

# *Intel® MPI Library Reference Manual*

## *Disclaimer and Legal Information*

---

Information in this document is provided in connection with Intel products. No license, express or implied, by estoppel or otherwise, to any intellectual property rights is granted by this document. EXCEPT AS PROVIDED IN INTEL'S TERMS AND CONDITIONS OF SALE FOR SUCH PRODUCTS, INTEL ASSUMES NO LIABILITY WHATSOEVER, AND INTEL DISCLAIMS ANY EXPRESS OR IMPLIED WARRANTY, RELATING TO SALE AND/OR USE OF INTEL PRODUCTS INCLUDING LIABILITY OR WARRANTIES RELATING TO FITNESS FOR A PARTICULAR PURPOSE, MERCHANTABILITY, OR INFRINGEMENT OF ANY PATENT, COPYRIGHT OR OTHER INTELLECTUAL PROPERTY RIGHT. Intel products are not intended for use in medical, life saving, life sustaining, critical control or safety systems, or in nuclear facility application. Intel may make changes to specifications and product descriptions at any time, without notice.

Designers must not rely on the absence or characteristics of any features or instructions marked "reserved" or "undefined." Intel reserves these for future definition and shall have no responsibility whatsoever for conflicts or incompatibilities arising from future changes to them.

The software described in this Intel® MPI Library Reference Manual may contain software defects which may cause the product to deviate from published specifications. Current characterized software defects are available on request.

This Intel® MPI Library Reference Manual, as well as the software described in it is furnished under license and may only be used or copied in accordance with the terms of the license. The information in this manual is furnished for informational use only, is subject to change without notice, and should not be construed as a commitment by Intel Corporation. Intel Corporation assumes no responsibility or liability for any errors or inaccuracies that may appear in this document or any software that may be provided in association with this document.

Except as permitted by such license, no part of this document may be reproduced, stored in a retrieval system, or transmitted in any form or by any means without the express written consent of Intel Corporation.

Developers must not rely on the absence or characteristics of any features or instructions marked "reserved" or "undefined." Improper use of reserved or undefined features or instructions may cause unpredictable behavior or failure in developer's software code when running on an Intel processor. Intel reserves these features or instructions for future definition and shall have no responsibility whatsoever for conflicts or incompatibilities arising from their unauthorized use.

Intel, the Intel logo, Intel Inside, the Intel Inside logo, Pentium, Itanium, Intel Xeon, Celeron, Intel SpeedStep, Intel Centrino, Intel NetBurst, Intel NetStructure, VTune, MMX, the MMX logo, Dialogic, i386, i486, iCOMP, Intel386, Intel486, Intel740, IntelDX2, IntelDX4 and IntelSX2 are trademarks or registered trademarks of Intel Corporation or its subsidiaries in the United States and other countries.

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

Copyright © 2004-2006 Intel Corporation

## **MPI Legal Notices**

Intel® MPI Library is based in part on the MPICH2\* implementation of MPI from Argonne National Laboratory\* (ANL).

Intel® MPI Library is also based in part on InfiniBand Architecture\* RDMA drivers from MVAPICH2\* from Ohio State University's Network-Based Computing Laboratory.

# Contents

---

<b>Intel® MPI Library Reference Manual</b>	<b>1</b>
<b>Disclaimer and Legal Information</b>	<b>2</b>
MPI Legal Notices .....	3
<b>Contents</b>	<b>4</b>
<b>Overview</b>	<b>5</b>
<b>Command Reference</b>	<b>6</b>
Compiler Commands .....	6
Compiler Command Options .....	7
Configuration Files.....	8
Environment Variables .....	8
Job Startup Commands .....	8
Global Options.....	9
Local Options.....	11
Configuration Files.....	12
Environment Variables .....	12
MPD Daemon Commands .....	15
Syntax.....	16
Arguments .....	16
Description.....	16
Configuration Files.....	19
Environment Variables .....	20
Simplified Job Startup Command.....	21
<b>Tuning Reference</b>	<b>23</b>
Process Pinning .....	23
Device Control.....	24
RDMA and RDSSM Device Control .....	25
Collective Operation Control .....	31
Miscellaneous .....	35

# Overview

---

The Intel® MPI Library is a multi-fabric message passing library that implements the Message Passing Interface, v2 (MPI-2) specification. It enables you to switch interconnection fabrics without re-linking.

The library is included in the following kits:

- *Intel MPI Library Runtime Environment* has the tools you need to run programs including MPD daemons and supporting utilities, shared (.so) libraries, Release Notes, a Getting Started Guide, and a Reference Manual.
- *Intel MPI Library Development Kit* includes all of the Runtime Environment components plus compilation tools including compiler commands such as `mpicc`, include files and modules, static (.a) libraries, debug libraries, trace libraries, and test codes.

The goal of this *Reference Manual* is to provide you with a complete command and tuning reference for the Intel MPI Library.

# Command Reference

## Compiler Commands

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

Compiler Command	Underlying Compiler	Supported Language(s)	Supported ABI(s)
<b>GNU* compilers</b>			
<code>mpicc</code>	<code>gcc, cc</code>	C	32/64 bit
<code>mpicxx</code>	<code>g++ v3.x</code>	C/C++	32/64 bit
<code>mpicxx2</code>	<code>g++ v2.x</code>	C/C++	32/64 bit
<code>mpif77</code>	<code>g77</code>	F77	32/64 bit
<code>mpif90</code>	<code>gfortran</code>	F95	32/64 bit
<b>Intel® compilers version 8.0, 8.1, 9.0 or 9.1</b>			
<code>mpiicc</code>	<code>icc</code>	C	32/64 bit
<code>mpiicpc</code>	<code>icpc</code>	C++	32/64 bit
<code>mpiifort</code>	<code>ifort</code>	F77/F95	32/64 bit
<b>Intel® compilers version 7.1</b>			
<code>mpiicc7</code>	<code>icc</code>	C	32 bit
<code>mpiicpc7</code>	<code>icpc</code>	C++	32 bit
<code>mpiifc</code>	<code>ifc</code>	F77/F90	32 bit
<code>mpiecc</code>	<code>ecc</code>	C	64 bit
<code>mpiecpc</code>	<code>ecpc</code>	C++	64 bit
<code>mpiefc</code>	<code>efc</code>	F77/F90	64 bit

## NOTES

- Compiler commands are only available in the Intel® MPI Library Development Kit.
- Compiler commands are in the `<installdir>/bin` directory. For Intel® EM64T, 64-bit-enabled compiler commands are in the `<installdir>/bin64` directory and 32-bit compiler commands are in the `<installdir>/bin` directory.
- Ensure that the corresponding underlying compilers (32-bit or 64-bit, as appropriate) are already in your `PATH`.
- To port existing, MPI-enabled applications to Intel MPI Library, recompile all sources.
- To compile and link without using the `mpicc` and related commands, run the appropriate command with the `-show` option added. The output will indicate the correct flags, options, includes, defines, and libraries to add to the compile and link lines. For example, use the following command to show the required compile flags, options, and then include paths for compiling source files:

```
$ mpicc -show -c test.c
```

- Use the following command to show the required link flags, options, and libraries for linking object files:

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

## Compiler Command Options

### -show

Use this option to display the compilation and linkage commands without actually running them. This is useful for debugging, for submitting support issues, or for determining compile and link options for complex builds.

### -echo

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

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

Use this option to set the path/name of the underlying compiler to be used.

### -g or -debug

Use the `-g` or `-debug` options to compile program in debug mode, and link the resulting executable against the debugging versions of the libraries. See also `I_MPI_DEBUG`, in Section [Environment variables](#), for information on how to use additional debug features with `-g` builds.

### -O

Use this option to enable optimization. If `-g` is used, `-O` is not implied. Specify `-O` explicitly if you want to enable optimization.

### -t or -trace

Use the `-t` or `-trace` option to link the resulting executable against the Intel® Trace Collector. Use the `-t=log` or `-trace=log` options to link the resulting executable against the logging versions of Intel MPI libraries and the Intel Trace Collector.

Include the installation path of the Intel Trace Collector into the `VT_ROOT` environment variable to use this option.

### -static\_mpi

Use this option to link the main `libmpi` library statically. This option does not affect the default linkage method for other libraries.

### -dynamic\_log

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

Include the `$VT_ROOT/slib` element into the `LD_LIBRARY_PATH` environment variable to run the resulting programs.

### -nocompchk

Use this option to disable compiler setup checks and to speed up compilation in some cases. By default, each compiler command performs checks to ensure that the appropriate underlying compiler is set up correctly.

## **-gcc-version**

Use this option for compiler drivers mpicxx and mpiicpc to link an application for running in a particular GNU\* C++ environment.

Use `-gcc-version=3` to build an application compatible to GNU\* C++ version up to 3.3.

Use `-gcc-version=4` to build an application compatible to GNU\* C++ version 3.4 or higher.

A library compatible with the detected version of the GNU\*C++ compiler is used by default.

## **Configuration Files**

You can create compiler configuration files using the following file naming convention:

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

where:

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

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

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

Source this file, if it exists, prior to compiling or linking to enable changes to the environment on a per-compiler-command basis.

## **Environment Variables**

### **MPICH\_{CC,CXX,F77,F90}=<compiler>**

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

### **CFLAGS=<flags>**

Add additional `CFLAGS` to be used in compile and/or link steps.

### **LDFLAGS=<flags>**

Set additional `LDFLAGS` to be used in the link step.

### **VT\_ROOT=<path>**

Set Intel Trace Collector installation directory path.

### **IDB\_HOME=<path>**

Set Intel® Debugger installation directory path.

## **Job Startup Commands**

### **mpiexec**

#### **Syntax**

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

or

```
mpiexec <g-options> <l-options> <executable> : \  
<l-options> <executable> ...
```

or

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



## Arguments

<code>&lt;g-options&gt;</code>	Global options that apply to all MPI processes
<code>&lt;l-options&gt;</code>	Local options that apply to a single arg-set
<code>&lt;executable&gt;</code>	<code>./a.out</code> , or <code>path/name</code> of executable, compiled with <code>mpicc</code> or related command
<code>&lt;file&gt;</code>	File with command-line options (see below)

## Description

In the first form, run the specified `<executable>` with the specified options. All the global and/or local options apply to all MPI processes. A single arg-set is assumed.

In the second form, divide the command line into multiple arg-sets, separated by colon characters. All the global options apply to all MPI processes, but the various local options and the `<executable>` that is executed can be specified separately for each arg-set.

In the third form, read the command line from the specified `<file>`. For a command with a single arg-set, the entire command should be specified on a single line in `<file>`. For a command with multiple arg-sets, each arg-set should be specified on a single, separate line in `<file>`. Global options should always appear at the beginning of the first line in `<file>`.

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

If `"."` is not in the `PATH` on all nodes in the cluster, specify the `<executable>` as `./a.out` rather than `a.out`.

## Global Options

### -version or -V

Use this option to output Intel MPI Library version information.

### -nolocal

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

### -perhost <# of processes>

Use this option to place the indicated number of consecutive MPI processes on every host.

The `mpiexec` command controls how the ranks of the processes are allocated to the nodes in the cluster. By default, `mpiexec` uses round-robin assignment of ranks to nodes. This placement algorithm may not be the best choice for your application, particularly for clusters with SMP nodes.

To change this default behavior, set the number of processes per host using the `-perhost` option, and set the total number of processes by using the `-n` option (see [Local Options](#)). Then the first `<# of processes>` indicated by the `-perhost` option will be run on the first host, the next `<# of processes>` on the next host, and so on.

This is shorthand for using the multiple arg-sets that run the same number of processes on each indicated host. The `-perhost` option does not make sense for the second form of the `mpiexec` command.

**-machinefile <machine file>**

Use this option to place the ranks of processes in compliance with machine file. The *<machine file>* has a list of host names one per line. Use a short host name or a fully qualified domain name. Repeat the same host name in compliance with quantity of starting processes on it. You can use the following format not to repeat the same host name: *<host name>:<number of processes>*. Set a # symbol at the beginning of line to comment it.

**-genv <ENVVAR> <value>**

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

**-genvnone**

Use this option to not propagate any environment variables to any MPI processes. The default is to propagate the entire environment from which `mpiexec` was called.

**-g<l-option>**

Use this option to apply the named local option *<l-option>* globally. See also Section [Local Options](#) for local options.

**-tv**

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

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

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

**-idb**

Use this option to run the *<executable>* under the Intel Debugger. For example:

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

Include the installation path of the Intel Debugger into the `IDB_HOME` environment variable to use this option.

**-gdb**

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

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

**-gdba**

Use this option to attach GNU\* debugger to existing job. For example:

```
$ mpiexec -gdba jobid
```

Obtain `jobid` by issuing `mpdlistjobs` command.

**-l**

Use this option to line labels by MPI rank

**-s <spec>**

Use this option to direct `stdin` to specified ranks. Use "all" as a *<spec>* value to specify all processes or `1,3,5` to specify exact list of processes or `2-4,6` to specify range of processes. The default value is `0`.

### **-ifhn <interface-hostname>**

Use this option to specify network interface which will be used for data flow. The **<interface-hostname>** should be IP address or hostname associated with alternate network interface

### **-m**

Use this option to merge output lines (default with gdb).

### **-a <alias name>**

Use this option to assign this alias **<alias name>** to the job.

### **-ecfn <filename>**

Use this option to output xml exit codes to file **<filename>**.

## **Local Options**

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

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

### **-env <ENVVAR> <value>**

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

### **-host <nodename>**

Use this option to specify the particular **<nodename>** on which the MPI processes for the current arg-set are to be run.

### **-path <directory>**

Use this option to specify the path to find the **<executable>** that is to be executed for the current arg-set.

### **-wdir <directory>**

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

### **-umask <umask>**

Use this option to perform umask for remote process.

### **-envall**

Use this option to pass all environment variables in current environment.

### **-envnone**

Use this option not to pass environment variables.

### **-envlist <list of env var names>**

Use this option to pass a list of environment variables with current values.

**-configfile <filename>**

Get command line options from the file `filename`.

**Configuration Files**

Create `mpiexec` configuration files using the following file naming convention:

```
<installdir>/etc/mpiexec.conf
```

```
$HOME/.mpiexec.conf
```

```
$PWD/mpiexec.conf
```

**Syntax**

The format of the `mpiexec.conf` files is free-format text containing default `mpiexec` command-line options. Blank lines and lines that start with a '#' character in the very first column of the line are ignored.

**Description**

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

1. System-wide `<installdir>/etc/mpiexec.conf` (if any)
2. User-specific `$HOME/.mpiexec.conf` (if any)
3. Session-specific `$PWD/mpiexec.conf` (if any)

This applies to all forms of the `mpiexec` command.

Use the `mpiexec.conf` files to specify the default options you will apply to all `mpiexec` commands. For example, to specify a default device, add the following to the respective `mpiexec.conf` file:

```
-genv I_MPI_DEVICE <device>
```

**Environment Variables****MPIEXEC\_TIMEOUT**

Set the `mpiexec` timeout.

**Syntax**

```
MPIEXEC_TIMEOUT=<timeout>
```

**Arguments**

<code>&lt;timeout&gt;</code>	Defines <code>mpiexec</code> timeout period in seconds
<code>&gt; 0</code>	There is no default timeout value

**Description**

Set this variable to make `mpiexec` terminate the job `<timeout>` seconds after its launch.

**I\_MPI\_DEVICE**

Select the particular network fabric to be used.

**Syntax**

```
I_MPI_DEVICE=<device>[:<provider>]
```

**Arguments**

<code>&lt;device&gt;</code>	One of { <code>sock</code> , <code>shm</code> , <code>ssm</code> }
<code>sock</code>	TCP/Ethernet/sockets
<code>shm</code>	shared-memory only (no sockets)
<code>ssm</code>	Combined TCP + shared memory (for clusters with SMP nodes)

<code>&lt;device&gt;</code>	One of { <code>rdma</code> , <code>rdssm</code> }
<code>&lt;provider&gt;</code>	Optional DAPL* provider name
<code>rdma</code>	RDMA-capable network fabrics including InfiniBand*, Myrinet* (via DAPL*)
<code>rdssm</code>	Combined TCP + shared memory + DAPL* (for clusters with SMP nodes and RDMA-capable network fabrics)

**Description**

Set this variable to select specific fabric combination. If the `I_MPI_DEVICE` variable is not defined, Intel MPI Library selects the most appropriate fabric combination automatically.

For example, to select the shared-memory as fabric, use the following command:

```
$ mpiexec -n <#ranks> -env I_MPI_DEVICE shm <executable>
```

Use the `<provider>` specification only for the {`rdma`, `rdssm`} devices. For these devices, if `<provider>` is not specified, the first DAPL\* provider in `/etc/dat.conf` is used. If the `<provider>` is set to `none`, the `rdssm` device establishes sockets connections between the nodes without trying to establish DAPL\* connections first.

**NOTES**

- If you build the MPI program using `mpicc -g`, the debug-enabled version of the library will be used
- If you build the MPI program using `mpicc -t=log`, the trace-enabled version of the library will be used.
- The debug-enabled and trace-enabled versions of the library are only available when you use the Intel® MPI Library Development Kit.

**I\_MPI\_FALLBACK\_DEVICE**

Control fallback upon the available fabric. It is valid only for `rdssm` and `rdma` modes.

**Syntax**

```
I_MPI_FALLBACK_DEVICE=<arg>
```

**Arguments**

<code>&lt;arg&gt;</code>	Binary indicator
<code>enable, yes, on, 1</code>	Fall back upon the <code>ssm</code> fabric if initialization of DAPL* fabric fails. This is the default value.
<code>disable, no, off, 0</code>	Terminate the job if the fabric selected by the <code>I_MPI_DEVICE</code> environment variable cannot be initialized.

## Description

Set this variable to control fallback upon the available fabric.

If the `I_MPI_FALLBACK_DEVICE` is set to `enable` and an attempt to initialize specified fabric fails, the library falls back upon the shared memory and/or socket fabrics. The exact combination depends on number of processes started per node. This device ensures that the job will run but it may not provide the highest possible performance for the given cluster configuration.

If the `I_MPI_FALLBACK_DEVICE` is set to `disable` and an attempt to initialize specified fabric fails, the library terminates the MPI job.

## I\_MPI\_DEBUG

Print out debugging information when an MPI program starts running.

### Syntax

`I_MPI_DEBUG=<level>`

### Arguments

<code>&lt;level&gt;</code>	Indicates level of debug information provided
<code>(unset)</code>	Print no debugging information
<code>1</code>	Print warnings if the specified <code>I_MPI_DEVICE</code> could not be used for some reason
<code>2</code>	Use to positively confirm which <code>I_MPI_DEVICE</code> was used
<code>&gt; 2</code>	Add extra levels of debug information

### Description

Set this variable to control output of the debugging information.

The `I_MPI_DEBUG` mechanism augments the `MPICH_DBG_OUTPUT` debug mechanism from `MPICH2*`. `I_MPI_DEBUG` overrides and implies `MPICH_DBG_OUTPUT=stdout`.

Compiling with `mpicc -g` causes considerable amounts of additional debug information to be printed.

In order to simplify process identification add the '+' or '-' sign in front of the numerical value for `I_MPI_DEBUG`. This setting makes debug output lines prepended with MPI process rank, UNIX process pid and host name as defined at the process launch time. For example:

```
$ mpiexec -n <# of processes> -env I_MPI_DEBUG +2 ./a.out
I_MPI: [rank#pid@hostname]Debug message
```

## TOTALVIEW

Select the particular TotalView\* executable file to use.

### Syntax

`TOTALVIEW=<path>`

### Arguments

<code>&lt;path&gt;</code>	Path/name of the TotalView* executable file instead of the default <code>totalview</code>
---------------------------	---

### Description

Set this variable to or select a particular TotalView\* executable file.

## MPD Daemon Commands

### mpdboot

#### Syntax

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

or

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

#### Arguments

-h, --help	Display help message
-d, --debug	Print debug information
-v, --verbose	Print extra verbose information. Show the rshcmd attempts
-n <#nodes> --totalnum=<#nodes>	Number of nodes in mpd.hosts on which daemons start
-r <rshcmd> --rsh=<rshcmd>	Specify remote shell to start daemons and jobs. The rsh is default value
-f <hostsfile> --file=<hostsfile>	Path/name of file that has the list of machine names on which daemons start.
-l	Remove a restriction of starting only one mpd per machine
-m <mpdcmd> --mpd=<mpdcms>	Specify the full path name of mpd on the remote hosts
-s, --shell	Specify shell
-u <user> --user=<user>	Specify user
--locons	Do not create local MPD consoles
--remcons	Do not create remote MPD consoles
--ncpus=<ncpus>	Indicate how many processors to use on the local machine (other nodes are listed in the hosts file)
-o, --ordered	Start all the mpd daemons exactly in the order as in the mpd.hosts file

#### Description

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

The **mpd** daemons are started using the **rsh** command by default. If the **rsh** connectivity is not enabled, use the **-r ssh** option to switch over to the **ssh**. Make sure that all nodes of the cluster

can connect to each other via `rsh` command without password or, if the `-r ssh` option is used, via `ssh` command without password.

## NOTES

- The `mpdboot` command will spawn a MPD daemon on the host machine, even if the machine name is not listed in `mpd.hosts` file.

## mpd

### Syntax

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

### Arguments

[ --help ]	Display help message
[ -h <host> -p <portnum> ] [ --host=<host> --port=<portnum> ]	Specify host and port to be used for entering an existing ring. <code>--host</code> and <code>-port</code> must be specified together
[ -n ] [ --noconsole ]	Do not create console at startup.
[ -t ] [ --trace ]	Print a lot of trace information
[ -e ] [ --echo ]	Print port number at startup which other mpd may connect
[ -d ] [ --daemon ]	Start mpd in daemon mode
[ --ifhn=<interface/hostname> ]	Specify an interface/hostname for the host
[ -l <listenport> ] [ --listenport=<listenport> ]	Specify a port for this mpd to listen on.

### Description

MPD is a process management system for starting parallel jobs. Before running a job start `mpd` daemons on each host and connect them into a ring. Long parameter names may be abbreviated to their first letters by using only one hyphen and no equal sign:

`mpd -h masterhost -p 4268 -n` is equivalent to `mpd --host=masterhost --port=4268 --noconsole`

A file named `.mpd.conf` file must be present in the user's home directory with read and write access only for the user, and must contain at least a line with `secretword=<secretword>`. Install `mpd.conf` in the `/etc` directory to run `mpd` as root.



## mpdtrace

Determine whether `mpd` is running.

### Syntax

```
mpdtrace [ -l ]
```

### Arguments

<code>-l</code>	Show MPD identifiers instead of the hostnames
-----------------	---

### Description

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

## mpdlistjobs

List running processes of jobs.

### Syntax

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

or

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

### Arguments

<code>-u &lt;username&gt;</code> <code>--user=&lt;username&gt;</code>	List jobs of particular user
<code>-a &lt;jobalias&gt;</code> <code>--alias=&lt;jobalias&gt;</code>	List information about particular job specified by jobalias
<code>-j &lt;jobid&gt;</code> <code>--jobid=&lt;jobid&gt;</code>	List information about particular job specified by jobid

### Description

Use this command to list running processes of jobs. All jobs are displayed by default.

## mpdkilljobs

Kill the job.

### Syntax

```
mpdkilljobs [ jobnum ] [ -a <jobalias> ]
```

### Arguments

<code>jobnum</code>	Kill job specified by <code>jobnum</code>
<code>-a &lt;jobalias&gt;</code>	Kill job specified by <code>&lt;jobalias&gt;</code>

### Description

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

## mpdringtest

Determine how much time is required for a ring loading.

### Syntax

```
mpdringtest [ number of loops ]
```

### Arguments

<code>number of loops</code>	Number of loops
------------------------------	-----------------

### Description

Use this command to test how much it takes for a message to circle the ring.

## mpdexit

Shut down a single `mpd`.

### Syntax

```
mpdexit mpdid
```

### Arguments

<code>mpdid</code>	Specify mpd daemon to kill
--------------------	----------------------------

### Description

Use this command to cause a single mpd to exit. Use mpdid obtained via mpdtrace -l command.

## mpdallexit

Shut down all `mpd` daemons on all nodes.

### Arguments

No arguments

### Description

Use this command to shutdown all `mpd` rings.

## mpdcleanup

### Syntax

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

or

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

### Arguments

<code>-f &lt;hostsfile&gt;</code> <code>--file=&lt;hostsfile&gt;</code>	Specify the file of machines to cleanup
<code>-r &lt;rshcmd&gt;</code> <code>--rsh=&lt;rshcmd&gt;</code>	Specify remote shell to use
<code>-u &lt;user&gt;</code> <code>--user=&lt;user&gt;</code>	Specify user

<code>-c &lt;cleancmd&gt;</code>	Specify command to use for removing UNIX* socket
<code>--clean=&lt;cleancmd&gt;</code>	

**Description**

Use this command to remove the UNIX\* socket on local and remote machines.

**mpdsigjob**

Deliver a signal to the application process of a job.

**Syntax**

```
mpdsigjob sigtype [-j jobid | -a jobalias ] [-s | -g ]
```

**Arguments**

<code>sygtype</code>	Specify signal to send
<code>-a jobalias</code>	Send a signal to job specified by <code>jobalias</code>
<code>-j jobid</code>	Send signal to job specified by <code>jobid</code>
<code>-s</code>	Delivery signal to the single user process
<code>-g</code>	Delivery signal to the group of processes. It is default behavior.

**Description**

Use this command to deliver a specific signal to the application processes of a job. Specified signal is the first argument. Specify only one of `-j` or `-a` options.

**mpdhelp****Syntax**

```
mpdhelp
```

**Arguments**

No arguments

**Description**

Use this command to get short help about mpd commands.

**Configuration Files****\$HOME/.mpd.conf**

This file has the `mpd` daemon password. Use it to control access to the daemons by various Intel MPI Library users.

**Syntax**

The file has a single line:

```
secretword=<mpd password>
```

or

```
MPD_SECRETWORD=<mpd password>
```

**Description**

An arbitrary `<mpd password>` string only controls access to the `mpd` daemons by various cluster users. Do not use any Linux\* login password here.

Place the `$HOME/.mpd.conf` file on a network-mounted file system, or replicate this file so that it is accessible as `$HOME/.mpd.conf` on all nodes in the cluster.

When `mpdboot` is executed by some non-root `<user>`, this file should have owner set to `<user>`, group set to `<<user>'s group>`, and mode set to `600` (user read and write privileges only).

## NOTES

- The `MPD_SECRETWORD` is a synonym for `secretword`.

## mpd.hosts

This file has the list of node machine names which the `mpdboot` command uses.

Ensure that this file only needs to be accessible by the user who runs `mpdboot` on the node/machine where the `mpdboot` command is actually invoked.

### Syntax

The format of the `mpd.hosts` file is a list of machine names, one name per line. Blank lines, and lines that start with a '#' character in the very first column of the line, are ignored.

## Environment Variables

### PATH

Make the `PATH` settings required for `mpdboot` and other `mpd` daemon commands.

## NOTES

- The `<installdir>/bin` directory (`<installdir>/bin64` directory for Intel® EM64T 64-bit mode) and the path to Python\* version 2.2 or higher should be in the `PATH` in order for `mpd` daemon commands to succeed.

## MPD\_CON\_EXT

Set unique name of the `mpd` console file.

### Syntax

`MPD_CON_EXT=<tag>`

### Arguments

<code>&lt;tag&gt;</code>	Unique MPD identifier
--------------------------	-----------------------

### Description

Set this variable to different unique values to allow several `mpd` rings to co-exist. Once this variable is set you can start one `mpd` ring and work with it without affect to other available `mpd` rings. Each `mpd` ring is associated with one `MPD_CON_EXT` value. Set the appropriate `MPD_CON_EXT` value to work with particular `mpd` ring.

Normally, every new `mpd` ring totally replaces the older one.

See section [Simplified Job Startup Command](#) to learn about an easier way to run several Intel MPI Library jobs at once.

## I\_MPI\_MPD\_CONF

Set the path/name of the `mpd` configuration file.

## Syntax

`I_MPI_MPD_CONF=<path/name>`

## Arguments

<code>&lt;path/name&gt;</code>	Absolute path of the MPD configuration file
--------------------------------	---

## Description

Set this variable to define the absolute path of the file that will be used by the `mpdboot` script instead of the default value `${HOME}/.mpd.conf`.

## I\_MPI\_MPD\_CONNECTION\_TIMEOUT

Set the `mpd` connection timeout.

## Syntax

`I_MPI_MPD_CONNECTION_TIMEOUT=<timeout>`

## Arguments

<code>&lt;timeout&gt;</code>	Defines MPD connection timeout period in seconds
<code>&gt; 0</code>	The default <i>timeout</i> value is equal to 20 seconds

## Description

Set this variable to make `mpd` terminate the job if another `mpd` cannot be connected to in at most `<timeout>` seconds.

## Simplified Job Startup Command

### `mpirun`

## Syntax

`mpirun [ <mpdboot options> ] <mpiexec options>`

## Arguments

<code>&lt;mpdboot options&gt;</code>	<code>mpdboot</code> options as described in the <code>mpdboot</code> section above, except <code>-n</code>
<code>&lt;mpiexec options&gt;</code>	<code>mpiexec</code> options as described in the <code>mpiexec</code> section above

## Description

Use this command to start an independent ring of `mpd` daemons, launch an MPI job, and shut down the `mpd` ring upon the job termination.

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

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

The set of hosts is defined by the following rules checked in order:

1. All host names from the `mpdboot` host file (either `mpd.hosts` or the file specified by the `-f` option).

2. All host names returned by the `mpdtrace` command, in case there is an `mpd` ring running.
3. Local host (a warning is issued in this case).

The `mpirun` command also detects if the MPI job is submitted in a session allocated using a job scheduler like Torque\*, PBS Pro\*, LSF\* or Fujitsu\* NQS\*. In this case, the `mpirun` command extracts the host list from the respective environment and uses these nodes fully automatically according to the above scheme.

In other words, if you work under one of the aforementioned job schedulers, you don't have to create the `mpd.hosts` file yourself. Just allocate the session you need using the particular job scheduler installed on your system, and use the `mpirun` command inside this session to run your MPI job.

See the product *Release Notes* for a complete list of the supported job schedulers.

# Tuning Reference

---

The Intel MPI Library provides many environment variables that can be used to influence program behavior and performance at run time. These variables are described below.

## Process Pinning

### I\_MPI\_PIN\_MODE

### I\_MPI\_PIN\_PROCS

Pin processes to the CPUs to prevent undesired process migration.

#### Syntax

`I_MPI_PIN_MODE=<pinmode>`

`I_MPI_PIN_PROCS=<proclist>`

#### Arguments

<code>&lt;pinmode&gt;</code>	Selects CPU pinning mode
<code>mpd</code>	Pin processes inside MPD
<code>lib</code>	Pin processes inside MPI library
<code>enable, yes, on, 1</code>	Pin processes inside MPD
<code>disable, no, off, 0</code>	Do not pin processes

<code>&lt;proclist&gt;</code>	Defines process to CPU map
<code>all</code>	Use all CPUs in order
<code>allcores</code>	Use all CPU cores in order. This is default value
<code>n</code>	Use only CPU number n (0,1, ... , total number of CPUs - 1)
<code>m-n</code>	Use CPUs from m to n
<code>k,l-m,n</code>	Use CPUs k, l thru m, and n

#### Description

Set these variables to enable and control process pinning.

Set the variable `I_MPI_PIN_MODE` to `lib` to make the Intel MPI Library pin the processes. Set the `I_MPI_PIN_PROCS` variable to define the set of processors.

Set the variable `I_MPI_PIN_MODE` to `mpd` to make `mpd` daemon pin processes via system specific means if they are available. Set the `I_MPI_PIN_PROCS` variable to define the set of processors. It may allow memory co-location to be performed in addition to the process pinning.

If only the variable `I_MPI_PIN_PROCS` is defined, the `I_MPI_PIN_MODE` value `mpd` is assumed. If only the variable `I_MPI_PIN_MODE` is defined, the `I_MPI_PIN_PROCS` value `allcores` is assumed.

Process pinning is performed if the number of CPUs on a node is not less than number of processes and if the operating system provides the necessary kernel interfaces.

If no CPU set is defined in the system, the number and order of the processors corresponds to the output of the `cat /proc/cpuinfo` command. If a CPU set is defined in the system, the `I_MPI_PIN_PROCS` value refers to the logical processors enabled in the current process set.

This variable does not influence the process placement that is controlled by the `mpdboot` and `mpiexec` commands. However, when this variable is defined and a process is placed upon the node, this process is bound to the next CPU out of the specified set.

Note that every host can be made to use their own value of an environment variable, or use a global value.

## Device Control

### `I_MPI_SPIN_COUNT`

Control the spin wait mode.

#### Syntax

`I_MPI_SPIN_COUNT=<scout>`

#### Arguments

<code>&lt;scout&gt;</code>	Defines the spin count for loop of polling fabric(s) in spin wait mode
<code>&gt; 0</code>	The default <code>&lt;scout&gt;</code> value is equal to 1 for sock, shm and ssm devices and equal to 250 for rdma and rdssm devices.

#### Description

Change the spin count for loop of polling fabric(s) before freeing processor in case absent messages for processing.

### `I_MPI_EAGER_THRESHOLD`

Change the eager/rendezvous cutover point for all devices.

#### Syntax

`I_MPI_EAGER_THRESHOLD=<nbytes>`

#### Arguments

<code>&lt;nbytes&gt;</code>	Defines eager/rendezvous cutover point
<code>&gt; 0</code>	The default <code>nbytes</code> value is equal to 262144

#### Description

Set this variable to control the point-to-point protocol switchover point.

There are eager and rendezvous protocols for data transferred by the library. Messages shorter than or equal in size to `<nbytes>` are sent eagerly. Larger messages are sent by using more memory efficient rendezvous protocol.



## I\_MPI\_INTRANODE\_EAGER\_THRESHOLD

Change the eager/rendezvous cutover point for intra-node communication mode.

### Syntax

`I_MPI_INTRANODE_EAGER_THRESHOLD=<nbytes>`

### Arguments

<code>&lt;nbytes&gt;</code>	Defines threshold for DAPL* intra-node communication
<code>&gt; 0</code>	The default <code>nbytes</code> value is equal to 262144

### Description

Set this variable to change threshold for intra-node communication mode.

There are eager and rendezvous protocols for data transferred by the library within the node. Messages shorter than or equal in size to `<nbytes>` are sent eagerly. Larger messages are sent by using more memory efficient rendezvous protocol.

If `I_MPI_INTRANODE_EAGER_THRESHOLD` is not set, the value of `I_MPI_EAGER_THRESHOLD` is used.

## RDMA and RDSSM Device Control

### RDMA\_IBA\_EAGER\_THRESHOLD

Change the eager/rendezvous cutover point.

### Syntax

`RDMA_IBA_EAGER_THRESHOLD=<nbytes>`

### Arguments

<code>&lt;nbytes&gt;</code>	Defines eager/rendezvous cutover point
<code>&gt; 0</code>	The default <code>nbytes</code> value is equal to 16512

### Description

Set this variable to control low level point-to-point protocol switchover point.

There are low level eager and rendezvous protocols for data transferred by the `rdma` and `rdssm` devices. Messages shorter than or equal in size to `<nbytes>` are sent eagerly through internal pre-registered buffers. Larger messages are sent by using more memory efficient rendezvous protocol.

### NOTES

- *This variable also determines the size of every pre-registered buffer. The higher it is, the more memory will be used for every established connection.*

### NUM\_RDMA\_BUFFER

Change the number of internal pre-registered buffers for each pair in a process group.

### Syntax

`NUM_RDMA_BUFFER=<nbuf>`

**Arguments**

<i>&lt;nbuf&gt;</i>	Defines the number of buffers for each pair in a process group
<i>&gt; 0</i>	The default <i>nbuf</i> value ranges between 8 and 40 depending on the cluster size and platform

**Description**

Set this variable to change the number of internal pre-registered buffers for each pair in a process group.

**NOTES**

- *The more pre-registered buffers are available, the more memory will be used for every established connection.*

**I\_MPI\_RDMA\_VBUF\_TOTAL\_SIZE**

Change the size of internal pre-registered buffers for each pair in a process group.

**Syntax**

*I\_MPI\_RDMA\_VBUF\_TOTAL\_SIZE=<nbytes>*

**Arguments**

<i>&lt;nbytes&gt;</i>	Defines the size of pre-registered buffers
<i>&gt; 0</i>	The default <i>nbytes</i> value is equal to 16640

**Description**

Setting the value *<nbytes>* for this environment variable directs the rdma and rdssm devices to set the size of internal pre-registered buffer for each pair in a process group according to specified value. The actual size calculated by adjusting *<nbytes>* for alignment buffer to optimal value.

**I\_MPI\_RDMA\_TRANSLATION\_CACHE**

Turn on/off the mode of using a registration cache.

**Syntax**

*I\_MPI\_RDMA\_TRANSLATION\_CACHE=<arg>*

**Arguments**

<i>&lt;arg&gt;</i>	Binary indicator
<i>enable, yes, on, 1</i>	Turn the memory registration cache on. This is the default value
<i>disable, no, off, 0</i>	Turn the memory registration cache off

**Description**

Set this variable to turn the memory registration cache on or off.

The cache substantially increases performance but may lead to correctness issues in certain rare situations. See the product *Release Notes* for further details.

## I\_MPI\_DAPL\_IP\_ADDR

## I\_MPI\_DAPL\_HOST

## I\_MPI\_DAPL\_HOST\_SUFFIX

Specify the Interface Adapter (IA) address.

### Syntax

`I_MPI_DAPL_IP_ADDR=<ipaddr>`

`I_MPI_DAPL_HOST=<hostname>`

`I_MPI_DAPL_HOST_SUFFIX=<hostsuff>`

### Arguments

<code>&lt;ipaddr&gt;</code>	Defines the IA address as an explicit IP address. The value <i>ipaddr</i> should have IP address of the host in the usual convention
<code>&lt;hostname&gt;</code>	Defines the IA address using a host name
<code>&lt;hostsuff&gt;</code>	Provides explicit hostname suffix that is prepended to the host name.

### Description

Set the `I_MPI_DAPL_IP_ADDR`, `I_MPI_DAPL_HOST`, or `I_MPI_DAPL_HOST_SUFFIX` variables to control the identity of the Interface Adapter (IA).

### NOTES

- *If none of these three variables is set, the IA address is determined automatically. This is the recommended mode of operation.*

## I\_MPI\_DAPL\_PORT

Specify the PSP (Public Service Point) value.

### Syntax

`I_MPI_DAPL_PORT=<port>`

### Arguments

<code>&lt;port&gt;</code>	Defines the port value
Between 1024 and 65536	The value of <i>port</i> must be an integer number between 1024 and 65536

### Description

Set this variable to specify the PSP value.

### NOTES

- *If this variable is not defined, the PSP port value is calculated automatically. This is the recommended mode of operation.*

## I\_MPI\_USE\_RENDEZVOUS\_RDMA\_WRITE

Turn on/off the use of rendezvous RDMA Write protocol instead of the default RDMA Read protocol.

### Syntax

`I_MPI_USE_RENDEZVOUS_RDMA_WRITE=<arg>`

### Arguments

<code>&lt;arg&gt;</code>	Binary indicator
<code>enable, yes, on, 1</code>	Turn the RDMA Write rendezvous protocol on
<code>disable, no, off, 0</code>	Turn the RDMA Write rendezvous protocol off. This is the default value

### Description

Set this variable to select RDMA Write based rendezvous protocol.

Certain DAPL\* providers have slow RDMA Read implementation on certain platforms. Switching on the rendezvous protocol based on RDMA Write operation may increase performance in these cases.

## I\_MPI\_RDMA\_USE\_EVD\_FALLBACK

Turn on/off the Event Dispatcher (EVD) based polling fallback path.

### Syntax

`I_MPI_RDMA_USE_EVD_FALLBACK=<arg>`

### Arguments

<code>&lt;arg&gt;</code>	Binary indicator
<code>enable, yes, on, 1</code>	Turn the EVD based fallback on
<code>disable, no, off, 0</code>	Turn the EVD based fallback off. This is the default value

### Description

Set this variable to use DAPL\* Event Dispatcher (EVD) for detecting incoming messages.

Use this method instead of the default method of buffer polling if the DAPL\* provider does not guarantee the delivery of the transmitted data in order from low to high addresses.

## NOTES

- *Note that the EVD path is typically substantially slower than the default algorithm.*

## I\_MPI\_USE\_DYNAMIC\_CONNECTIONS

Turn on/off the dynamic connection establishment.

### Syntax

`I_MPI_USE_DYNAMIC_CONNECTIONS=<arg>`

### Arguments

<code>&lt;arg&gt;</code>	Binary indicator
<code>enable, yes, on, 1</code>	Turn the dynamic connection establishment on. This is the default value.
<code>disable, no, off, 0</code>	Turn the dynamic connection establishment off.

### Description

Set this variable to control dynamic connection establishment.

If this mode is enabled, connections are established upon first communication between each pair of processes. This is the default behavior. All connections are established upfront if this variable is off.

## I\_MPI\_DYNAMIC\_CONNECTIONS\_MODE

Choose the algorithm of dynamic establishment of the DAPL\* connections.

### Syntax

`I_MPI_DYNAMIC_CONNECTION_MODE=<arg>`

### Arguments

<code>&lt;arg&gt;</code>	Mode selector
<code>reject</code>	Deny one of simultaneous connection requests. This is the default value
<code>disconnect</code>	Deny one of simultaneous connection requests after both connections established

### Description

Set this variable to choose the algorithm for handling dynamically established connections for DAPL\* capable fabrics.

In the `reject` mode one of the requests is rejected if two processes initiate the connection simultaneously. In the `disconnect` mode both connections are established, but then one is disconnected. The `disconnect` mode is provided to avoid a bug in some broken providers.

## I\_MPI\_DAPL\_CONNECTION\_TIMEOUT

Specify DAPL\* connection timeout.

### Syntax

`I_MPI_DAPL_CONNECTION_TIMEOUT=<value>`

### Arguments

<code>&lt;value&gt;</code>	Defines DAPL* connection timeout value in microseconds
<code>&gt; 0</code>	Default value is infinite

**Description**

Set this variable to specify timeout for DAPL\* connection establishment operations.

**NOTES**

- *If this variable is not defined, infinite timeout is used. This is the recommended mode of operation.*

**I\_MPI\_TWO\_PHASE\_BUF\_ENLARGEMENT**

Turn on/off the mode of using two-phase buffer enlargement.

**Syntax**

`I_MPI_TWO_PHASE_BUF_ENLARGEMENT=<arg>`

**Arguments**

<code>&lt;arg&gt;</code>	Binary indicator
<code>enable, yes, on, 1</code>	Turn the mode of using two-pharse buffer enlargement on
<code>disable, no, off, 0</code>	Turn the mode of using two-pharse buffer enlargement off. This is the default value

**Description**

Set this variable to turn on/off the mode of using two- phase buffer enlargement.

If this mode is turned on, small size internal pre-registered RDMA buffers are allocated and enlarged later if the data to transfer exceeds some threshold. Two-phase buffer enlargement is off by default.

**I\_MPI\_RDMA\_SHORT\_BUF\_THRESHOLD**

Change threshold for two-pharse buffer enlargement mode.

**Syntax**

`I_MPI_RDMA_SHORT_BUF_THRESHOLD=<nbytes>`

**Arguments**

<code>&lt;nbytes&gt;</code>	Defines threshold for starting RDMA buffers enlargement
<code>&gt; 0</code>	The default <i>nbytes</i> value is equal to 580

**Description**

Set this variable to change threshold for two-phase buffer enlargement mode. This variable is used only if `I_MPI_TWO_PHASE_BUF_ENLARGEMENT` is set to `enable`.

**I\_MPI\_USE\_DAPL\_INTRANODE**

Turn on/off the DAPL\* intra-node communication mode.

**Syntax**

`I_MPI_USE_DAPL_INTRANODE=<arg>`

**Arguments**

<code>&lt;arg&gt;</code>	Binary indicator
<code>enable, yes, on, 1</code>	Turn the DAPL* intra-node communication on
<code>disable, no, off, 0</code>	Turn the DAPL* intra-node communication off. This is the default value

**Description**

Set this variable to specify the intra-node communication mode for the universal device. If DAPL\* intra-node communication mode is turned on then small messages are transferred using shared memory and large ones via the DAPL\* layer. The threshold value for selecting communication layer of each message is determined by setting the `I_MPI_INTRANODE_EAGER_THRESHOLD` variable. It works only if neither shared memory nor DAPL\* layer were turned off by setting `I_MPI_DEVICE` environment variable.

**I\_MPI\_CONN\_EVD\_QLEN**

Define event queue size of DAPL\* event dispatcher.

**Syntax**

`I_MPI_CONN_EVD_QLEN=<size>`

**Arguments**

<code>&lt;size&gt;</code>	Defines length of event queue
<code>&gt; 0</code>	The default value queried from DAPL provider.

**Description**

Set this variable to define event queue size of DAPL event dispatcher. If this variable is set, the minimum value between size and the value obtained from the provider is used as the size of the event queue. The provider is required to provide a queue size that is at least equal to the calculated value, but the provider is free to provide a larger queue size.

**Collective Operation Control****I\_MPI\_FAST\_COLLECTIVES**

Turn on/off the optimization of the collective operations.

**Syntax**

`I_MPI_FAST_COLLECTIVES=<arg>`

**Arguments**

<code>&lt;arg&gt;</code>	Binary indicator
<code>enable, yes, on, 1</code>	Turn the collective optimizations on
<code>disable, no, off, 0</code>	Turn the collective optimizations off. This is the default value

**Description**

Set this variable to controls optimization level of the collective operations.

The character of optimization depends upon internal package settings. All collective optimizations are turned off by default.

**NOTES**

- If `I_MPI_FAST_COLLECTIVES` is turned on, then all other settings related to the collective operations (see `I_MPI_BCAST_NUM_PROCS`, `I_MPI_BCAST_MSG`, and so on) are not observed directly, because more suitable algorithms are chosen automatically in this case.
- Some optimizations of the collective operations may lead to violation of the MPI recommendation regarding the order of execution of the collective operations. Therefore results obtained in two different runs may differ depending on the process layout with respect to the processors and certain other factors.
- Some optimizations controlled by this variable may have an experimental character. In case of failure, turn the collective optimizations off and repeat the run.

**I\_MPI\_BCAST\_NUM\_PROCS****I\_MPI\_BCAST\_MSG**

Control `MPI_Bcast` algorithm thresholds.

**Syntax**

`I_MPI_BCAST_NUM_PROCS=<nproc>`

`I_MPI_BCAST_MSG=<nbytes1,nbytes2>`

**Arguments**

<code>&lt;nproc&gt;</code>	Defines the <code>MPI_Bcast</code> number of processes algorithm threshold
<code>&gt; 0</code>	The default value is 8

<code>&lt;nbytes1,nbytes2&gt;</code>	Defines the <code>MPI_Bcast</code> buffer size algorithm thresholds in bytes
<code>&gt; 0</code> <code>nbytes2 &gt;= nbytes1</code>	The default value is 12288,524288

**Description**

Set these variables to control selection of the `MPI_Bcast` algorithms according to the following scheme:

1. The first algorithm is selected if the message size is below `<nbytes1>`, or the number of processes in the operation is below `<nproc>`.
2. The second algorithm is selected if the message size lies between `<nbytes1>` and `<nbytes2>`, and the number of processes in the operation is a power of two.
3. The third algorithm is selected otherwise.

**I\_MPI\_ALLTOALL\_NUM\_PROCS****I\_MPI\_ALLTOALL\_MSG**

Control `MPI_Alltoall` algorithm thresholds.

**Syntax**

`I_MPI_ALLTOALL_NUM_PROCS=<nproc>`

`I_MPI_ALLTOALL_MSG=<nbytes1,nbytes2>`



**Arguments**

<code>&lt;nproc&gt;</code>	Defines the MPI_Alltoall number of processes algorithm thresholds
<code>&gt; 0</code>	The default value is 8

<code>&lt;nbytes1,nbytes2&gt;</code>	Defines the MPI_Alltoall buffer size algorithm thresholds in bytes
<code>&gt; 0</code> <code>nbytes2 &gt;= nbytes1</code>	The default value is 256,32768

**Description**

Set these variables to control selection of the `MPI_Alltoall` algorithms according to the following scheme:

1. The first algorithm is selected if the message size is below `<nbytes1>`, and the number of processes in the operation is not less than `<nproc>`.
2. The second algorithm is selected if the message size lies between `<nbytes1>` and `<nbytes2>`, or if the message size lies below `<nbytes1>` and the number of processes in the operation is less than `<nproc>`.
3. The third algorithm is selected otherwise.

**I\_MPI\_ALLGATHER\_MSG**

Control `MPI_Allgather` algorithm thresholds.

**Syntax**

`I_MPI_ALLGATHER_MSG=<nbytes1,nbytes2>`

**Arguments**

<code>&lt;nbytes1,nbytes2&gt;</code>	Defines the MPI_Allgather buffer size algorithm thresholds in bytes
<code>&gt; 0</code> <code>nbytes2 &gt;= nbytes1</code>	The default value is 81920,524288

**Description**

Set this variable to control selection of the `MPI_Allgather` algorithms according to the following scheme:

1. The first algorithm is selected if the message size lies below `<nbytes2>` and the number of processes in the operation is a power of two.
2. The second algorithm is selected if the message size lies below `<nbytes1>` and number of processes in the operation is not a power of two.
3. The third algorithm is selected otherwise.

**I\_MPI\_ALLREDUCE\_MSG**

Control `MPI_Allreduce` algorithm thresholds.

**Syntax**

`I_MPI_ALLREDUCE_MSG=<nbytes>`

**Arguments**

<i>&lt;nbytes&gt;</i>	Defines the MPI_Allreduce buffer size algorithm threshold in bytes
<i>&gt; 0</i>	The default value is 2048

**Description**

Set this variable to control selection of the **MPI\_Allreduce** algorithms according to the following scheme:

1. The first algorithm is selected if the message size lies below *<nbytes>*, or the reduction operation is user-defined, or the count argument is less than the nearest power of two less than or equal to the number of processes.
2. The second algorithm is selected otherwise.

**I\_MPI\_REDUCE\_MSG**

Control **MPI\_Reduce** algorithm thresholds.

**Syntax**

*I\_MPI\_REDUCE\_MSG=<nbytes>*

**Arguments**

<i>&lt;nbytes&gt;</i>	Defines the MPI_Reduce buffer size protocol threshold in bytes
<i>&gt; 0</i>	The default value is 2048

**Description**

Set this variable to control selection of the **MPI\_Reduce** algorithms according to the following scheme:

1. The first algorithm is selected if the message size lies above *<nbytes>*, the reduction operation is not user defined, and the count argument is not less than the nearest power of two less than or equal to the number of processes.
2. The second algorithm is selected otherwise.

**I\_MPI\_SCATTER\_MSG**

Control **MPI\_Scatter** algorithm thresholds.

**Syntax**

*I\_MPI\_SCATTER\_MSG=<nbytes>*

**Arguments**

<i>&lt;nbytes&gt;</i>	Defines the MPI_Scatter buffer size algorithm threshold in bytes
<i>&gt; 0</i>	The default value is 2048

**Description**

Set this variable to control selection of the **MPI\_Scatter** algorithms according to the following scheme:

1. The first algorithm is selected on intercommunicators if the message size lies above *<nbytes>*.
2. The second algorithm is selected otherwise.

## I\_MPI\_GATHER\_MSG

Control `MPI_Gather` algorithm thresholds.

### Syntax

`I_MPI_GATHER_MSG=<nbytes>`

### Arguments

<code>&lt;nbytes&gt;</code>	Defines the <code>MPI_Gather</code> buffer size algorithm threshold in bytes
<code>&gt; 0</code>	The default value is 2048

### Description

Set this variable to control selection of the `MPI_Gather` algorithms according to the following scheme:

1. The first algorithm is selected on intercommunicators if the message size lies above `<nbytes>`.
2. The second algorithm is selected otherwise.

## I\_MPI\_REDS CAT\_MSG

Control `MPI_Reduce_scatter` algorithm thresholds.

### Syntax

`I_MPI_REDS CAT_MSG=<nbytes1,nbytes2>`

### Arguments

<code>&lt;nbytes1,nbytes2&gt;</code>	Defines the <code>MPI_Reduce_scatter</code> buffer size algorithm threshold in bytes
<code>&gt; 0</code> <code>nbytes2 &gt;= nbytes1</code>	The default value is 512,524288

### Description

Set this variable to control selection of the `MPI_Reduce_scatter` algorithms according to the following scheme:

1. The first algorithm is selected if the reduction operation is commutative and the message size lies below `<nbytes2>`.
2. The second algorithm is selected if the reduction operation is commutative and message size lies above `<nbytes2>`, or if the reduction operation is not commutative and the message size lies above `<nbytes1>`.
3. The third algorithm is selected otherwise.

## Miscellaneous

## I\_MPI\_TIMER\_KIND

Select the timer used by the `MPI_Wtime` and `MPI_Wtick` calls.

### Syntax

`I_MPI_TIMER_KIND=<timename>`

**Arguments**

<code>&lt;timername&gt;</code>	Defines timer type
<code>gettimeofday</code>	MPI_Wtime and MPI_Wtick functions will work through the function <code>gettimeofday(2)</code> . This is a default value.
<code>rdtsc</code>	MPI_Wtime and MPI_Wtick functions will use the high resolution RDTSC timer

**Description**

Set this variable to select either the ordinary or RDTSC timer.

**NOTES**

- The resolution of the default `gettimeofday(2)` timer may be insufficient on certain platforms.