | Multi-threaded optimized library. This is the default value |
| debug | Multi-threaded debug library |

Description

Use this variable to set an argument for the `mpivars.[c]sh` script. This script establishes the Intel® MPI Library environment and enables you to specify the appropriate library configuration. To ensure that the desired configuration is set, check the `LD_LIBRARY_PATH` variable.

Example

```
> export I_MPI_LIBRARY_KIND=debug
```

Setting this variable is equivalent to passing an argument directly to the `mpivars.[c]sh` script:

Example

```
> <installdir>\intel64\bin\mpivars.bat release
```

4. Miscellaneous

4.1. User Authorization

Intel® MPI Library supports several authentication methods under Windows* OS:

I_MPI_AUTH_METHOD

Select a user authorization method.

Syntax

I_MPI_AUTH_METHOD=<method>

Arguments

| | |
|-------------|--|
| <method> | Define the authorization method |
| password | Use the password-based authorization. This is the default value. |
| delegate | Use the domain-based authorization with delegation ability. |
| impersonate | Use the limited domain-based authorization. You will not be able to open files on remote machines or access mapped network drives. |

Description

Set this environment variable to select a desired authorization method. If this environment variable is not defined, `mpiexec` uses the password-based authorization method by default. Alternatively, you can change the default behavior by using the `-delegate` or `-impersonate` options.

For more details, see the *Developer Guide*, section *Installation and Prerequisites > User Authorization*.