## ARGONNE NATIONAL LABORATORY 9700 South Cass Avenue Argonne, IL 60439

ANL/MCS-TM-234

# Users Guide for ROMIO: A High-Performance, Portable MPI-IO Implementation

by

Rajeev Thakur, Robert Ross, Ewing Lusk, William Gropp, Robert Latham

Mathematics and Computer Science Division

Technical Memorandum No. 234

Revised May 2004, November 2007, April 2010

This work was supported by the Mathematical, Information, and Computational Sciences Division subprogram of the Office of Advanced Scientific Computing Research, U.S. Department of Energy, under Contract W-31-109-Eng-38; and by the Scalable I/O Initiative, a multiagency project funded by the Defense Advanced Research Projects Agency (Contract DABT63-94-C-0049), the Department of Energy, the National Aeronautics and Space Administration, and the National Science Foundation.

## Contents

Al	ostract	1							
1	Introduction								
2	Major Changes in This Version	1							
3	General Information         3.1 ROMIO Optimizations         3.2 Hints          3.2.1 Hints for PFS          3.2.2 Hints for XFS          3.2.3 Hints for PVFS (v1)          3.2.4 Hints for PVFS (v2)          3.2.5 Hints for Lustre          3.2.6 Hints for PANFS (Panasas)          3.2.7 Systemwide Hints          3.3 Using ROMIO on NFS          3.3.1 ROMIO, NFS, and Synchronization          3.4 Using testfs          3.5 ROMIO and MPI_FILE_SYNC          3.6 ROMIO and MPI_FILE_SET_SIZE	1 2 3 5 5 6 6 6 8 8 9 9 10 10							
4	Installation Instructions 4.1 Configuring for Linux and Large Files	<b>10</b> 11							
5	Testing ROMIO	<b>12</b>							
6	Compiling and Running MPI-IO Programs	<b>12</b>							
7	Limitations of This Version of ROMIO	13							
8	Usage Tips	13							
9	Reporting Bugs	14							
10	ROMIO Internals	14							
11	Learning MPI-IO	14							
	Major Changes in Previous Releases  12.1 Major Changes in Version 1.2.3	14 14 15 15 16							
100	2101 011000	10							

# Users Guide for ROMIO: A High-Performance, Portable MPI-IO Implementation

by

Rajeev Thakur, Robert Ross, Ewing Lusk, and William Gropp

#### Abstract

ROMIO is a high-performance, portable implementation of MPI-IO (the I/O chapter in the MPI Standard). This document describes how to install and use ROMIO version 1.2.4 on various machines.

## 1 Introduction

ROMIO<sup>1</sup> is a high-performance, portable implementation of MPI-IO (the I/O chapter in MPI [4]). This document describes how to install and use ROMIO version 1.2.4 on various machines.

## 2 Major Changes in This Version

- Added section describing ROMIO MPI\_FILE\_SYNC and MPI\_FILE\_CLOSE behavior to User's Guide
- Bug removed from PVFS ADIO implementation regarding resize operations
- Added support for PVFS listic operations (see Section 3.2)
- Added the following working hints: romio\_pvfs\_listio\_read, romio\_pvfs\_listio\_write

## 3 General Information

This version of ROMIO includes everything defined in the MPI I/O chapter except support for file interoperability and user-defined error handlers for files (§ 4.13.3). The subarray and distributed array datatype constructor functions from Chapter 4 (§ 4.14.4 & § 4.14.5) have been implemented. They are useful for accessing arrays stored in files. The functions MPI\_File\_f2c and MPI\_File\_c2f (§ 4.12.4) are also implemented. C, Fortran, and profiling interfaces are provided for all functions that have been implemented.

This version of ROMIO runs on at least the following machines: IBM SP; Intel Paragon; HP Exemplar; SGI Origin2000; Cray T3E; NEC SX-4; other symmetric multiprocessors from HP, SGI, DEC, Sun, and IBM; and networks of workstations (Sun, SGI, HP, IBM, DEC, Linux, and FreeBSD). Supported file systems are IBM PIOFS, Intel PFS, HP/Convex HFS, SGI XFS, NEC SFS, PVFS, NFS, NTFS, and any Unix file system (UFS).

This version of ROMIO is included in MPICH 1.2.4; an earlier version is included in at least the following MPI implementations: LAM, HP MPI, SGI MPI, and NEC MPI.

<sup>1</sup>http://www.mcs.anl.gov/romio

Note that proper I/O error codes and classes are returned and the status variable is filled only when used with MPICH revision 1.2.1 or later.

You can open files on multiple file systems in the same program. The only restriction is that the directory where the file is to be opened must be accessible from the process opening the file. For example, a process running on one workstation may not be able to access a directory on the local disk of another workstation, and therefore ROMIO will not be able to open a file in such a directory. NFS-mounted files can be accessed.

An MPI-IO file created by ROMIO is no different from any other file created by the underlying file system. Therefore, you may use any of the commands provided by the file system to access the file, for example, ls, mv, cp, rm, ftp.

Please read the limitations of this version of ROMIO that are listed in Section 7 of this document (e.g., restriction to homogeneous environments).

## 3.1 ROMIO Optimizations

ROMIO implements two I/O optimization techniques that in general result in improved performance for applications. The first of these is data sieving [2]. Data sieving is a technique for efficiently accessing noncontiguous regions of data in files when noncontiguous accesses are not provided as a file system primitive. The naive approach to accessing noncontiguous regions is to use a separate I/O call for each contiguous region in the file. This results in a large number of I/O operations, each of which is often for a very small amount of data. The added network cost of performing an I/O operation across the network, as in parallel I/O systems, is often high because of latency. Thus, this naive approach typically performs very poorly because of the overhead of multiple operations. In the data sieving technique, a number of noncontiguous regions are accessed by reading a block of data containing all of the regions, including the unwanted data between them (called "holes"). The regions of interest are then extracted from this large block by the client. This technique has the advantage of a single I/O call, but additional data is read from the disk and passed across the network.

There are four hints that can be used to control the application of data sieving in ROMIO: ind\_rd\_buffer\_size, ind\_wr\_buffer\_size, romio\_ds\_read, and romio\_ds\_write. These are discussed in Section 3.2.

The second optimization is  $two-phase\ I/O\ [1]$ . Two-phase I/O, also called collective buffering, is an optimization that only applies to collective I/O operations. In two-phase I/O, the collection of independent I/O operations that make up the collective operation are analyzed to determine what data regions must be transferred (read or written). These regions are then split up amongst a set of aggregator processes that will actually interact with the file system. In the case of a read, these aggregators first read their regions from disk and redistribute the data to the final locations, while in the case of a write, data is first collected from the processes before being written to disk by the aggregators.

There are five hints that can be used to control the application of two-phase I/O: cb\_config\_list, cb\_nodes, cb\_buffer\_size, romio\_cb\_read, and romio\_cb\_write. These are discussed in Subsection 3.2.

## 3.2 Hints

If ROMIO doesn't understand a hint, or if the value is invalid, the hint will be ignored. The values of hints being used by ROMIO for a file can be obtained at any time via MPI\_File\_get\_info.

The following hints control the data sieving optimization and are applicable to all file system types:

- ind\_rd\_buffer\_size Controls the size (in bytes) of the intermediate buffer used by ROMIO when performing data sieving during read operations. Default is 4194304 (4 Mbytes).
- ind\_wr\_buffer\_size Controls the size (in bytes) of the intermediate buffer used by ROMIO when performing data sieving during write operations. Default is 524288 (512 Kbytes).
- romio\_ds\_read Determines when ROMIO will choose to perform data sieving. Valid values are enable, disable, or automatic. Default value is automatic. In automatic mode ROMIO may choose to enable or disable data sieving based on heuristics.
- romio\_ds\_write Same as above, only for writes.

The following hints control the two-phase (collective buffering) optimization and are applicable to all file system types:

- cb\_buffer\_size Controls the size (in bytes) of the intermediate buffer used in two-phase collective I/O. If the amount of data that an aggregator will transfer is larger than this value, then multiple operations are used. The default is 4194304 (4 Mbytes).
- cb\_nodes Controls the maximum number of aggregators to be used. By default this is set to the number of unique hosts in the communicator used when opening the file.
- romio\_cb\_read Controls when collective buffering is applied to collective read operations. Valid values are enable, disable, and automatic. Default is automatic. When enabled, all collective reads will use collective buffering. When disabled, all collective reads will be serviced with individual operations by each process. When set to automatic, ROMIO will use heuristics to determine when to enable the optimization.
- romio\_cb\_write Controls when collective buffering is applied to collective write operations. Valid values are enable, disable, and automatic. Default is automatic. See the description of romio\_cb\_read for an explanation of the values.
- romio\_no\_indep\_rw This hint controls when "deferred open" is used. When set to true, ROMIO will make an effort to avoid performing any file operation on non-aggregator nodes. The application is expected to use only collective operations. This is discussed in further detail below.
- cb\_config\_list Provides explicit control over aggregators. This is discussed in further detail below.

For some systems configurations, more control is needed to specify which hardware resources (processors or nodes in an SMP) are preferred for collective I/O, either for performance reasons

or because only certain resources have access to storage. The additional MPLInfo key name cb\_config\_list specifies a comma-separated list of strings, each string specifying a particular node and an optional limit on the number of processes to be used for collective buffering on this node.

This refers to the same processes that **cb\_nodes** refers to, but specifies the available nodes more precisely.

The format of the value of cb\_config\_list is given by the following BNF:

The value hostname identifies a processor. This name must match the name returned by MPI\_Get\_processor\_name <sup>2</sup> for the specified hardware. The value \* as a hostname matches all processors. The value of maxprocesses may be any nonnegative integer (zero is allowed).

The value maxprocesses specifies the maximum number of processes that may be used for collective buffering on the specified host. If no value is specified, the value one is assumed. If \* is specified for the number of processes, then all MPI processes with this same hostname will be used..

Leftmost components of the info value take precedence.

Note: Matching of processor names to cb\_config\_list entries is performed with string matching functions and is independent of the listing of machines that the user provides to mpirun/mpiexec. In other words, listing the same machine multiple times in the list of hosts to run on will not cause a \*:1 to assign the same host four aggregators, because the matching code will see that the processor name is the same for all four and will assign exactly one aggregator to the processor.

The value of this info key must be the same for all processes (i.e., the call is collective and each process must receive the same hint value for these collective buffering hints). Further, in the ROMIO implementation the hint is only recognized at MPI\_File\_open time.

The set of hints used with a file is available through the routine MPI\_File\_get\_info, as documented in the MPI standard. As an additional feature in the ROMIO implementation, wildcards will be expanded to indicate the precise configuration used with the file, with the hostnames in the rank order used for the collective buffering algorithm (this is not implemented at this time).

Here are some examples of how this hint might be used:

- \*:1 One process per hostname (i.e., one process per node)
- box12:30,\*:0 Thirty processes on one machine, namely box12, and none anywhere else.
- n01,n11,n21,n31,n41 One process on each of these specific nodes only.

When the values specified by cb\_config\_list conflict with other hints (e.g., the number of collective buffering nodes specified by cb\_nodes), the implementation is encouraged to take the

<sup>&</sup>lt;sup>2</sup>The MPI standard requires that the output from this routine identify a particular piece of hardware; some MPI implementations may not conform to this requirement. MPICH does conform to the MPI standard.

minimum of the two values. In other words, if cb\_config\_list specifies ten processors on which I/O should be performed, but cb\_nodes specifies a smaller number, then an implementation is encouraged to use only cb\_nodes total aggregators. If cb\_config\_list specifies fewer processes than cb\_nodes, no more than the number in cb\_config\_list should be used.

The implementation is also encouraged to assign processes in the order that they are listed in cb\_config\_list.

The following hint controls the deferred open feature of romio and are also applicable to all file system types:

• romio\_no\_indep\_rw - If the application plans on performing only collective operations and this hint is set to "true", then ROMIO can have just the aggregators open a file. The cb\_config\_list and cb\_nodes hints can be given to further control which nodes are aggregators.

For PVFS, PIOFS, and PFS:

- striping\_factor Controls the number of I/O devices to stripe across. The default is file system dependent, but for PVFS it is -1, indicating that the file should be striped across all I/O devices.
- striping\_unit Controls the striping unit (in bytes). For PVFS the default will be the PVFS file system default strip size.
- start\_iodevice Determines what I/O device data will first be written to. This is a number in the range of 0 ... striping\_factor 1.

#### 3.2.1 Hints for PFS

• pfs\_svr\_buf - Turns on PFS server buffering. Valid values are true and false. Default is false.

#### 3.2.2 Hints for XFS

For XFS control is provided for the direct I/O optimization:

- direct\_read Controls direct I/O for reads. Valid values are true and false. Default is false.
- direct\_write Controls direct I/O for writes. Valid values are true and false. Default is false.

## 3.2.3 Hints for PVFS (v1)

For PVFS control is provided for the use of the listic interface. This interface to PVFS allows for a collection of noncontiguous regions to be requested (for reading or writing) with a single operation. This can result in substantially higher performance when accessing noncontiguous regions. Support for these operations in PVFS exists after version 1.5.4, but has not been heavily tested, so use of the interface is disabled in ROMIO by default at this time. The hints to control listic use are:

- romio\_pvfs\_listio\_read Controls use of listio for reads. Valid values are enable, disable, and automatic. Default is disable.
- romio\_pvfs\_listio\_write Controls use of listio for writes. Valid values are enable, disable, and automatic. Default is disable.

## 3.2.4 Hints for PVFS (v2)

The PVFS v2 file system has many tuning parameters.

• dtype i/o

#### 3.2.5 Hints for Lustre

• romio\_lustre\_co\_ratio

In stripe-contiguous IO pattern, each OST will be accessed by a group of IO clients. CO means \*C\*lient/\*O\*ST ratio, or the max. number of IO clients for each OST. CO=1 by default.

#### • romio\_lustre\_coll\_threshold

We won't do collective I/O if this hint is set and the IO request size is bigger than this value. That's because when the request size is big, the collective communication overhead increases and the benefits from collective I/O becomes limited. A value of 0 means always perform collective I/O

#### • romio\_lustre\_cb\_ds\_threshold

ROMIO can optimize collective I/O with a version of data sieving. If the I/O request is smaller than this hint's value, though, ROMIO will not try to apply the data sieving optimization.

#### • romio\_lustre\_ds\_in\_coll

Collective IO will apply read-modify-write to deal with non-contiguous data by default. However, it will introduce some overhead(IO operation and locking). The Lustre developers have run tests where data sieving showed bad collective write performance for some kinds of workloads. So, to avoid this, we define the romio\_lustre\_ds\_in\_coll hint to disable the read-modify-write step in collective I/O. This optimization is distinct from the one in independent I/O (controlled by romio\_ds\_read and romio\_ds\_write).

#### 3.2.6 Hints for PANFS (Panasas)

PanFS allows users to specify the layout of a file at file-creation time. Layout information includes the number of StorageBlades (SB) across which the data is stored, the number of SBs across which a parity stripe is written, and the number of consecutive stripes that are placed on the same set of SBs. The panfs\_layout\_\* hints are only used if supplied at file-creation time.

- panfs\_layout\_type Specifies the layout of a file: 2 = RAID0 3 = RAID5 Parity Stripes
- panfs\_layout\_stripe\_unit The size of the stripe unit in bytes

- panfs\_layout\_total\_num\_comps The total number of StorageBlades a file is striped across.
- panfs\_layout\_parity\_stripe\_width If the layout type is RAID5 Parity Stripes, this hint specifies the number of StorageBlades in a parity stripe.
- panfs\_layout\_parity\_stripe\_depth If the layout type is RAID5 Parity Stripes, this hint specifies the number of contiguous parity stripes written across the same set of SBs.
- panfs\_layout\_visit\_policy If the layout type is RAID5 Parity Stripes, the policy used to determine the parity stripe a given file offset is written to: 1 = Round Robin

PanFS supports the "concurrent write" (CW) mode, where groups of cooperating clients can disable the PanFS consistency mechanisms and use their own consistency protocol. Clients participating in concurrent write mode use application specific information to improve performance while maintaining file consistency. All clients accessing the file(s) must enable concurrent write mode. If any client does not enable concurrent write mode, then the PanFS consistency protocol will be invoked. Once a file is opened in CW mode on a machine, attempts to open a file in non-CW mode will fail with EACCES. If a file is already opened in non-CW mode, attempts to open the file in CW mode will fail with EACCES. The following hint is used to enable concurrent write mode.

• panfs\_concurrent\_write If set to 1 at file open time, the file is opened using the PanFS concurrent write mode flag. Concurrent write mode is not a persistent attribute of the file.

Below is an example PanFS layout using the following parameters:

```
= 3
- panfs_layout_type
- panfs_layout_total_num_comps
                                   = 100
- panfs_layout_parity_stripe_width = 10
- panfs_layout_parity_stripe_depth = 8
- panfs_layout_visit_policy
```

Parity Stripe Group 1			Parity	Stripe Gr	oup 2	Parity	rity Stripe Group 10		
SB1	SB2	SB10	SB11	SB12	SB20	SB91	SB92	. SB100	
D1	D2	D10	D91	D92	D100	D181	D182 .	D190	
D11	D12	D20	D101	D102	D110	D191	D192	D193	
D21	D22	D30							
D31	D32	D40							
D41	D42	D50							
D51	D52	D60							
D61	D62	D70							
D71	D72	D80							
D81	D82	D90	D171	D172	D180	D261	D262	D270	
D271	D272	D273							

### 3.2.7 Systemwide Hints

A site administrator with knowledge of the storage and networking capabilities of a machine might be able to come up with a set of hint values that work better for that machine than the ROMIO default values. As an extention to the standard, ROMIO will consult a "hints file". This file provides an additional mechanism for setting MPI-IO hints, albeit in a ROMIO-specific manner. The hints file contains a list of hints and their values. ROMIO will use these initial hint settings, though programs are free to override any of them.

The format of the hints file is a list of hints and their values, one per line. A # character in the first column indicates a comment, and ROMIO will ignore the entire line. Here's an example:

```
# this is a comment describing the following setting
cb_nodes 32
# these nodes happen to have the best connection to storage
cb_config_list n01,n11,n21,n31,n41
```

ROMIO will look for these hints in the file /etc/romio-hints. A user can set the environment variable ROMIO HINTS to the name of a file which ROMIO will use instead.

## 3.3 Using ROMIO on NFS

It is worth first mentioning that in no way do we encourage the use of ROMIO on NFS volumes. NFS is not a high-performance protocol, nor are NFS servers typically very good at handling the types of concurrent access seen from MPI-IO applications. Nevertheless, NFS is a very popular mechanism for providing access to a shared space, and ROMIO does support MPI-IO to NFS volumes, provided that they are configured properly.

To use ROMIO on NFS, file locking with fcntl must work correctly on the NFS installation. On some installations, fcntl locks don't work. To get them to work, you need to use Version 3 of NFS, ensure that the lockd daemon is running on all the machines, and have the system administrator mount the NFS file system with the "noac" option (no attribute caching). Turning off attribute caching may reduce performance, but it is necessary for correct behavior.

The following are some instructions we received from Ian Wells of HP for setting the noac option on NFS. We have not tried them ourselves. We are including them here because you may find them useful. Note that some of the steps may be specific to HP systems, and you may need root permission to execute some of the commands.

```
>1. first confirm you are running nfs version 3
>rpcnfo -p 'hostname' | grep nfs
>ie
>
     goedel >rpcinfo -p goedel | grep nfs
>
     100003
               2
                          2049 nfs
                    udp
     100003
               3
                   udp
                          2049 nfs
>
>
```

```
>2. then edit /etc/fstab for each nfs directory read/written by MPIO
    on each machine used for multihost MPIO.
>
>
    Here is an example of a correct fstab entry for /epm1:
>
>
    ie grep epm1 /etc/fstab
>
>
       ROOOOT 11>grep epm1 /etc/fstab
>
       gershwin:/epm1 /rmt/gershwin/epm1 nfs bg,intr,noac 0 0
>
    if the noac option is not present, add it
    and then remount this directory
>
>
    on each of the machines that will be used to share MPIO files
>
>ie
>
>ROOOOT >umount /rmt/gershwin/epm1
>ROOOOT >mount /rmt/gershwin/epm1
>
>3. Confirm that the directory is mounted noac:
>ROOOOT >grep gershwin /etc/mnttab
>gershwin:/epm1 /rmt/gershwin/epm1 nfs
>noac,acregmin=0,acregmax=0,acdirmin=0,acdirmax=0 0 0 899911504
```

#### 3.3.1 ROMIO, NFS, and Synchronization

NFS has a "sync" option that specifies that the server should put data on the disk before replying that an operation is complete. This means that the actual I/O cost on the server side cannot be hidden with caching, etc. when this option is selected.

In the "async" mode the server can get the data into a buffer (and perhaps put it in the write queue; this depends on the implementation) and reply right away. Obviously if the server were to go down after the reply was sent but before the data was written, the system would be in a strange state, which is why so many articles suggest the "sync" option.

Some systems default to "sync", while others default to "async", and the default can change from version to version of the NFS software. If you find that access to an NFS volume through MPI-IO is particularly slow, this is one thing to check out.

#### 3.4 Using testfs

The testfs ADIO implementation provides a harness for testing components of ROMIO or discovering the underlying I/O access patterns of an application. When testfs is specified as the file system type, no actual files will be opened. Instead debugging information will be displayed on the processes opening the file. Subsequent I/O operations on this testfs file will provide additional debugging information.

The intention of the testfs implementation is that it serve as a starting point for further instrumentation when debugging new features or applications. As such it is expected that users will want to modify the ADIO implementation in order to get the specific output they desire.

#### 3.5 ROMIO and MPI\_FILE\_SYNC

The MPI specification notes that a call to MPI\_FILE\_SYNC "causes all previous writes to fh by the calling process to be transferred to the storage device." Likewise, calls to MPI\_FILE\_CLOSE have this same semantic. Further, "if all processes have made updates to the storage device, then all such updates become visible to subsequent reads of fh by the calling process."

The intended use of MPI\_FILE\_SYNC is to allow all processes in the communicator used to open the file to see changes made to the file by each other (the second part of the specification). The definition of "storage device" in the specification is vague, and it isn't necessarily the case that calling MPI\_FILE\_SYNC will force data out to permanent storage.

Since users often use MPI\_FILE\_SYNC to attempt to force data out to permanent storage (i.e. disk), the ROMIO implementation of this call enforces stronger semantics for most underlying file systems by calling the appropriate file sync operation when MPI\_FILE\_SYNC is called (e.g. fsync). However, it is still unwise to assume that the data has all made it to disk because some file systems (e.g. NFS) may not force data to disk when a client system makes a sync call.

For performance reasons we do *not* make this same file system call at MPI\_FILE\_CLOSE time. At close time ROMIO ensures any data has been written out to the "storage device" (file system) as defined in the standard, but does not try to push the data beyond this and into physical storage. Users should call MPI\_FILE\_SYNC before the close if they wish to encourage the underlying file system to push data to permanent storage.

#### 3.6 ROMIO and MPI\_FILE\_SET\_SIZE

MPI\_FILE\_SET\_SIZE is a collective routine used to resize a file. It is important to remember that a MPI-IO routine being collective does not imply that the routine synchronizes the calling processes in any way (unless this is specified explicitly).

As of 1.2.4, ROMIO implements MPI\_FILE\_SET\_SIZE by calling ftruncate from all processes. Since different processes may call the function at different times, it means that unless external synchronization is used, a resize operation mixed in with writes or reads could have unexpected results.

In short, if synchronization after a set size is needed, the user should add a barrier or similar operation to ensure the set size has completed.

## 4 Installation Instructions

Since ROMIO is included in MPICH, LAM, HP MPI, SGI MPI, and NEC MPI, you don't need to install it separately if you are using any of these MPI implementations. If you are using some other MPI, you can configure and build ROMIO as follows:

Untar the tar file as

```
gunzip -c romio.tar.gz | tar xvf -
```

```
or
```

```
zcat romio.tar.Z | tar xvf -
then

cd romio
   ./configure
   make
```

Some example programs and a Makefile are provided in the romio/test directory. Run the examples as you would run any MPI program. Each program takes the filename as a command-line argument "-fname filename".

The configure script by default configures ROMIO for the file systems most likely to be used on the given machine. If you wish, you can explicitly specify the file systems by using the "-file\_system" option to configure. Multiple file systems can be specified by using '+' as a separator, e.g.,

If you intend to leave ROMIO where you built it, you should *not* install it; make install is used only to move the necessary parts of a built ROMIO to another location. The installed copy will have the include files, libraries, man pages, and a few other odds and ends, but not the whole source tree. It will have a test directory for testing the installation and a location-independent Makefile built during installation, which users can copy and modify to compile and link against the installed copy.

```
To rebuild ROMIO with a different set of configure options, do make distclean
```

to clean everything, including the Makefiles created by configure. Then run configure again with the new options, followed by make.

## 4.1 Configuring for Linux and Large Files

32-bit systems running linux kernel version 2.4.0 or newer and glibc version 2.2.0 or newer can support files greater than 2 GBytes in size. This support is currently automaticly detected and enabled. We document the manual steps should the automatic detection not work for some reason.

The two macros \_FILE\_OFFSET\_BITS=64 and \_LARGEFILE64\_SOURCE tell gnu libc it's ok to support large files on 32 bit platforms. The former changes the size of off\_t (no need to change source. might affect interoperability with libraries compiled with a different size of off\_t). The latter exposes the gnu libc functions open64(), write64(), read64(), etc. ROMIO does not make use of the 64 bit system calls directly at this time, but we add this flag for good measure.

If your linux system is relatively new, there is an excellent chance it is running kernel 2.4.0 or newer and glibc-2.2.0 or newer. Add the string

```
"-D_FILE_OFFSET_BITS=64 -D_LARGEFILE64_SOURCE"
```

to your CFLAGS environment variable before runnint ./configure

## 5 Testing ROMIO

To test if the installation works, do make testing

in the romio/test directory. This calls a script that runs the test programs and compares the results with what they should be. By default, make testing causes the test programs to create files in the current directory and use whatever file system that corresponds to. To test with other file systems, you need to specify a filename in a directory corresponding to that file system as follows:

make testing TESTARGS="-fname=/foo/piofs/test"

## 6 Compiling and Running MPI-IO Programs

If ROMIO is not already included in the MPI implementation, you need to include the file mpio.h for C or mpiof.h for Fortran in your MPI-IO program.

Note that on HP machines running HPUX and on NEC SX-4, you need to compile Fortran programs with mpif90, because mpif77 does not support 8-byte integers.

```
With MPICH, HP MPI, or NEC MPI, you can compile MPI-IO programs as mpicc foo.c or mpif77 foo.f or mpif90 foo.f
```

As mentioned above, mpif90 is preferred over mpif77 on HPUX and NEC because the f77 compilers on those machines do not support 8-byte integers.

With SGI MPI, you can compile MPI-IO programs as

```
cc foo.c -lmpi
or
     f77 foo.f -lmpi
or
     f90 foo.f -lmpi

With LAM, you can compile MPI-IO programs as
     hcc foo.c -lmpi
or
     hf77 foo.f -lmpi
```

If you have built ROMIO with some other MPI implementation, you can compile MPI-IO programs by explicitly giving the path to the include file mpio.h or mpiof.h and explicitly specifying the path to the library libmpio.a, which is located in \$(ROMIO\_HOME)/lib/\$(ARCH)/libmpio.a.

Run the program as you would run any MPI program on the machine. If you use mpirun, make sure you use the correct mpirun for the MPI implementation you are using. For example, if you are using MPICH on an SGI machine, make sure that you use MPICH's mpirun and not SGI's mpirun.

## 7 Limitations of This Version of ROMIO

- When used with any MPI implementation other than MPICH revision 1.2.1 or later, the status argument is not filled in any MPI-IO function. Consequently, MPI\_Get\_count and MPI\_Get\_elements will not work when passed the status object from an MPI-IO operation.
- Additionally, when used with any MPI implementation other than MPICH revision 1.2.1 or later, all MPI-IO functions return only two possible error codes—MPI\_SUCCESS on success and MPI\_ERR\_UNKNOWN on failure.
- This version works only on a homogeneous cluster of machines, and only the "native" file data representation is supported.
- Shared file pointers are not supported on PVFS and IBM PIOFS file systems because they don't support fcntl file locks, and ROMIO uses that feature to implement shared file pointers.
- On HP machines running HPUX and on NEC SX-4, you need to compile Fortran programs with mpif90 instead of mpif77, because the f77 compilers on these machines don't support 8-byte integers.
- The file-open mode MPI\_MODE\_EXCL does not work on Intel PFS file system, due to a bug in PFS.

# 8 Usage Tips

When using ROMIO with SGI MPI, you may sometimes get an error message from SGI MPI: "MPI has run out of internal datatype entries. Please set the environment variable MPI\_TYPE\_MAX for additional space." If you get this error message, add the following line to your .cshrc file:

```
setenv MPI_TYPE_MAX 65536
```

Use a larger number if you still get the error message.

- If a Fortran program uses a file handle created using ROMIO's C interface, or vice versa, you must use the functions MPI\_File\_c2f or MPI\_File\_f2c (see § 4.12.4 in [4]). Such a situation occurs, for example, if a Fortran program uses an I/O library written in C with MPI-IO calls. Similar functions MPIO\_Request\_f2c and MPIO\_Request\_c2f are also provided.
- For Fortran programs on the Intel Paragon, you may need to provide the complete path to mpif.h in the include statement, e.g.,

```
include '/usr/local/mpich/include/mpif.h' instead of
```

```
include 'mpif.h'
```

This is because the -I option to the Paragon Fortran compiler if77 doesn't work correctly.

It always looks in the default directories first and, therefore, picks up Intel's mpif.h, which is actually the mpif.h of an older version of MPICH.

## 9 Reporting Bugs

If you have trouble, first check the users guide. Then check if there is a list of known bugs and patches on the ROMIO web page at http://www.mcs.anl.gov/romio. Finally, if you still have problems, send a detailed message containing:

- the type of system (often uname -a),
- the output of configure,
- the output of make, and
- any programs or tests

to romio-maint@mcs.anl.gov.

## 10 ROMIO Internals

A key component of ROMIO that enables such a portable MPI-IO implementation is an internal abstract I/O device layer called ADIO [5]. Most users of ROMIO will not need to deal with the ADIO layer at all. However, ADIO is useful to those who want to port ROMIO to some other file system. The ROMIO source code and the ADIO paper [5] will help you get started.

MPI-IO implementation issues are discussed in [6]. All ROMIO-related papers are available online at http://www.mcs.anl.gov/romio.

# 11 Learning MPI-IO

The book *Using MPI-2: Advanced Features of the Message-Passing Interface* [3], published by MIT Press, provides a tutorial introduction to all aspects of MPI-2, including parallel I/O. It has lots of example programs. See http://www.mcs.anl.gov/mpi/usingmpi2 for further information about the book.

# 12 Major Changes in Previous Releases

#### 12.1 Major Changes in Version 1.2.3

- Added explicit control over aggregators for collective operations (see description of cb\_config\_list).
- Added the following working hints: cb\_config\_list, romio\_cb\_read, romio\_cb\_write, romio\_ds\_read. These additional hints have been added but are currently ignored by the implementation: romio\_ds\_write, romio\_no\_indep\_rw.
- Added NTFS ADIO implementation.
- Added testfs ADIO implementation for use in debugging.
- Added delete function to ADIO interface so that file systems that need to use their own delete function may do so (e.g. PVFS).

• Changed version numbering to match version number of MPICH release.

### 12.2 Major Changes in Version 1.0.3

- When used with MPICH 1.2.1, the MPI-IO functions return proper error codes and classes, and the status object is filled in.
- On SGI's XFS file system, ROMIO can use direct I/O even if the user's request does not meet the various restrictions needed to use direct I/O. ROMIO does this by doing part of the request with buffered I/O (until all the restrictions are met) and doing the rest with direct I/O. (This feature hasn't been tested rigorously. Please check for errors.)

By default, ROMIO will use only buffered I/O. Direct I/O can be enabled either by setting the environment variables MPIO\_DIRECT\_READ and/or MPIO\_DIRECT\_WRITE to TRUE, or on a per-file basis by using the info keys direct\_read and direct\_write.

Direct I/O will result in higher performance only if you are accessing a high-bandwidth disk system. Otherwise, buffered I/O is better and is therefore used as the default.

• Miscellaneous bug fixes.

### 12.3 Major Changes in Version 1.0.2

- Implemented the shared file pointer functions and split collective I/O functions. Therefore, the main components of the MPI I/O chapter not yet implemented are file interoperability and error handling.
- Added support for using "direct I/O" on SGI's XFS file system. Direct I/O is an optional feature of XFS in which data is moved directly between the user's buffer and the storage devices, bypassing the file-system cache. This can improve performance significantly on systems with high disk bandwidth. Without high disk bandwidth, regular I/O (that uses the file-system cache) perfoms better. ROMIO, therefore, does not use direct I/O by default. The user can turn on direct I/O (separately for reading and writing) either by using environment variables or by using MPI's hints mechanism (info). To use the environment-variables method, do

setenv MPIO\_DIRECT\_READ TRUE setenv MPIO\_DIRECT\_WRITE TRUE

To use the hints method, the two keys are direct\_read and direct\_write. By default their values are false. To turn on direct I/O, set the values to true. The environment variables have priority over the info keys. In other words, if the environment variables are set to TRUE, direct I/O will be used even if the info keys say false, and vice versa. Note that direct I/O must be turned on separately for reading and writing. The environment-variables method assumes that the environment variables can be read by each process in the MPI job. This is not guaranteed by the MPI Standard, but it works with SGI's MPI and the ch\_shmem device of MPICH.

- Added support (new ADIO device, ad\_pvfs) for the PVFS parallel file system for Linux clusters, developed at Clemson University (see http://www.parl.clemson.edu/pvfs). To use it, you must first install PVFS and then when configuring ROMIO, specify -file\_system=pvfs in addition to any other options to configure. (As usual, you can configure for multiple file systems by using "+"; for example, -file\_system=pvfs+ufs+nfs.) You will need to specify the path to the PVFS include files via the -cflags option to configure, for example, configure -cflags=-I/usr/pvfs/include. You will also need to specify the full path name of the PVFS library. The best way to do this is via the -lib option to MPICH's configure script (assuming you are using ROMIO from within MPICH).
- Uses weak symbols (where available) for building the profiling version, i.e., the PMPI routines. As a result, the size of the library is reduced considerably.
- The Makefiles use *virtual paths* if supported by the make utility. GNU make supports it, for example. This feature allows you to untar the distribution in some directory, say a slow NFS directory, and compile the library (create the .o files) in another directory, say on a faster local disk. For example, if the tar file has been untarred in an NFS directory called /home/thakur/romio, one can compile it in a different directory, say /tmp/thakur, as follows:

cd /tmp/thakur
/home/thakur/romio/configure
make

The .o files will be created in /tmp/thakur; the library will be created in /home/thakur/romio/lib/\$ARCH/libmpio.a. This method works only if the make utility supports *virtual paths*. If the default make utility does not, you can install GNU make which does, and specify it to configure as

/home/thakur/romio/configure -make=/usr/gnu/bin/gmake (or whatever)

- Lots of miscellaneous bug fixes and other enhancements.
- This version is included in MPICH 1.2.0. If you are using MPICH, you need not download ROMIO separately; it gets built as part of MPICH. The previous version of ROMIO is included in LAM, HP MPI, SGI MPI, and NEC MPI. NEC has also implemented the MPI-IO functions missing in ROMIO, and therefore NEC MPI has a complete implementation of MPI-IO.

#### 12.4 Major Changes in Version 1.0.1

- This version is included in MPICH 1.1.1 and HP MPI 1.4.
- Added support for NEC SX-4 and created a new device ad\_sfs for NEC SFS file system.
- New devices ad\_hfs for HP HFS file system and ad\_xfs for SGI XFS file system.
- Users no longer need to prefix the filename with the type of file system; ROMIO determines the file-system type on its own.

- Added support for 64-bit file sizes on IBM PIOFS, SGI XFS, HP HFS, and NEC SFS file systems.
- MPI\_Offset is an 8-byte integer on machines that support 8-byte integers. It is of type long long in C and integer\*8 in Fortran. With a Fortran 90 compiler, you can use either integer\*8 or integer(kind=MPI\_OFFSET\_KIND). If you printf an MPI\_Offset in C, remember to use %11d or %1d as required by your compiler. (See what is used in the test program romio/test/misc.c). On some machines, ROMIO detects at configure time that long long is either not supported by the C compiler or it doesn't work properly. In such cases, configure sets MPI\_Offset to long in C and integer in Fortran. This happens on Intel Paragon, Sun4, and FreeBSD.
- Added support for passing hints to the implementation via the MPI\_Info parameter. ROMIO understands the following hints (keys in MPI\_Info object): cb\_buffer\_size, cb\_nodes, ind\_rd\_buffer\_size, ind\_wr\_buffer\_size (on all but IBM PIOFS), striping\_factor (on PFS and PIOFS), striping\_unit (on PFS and PIOFS), start\_iodevice (on PFS and PIOFS), and pfs\_svr\_buf (on PFS only).

## References

- [1] Rajesh Bordawekar, Juan Miguel del Rosario, and Alok Choudhary. Design and evaluation of primitives for parallel I/O. In *Proceedings of Supercomputing '93*, pages 452–461, Portland, OR, 1993. IEEE Computer Society Press.
- [2] Alok Choudhary, Rajesh Bordawekar, Michael Harry, Rakesh Krishnaiyer, Ravi Ponnusamy, Tarvinder Singh, and Rajeev Thakur. PASSION: parallel and scalable software for inputoutput. Technical Report SCCS-636, ECE Dept., NPAC and CASE Center, Syracuse University, September 1994.
- [3] William Gropp, Ewing Lusk, and Rajeev Thakur. *Using MPI-2: Advanced Features of the Message-Passing Interface*. MIT Press, Cambridge, MA, 1999.
- [4] Message Passing Interface Forum. MPI-2: Extensions to the Message-Passing Interface, July 1997. http://www.mpi-forum.org/docs/docs.html.
- [5] Rajeev Thakur, William Gropp, and Ewing Lusk. An abstract-device interface for implementing portable parallel-I/O interfaces. In *Proceedings of the 6th Symposium on the Frontiers of Massively Parallel Computation*, pages 180–187. IEEE Computer Society Press, October 1996.
- [6] Rajeev Thakur, William Gropp, and Ewing Lusk. On implementing MPI-IO portably and with high performance. In *Proceedings of the 6th Workshop on I/O in Parallel and Distributed Systems*, pages 23–32. ACM Press, May 1999.