
IOR Documentation

Release 3.4.0+dev

IOR

Sep 28, 2022

Contents

1	Introduction	3
2	Install	5
2.1	Building	5
3	First Steps with IOR	7
3.1	Running IOR	7
3.2	Getting Started with IOR	7
3.3	Effect of Page Cache on Benchmarking	9
3.4	Defeating Page Cache	11
3.5	Corollary	13
4	Options	15
4.1	Command line options	15
4.2	Directive Options	17
4.3	Verbosity levels	21
4.4	Incompressible notes	21
5	Scripting	23
6	Version Compatibility	25
7	Frequently Asked Questions	27
8	Doxygen	31
9	Continuous Integration	33
10	Release Process	35
10.1	General release process	35
10.2	Building a release of IOR	35
10.3	Feature freezing for a new release	36
10.4	Creating a new release candidate	36
10.5	Applying patches to a new microrelease	37
10.6	Creating a new release	37
11	Changes in IOR	39
11.1	Version 3.4.0+dev	39

11.2	Version 3.3.0	39
11.3	Version 3.2.1	40
11.4	Version 3.2.0	40
11.5	Version 3.0.2	41
11.6	Version 3.0.0	41
11.7	Version 2.10.3	42
11.8	Version 2.10.2	42
11.9	Version 2.10.1	42
11.10	Version 2.10.0.1	43
11.11	Version 2.9.6.1	43
11.12	Version 2.9.5	43
11.13	Version 2.9.4	43
11.14	Version 2.9.3	43
11.15	Version 2.9.2	44
11.16	Version 2.9.1	44
11.17	Version 2.9.0	44

IOR is a parallel IO benchmark that can be used to test the performance of parallel storage systems using various interfaces and access patterns. The IOR repository also includes the mdtest benchmark which specifically tests the peak metadata rates of storage systems under different directory structures. Both benchmarks use a common parallel I/O abstraction backend and rely on MPI for synchronization.

This documentation consists of two parts.

User documentation includes installation instructions (*Install*), a beginner's tutorial (*First Steps with IOR*), and information about IOR's runtime *Options*.

Developer documentation consists of code documentation generated with Doxygen and some notes about the contiguous integration with Travis.

Many aspects of both IOR/mdtest user and developer documentation are incomplete, and contributors are encouraged to comment the code directly or expand upon this documentation.

CHAPTER 1

Introduction

IOR is a parallel IO benchmark that can be used to test the performance of parallel storage systems using various interfaces and access patterns. The IOR repository also includes the mdtest benchmark which specifically tests the peak metadata rates of storage systems under different directory structures. Both benchmarks use a common parallel I/O abstraction backend and rely on MPI for synchronization.

This documentation consists of two parts.

User documentation includes installation instructions (*Install*), a beginner's tutorial (*First Steps with IOR*), and information about IOR's runtime *Options*.

Developer documentation consists of code documentation generated with Doxygen and some notes about the contiguous integration with Travis.

Many aspects of both IOR/mdtest user and developer documentation are incomplete, and contributors are encouraged to comment the code directly or expand upon this documentation.

2.1 Building

0. If `configure` is missing from the top level directory, you probably retrieved this code directly from the repository. Run `./bootstrap`.

If your versions of the autotools are not new enough to run this script, download and official tarball in which the `configure` script is already provided.

1. Run `./configure`
See `./configure --help` for configuration options.
2. Run `make`
3. Optionally, run `make install`. The installation prefix can be changed as an option to the `configure` script.

First Steps with IOR

This is a short tutorial for the basic usage of IOR and some tips on how to use IOR to handle caching effects as these are very likely to affect your measurements.

3.1 Running IOR

There are two ways of running IOR:

- 1) Command line with arguments – executable followed by command line options.

```
$ ./IOR -w -r -o filename
```

This performs a write and a read to the file ‘filename’.

- 2) Command line with scripts – any arguments on the command line will establish the default for the test run, but a script may be used in conjunction with this for varying specific tests during an execution of the code. Only arguments before the script will be used!

```
$ ./IOR -W -f script
```

This defaults all tests in ‘script’ to use write data checking.

In this tutorial the first one is used as it is much easier to toy around with an get to know IOR. The second option thought is much more useful to safe benchmark setups to rerun later or to test many different cases.

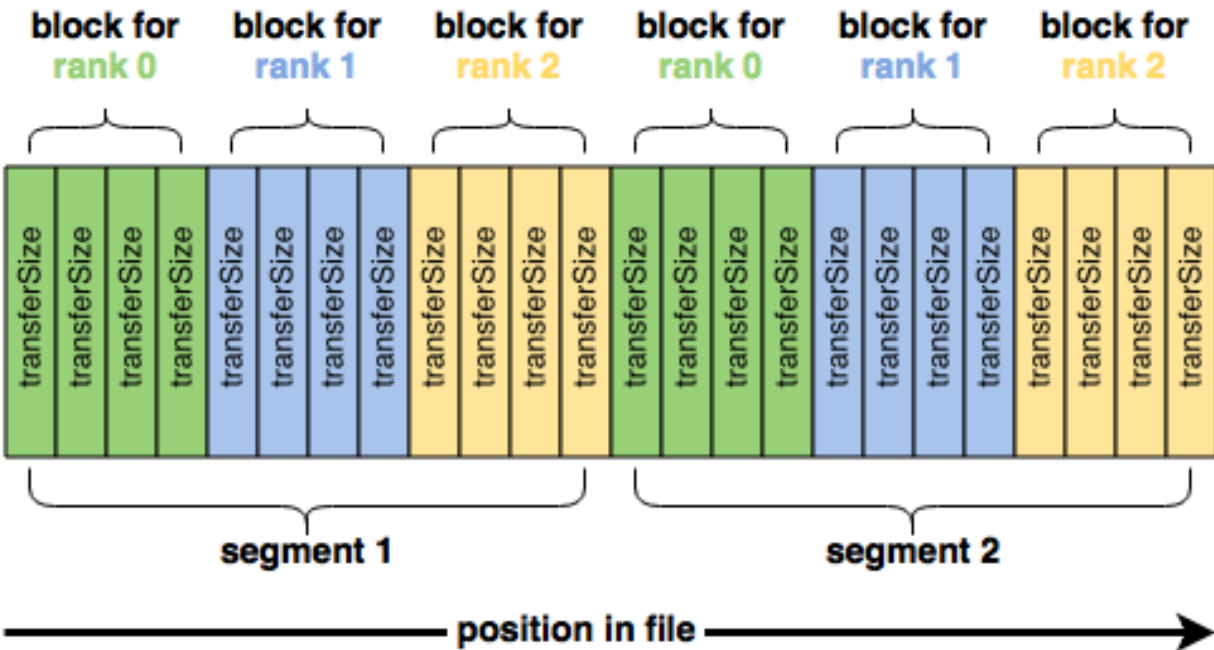
3.2 Getting Started with IOR

IOR writes data sequentially with the following parameters:

- `blockSize (-b)`
- `transferSize (-t)`

- segmentCount (-s)
- numTasks (-n)

which are best illustrated with a diagram:



These four parameters are all you need to get started with IOR. However, naively running IOR usually gives disappointing results. For example, if we run a four-node IOR test that writes a total of 16 GiB:

```
$ mpirun -n 64 ./ior -t 1m -b 16m -s 16
...
access bw(MiB/s) block(KiB) xfer(KiB) open(s) wr/rd(s) close(s) total(s) iter
-----
write 427.36 16384 1024.00 0.107961 38.34 32.48 38.34 2
read 239.08 16384 1024.00 0.005789 68.53 65.53 68.53 2
remove - - - - - - 0.534400 2
```

we can only get a couple hundred megabytes per second out of a Lustre file system that should be capable of a lot more.

Switching from writing to a single-shared file to one file per process using the `-F` (`filePerProcess=1`) option changes the performance dramatically:

```
$ mpirun -n 64 ./ior -t 1m -b 16m -s 16 -F
...
access bw(MiB/s) block(KiB) xfer(KiB) open(s) wr/rd(s) close(s) total(s) iter
-----
write 33645 16384 1024.00 0.007693 0.486249 0.195494 0.486972 1
read 149473 16384 1024.00 0.004936 0.108627 0.016479 0.109612 1
remove - - - - - - 6.08 1
```

This is in large part because letting each MPI process work on its own file cuts out any contention that would arise because of file locking.

However, the performance difference between our naive test and the file-per-process test is a bit extreme. In fact, the only way that 146 GB/sec read rate could be achievable on Lustre is if each of the four compute nodes had over 45 GB/sec of network bandwidth to Lustre—that is, a 400 Gbit link on every compute and storage node.

3.3 Effect of Page Cache on Benchmarking

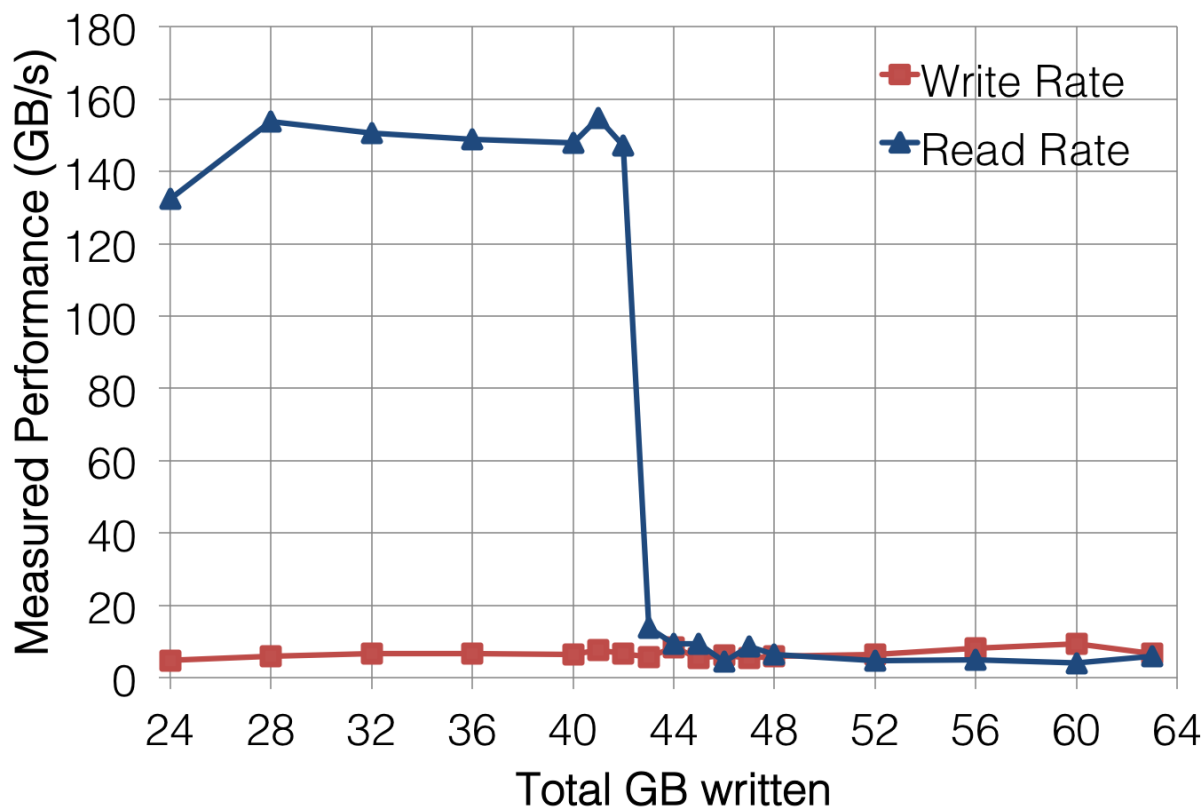
What's really happening is that the data being read by IOR isn't actually coming from Lustre; rather, files' contents are already cached, and IOR is able to read them directly out of each compute node's DRAM. The data wound up getting cached during the write phase of IOR as a result of Linux (and Lustre) using a write-back cache to buffer I/O, so that instead of IOR writing and reading data directly to Lustre, it's actually mostly talking to the memory on each compute node.

To be more specific, although each IOR process thinks it is writing to a file on Lustre and then reading back the contents of that file from Lustre, it is actually

- 1) writing data to a copy of the file that is cached in memory. If there is no copy of the file cached in memory before this write, the parts being modified are loaded into memory first.
- 2) those parts of the file in memory (called "pages") that are now different from what's on Lustre are marked as being "dirty"
- 3) the write() call completes and IOR continues on, even though the written data still hasn't been committed to Lustre
- 4) independent of IOR, the OS kernel continually scans the file cache for files who have been updated in memory but not on Lustre ("dirty pages"), and then commits the cached modifications to Lustre
- 5) dirty pages are declared non-dirty since they are now in sync with what's on disk, but they remain in memory

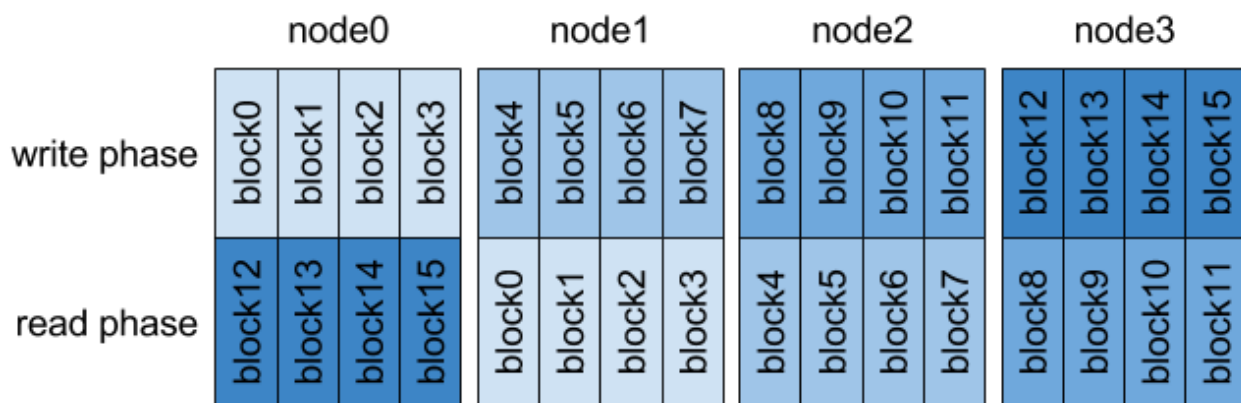
Then when the read phase of IOR follows the write phase, IOR is able to just retrieve the file's contents from memory instead of having to communicate with Lustre over the network.

There are a couple of ways to measure the read performance of the underlying Lustre file system. The most crude way is to simply write more data than will fit into the total page cache so that by the time the write phase has completed, the beginning of the file has already been evicted from cache. For example, increasing the number of segments (`-s`) to write more data reveals the point at which the nodes' page cache on my test system runs over very clearly:



However, this can make running IOR on systems with a lot of on-node memory take forever.

A better option would be to get the MPI processes on each node to only read data that they didn't write. For example, on a four-process-per-node test, shifting the mapping of MPI processes to blocks by four makes each node N read the data written by node N-1.



Since page cache is not shared between compute nodes, shifting tasks this way ensures that each MPI process is reading data it did not write.

IOR provides the `-C` option (`reorderTasks`) to do this, and it forces each MPI process to read the data written by its neighboring node. Running IOR with this option gives much more credible read performance:

```
$ mpirun -n 64 ./ior -t 1m -b 16m -s 16 -F -C
...
```

(continues on next page)

(continued from previous page)

access	bw (MiB/s)	block (KiB)	xfer (KiB)	open (s)	wr/rd (s)	close (s)	total (s)	iter
write	41326	16384	1024.00	0.005756	0.395859	0.095360	0.396453	0
read	3310.00	16384	1024.00	0.011786	4.95	4.20	4.95	1
remove	-	-	-	-	-	-	0.237291	1

But now it should seem obvious that the write performance is also ridiculously high. And again, this is due to the page cache, which signals to IOR that writes are complete when they have been committed to memory rather than the underlying Lustre file system.

To work around the effects of the page cache on write performance, we can issue an `fsync()` call immediately after all of the write(s) return to force the dirty pages we just wrote to flush out to Lustre. Including the time it takes for `fsync()` to finish gives us a measure of how long it takes for our data to write to the page cache and for the page cache to write back to Lustre.

IOR provides another convenient option, `-e (fsync)`, to do just this. And, once again, using this option changes our performance measurement quite a bit:

```
$ mpirun -n 64 ./ior -t lm -b 16m -s 16 -F -C -e
...
access bw (MiB/s) block (KiB) xfer (KiB) open (s) wr/rd (s) close (s) total (s) iter
-----
write 2937.89 16384 1024.00 0.011841 5.56 4.93 5.58 0
read 2712.55 16384 1024.00 0.005214 6.04 5.08 6.04 3
remove - - - - - - - 0.037706 0
```

and we finally have a believable bandwidth measurement for our file system.

3.4 Defeating Page Cache

Since IOR is specifically designed to benchmark I/O, it provides these options that make it as easy as possible to ensure that you are actually measuring the performance of your file system and not your compute nodes' memory. That being said, the I/O patterns it generates are designed to demonstrate peak performance, not reflect what a real application might be trying to do, and as a result, there are plenty of cases where measuring I/O performance with IOR is not always the best choice. There are several ways in which we can get clever and defeat page cache in a more general sense to get meaningful performance numbers.

When measuring write performance, bypassing page cache is actually quite simple; opening a file with the `O_DIRECT` flag going directly to disk. In addition, the `fsync()` call can be inserted into applications, as is done with IOR's `-e` option.

Measuring read performance is a lot trickier. If you are fortunate enough to have root access on a test system, you can force the Linux kernel to empty out its page cache by doing

```
# echo 1 > /proc/sys/vm/drop_caches
```

and in fact, this is often good practice before running any benchmark (e.g., Linpack) because it ensures that you aren't losing performance to the kernel trying to evict pages as your benchmark application starts allocating memory for its own use.

Unfortunately, many of us do not have root on our systems, so we have to get even more clever. As it turns out, there is a way to pass a hint to the kernel that a file is no longer needed in page cache:

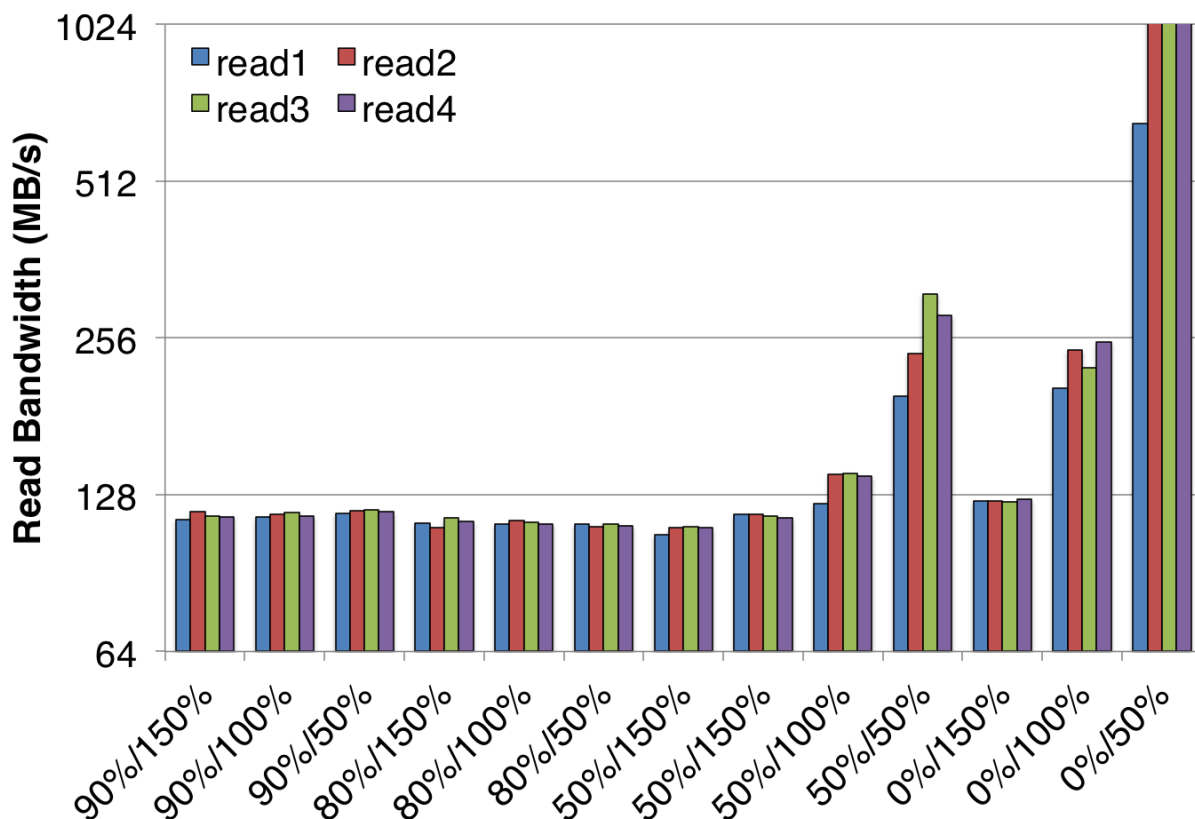
```

#define _XOPEN_SOURCE 600
#include <unistd.h>
#include <fcntl.h>
int main(int argc, char *argv[]) {
    int fd;
    fd = open(argv[1], O_RDONLY);
    fdatasync(fd);
    posix_fadvise(fd, 0, 0, POSIX_FADV_DONTNEED);
    close(fd);
    return 0;
}

```

The effect of passing `POSIX_FADV_DONTNEED` using `posix_fadvise()` is usually that all pages belonging to that file are evicted from page cache in Linux. However, this is just a hint—not a guarantee—and the kernel evicts these pages asynchronously, so it may take a second or two for pages to actually leave page cache. Fortunately, Linux also provides a way to probe pages in a file to see if they are resident in memory.

Finally, it's often easiest to just limit the amount of memory available for page cache. Because application memory always takes precedence over cache memory, simply allocating most of the memory on a node will force most of the cached pages to be evicted. Newer versions of IOR provide the `memoryPerNode` option that do just that, and the effects are what one would expect:



The above diagram shows the measured bandwidth from a single node with 128 GiB of total DRAM. The first percent on each x-label is the amount of this 128 GiB that was reserved by the benchmark as application memory, and the second percent is the total write volume. For example, the “50%/150%” data points correspond to 50% of the node memory (64 GiB) being allocated for the application, and a total of 192 GiB of data being read.

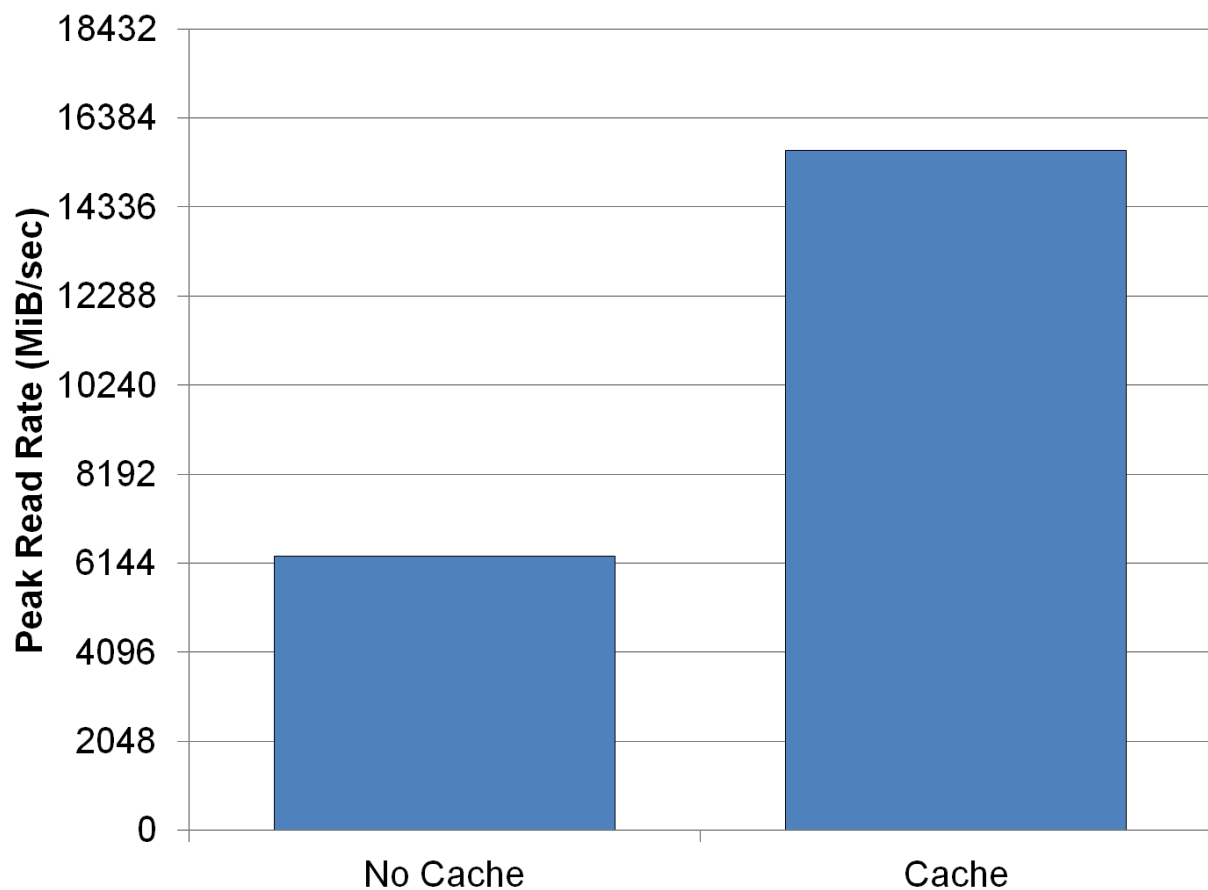
This benchmark was run on a single spinning disk which is not capable of more than 130 MB/sec, so the conditions

that showed performance higher than this were benefiting from some pages being served from cache. And this makes perfect sense given that the anomalously high performance measurements were obtained when there was plenty of memory to cache relative to the amount of data being read.

3.5 Corollary

Measuring I/O performance is a bit trickier than CPU performance in large part due to the effects of page caching. That being said, page cache exists for a reason, and there are many cases where an application's I/O performance really is best represented by a benchmark that heavily utilizes cache.

For example, the BLAST bioinformatics application re-reads all of its input data twice; the first time initializes data structures, and the second time fills them up. Because the first read caches each page and allows the second read to come out of cache rather than the file system, running this I/O pattern with page cache disabled causes it to be about 2x slower:



Thus, letting the page cache do its thing is often the most realistic way to benchmark with realistic application I/O patterns. Once you know how page cache might be affecting your measurements, you stand a good chance of being able to reason about what the most meaningful performance metrics are.

IOR provides many options, in fact there are now more than there are one letter flags in the alphabet. For this and to run IOR by a config script, there are some options which are only available via directives. When both script and command line options are in use, command line options set in front of -f are the defaults which may be overridden by the script. Directives can also be set from the command line via “-O” option. In combination with a script they behave like the normal command line options. But directives and normal parameters override each other, so the last one executed.

4.1 Command line options

These options are to be used on the command line (e.g., ./ior -a POSIX -b 4K).

-a S	api – API for I/O [POSIX MPIIO HDF5 HDFS S3 S3_EMU NCMPI RADOS]
-A N	refNum – user reference number to include in long summary
-b N	blockSize – contiguous bytes to write per task (e.g.: 8, 4k, 2m, 1g)
-c	collective – collective I/O
-C	reorderTasksConstant – changes task ordering to n+1 ordering for readback
-d N	interTestDelay – delay between reps in seconds
-D N	deadlineForStonewalling – seconds before stopping write or read phase
-e	fsync – perform fsync upon POSIX write close
-E	useExistingTestFile – do not remove test file before write access
-f S	scriptFile – test script name
-F	filePerProc – file-per-process
-g	intraTestBarriers – use barriers between open, write/read, and close

-G N	setTimeStampSignature – set value for time stamp signature
-h	showHelp – displays options and help
-H	showHints – show hints
-i N	repetitions – number of repetitions of test
-I	individualDataSets – datasets not shared by all procs [not working]
-j N	outlierThreshold – warn on outlier N seconds from mean
-J N	setAlignment – HDF5 alignment in bytes (e.g.: 8, 4k, 2m, 1g)
-k	keepFile – don't remove the test file(s) on program exit
-K	keepFileWithError – keep error-filled file(s) after data-checking
-l	data packet type– type of packet that will be created [offsetlincompressibletimestamploilt]
-m	multiFile – use number of reps (-i) for multiple file count
-M N	memoryPerNode – hog memory on the node (e.g.: 2g, 75%)
-n	noFill – no fill in HDF5 file creation
-N N	numTasks – number of tasks that should participate in the test
-o S	testFile – full name for test
-O S	string of IOR directives (e.g. -O checkRead=1,GPUid=2)
-p	preallocate – preallocate file size
-P	useSharedFilePointer – use shared file pointer [not working]
-q	quitOnError – during file error-checking, abort on error
-Q N	taskPerNodeOffset for read tests use with -C & -Z options (-C constant N, -Z at least N) [!HDF5]
-r	readFile – read existing file
-R	checkRead – check read after read
-s N	segmentCount – number of segments
-S	useStridedDatatype – put strided access into datatype [not working]
-t N	transferSize – size of transfer in bytes (e.g.: 8, 4k, 2m, 1g)
-T N	maxTimeDuration – max time in minutes to run tests
-u	uniqueDir – use unique directory name for each file-per-process
-U S	hintsFileName – full name for hints file
-v	verbose – output information (repeating flag increases level)
-V	useFileView – use MPI_File_set_view
-w	writeFile – write file
-W	checkWrite – check read after write
-x	singleXferAttempt – do not retry transfer if incomplete
-X N	reorderTasksRandomSeed – random seed for -Z option
-Y	fsyncPerWrite – perform fsync after each POSIX write

-z	randomOffset – access is to random, not sequential, offsets within a file
-Z	reorderTasksRandom – changes task ordering to random ordering for readback

- S is a string, N is an integer number.
- For transfer and block sizes, the case-insensitive K, M, and G suffices are recognized. I.e., ‘4k’ or ‘4K’ is accepted as 4096.

Various options are only valid for specific modules, you can see details when running \$./ior -h These options are typically prefixed with the module name, an example is: -posix.odirect

4.2 Directive Options

For all true/false options below, [1]=true, [0]=false. All options are case-insensitive.

4.2.1 GENERAL

- refNum - user supplied reference number, included in long summary (default: 0)
- api - must be set to one of POSIX, MPIIO, HDF5, HDFS, S3, S3_EMCM, NCMPI, IME, MMAP, or RAODS depending on test (default: POSIX)
- testFile - name of the output file [testFile]. With filePerProc set, the tasks can round robin across multiple file names via -o S@S@S. If only a single file name is specified in this case, IOR appends the MPI rank to the end of each file generated (e.g., testFile.00000059) (default: testFile)
- hintsFileName - name of the hints file (default: none)
- repetitions - number of times to run each test (default: 1)
- multiFile - creates multiple files for single-shared-file or file-per-process modes for each iteration (default: 0)
- reorderTasksConstant - reorders tasks by a constant node offset for writing/reading neighbor’s data from different nodes (default: 0)
- taskPerNodeOffset - for read tests. Use with -C and -Z options. With reorderTasks, constant N. With reordertasksrandom, >= N (default: 1)
- reorderTasksRandom - reorders tasks to random ordering for read tests (default: 0)
- **reorderTasksRandomSeed - random seed for reordertasksrandom option. (default: 0)**
 - When > 0, use the same seed for all iterations
 - When < 0, different seed for each iteration
- quitOnError - upon error encountered on checkWrite or checkRead, display current error and then stop execution. Otherwise, count errors and continue (default: 0)
- numTasks - number of tasks that should participate in the test. 0 denotes all tasks. (default: 0)
- interTestDelay - time (in seconds) to delay before beginning a write or read phase in a series of tests This does not delay before check-write or check-read phases. (default: 0)
- outlierThreshold - gives warning if any task is more than this number of seconds from the mean of all participating tasks. The warning includes the offending task, its timers (start, elapsed create, elapsed transfer,

elapsed close, end), and the mean and standard deviation for all tasks. When zero, disable this feature. (default: 0)

- `intraTestBarriers` - use barrier between open, write/read, and close phases (default: 0)
- `uniqueDir` - create and use unique directory for each file-per-process (default: 0)
- `writeFile` - write file(s), first deleting any existing file. The defaults for `writeFile` and `readFile` are set such that if there is not at least one of `-w`, `-r`, `-W`, or `-R`, `-w` and `-r` are enabled. If either `writeFile` or `readFile` are explicitly enabled, though, its complement is *not* also implicitly enabled.
- `readFile` - reads existing file(s) as specified by the `testFile` option. The defaults for `writeFile` and `readFile` are set such that if there is not at least one of `-w`, `-r`, `-W`, or `-R`, `-w` and `-r` are enabled. If either `writeFile` or `readFile` are explicitly enabled, though, its complement is *not* also implicitly enabled.
- `filePerProc` - have each MPI process perform I/O to a unique file (default: 0)
- `checkWrite` - read data back and check for errors against known pattern. Can be used independently of `writeFile`. Data checking is not timed and does not affect other performance timings. All errors detected are tallied and returned as the program exit code unless `quitOnError` is set. (default: 0)
- `checkRead` - re-read data and check for errors between reads. Can be used independently of `readFile`. Data checking is not timed and does not affect other performance timings. All errors detected are tallied and returned as the program exit code unless `quitOnError` is set. (default: 0)
- `keepFile` - do not remove test file(s) on program exit (default: 0)
- `keepFileWithError` - do not delete any files containing errors if detected during read-check or write-check phases. (default: 0)
- `useExistingTestFile` - do not remove test file(s) before write phase (default: 0)
- `segmentCount` - number of segments in file, where a segment is a contiguous chunk of data accessed by multiple clients each writing/reading their own contiguous data (blocks). The exact semantics of segments depend on the API used; for example, HDF5 repeats the pattern of an entire shared dataset. (default: 1)
- `blockSize` - size (in bytes) of a contiguous chunk of data accessed by a single client. It is comprised of one or more transfers (default: 1048576)
- `transferSize` - size (in bytes) of a single data buffer to be transferred in a single I/O call (default: 262144)
- `verbose` - output more information about what IOR is doing. Can be set to levels 0-5; repeating the `-v` flag will increase verbosity level. (default: 0)
- `setTimeStampSignature` - Value to use for the time stamp signature. Used to rerun tests with the exact data pattern by setting data signature to contain positive integer value as timestamp to be written in data file; if set to 0, is disabled (default: 0)
- `showHelp` - display options and help (default: 0)
- `storeFileOffset` - use file offset as stored signature when writing file. This will affect performance measurements (default: 0)
- `memoryPerNode` - allocate memory on each node to simulate real application memory usage or restrict page cache size. Accepts a percentage of node memory (e.g. 50%) on systems that support `sysconf(_SC_PHYS_PAGES)` or a size. Allocation will be split between tasks that share the node. (default: 0)
- `memoryPerTask` - allocate specified amount of memory (in bytes) per task to simulate real application memory usage. (default: 0)
- `maxTimeDuration` - max time (in minutes) to run all tests. Any current read/write phase is not interrupted; only future I/O phases are cancelled once this time is exceeded. Value of zero unsets disables. (default: 0)

- `deadlineForStonewalling` - seconds before stopping write or read phase. Used for measuring the amount of data moved in a fixed time. After the barrier, each task starts its own timer, begins moving data, and then stops moving data at a pre-arranged time. Instead of measuring the amount of time to move a fixed amount of data, this option measures the amount of data moved in a fixed amount of time. The objective is to prevent straggling tasks from skewing the performance. This option is incompatible with read-check and write-check modes. Value of zero unsets this option. (default: 0)
- `randomOffset` - randomize access offsets within test file(s). Currently incompatible with `checkRead`, `storeFileOffset`, `MPIIO` collective and `useFileView`, and HDF5 and NCMPI APIs. (default: 0)
- `summaryAlways` - Always print the long summary for each test even if the job is interrupted. (default: 0)

4.2.2 POSIX-ONLY

- `useO_DIRECT` - use direct I/O for POSIX, bypassing I/O buffers (default: 0)
- `singleXferAttempt` - do not continue to retry transfer entire buffer until it is transferred. When performing a `write()` or `read()` in POSIX, there is no guarantee that the entire requested size of the buffer will be transferred; this flag keeps the retrying a single transfer until it completes or returns an error (default: 0)
- `fsyncPerWrite` - perform `fsync` after each POSIX write (default: 0)
- `fsync` - perform `fsync` after POSIX file close (default: 0)

4.2.3 MPIIO-ONLY

- `preallocate` - preallocate the entire file before writing (default: 0)
- `useFileView` - use an MPI datatype for setting the file view option to use individual file pointers. Default IOR uses explicit file pointers. (default: 0)
- `useSharedFilePointer` - use a shared file pointer. Default IOR uses explicit file pointers. (default: 0)
- `useStridedDatatype` - create a datatype (max=2GB) for strided access; akin to `MULTIBLOCK_REGION_SIZE` (default: 0)

4.2.4 HDF5-ONLY

- `individualDataSets` - within a single file, each task will access its own dataset. Default IOR creates a dataset the size of `numTasks * blockSize` to be accessed by all tasks (default: 0)
- `noFill` - do not pre-fill data in HDF5 file creation (default: 0)
- `setAlignment` - set the HDF5 alignment in bytes (e.g.: 8, 4k, 2m, 1g) (default: 1)
- `hdf5.collectiveMetadata` - enable HDF5 collective metadata (available since HDF5-1.10.0)

4.2.5 MPIIO-, HDF5-, AND NCMPI-ONLY

- `collective` - uses collective operations for access (default: 0)
- `showHints` - show hint/value pairs attached to open file (default: 0)

4.2.6 LUSTRE-SPECIFIC

- POSIX-ONLY:
 - `--posix.lustre.stripecount` - set the Lustre stripe count for the test file(s) (default: 0)
 - `--posix.lustre.stripesize` - set the Lustre stripe size for the test file(s) (default: 0)
 - `--posix.lustre.startost` - set the starting OST for the test file(s) (default: -1)
 - `--posix.lustre.ignorelocks` - disable Lustre range locking (default: 0)
- MPIIO-, HDF5-, AND NCMPI-ONLY:
 - ROMIO-based IO (see [here](#)):
 - * requires setting the environment variable `ROMIO_FSTYPE_FORCE=LUSTRE`: (or similar for specific MPIs) to enable ROMIO's Lustre ADIO
 - * `IOR_HINT__MPI__striping_factor` - set the Lustre stripe count for the test file(s) (default: -1)
 - * `IOR_HINT__MPI__striping_unit` - set the Lustre stripe size for the test file(s) (default: 0)
 - * `IOR_HINT__MPI__romio_lustre_start_iodevice` - set the starting OST for the test file(s) (default: -1)
 - OMPIO-based IO (see [here](#))
 - * not setting either of the environment variables below causes a fatal *Floating point exception: Integer divide-by-zero* error
 - * execution with either of the environment variables causes this message *ior: setstripe error for 'test-file': stripe already set* which can safely be ignored as OMPIO tries to modify the stripe settings twice although the first time succeeds
 - * `OMPI_MCA_fs_lustre_stripe_width / IOR_HINT__MPI__stripe_width` - set the Lustre stripe count for the test file(s) (default: 0)
 - * `OMPI_MCA_fs_lustre_stripe_size / IOR_HINT__MPI__stripe_size` - set the Lustre stripe size for the test file(s) (default: 0)

4.2.7 GPFS-SPECIFIC (POSIX-ONLY)

- `--posix.gpfs.hintaccess` - use `gpfs_fcntl` hints to pre-declare accesses
- `--posix.gpfs.releasetoken` - immediately after opening or creating file, release all locks. Might help mitigate lock-revocation traffic when many processes write/read to same file.
- `--posix.gpfs.finegrainwritesharing` - This hint optimizes the performance of small strided writes to a shared file from a parallel application
- `--posix.gpfs.finegrainreadsharing` - This hint optimizes the performance of small strided reads from a shared file from a parallel application

4.2.8 BeeGFS-SPECIFIC (POSIX-ONLY):

- `--posix.beegfs.NumTargets` - set the number of storage targets to use
- `--posix.beegfs.ChunkSize` - set the striping chunk size. Must be a power of two, and greater than 64kiB, (e.g.: 256k, 1M, ...)

4.3 Verbosity levels

The verbosity of output for IOR can be set with `-v`. Increasing the number of `-v` instances on a command line sets the verbosity higher.

Here is an overview of the information shown for different verbosity levels:

Level	Behavior
0	default; only bare essentials shown
1	max clock deviation, participating tasks, free space, access pattern, commence/verify access notification with time
2	rank/hostname, machine name, timer used, individual repetition performance results, timestamp used for data signature
3	full test details, transfer block/offset compared, individual data checking errors, environment variables, task writing/reading file name, all test operation times
4	task id and offset for each transfer
5	each 8-byte data signature comparison (WARNING: more data to STDOUT than stored in file, use carefully)

4.4 Incompressible notes

Please note that incompressibility is a factor of how large a block compression algorithm uses. The incompressible buffer is filled only once before write times, so if the compression algorithm takes in blocks larger than the transfer size, there will be compression. Below are some baselines for zip, gzip, and bzip.

- 1) zip: For zipped files, a transfer size of 1k is sufficient.
- 2) gzip: For gzipped files, a transfer size of 1k is sufficient.
- 3) bzip2: For bzipped files a transfer size of 1k is insufficient (~50% compressed). To avoid compression a transfer size of greater than the bzip block size is required (default = 900KB). I suggest a transfer size of greater than 1MB to avoid bzip2 compression.

Be aware of the block size your compression algorithm will look at, and adjust the transfer size accordingly.

IOR can use an input script with the command line using the `-f` option. **Any options on the command line set before the ‘-f’ option is given will be considered the default settings for running the script.** For example,

```
mpirun ./ior -W -f script
```

will run all tests in the script with an implicit `-W`. The script itself can override these settings and may be set to run many different tests of IOR under a single execution, and it is important to note that **any command-line options specified after “-f” will not be applied to the runs dictated by the script.** For example,

```
mpirun ./ior -f script -W
```

will *not* run any tests with the implicit `-W` since that argument does not get applied until after the `-f` option (and its constituent runs) are complete.

Input scripts are specified using the long-form option names that correspond to each command-line option. In addition to long-form options,

- IOR START and IOR END mark the beginning and end of the script
- RUN dispatches the test using all of the options specified before it
- All previous set parameter stay set for the next test. They are not reset to the default! For default the must be rest manually.
- White space is ignored in script, as are comments starting with #.
- Not all test parameters need be set.

An example of a script:

```
IOR START
  api=[POSIX|MPIO|HDF5|HDFS|S3|S3_EMC|NCMPI|RADOS]
  testFile=testFile
  hintsFileName=hintsFile
  repetitions=8
  multiFile=0
```

(continues on next page)

(continued from previous page)

```
interTestDelay=5
readFile=1
writeFile=1
filePerProc=0
checkWrite=0
checkRead=0
keepFile=1
quitOnError=0
segmentCount=1
blockSize=32k
outlierThreshold=0
setAlignment=1
transferSize=32
singleXferAttempt=0
individualDataSets=0
verbose=0
numTasks=32
collective=1
preallocate=0
useFileView=0
keepFileWithError=0
setTimeStampSignature=0
useSharedFilePointer=0
useStridedDatatype=0
uniqueDir=0
fsync=0
storeFileOffset=0
maxTimeDuration=60
deadlineForStonewalling=0
useExistingTestFile=0
useO_DIRECT=0
showHints=0
showHelp=0
RUN
# additional tests are optional
    transferSize=64
    blockSize=64k
    segmentcount=2
RUN
    transferSize=4K
    blockSize=1M
    segmentcount=1024
RUN
IOR STOP
```

Version Compatibility

IOR has a long history and only IOR version 3 is currently supported. However, there are many forks of IOR based on earlier versions, and the following incompatibilities are known to exist between major versions.

- 1) IOR version 1 (c. 1996-2002) and IOR version 2 (c. 2003-present) are incompatible. Input decks from one will not work on the other. As version 1 is not included in this release, this shouldn't be case for concern. All subsequent compatibility issues are for IOR version 2.
- 2) IOR versions prior to release 2.8 provided data size and rates in powers of two. E.g., 1 MB/sec referred to 1,048,576 bytes per second. With the IOR release 2.8 and later versions, MB is now defined as 1,000,000 bytes and MiB is 1,048,576 bytes.
- 3) In IOR versions 2.5.3 to 2.8.7, IOR could be run without any command line options. This assumed that if both write and read options (-w -r) were omitted, the run with them both set as default. Later, it became clear that in certain cases (data checking, e.g.) this caused difficulties. In IOR versions 2.8.8 and later, if not one of the -w -r -W or -R options is set, then -w and -r are set implicitly.
- 4) IOR version 3 (Jan 2012-present) has changed the output of IOR somewhat, and the "testNum" option was renamed "refNum".

Frequently Asked Questions

HOW DO I PERFORM MULTIPLE DATA CHECKS ON AN EXISTING FILE?

Use this command line: `IOR -k -E -W -i 5 -o file`

-k keeps the file after the access rather than deleting it -E uses the existing file rather than truncating it first -W performs the writecheck -i number of iterations of checking -o filename

On versions of IOR prior to 2.8.8, you need the -r flag also, otherwise you'll first overwrite the existing file. (In earlier versions, omitting -w and -r implied using both. This semantic has been subsequently altered to be omitting -w, -r, -W, and -R implied using both -w and -r.)

If you're running new tests to create a file and want repeat data checking on this file multiple times, there is an undocumented option for this. It's -O multiReRead=1, and you'd need to have an IOR version compiled with the USE_UNDOC_OPT=1 (in iordef.h). The command line would look like this:

`IOR -k -E -w -W -i 5 -o file -O multiReRead=1`

For the first iteration, the file would be written (w/o data checking). Then for any additional iterations (four, in this example) the file would be reread for whatever data checking option is used.

HOW DOES IOR CALCULATE PERFORMANCE?

IOR performs get a time stamp START, then has all participating tasks open a shared or independent file, transfer data, close the file(s), and then get a STOP time. A `stat()` or `MPI_File_get_size()` is performed on the file(s) and compared against the aggregate amount of data transferred. If this value does not match, a warning is issued and the amount of data transferred as calculated from `write()`, e.g., return codes is used. The calculated bandwidth is the amount of data transferred divided by the elapsed STOP-minus-START time.

IOR also gets time stamps to report the open, transfer, and close times. Each of these times is based on the earliest start time for any task and the latest stop time for any task. Without using barriers between these operations (-g), the sum of the open, transfer, and close times may not equal the elapsed time from the first open to the last close.

HOW DO I ACCESS MULTIPLE FILE SYSTEMS IN IOR?

It is possible when using the `filePerProc` option to have tasks round-robin across multiple file names. Rather than use a single file name '-o file', additional names '-o file1@file2@file3' may be used. In this

case, a file per process would have three different file names (which may be full path names) to access. The '@' delimiter is arbitrary, and may be set in the FILENAME_DELIMITER definition in iordef.h.

Note that this option of multiple filenames only works with the filePerProc -F option. This will not work for shared files.

HOW DO I BALANCE LOAD ACROSS MULTIPLE FILE SYSTEMS?

As for the balancing of files per file system where different file systems offer different performance, additional instances of the same destination path can generally achieve good balance.

For example, with FS1 getting 50% better performance than FS2, set the '-o' flag such that there are additional instances of the FS1 directory. In this case, '-o FS1/file@FS1/file@FS1/file@FS2/file@FS2/file' should adjust for the performance difference and balance accordingly.

HOW DO I USE STONEWALLING?

To use stonewalling (-D), it's generally best to separate write testing from read testing. Start with writing a file with '-D 0' (stonewalling disabled) to determine how long the file takes to be written. If it takes 10 seconds for the data transfer, run again with a shorter duration, '-D 7' e.g., to stop before the file would be completed without stonewalling. For reading, it's best to create a full file (not an incompletely written file from a stonewalling run) and then run with stonewalling set on this preexisting file. If a write and read test are performed in the same run with stonewalling, it's likely that the read will encounter an error upon hitting the EOF. Separating the runs can correct for this. E.g.,

```
IOR -w -k -o file -D 10 # write and keep file, stonewall after 10 seconds
IOR -r -E -o file -D 7 # read existing file, stonewall after 7 seconds
```

Also, when running multiple iterations of a read-only stonewall test, it may be necessary to set the -D value high enough so that each iteration is not reading from cache. Otherwise, in some cases, the first iteration may show 100 MB/s, the next 200 MB/s, the third 300 MB/s. Each of these tests is actually reading the same amount from disk in the allotted time, but they are also reading the cached data from the previous test each time to get the increased performance. Setting -D high enough so that the cache is overfilled will prevent this.

HOW DO I BYPASS CACHING WHEN READING BACK A FILE I'VE JUST WRITTEN?

One issue with testing file systems is handling cached data. When a file is written, that data may be stored locally on the node writing the file. When the same node attempts to read the data back from the file system either for performance or data integrity checking, it may be reading from its own cache rather than from the file system.

The reorderTasksConstant '-C' option attempts to address this by having a different node read back data than wrote it. For example, node N writes the data to file, node N+1 reads back the data for read performance, node N+2 reads back the data for write data checking, and node N+3 reads the data for read data checking, comparing this with the reread data from node N+4. The objective is to make sure on file access that the data is not being read from cached data.

Node 0: writes data Node 1: reads data Node 2: reads written data for write checking Node 3:
reads written data for read checking Node 4: reads written data for read checking, comparing
with Node 3

The algorithm for skipping from N to N+1, e.g., expects consecutive task numbers on nodes (block assignment), not those assigned round robin (cyclic assignment). For example, a test running 6 tasks on 3 nodes would expect tasks 0,1 on node 0; tasks 2,3 on node 1; and tasks 4,5 on node 2. Were the assignment for tasks-to-node in round robin fashion, there would be tasks 0,3 on node 0; tasks 1,4 on node 1; and tasks 2,5 on node 2. In this case, there would be no expectation that a task would not be reading from data cached on a node.

HOW DO I USE HINTS?

It is possible to pass hints to the I/O library or file system layers following this form:

```
'setenv IOR_HINT__<layer>__<hint> <value>'
```

For example:: 'setenv IOR_HINT__MPI__IBM_largeblock_io true' 'setenv IOR_HINT__GPFS__important_hint true'

or, in a file in the form:: 'IOR_HINT__<layer>__<hint>=<value>'

Note that hints to MPI from the HDF5 or NCMPI layers are of the form:: 'setenv IOR_HINT__MPI__<hint> <value>'

HOW DO I EXPLICITLY SET THE FILE DATA SIGNATURE?

The data signature for a transfer contains the MPI task number, transfer- buffer offset, and also timestamp for the start of iteration. As IOR works with 8-byte long long ints, the even-numbered long longs written contain a 32-bit MPI task number and a 32-bit timestamp. The odd-numbered long longs contain a 64-bit transferbuffer offset (or file offset if the '-l' storeFileOffset option is used). To set the timestamp value, use '-G' or setTimeStampSignature.

HOW DO I EASILY CHECK OR CHANGE A BYTE IN AN OUTPUT DATA FILE?

There is a simple utility IOR/src/C/cbif/cbif.c that may be built. This is a stand-alone, serial application called cbif (Change Byte In File). The utility allows a file offset to be checked, returning the data at that location in IOR's data check format. It also allows a byte at that location to be changed.

HOW DO I CORRECT FOR CLOCK SKEW BETWEEN NODES IN A CLUSTER?

To correct for clock skew between nodes, IOR compares times between nodes, then broadcasts the root node's timestamp so all nodes can adjust by the difference. To see an egregious outlier, use the '-j' option. Be sure to set this value high enough to only show a node outside a certain time from the mean.

CHAPTER 8

Doxygen

Click [here](#) for doxygen.

This documentation utilities doxygen for parsing the c code. Therefore a doxygen instances is created in the background anyway. This might be helpful as doxygen produces nice call graphs.

Continuous Integration

Continuous Integration is used for basic sanity checking. Travis-CI provides free CI for open source github projects and is configured via a `.travis.yml`.

For now this is set up to compile IOR on a ubuntu 14.04 machine with gcc 4.8, openmpi and hdf5 for the backends. This is a pretty basic check and should be advance over time. Nevertheless this should detect major errors early as they are shown in pull requests.

10.1 General release process

The versioning for IOR is encoded in the `META` file in the root of the repository. The nomenclature is

- 3.2.0 designates a proper release
- 3.2.0rc1 designates the first release candidate in preparation for the 3.2.0 release
- 3.2.0+dev indicates development towards 3.2.0 prior to a feature freeze
- 3.2.0rc1+dev indicates development towards 3.2.0's first release candidate after a feature freeze

10.2 Building a release of IOR

To build a new version of IOR, e.g., from the 3.2 release branch:

```
$ docker run -it ubuntu bash
$ apt-get update
$ apt-get install -y git automake autoconf make gcc mpich
$ git clone -b 3.2 https://github.com/hpc/ior
$ cd ior
$ ./travis-build.sh
```

Alternatively you can build an arbitrary branch in Docker using a bind mount. This will be wrapped into a build-release Dockerfile in the future:

```
$ docker run -it --mount type=bind,source=$PWD,target=/ior ubuntu
$ apt-get update
$ apt-get install -y git automake autoconf make gcc mpich
$ ./travis-build.sh
```

10.3 Feature freezing for a new release

1. Branch *major.minor* from the commit at which the feature freeze should take effect.
2. Append the “rc1+dev” designator to the Version field in the META file, and update the NEWS file to have this new version as the topmost heading
3. Commit and push this new branch 2. Update the `Version:` field in META *of the main branch* to be the *next* release version, not the one whose features have just been frozen, and update the NEWS file as you did in step 2.

For example, to feature-freeze for version 3.2:

```
$ git checkout 11469ac
$ git checkout -B 3.2
$ vim META # update the ``Version:`` field to 3.2.0rc1+dev
$ vim NEWS # update the topmost version number to 3.2.0rc1+dev
$ git add NEWS META
$ git commit -m "Update version for feature freeze"
$ git push upstream 3.2
$ git checkout main
$ vim META # update the ``Version:`` field to 3.3.0+dev
$ vim NEWS # update the topmost version number to 3.3.0+dev
$ git add NEWS META
$ git commit -m "Update version number"
$ git push upstream main
```

10.4 Creating a new release candidate

1. Check out the appropriate commit from the *major.minor* branch
2. Disable the `check-news` option in `AM_INIT_AUTOMAKE` inside `configure.ac`
3. Remove the “+dev” designator from the Version field in META
4. Build a release package as described above
5. Revert the change from #2 (it was just required to build a non-release tarball) 5. Tag and commit the updated META so one can easily recompile this rc from git 6. Update the “rcX” number and add “+dev” back to the `Version:` field in

META. This will allow anyone playing with the tip of this branch to see that this the state is in preparation of the next rc, but is unreleased because of +dev.

7. Commit

For example to release 3.2.0rc1:

```
$ git checkout 3.2
$ # edit configure.ac and remove the check-news option
$ # remove +dev from the Version field in META (Version: 3.2.0rc1)
$ # build
$ git checkout configure.ac
$ git add META
$ git commit -m "Release candidate for 3.2.0rc1"
$ git tag 3.2.0rc1
$ # uptick rc number and re-add +dev to META (Version: 3.2.0rc2+dev)
```

(continues on next page)

(continued from previous page)

```
$ git add META # should contain Version: 3.2.0rc2+dev
$ git commit -m "Uptick version after release"
$ git push && git push --tags
```

10.5 Applying patches to a new microrelease

If a released version 3.2.0 has bugs, cherry-pick the fixes from main into the 3.2 branch:

```
$ git checkout 3.2
$ git cherry-pick cb40c99
$ git cherry-pick aafdf89
$ git push upstream 3.2
```

Once you’ve accumulated enough bugs, move on to issuing a new release below.

10.6 Creating a new release

This is a two-phase process because we need to ensure that NEWS in main contains a full history of releases, and we achieve this by always merging changes from main into a release branch.

1. Check out main
2. Ensure that the latest release notes for this release are reflected in NEWS
3. Commit that to main

Then work on the release branch:

1. Check out the relevant *major.minor* branch
2. Remove any “rcX” and “+dev” from the Version field in META
3. Cherry-pick your NEWS update commit from main into this release branch. Resolve conflicts and get rid of news that reflect future releases.
4. Build a release package as described above
5. Tag and commit the updated NEWS and META so one can easily recompile this release from git
6. Update the Version field to the next rc version and re-add “+dev”
7. Commit
8. Create the major.minor.micro release on GitHub from the associated tag

For example to release 3.2.0:

```
$ git checkout main
$ vim NEWS # add release notes from ``git log --oneline 3.2.0rc1..``
$ git commit
```

Let’s say the above generated commit abc345e on main. Then:

```
$ git checkout 3.2
$ vim META # 3.2.0rc2+dev -> 3.2.0
$ git cherry-pick abc345e
```

(continues on next page)

(continued from previous page)

```
$ vim NEWS # resolve conflicts, delete stuff for e.g., 3.4
$ # build
$ git add NEWS META
$ git commit -m "Release v3.2.0"
$ git tag 3.2.0
$ vim META # 3.2.0 -> 3.2.1rc1+dev
# vim NEWS # add a placeholder for 3.2.1rc2+dev so automake is happy
$ git add NEWS META
$ git commit -m "Uptick version after release"
```

Then push your main and your release branch and also push tags:

```
$ git checkout main && git push && git push --tags
$ git checkout 3.2 && git push && git push --tags
```

11.1 Version 3.4.0+dev

New major features:

New minor features:

Bugfixes:

11.2 Version 3.3.0

New major features:

- Add CephFS AIORI (Mark Nelson)
- Add Gfarm AIORI (Osamu Tatebe)
- Add DAOS AIORI (Mohamad Chaarawi)
- Add DAOS DFS AIORI (Mohamad Chaarawi)
- -B option has been replaced with `-posix.odirect`

New minor features:

- Display outlier host names (Jean-Yves Vet)
- Enable global default dir layout for subdirs in Lustre (Petros Koutoupis)
- Removed pound signs (#) from mdtest output file names (Julian Kunkel)
- Print I/O hints from NCMPi (Wei-keng Liao)
- Add mknod support to mdtest (Gu Zheng)
- Refactor AIORI-specific options (Julian Kunkel)
- Enable IME native backend for mdtest (Jean-Yves Vet)

- Enable mkdir/rmdir to IME AIORI (Jean-Yves Vet)
- Add HDF5 collective metadata option (Rob Latham)
- Add support for sync to AIORIs (Julian Kunkel)

General user improvements and bug fixes:

- Allocate aligned buffers to support DirectIO for BeeGFS (Sven Breuner)
- Added IOPS and latency results to json output (Robert LeBlanc)
- Fixed case where numTasks is not evenly divisible by tasksPerNode (J. Schwartz)
- Fix several memory leaks and buffer alignment problems (J. Schwartz, Axel Huebl, Sylvain Didelot)
- Add mdtest data verification (Julian Kunkel)
- Clean up functionality of stonewall (Julian Kunkel)
- Fix checks for lustre_user.h (Andreas Dilger)
- Make write verification work without read test (Jean-Yves Vet)
- Documentation updates (Vaclav Hapla, Glenn Lockwood)
- Add more debugging support (J. Schwartz)

General developer improvements:

- Fix type casting errors (Vaclav Hapla)
- Add basic test infrastructure (Julian Kunkel, Glenn Lockwood)
- Conform to strict C99 (Glenn Lockwood)

Known issues:

- S3 and HDFS backends may not compile with new versions of respective libraries

11.3 Version 3.2.1

- Fixed a memory protection bug in mdtest (Julian Kunkel)
- Fixed correctness bugs in mdtest leaf-only mode (#147) (rmn1)
- Fixed bug where mdtest attempted to stat uncreated files (Julian Kunkel)

11.4 Version 3.2.0

New major features:

- mdtest now included as a frontend for the IOR aiori backend (Nathan Hjelm, LANL)
- Added mmap AIORI (Li Dongyang, DDN)
- Added RADOS AIORI (Shane Snyder, ANL)
- Added IME AIORI (Jean-Yves Vet, DDN)
- Added stonewalling for mdtest (Julian Kunkel, U Reading)

New minor features:

- Dropped support for PLFS AIORI (John Bent)

- Added stoneWallingWearOut functionality to allow stonewalling to ensure that each MPI rank writes the same amount of data and captures the effects of slow processes (Julian Kunkel, DKRZ)
- Added support for JSON output (Enno Zickler, U Hamburg; Julian Kunkel, U Reading)
- Added dummy AIORI (Julian Kunkel, U Reading)
- Added support for HDF5 collective metadata operations (Rob Latham, ANL)

General user improvements:

- BeeGFS parameter support (Oliver Steffen, ThinkParQ)
- Semantics of `-R` now compares to expected signature (`-G`) (Julian Kunkel, DKRZ)
- Improved macOS support for ncmpi (Vinson Leung)
- Added more complete documentation (Enno Zickler, U Hamburg)
- Assorted bugfixes and code refactoring (Adam Moody, LLNL; Julian Kunkel, U Reading; Enno Zickler, U Hamburg; Nathan Hjelm, LANL; Rob Latham, ANL; Jean-Yves Vet, DDN)
- More robust support for non-POSIX-backed AIORIs (Shane Snyder, ANL)

General developer improvements:

- Improvements to build process (Nathan Hjelm, LANL; Ian Kirker, UCL)
- Added continuous integration support (Enno Zickler, U Hamburg)
- Added better support for automated testing (Julian Kunkel, U Reading)
- Rewritten option handling to improve support for backend-specific options (Julian Kunkel, U Reading)

Known issues:

- `api=RADOS` cannot work correctly if specified in a config file (via `-f`) because `-ul/-cl/-p` cannot be specified (issue #98)
- `writeCheck` cannot be enabled for write-only tests using some AIORIs such as MPI-IO (pull request #89)

11.5 Version 3.0.2

- IOR and mdtest now share a common codebase. This will make it easier run performance benchmarks on new hardware.
- Note: this version was never properly released

11.6 Version 3.0.0

- Reorganization of the build system. Now uses autoconf/automake. N.B. Windows support is not included. Patches welcome.
- Much code refactoring.
- Final summary table is printed after all tests have finished.
- Error messages significantly improved.
- Drop all “undocumented changes”. If they are worth having, they need to be implemented well and documented.

11.7 Version 2.10.3

- bug 2962326 “Segmentation Fault When Summarizing Results” fixed
- **bug 2786285 “-Wrong number of parameters to function H5Dcreate” fixed** (NOTE: to compile for HDF5 1.6 libs use “-D H5_USE_16_API”)
- bug 1992514 “delay (-d) doesn’t work” fixed

Contributed by demyn@users.sourceforge.net - Ported to Windows. Required changes related to ‘long’ types, which on Windows

are always 32-bits, even on 64-bit systems. Missing system headers and functions account for most of the remaining changes. New files for Windows:

- IOR/ior.vcproj - Visual C project file
- IOR/src/C/win/getopt.{h,c} - GNU getopt() support

See updates in the USER_GUIDE for build instructions on Windows.

- Fixed bug in incrementing transferCount
- Fixed bugs in SummarizeResults with mismatched format specifiers
- Fixed inconsistencies between option names, -h output, and the USER_GUIDE.

11.8 Version 2.10.2

- Extend existing random I/O capabilities and enhance performance output statistics. (Hodson, 8/18/2008)

11.9 Version 2.10.1

- Added ‘-J’ setAlignment option for HDF5 alignment in bytes; default value is 1, which does not set alignment
- Changed how HDF5 and PnetCDF calculate performance – formerly each used the size of the stat()’ed file; changed it to be number of data bytes transferred. these library-generated files can have large headers and filler as well as sparse file content
- Found potential overflow error in cast – using IOR_offset_t, not int now
- Replaced HAVE_HDF5_NO_FILL macro to instead directly check if H5_VERS_MAJOR H5_VERS_MINOR are defined; if so, then must be HDF5-1.6.x or higher for no-fill usage.
- Corrected IOR_GetFileSize() function to point to HDF5 and NCMPI versions of IOR_GetFileSize() calls
- Changed the netcdf dataset from 1D array to 4D array, where the 4 dimensions are: [segment-Count][numTasks][numTransfers][transferSize] This patch from Wei-keng Liao allows for file sizes > 4GB (provided no single dimension is > 4GB).
- Finalized random-capability release
- Changed statvfs() to be for __sun (SunOS) only
- Retired Python GUI

11.10 Version 2.10.0.1

- Cleaned up WriteOrRead(), reducing much to smaller functions.
- Added random capability for transfer offset.
- Modified strtok(NULL, "trn") in ExtractHints() so no trailing characters
- Added capability to set hints in NCMPI

11.11 Version 2.9.6.1

- For 'pvfs2:' filename prefix, now skips DisplayFreeSpace(); formerly this caused a problem with statvfs()
- Changed gethostname() to MPI_Get_processor_name() so for certain cases when gethostname() would only return the frontend node name
- Added SunOS compiler settings for makefile
- Updated O_DIRECT usage for SunOS compliance
- Changed statfs() to instead use statvfs() for SunOS compliance
- Renamed compiler directive _USE_LUSTRE to _MANUALLY_SET_LUSTRE_STRIPING

11.12 Version 2.9.5

- Wall clock deviation time relabeled to be "Start time skew across all tasks".
- Added notification for "Using reorderTasks '-C' (expecting block, not cyclic, task assignment)"
- Corrected bug with read performance with stonewalling (was using full size, stat'ed file instead of bytes transferred).

11.13 Version 2.9.4

- Now using IOR_offset_t instead of int for tmpOffset in IOR.c:WriteOrRead(). Formerly, this would cause error in file(s) > 2GB for ReadCheck. The more-commonly-used WriteCheck option was not affected by this.

11.14 Version 2.9.3

- Changed FILE_DELIMITER from ':' to '@'.
- Any time skew between nodes is automatically adjusted so that all nodes are calibrated to the root node's time.
- Wall clock deviation time is still reported, but have removed the warning message that was generated for this. (Didn't add much new information, just annoyed folks.)
- The '-j' outlierThreshold option now is set to 0 (off) as default. To use this, set to a positive integer N, and any task who's time (for open, access, etc.) is not within N seconds of the mean of all the tasks will show up as an outlier.

11.15 Version 2.9.2

- Simple cleanup of error messages, format, etc.
- Changed error reporting so that with VERBOSITY=2 (-vv) any error encountered is displayed. Previously, this was with VERBOSITY=3, along with full test parameters, environment, and all access times of operations.
- Added deadlineForStonewalling option (-D). This option is used for measuring the amount of data moved in a fixed time. After the barrier, each task starts its own timer, begins moving data, and stops moving data at a prearranged time. Instead of measuring the amount of time to move a fixed amount of data, this option measures the amount of data moved in a fixed amount of time. The objective is to prevent tasks slow to complete from skewing the performance.

11.16 Version 2.9.1

- Updated test script to run through file1:file2 cases.
- Program exit code is now total numbers of errors (both writecheck and readcheck for all iterations), unless quitOnError (-q) is set.
- For various failure situations, replace abort with warning, including: - failed uname() for platform name now gives warning - failed unlink() of data file now gives warning - failed fsync() of data file now gives warning - failed open() of nonexistent script file now gives warning
- Changed labelling for error checking output to be (hopefully) more clear in details on errors for diagnostics.
- Another fix for -o file1:file2 option.
- Corrected bug in GetTestFileName() – now possible to handle -o file1:file2 cases for file-per-proc.

11.17 Version 2.9.0

- Improved checkRead option to reread data from different node (to avoid cache) and then compare both reads.
- Added outlierThreshold (-j) option to give warning if any task is more than this number of seconds from the mean of all participating tasks. If so, the task is identified, its time (start, elapsed create, elapsed transfer, elapsed close, or end) is reported, as is the mean and standard deviation for all tasks. The default for this is 5, i.e. any task not within 5 seconds of the mean for those times is displayed. This value can be set with outlierThreshold=<value> or -j <value>.
- Correct for clock skew between nodes - if skew greater than wallclock deviation threshold (WC_DEV_THRESHOLD) in seconds, then broadcast root node's timestamp and adjust by the difference. WC_DEV_THRESHOLD currently set to 5.
- Added a Users Guide.