

# mpiP: Lightweight, Scalable MPI Profiling

Version 3.4.1  
13 March 2014

Jeffrey Vetter  
[jsvetter@users.sourceforge.net](mailto:jsvetter@users.sourceforge.net)

Chris Chambreau  
[chcham@users.sourceforge.net](mailto:chcham@users.sourceforge.net)

---

## Contents

- [Introduction](#)
    - [Downloading](#)
    - [Contributing](#)
    - [New Features](#)
  - [Using mpiP \(summary\)](#)
    - [Supported Platforms](#)
  - [Configuring and Building mpiP](#)
    - [Linking Examples](#)
  - [Run-time Configuration](#)
  - [mpiP Output](#)
    - [Report Viewers](#)
  - [Controlling mpiP Profiling Scope](#)
  - [Caveats](#)
  - [List of Profiled Routines](#)
  - [List of Routines With Collected Sent Message Size Information](#)
  - [List of I/O Routines](#)
  - [List of RMA Routines](#)
  - [How to add MPI calls to profile](#)
  - [License](#)
- 

## Introduction

mpiP is a lightweight profiling library for MPI applications. Because it only collects statistical information about MPI functions, mpiP generates considerably less overhead and much less data than tracing tools. All the information captured by mpiP is task-local. It only uses communication during report generation, typically at the end of the experiment, to merge results from all of the tasks into one output file.

We have tested mpiP on a variety of C/C++/Fortran applications from 2 to 262144 processes, including a 262144-process run on the LLNL Sequoia BG/Q system.

Please send your comments, questions, and ideas for enhancements to [mpip-help@lists.sourceforge.net](mailto:mpip-help@lists.sourceforge.net). To receive mail regarding new mpiP releases, please subscribe to [mpip-announce@lists.sourceforge.net](mailto:mpip-announce@lists.sourceforge.net) (send e-mail with body "subscribe" to [mpip-announce-request@lists.sourceforge.net](mailto:mpip-announce-request@lists.sourceforge.net)). Please also consider subscribing to [mpip-users@lists.sourceforge.net](mailto:mpip-users@lists.sourceforge.net) to contribute and receive mpiP use and status information.

To learn more about performance analysis with mpiP, see Vetter, J.S. and M.O. McCracken, "[Statistical Scalability Analysis of Communication Operations in Distributed Applications](#)," Proc. ACM SIGPLAN Symp. on Principles and Practice of Parallel Programming (PPOPP), 2001.

## Downloading

You may download the current version of mpiP from <http://sourceforge.net/projects/mpip>.

## Contributing

We are constantly improving mpiP. Bug fixes and ports to new platforms are always welcome. Many thanks to the following contributors (chronological order):

- Michael McCracken (UCSD)
- Curt Janssen (Sandia National Laboratories)
- Mike Campbell (UIUC)
- Jim Brandt (Sandia National Laboratories)
- Philip Roth (Oak Ridge National Laboratory)
- Tushar Mohan (SiCortex)
- Philip Mucci (SiCortex)

- Karl Schulz (Texas Advanced Computing Center)

## New Features with Release 3.4.1

Release v3.4.1 addresses the following issue:

- Added de-activation of shared object source lookup when libbfd is not available.

Release v3.4 addresses the following issues:

- Compatibility with MPI-3.
- Histogram reporting for Point-to-point (-p) and Collective (-y) operation message sizes and communicators.
- Added a low-memory-use concise report format, with the ability to set the default report format and specify report formats at run time.
- Supports MPI call reporting (no call sites) with stack depth (-k) of 0.
- Configure can disable SO lookup functionality.

Release v3.3 addresses the following issues:

- Support for shared object source lookup with libbfd.
- Improved configuration process for recent versions of binutils and Cray XE6.
- Added "-z" MPIP run time flag to suppress report generation at Finalize.
- Corrected number of stack frames available when using glibc backtrace.

Release v3.2.1 addresses the following issue:

- Improved support for SLURM run-time instrumentation.

Release v3.2 addresses the following issues:

- Support for MPI RMA functions.
- Support for glibc backtrace.
- Default to MPI\_Wtime if platform-specific timers are not found.

Release v3.1.2 addresses the following issues:

- Better MPI support for Init\_thread, Testany, Testsome, Waitany, and Waitsome.
- Improved support for MIPS64-Linux.
- Added option to configure for generating weak Fortran symbols in the case of multiple Fortran mangling schemes in the application object files (--enable-fortranweak).
- Addressed various outstanding issues (see ChangeLog for more details).

Release v3.1.1 addresses the following issues:

- Revert to gettimeofday as default Linux timer.
- MIPS64-Linux stack walking support.
- Catamount dclock timer support.
- Greater install flexibility:
  - 'install' target only installs lib and doc files.
  - 'install-api', 'install-bin', 'install-all' targets provide additional install functionality.
  - New 'uninstall' target.

For more information, please see the ChangeLog in the distribution.

[Top](#)

---

## Using mpiP

Using mpiP is very simple. Because it gathers MPI information through the MPI profiling layer, mpiP is a link-time library. That is, you don't have to recompile your application to use mpiP. Note that you might have to recompile to include the '-g' option. This is important if you want mpiP to decode the PC to a source code filename and line number automatically. mpiP will work without -g, but mileage may vary.

To compile a simple program on an LLNL x86\_64-linux system where libunwind is installed, add the following libraries to your link command:

```
-L${mpiP_root}/lib -lmpiP -lm -lbfd -liberty -lunwind
```

For example, the new mpiP link command becomes

```
$ mpicc -g 1-hot-potato.o -o 1-hot-potato.exe -L${mpiP_root}/lib -lmpiP -lm -lbfd -liberty -lunwind
```

from

```
$ mpicc -g 1-hot-potato.o -o 1-hot-potato.exe
```

Make sure the mpiP library appears before the MPI library on your link line. The libraries (-lbfd -liberty ) provide

support for decoding the symbol information; they are part of GNU binutils.

Run your application. You can verify that mpiP is working by identifying the header and trailer in standard out.

```
mpiP:
mpiP: mpiP: mpiP V3.2.0 (Build Mar 10 2010/13:27:39)
mpiP: Direct questions and errors to mpiP-help@lists.sourceforge.net
mpiP:
mpiP:
mpiP: Storing mpiP output in [./1-hot-potato.exe.2.27872.1.mpiP].
mpiP:
```

By default, the output file is written to the current directory of the application. mpiP files are always much smaller than trace files, so writing them to this directory is safe.

## Supported Platforms

mpiP has been tested on several Linux, AIX, UNICOS and IBM BG systems. Please contact us with bug reports or questions regarding these platforms. The following table indicates platforms where mpiP was successfully run and any requirements for that platform.

Platform	OS	Compiler	MPI	binutils	Requirements
x86_64-Linux	2.6.18 CHAOS Kernel	Intel 9.1 PGI 7.0	MVAPICH 0.9.7	2.20.51	Example configure command: ./configure LDFLAGS=-L/usr/lib64 LIBS="-lbfd -liberty" --enable-collective-report-default --enable-demangling=GNU --with-cc=mpicc --with-cxx=mpiCC --with-f77=mpif77
IBM BG/Q	Driver V1R2M0	IBM XL 12.1		2.21.1	Source code lookup support requires zlib. Example configure command: CFLAGS=-I/bgsys/drivers/ppcfloor/toolchain/gnu/build-powerpc64-bgq-linux/binutils-2.21.1-build/bfd -I/bgsys/drivers/ppcfloor/toolchain/gnu/gcc-4.4.6/include -I/bgsys/drivers/ppcfloor/toolchain/gnu/gdb-7.1/include LIBS=-L/bgsys/drivers/ppcfloor/toolchain/gnu/build-powerpc64-bgq-linux/binutils-2.21.1-build/bfd -L/g0/chcham/ToolTesting/mpiP/bgqos_0/mpiP-3.3/libz/zlib-1.2.6 CC=mpixlc_r CXX=mpixlcxx_r F77=mpixlf77_r ./configure --enable-getarg
Cray XT3/XT4	Catamount 1.4.32	PGI 6.1-4 C/C/Fortran, from Cray PrgEnv-pgi module version 1.4.32	Cray XT3 Message Passing Toolkit (MPT) version 1.4.32	Binutils 2.17, built for x86_64- unknown- linux-gnu	--with-cc=cc --with-cxx=CC --with-f77=ftn --with-libs="-lpgf90 -lpgf90_rpm1 -lpgf902 -lpgf90rtl -lpgftnrtl" --build=x86_64-unknown-linux-gnu --host=x86_64-cray-catamount --target=x86_64-cray-catamount --enable-getarg \ --with-wtime \ --with-binutils-dir=<path to binutils-2.17 installation>
Cray X1E	UNICOS/mp 3.1.16	Cray Standard C 5.5.0.5, Cray Fortran 5.5.0.5	Cray Message Passing Toolkit (MPT) 2.4.0.7	from Cray Open Software module 3.6	Requires libelf/libdwf. Example configure flags: --disable-libunwind --enable-dwarf --disable-demangle --enable-getarg --with-cc=cc --with-cxx=CC --with-f77=ftn Python on an alternative system may be needed to "make wrappers.c", due to missing socket module.

[Top](#)

## Configuring and Building mpiP

### Configuring mpiP

Currently, mpiP requires a compatible GNU [binutils](#) installation for source lookup and demangling features. Alternatively, [libelf](#) and [libdwf](#) can be used for source lookup. The binutils installation location may need to be specified with either the --with-binutils-dir option or with the --with-include and --with-ldflags configure flags. It is likely that the compilers will need to be identified as well, with the --with-cc, --with-cxx, and --with-f77 flags. Use CFLAGS and FFLAGS variables to specify compiler options, as in CFLAGS="-O3" ./configure.

There are many configuration options available. Please use ./configure --help to list all of these options. Additional description are provided for the following options:

Flag	Effect	Description
<b>Reporting Options</b>		
--enable-demangling=[type]	Specify demangling support.	If the GNU option is specified, demangling is applied to each symbol by default using the libliberty implementation. For the IBM option, demangling support is implemented in the library libmpiPdmg.a. Use GNU for the Intel compiler.

--disable-mpi-io	Disable MPI-I/O reporting.	Useful for generating an mpiP library without MPI I/O for MPI implementations such as Quadrics MPI that has a separate MPI I/O library.
--enable-collective-report-default	Report data is aggregated on a per-callsite basis	By default, mpiP aggregates all process data at a single process which generates the report. Enabling this feature causes mpiP report generation to default to aggregating callsite data only for each individual callsite being reported. This dramatically reduces the memory requirements for large runs of applications that make many MPI calls. See run-time flags -l and -r to modify report generation behavior.
<b>Stack Trace Options</b>		
--enable-stackdepth=[depth]	Specify maximum stacktrace depth (default is 8).	Stacktraces with larger than 8 levels are sometime useful for some applications.
--disable-libunwind	Do not use libunwind to generate stack traces.	Currently, libunwind seems useful on IA64-Linux and x86-Linux platforms, although it can conflict with the libunwind.a provided with the Intel compiler.
<b>Address Lookup Options</b>		
--enable-dwarf	Use libdwarf/libelf for source lookup.	libdwarf and libelf can be used for address-to-source translation as an alternative to binutils libbfd.
--disable-bfd	Do not use GNU binutils libbfd for source lookup.	Binutils is not always available or compatible.
<b>Timing Options</b>		
--with-gettimeofday	Use gettimeofday for timing.	Use the gettimeofday call for timing instead of the default platform timer.
--with-wtime	Use MPI_Wtime for timing	Use the MPI_Wtime call for timing instead of the default platform timer.
--with-clock_gettime	Use clock_gettime for timing.	Use the clock_gettime monotonic timer for timing instead of the default platform timer.
--with-dclock	Use Catamount dclock for timing.	Use the dclock timer for timing on Catamount systems instead of the default platform timer.
--enable-check-time	Enable AIX check for negative time values.	Activate IBM timing debugging code.
<b>Fortran-related Options</b>		
--enable-getarg	Use getarg to get Fortran command line args.	This is used on UNICOS to provide access to the command line for Fortran applications.
--disable-fortran-xlate	Do not translate Fortran opaque objects.	Opaque object translation is not necessary on some platforms, but necessary for Fortran applications on some 64-bit platforms.
--enable-fortranweak	Generate weak symbols for additional Fortran symbol name styles.	If application objects have been created from compilers with different Fortran symbol name styles, it may be necessary to generate weak symbols to capture all MPI calls.

## Build targets

Command	Effect
---------	--------

make	Build mpiP library or libraries for MPI profiling
make install	Install bin, include, lib, and slib (if applicable) directories. The default install directory is the mpiP source directory. The installation location can be specified with the prefix variable as in <code>make prefix=[install directory] install</code> .
make shared	Make shared object version of library for runtime insertion (Linux). Support for runtime insertion on AIX and for MPI calls made within shared objects on Linux and AIX will be provided in a future release.
make API	Make standalone API library. See mpiP-API.c and mpiP-API.h for available features.
make check	Run mpiP dejagnu tests. Requires that runtest is available.
make add_binutils_objs	For convenience, add the binutils objects to the mpiP library. The binutils installation location must have been specified during configuration.
make add_libunwind_objs	For convenience, add the libunwind objects to the mpiP library.

## Example Application Link Commands

LLNL users can now use the `srun-mpip`, `poe-mpip`, and `poe-mpip-cxx` wrapper scripts to use mpiP without re-linking their application. AIX executables would still need to be linked with `-bnoobjreorder` for successful runtime address lookup. Additionally, all LLNL installations contain the appropriate binutils objects in the mpiP library, so the `-lbfd`, `-liberty`, and `-lintl` flags are not required. An example runtime script for `mpirun` is provided in the mpiP bin directory. Many of the following examples use LLNL-specific compile scripts.

OS	Compiler	Language	Example Link Command
AIX	Visual Age	C	<code>mpxlc -g -bnoobjreorder 1-hot-potato.c -o 1-hot-potato.exe -L/usr/local/tools/mpiP/lib -lmpiP -lbfd -liberty -lintl -lm</code>
		C++	<code>mpCC_r -g -bnoobjreorder 4-demangle.C -o 4-demangle.exe -L/usr/local/tools/mpiP/lib -lmpiPdmg -lbfd -liberty -lintl -lm</code>
		Fortran	<code>mpxlf -g -bnoobjreorder sweep-ops.f -o sweep-ops.exe -L/usr/local/tools/mpiP/lib -lmpiP -lbfd -liberty -lintl -lm</code>
Linux	Intel	C	<code>mpiicc -g 1-hot-potato.c -o 1-hot-potato.exe -L/usr/local/tools/mpiP/lib -lmpiP -lbfd -liberty -lm -lmpio</code>
		C++	<code>mpiicc -g 4-demangle.C -o 4-demangle.exe -L/usr/local/tools/mpiP/lib -lmpiP -lbfd -liberty -lm -lmpio</code>
		Fortran	<code>mpiifc -g sweep-ops.f -o sweep-ops.exe -L/usr/local/tools/mpiP/lib -lmpiP -lbfd -liberty -lm -lmpio</code>
	PGI	C	<code>mpipgcc -g 1-hot-potato.c -o 1-hot-potato.exe -L/usr/local/tools/mpiP/lib -lmpiP -lbfd -liberty -lm -lmpio</code>
		C++	<code>mpipgCC -g 4-demangle.C -o 4-demangle.exe -L/usr/local/tools/mpiP/lib -lmpiP -lbfd -liberty -lm -lmpio</code>
		Fortran	<code>mpipgf77 -g sweep-ops.f -o sweep-ops.exe -L/usr/local/tools/mpiP/lib -lmpiP -lbfd -liberty -lm -lmpio</code>
	GNU	C	<code>mpicc -g 1-hot-potato.c -o 1-hot-potato.exe -L/usr/local/tools/mpiP/lib -lmpiP -lbfd -liberty -lm</code>
		C++	<code>mpiCC -g 4-demangle.C -o 4-demangle.exe -L/usr/local/tools/mpiP/lib -lmpiP -lbfd -liberty -lm</code>
		Fortran	<code>mpif77 -g sweep-ops.f -o sweep-ops.exe -L/usr/local/tools/mpiP/lib -lmpiP -lbfd -liberty -lm</code>
Cray X1	Cray	C/C++/Fortran	Link with <code>-lmpiP -lbfd -liberty -ldwarf -lelf</code>
Cray XD1	GNU or PGI	C/C++/Fortran	Link with <code>[path_to_mpiP_install]/libmpiP.a -lbfd -liberty [mpich_libs]</code>

### Note:

- If source lookup is failing during report generation, the script `mpip-insert-src` can be used from a login node to translate addresses in the mpiP report to source information.
- Source lookup for callsites may fail with certain versions of binutils. If you are running into trouble, you may want to download a recent snapshot from <ftp://ftp.gnu.org/gnu/binutils/>.

[Top](#)

## Run-time Configuration of mpiP

mpiP has several configurable parameters that a user can set via the environment variable `MPIP`. Setting `MPIP` is done like command-line parameters: `"-t 10 -k 2"`. Additionally, a comma can be used to delimit multiple parameters, as in `"-t10,-k2"`. Currently, mpiP has several configurable parameters.

Option	Description	Default
--------	-------------	---------

-c	Generate concise version of report, omitting callsite process-specific detail.	
-d	Suppress printing of callsite detail sections.	
-e	Print report data using floating-point format.	
-f dir	Record output file in directory <dir>.	.
-g	Enable mpiP debug mode.	disabled
-k n	Sets callsite stack traceback depth to <n>.	1
-l	Use less memory to generate the report by using MPI collectives to generate callsite information on a callsite-by-callsite basis.	
-n	Do not truncate full pathname of filename in callsites.	
-o	Disable profiling at initialization. Application must enable profiling with MPI_Pcontrol().	
-p	Point-to-point histogram reporting on message size and communicator used.	
-r	Generate the report by aggregating data at a single task.	default
-s n	Set hash table size to <n>.	256
-t x	Set print threshold for report, where <x> is the MPI percentage of time for each callsite.	0.0
-v	Generates both concise and verbose report output.	
-x exe	Specify the full path to the executable.	
-y	Collective histogram reporting on message size and communicator used.	
-z	Suppress printing of the report at MPI_Finalize.	

For example, to set the callsite stack walking depth to 2 and the report print threshold to 10%, you simply need to define the mpiP string in your environment, as in any of the following examples:

```
$ export MPIP="-t 10.0 -k 2" (bash)
```

```
$ export MPIP=-t10.0,-k2 (bash)
```

```
$ setenv MPIP "-t 10.0 -k 2" (csh)
```

mpiP prints a message at initialization if it successfully finds this MPIP variable.

[Top](#)

## mpiP Output

Here is some sample output from mpiP with an application that has 4 MPI calls. It is broken down by sections below. Here also is the experiment setup. **Note that MPIP does not capture information about ALL MPI calls.** Local calls, such as `MPI_Comm_size`, are omitted from the profiling library measurement to reduce perturbation and mpiP output.

### The test code:

```
sleeptime = 10;
MPI_Init (&argc, &argv);
MPI_Comm_size (comm, &nprocs);
MPI_Comm_rank (comm, &rank);
MPI_Barrier (comm);
if (rank == 0)
{
    sleep (sleeptime);          /* slacker! delaying everyone else */
}
MPI_Barrier (comm);
MPI_Finalize ();
```

### The code was compiled with:

```
$ mpicc -g -DAIX 9-test-mpip-time.c -o 9-test-mpip-time.exe \
-L.. -L/g/g2/vetter/AIX/lib -lmpiP -lbfd -liberty -lintl -lm
```

### Environment variables were set as:

```
$ export MPIP="-t 10.0"
```

### The example was executed on MCR like this:

```
$ srun -n 4 -ppdebug ./9-test-mpip-time.exe
```

### This experiment produced an output file that we can now analyze:

```
./9-test-mpip-time.exe.4.25972.1.mpiP
```

Header information provides basic information about your performance experiment.

```
@ mpiP
@ Command : /g/g0/chcham/mpiP/devo/testing/./9-test-mpip-time.exe
@ Version : 2.8.2
```

```

@ MPIP Build date : Jan 10 2005, 15:15:47
@ Start time : 2005 01 10 16:01:32
@ Stop time : 2005 01 10 16:01:42
@ Timer Used : gettimeofday
@ MPIP env var : -t 10.0
@ Collector Rank : 0
@ Collector PID : 25972
@ Final Output Dir : .
@ MPI Task Assignment : 0 mcr88
@ MPI Task Assignment : 1 mcr88
@ MPI Task Assignment : 2 mcr89
@ MPI Task Assignment : 3 mcr89

```

This next section provides an overview of the application's time in MPI. Apptime is the wall-clock time from the end of MPI\_Init until the beginning of MPI\_Finalize. MPI\_Time is the wall-clock time for all the MPI calls contained within Apptime. MPI% shows the ratio of this MPI\_Time to Apptime. The asterisk (\*) is the aggregate line for the entire application.

```

-----
@--- MPI Time (seconds) -----
-----
Task   AppTime  MPITime  MPI%
  0      10    0.000243  0.00
  1      10      10    99.92
  2      10      10    99.92
  3      10      10    99.92
  *      40      30    74.94

```

The callsite section identifies all the MPI callsites within the application. The first number is the callsite ID for this mpiP file. The next column shows the type of MPI call (w/o the MPI\_ prefix). The name of the function that contains this MPI call is next, followed by the file name and line number. Finally, the last column shows the PC, or program counter, for that MPI callsite. Note that the default setting for callsite stack walk depth is 1. Other settings will enumerate callsites by the entire stack trace rather than the single callsite alone.

```

-----
@--- Callsites: 2 -----
-----
ID Lev File/Address      Line Parent_Funct      MPI_Call
  1  0 9-test-mpip-time.c    52 main                Barrier
  2  0 9-test-mpip-time.c    61 main                Barrier

```

The aggregate time section is a very quick overview of the top twenty MPI callsites that consume the most aggregate time in your application. Call identifies the type of MPI function. Site provides the callsite ID (as listed in the callsite section). Time is the aggregate time for that callsite in milliseconds. The next two columns show the ratio of that aggregate time to the total application time and to the total MPI time, respectively. The COV column indicates the variation in times of individual processes for this callsite by presenting the coefficient of variation as calculated from the individual process times. A larger value indicates more variation between the process times.

```

-----
@--- Aggregate Time (top twenty, descending, milliseconds) -----
-----
Call           Site      Time      App%      MPI%      COV
Barrier        2         3e+04    75.00    100.00    0.67
Barrier        1          0.405    0.00     0.00     0.59

```

The next section is similar to the aggregate time section, although it reports on the top 20 callsites for total sent message sizes. For example:

```

-----
@--- Aggregate Sent Message Size (top twenty, descending, bytes) -----
-----
Call           Site      Count      Total      Avrg      MPI%
Send           7         320      1.92e+06   6e+03     99.96
Bcast          1         12        336        28        0.02

```

The final sections are the ad nauseum listing of the statistics for each callsite across all tasks, followed by an aggregate line (indicated by an asterisk in the Rank column). The first section is for operation time followed by the section for message sizes.

```

-----
@--- Callsite Time statistics (all, milliseconds): 8 -----
-----
Name           Site Rank  Count      Max      Mean      Min      App%  MPI%
Barrier        1    0     1     0.107    0.107    0.107    0.00  44.03
Barrier        1    *     4     0.174    0.137    0.107    0.00  0.00

Barrier        2    0     1     0.136    0.136    0.136    0.00  55.97
Barrier        2    1     1     1e+04    1e+04    1e+04    99.92 100.00
Barrier        2    2     1     1e+04    1e+04    1e+04    99.92 100.00
Barrier        2    3     1     1e+04    1e+04    1e+04    99.92 100.00
Barrier        2    *     4     1e+04    7.5e+03  0.136    74.94 100.00

```

Remember that we configured MPIP to not print lines where MPI% was less than 10%. All aggregate lines are printed regardless of the configuration settings.

Column	Description
Name	Name of the MPI function at that callsite.

Site	Callsite ID as listed in the callsite section above.
Rank	Task rank in MPI_COMM_WORLD.
Count	Number of times this call was executed.
Max	Maximum wall-clock time for one call.
Mean	Arithmetic mean of the wall-clock time for one call.
Min	Minimum wall-clock time for one call.
App%	Ratio of time for this call to the overall application time for each task.
MPI%	Ratio of time for this call to the overall MPI time for each task.

The aggregate result for each call has the same measurement meaning; however, the statistics are gathered across all tasks and compared with the aggregate application and MPI times.

The section for sent message sizes has a similar format:

```
-----
@--- Callsite Message Sent statistics (all, sent bytes) -----
-----
Name           Site Rank  Count    Max    Mean    Min    Sum
Send           5     0     80    6000   6000   6000  4.8e+05
Send           5     1     80    6000   6000   6000  4.8e+05
Send           5     2     80    6000   6000   6000  4.8e+05
Send           5     3     80    6000   6000   6000  4.8e+05
Send           5     *    320   6000   6000   6000  1.92e+06
```

Column	Description
Name	Name of the MPI function at that callsite.
Site	Callsite ID as listed in the callsite section above.
Rank	Task rank in MPI_COMM_WORLD.
Count	Number of times this call was executed.
Max	Maximum sent message size in bytes for one call.
Mean	Arithmetic mean of the sent message sizes in bytes for one call.
Min	Minimum sent message size in bytes for one call.
Sum	Total of all message sizes for this operation and callsite.

The format of MPI I/O report section is very similar to the sent message sizes section:

```
-----
@--- Callsite I/O statistics (all, I/O bytes) -----
-----
Name           Site Rank  Count    Max    Mean    Min    Sum
File_read     1     0     20     64     64     64    1280
File_read     1     1     20     64     64     64    1280
File_read     1     *     40     64     64     64    2560
```

## Report Viewers

- The [Tool Gear](#) project has a Qt mpiP viewer. LLNL users can run this as mpipview.

[Top](#)

## Controlling the Scope of mpiP Profiling in your Application

In mpiP, you can limit the scope of profiling measurements to specific regions of your code using the `MPI_Pcontrol(int level)` subroutine. A value of zero disables mpiP profiling, while any nonzero value enables profiling. To disable profiling initially at `MPI_Init`, use the `-o` configuration option. mpiP will only record information about MPI commands encountered between activation and deactivation. There is no limit to the number of times that an application can activate profiling during execution.

For example, in your application you can capture the MPI activity for timestep 5 only using `Pcontrol`. Remember to set the mpiP environment variable to include `-o` when using this feature.

```
for(i=1; i < 10; i++)
{
  switch(i)
  {
    case 5:
      MPI_Pcontrol(1);
      break;
    case 6:
      MPI_Pcontrol(0);
      break;
    default:
      break;
  }
}
```

```

}
/* ... compute and communicate for one timestep ... */
}

```

## Arbitrary Report Generation

You can also generate arbitrary reports by making calls to `MPI_Pcontrol()` with an argument of 3 or 4 (see table below). The first report generated will have the default report filename. Subsequent report files will have an index number included, such as `sweep3d.mpi.4.7371.1.mpiP`, `sweep3d.mpi.4.7371.2.mpiP`, etc. The final report will still be generated during `MPI_Finalize`. **NOTE:** In the current release, callsite IDs will not be consistent between reports. Comparison of callsite data between reports must be done by source location and callstack.

`MPI_Pcontrol` features should be fully functional for C/C++ as well as Fortran.

Pcontrol Argument	Behavior
0	Disable profiling.
1	Enable Profiling.
2	Reset all callsite data.
3	Generate verbose report.
4	Generate concise report.

If you want to generate individual reports each time a section of code is executed, but don't want the profile data to accumulate, you can specify code to reset the profile data, profile, and then generate reports. For example:

```

for(i=1; i < 10; i++)
{
  switch(i)
  {
    case 5:
      MPI_Pcontrol(2); // make sure profile data is reset
      MPI_Pcontrol(1); // enable profiling
      break;
    case 6:
      MPI_Pcontrol(3); // generate verbose report
      MPI_Pcontrol(4); // generate concise report
      MPI_Pcontrol(0); // disable profiling
      break;
    default:
      break;
  }
  /* ... compute and communicate for one timestep ... */
}

```

[Top](#)

## Caveats

- If mpiP has problems with the source code translation, you might be able to decode the program counters on LLNL systems with some of the following techniques. You can use `instmap`, `addr2line`, or look at the assembler code itself.
- Compiler transformations like loop unrolling can sometimes make one source code line appear as many different PCs. You can verify this by looking at the assembler. In my experience, both `instmap` and `addr2line` do a pretty good job of mapping these transformed PCs into a file name and line number.
  - `instmap`—an IBM utility
  - `addr2line`—a gnu tool
  - look at the assembler listing, or with GNU's `objdump (-d -S)`
  - use Totalview or `gdb` to translate the PC
- There are known incompatibilities with certain binutils versions and recent versions of the IBM compilers. As of this release, a fix has not been incorporated into binutils, however, using the `-bnoobjreorder` option is a valid work-around.
- In one case, we encountered problems on IBM machines with source lookup of 64-bit Fortran applications. It appears that an incorrect compiler configuration file was being used, incorrectly matching debugging information and PC values. We addressed this by using the link flag `-bpT:0x100000000`.
- Issues when stack walking optimized applications:
  - Applications compiled with `gcc` may return incorrect parent functions; however, the file and line number information may be correct.
  - Applications compiled with the Intel compiler may not be able to identify parent stack frames.
- If you are calling MPI functions from within dynamically loaded objects, you may need to recompile the library as a shared object.
- We have encountered occasional negative report values on Linux and AIX systems. We will continue to investigate this issue, but it is possible that this behavior may be experienced with mpiP.

[Top](#)

---

## MPI Routines Profiled with mpiP

MPI\_Allgather  
MPI\_Allgatherv  
MPI\_Allreduce  
MPI\_Alltoall  
MPI\_Alltoallv  
MPI\_Attr\_delete  
MPI\_Attr\_get  
MPI\_Attr\_put  
MPI\_Barrier  
MPI\_Bcast  
MPI\_Bsend  
MPI\_Bsend\_init  
MPI\_Buffer\_attach  
MPI\_Buffer\_detach  
MPI\_Cancel  
MPI\_Cart\_coords  
MPI\_Cart\_create  
MPI\_Cart\_get  
MPI\_Cart\_map  
MPI\_Cart\_rank  
MPI\_Cart\_shift  
MPI\_Cart\_sub  
MPI\_Cartdim\_get  
MPI\_Comm\_create  
MPI\_Comm\_dup  
MPI\_Comm\_group  
MPI\_Comm\_remote\_group  
MPI\_Comm\_remote\_size  
MPI\_Comm\_split  
MPI\_Comm\_test\_inter  
MPI\_Dims\_create  
MPI\_Error\_class  
MPI\_File\_close  
MPI\_File\_open  
MPI\_File\_preallocate  
MPI\_File\_read  
MPI\_File\_read\_all  
MPI\_File\_read\_at  
MPI\_File\_seek  
MPI\_File\_set\_view  
MPI\_File\_write  
MPI\_File\_write\_all  
MPI\_File\_write\_at  
MPI\_Gather  
MPI\_Gatherv  
MPI\_Graph\_create  
MPI\_Graph\_get  
MPI\_Graph\_map  
MPI\_Graph\_neighbors  
MPI\_Graph\_neighbors\_count  
MPI\_Graphdims\_get  
MPI\_Group\_compare  
MPI\_Group\_difference  
MPI\_Group\_excl  
MPI\_Group\_free  
MPI\_Group\_incl  
MPI\_Group\_intersection  
MPI\_Group\_translate\_ranks  
MPI\_Group\_union  
MPI\_Ibsend  
MPI\_Intercomm\_create  
MPI\_Intercomm\_merge  
MPI\_Iprobe  
MPI\_Irecv  
MPI\_Irsend  
MPI\_Isend  
MPI\_Issend  
MPI\_Keyval\_create  
MPI\_Keyval\_free  
MPI\_Pack  
MPI\_Probe

MPI\_Recv  
MPI\_Recv\_init  
MPI\_Reduce  
MPI\_Reduce\_scatter  
MPI\_Request\_free  
MPI\_Rsend  
MPI\_Rsend\_init  
MPI\_Scan  
MPI\_Scatter  
MPI\_Scatterv  
MPI\_Send  
MPI\_Send\_init  
MPI\_Sendrecv  
MPI\_Sendrecv\_replace  
MPI\_Ssend  
MPI\_Ssend\_init  
MPI\_Start  
MPI\_Startall  
MPI\_Test  
MPI\_Testall  
MPI\_Testany  
MPI\_Testsome  
MPI\_Topo\_test  
MPI\_Type\_commit  
MPI\_Type\_free  
MPI\_Type\_get\_contents  
MPI\_Type\_get\_envelope  
MPI\_Unpack  
MPI\_Wait  
MPI\_Waitall  
MPI\_Waitany  
MPI\_Waitsome  
MPI\_Win\_complete  
MPI\_Win\_create  
MPI\_Win\_fence  
MPI\_Win\_free  
MPI\_Win\_get\_group  
MPI\_Win\_lock  
MPI\_Win\_post  
MPI\_Win\_start  
MPI\_Win\_test  
MPI\_Win\_unlock  
MPI\_Win\_wait

[Top](#)

---

## **MPI Routines For Which mpiP Gathers Sent Message Size Data**

MPI\_Allgather  
MPI\_Allgatherv  
MPI\_Allreduce  
MPI\_Alltoall  
MPI\_Bcast  
MPI\_Bsend  
MPI\_Gather  
MPI\_Gatherv  
MPI\_Ibsend  
MPI\_Irsend  
MPI\_Isend  
MPI\_Issend  
MPI\_Reduce  
MPI\_Rsend  
MPI\_Scan  
MPI\_Scatter  
MPI\_Send  
MPI\_Sendrecv  
MPI\_Sendrecv\_replace  
MPI\_Ssend

[Top](#)

---

## **MPI Routines For Which mpiP Gathers I/O Data**

MPI\_File\_close

MPI\_File\_open  
MPI\_File\_preallocate  
MPI\_File\_read  
MPI\_File\_read\_all  
MPI\_File\_read\_at  
MPI\_File\_seek  
MPI\_File\_set\_view  
MPI\_File\_write  
MPI\_File\_write\_all  
MPI\_File\_write\_at

[Top](#)

---

## MPI Routines For Which mpiP Gathers RMA Origin Data

MPI\_Accumulate  
MPI\_Get  
MPI\_Put

[Top](#)

---

## How to add MPI calls to profile

Here is an example of how to add MPI calls to mpiP, using the MPI\_Comm\_spawn call as an example:

1. Insert the appropriate call with appropriate arguments into the mpi.protos.txt.in file:

```
int MPI_Comm_spawn ( char *command, char *argv[], int maxprocs, MPI_Info info, int root,  
MPI_Comm comm, MPI_Comm *intercomm, int array_of_errcodes[] )
```

2. Configure mpiP or, if you have already configured mpiP, run ./config.status.
3. Currently, it is necessary to add entries for MPI opaque objects to the make-wrappers.py script. MPI\_Comm\_spawn has 3 arguments that are MPI opaque object which need to be added to make-wrappers.py dictionaries:

1. Add the following entries to the opaqueInArgDict:

```
("MPI_Comm_spawn", "info"): "MPI_Info",  
("MPI_Comm_spawn", "comm"): "MPI_Comm",
```

2. Add the following entry to the opaqueOutArgDict:

```
("MPI_Comm_spawn", "intercomm"): "MPI_Comm",
```

4. When you build mpiP, you should see an MPI\_Comm\_spawn wrapper in the generated wrappers.c file.

[Top](#)

---

## License

Copyright (c) 2006, The Regents of the University of California.  
Produced at the Lawrence Livermore National Laboratory  
Written by Jeffery Vetter and Christopher Chembreau.  
UCRL-CODE-223450.  
All rights reserved.

This file is part of mpiP. For details, see <http://mpip.sourceforge.net/>.

Redistribution and use in source and binary forms, with or without modification, are permitted provided that the following conditions are met:

\* Redistributions of source code must retain the above copyright notice, this list of conditions and the disclaimer below.

\* Redistributions in binary form must reproduce the above copyright notice, this list of conditions and the disclaimer (as noted below) in the documentation and/or other materials provided with the distribution.

\* Neither the name of the UC/LLNL nor the names of its contributors may be used to endorse or promote products derived from this software without specific prior written permission.

THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS

"AS IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE REGENTS OF THE UNIVERSITY OF CALIFORNIA, THE U.S. DEPARTMENT OF ENERGY OR CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

#### Additional BSD Notice

1. This notice is required to be provided under our contract with the U.S. Department of Energy (DOE). This work was produced at the University of California, Lawrence Livermore National Laboratory under Contract No. W-7405-ENG-48 with the DOE.

2. Neither the United States Government nor the University of California nor any of their employees, makes any warranty, express or implied, or assumes any liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately-owned rights.

3. Also, reference herein to any specific commercial products, process, or services by trade name, trademark, manufacturer or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or the University of California. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or the University of California, and shall not be used for advertising or product endorsement purposes.

[Top](#)

---

For further information please send mail to [mpip-help@lists.sourceforge.net](mailto:mpip-help@lists.sourceforge.net).