Trading Memory for Disk Using Parallel Access to Fast InfiniBand Disk Arrays for Large Computational Chemistry Applications

Kyle Schochenmaier, Troy Benjegerdes and Brett Bode
Ames Laboratory, U.S. DOE
Iowa State University
Scalable Computing Laboratory
Ames, Iowa, 50011
Email: {kschoche, troy, brett}@scl.ameslab.gov

1 Abstract

We present a novel approach for using high performance network attached parallel storage for out-ofcore computation. Our approach utilizes many parallel disks and storage controllers with a near 1:1 ratio of compute nodes to storage servers. This, when combined with 30 Gigabit 12X InfiniBand interconnects, allows remote storage subsystem bandwidth to reach the same performance level as local memory-cache file I/O performance. This combination allows computational chemistry application problem sizes which require more than 100GB of intermediate data to store this data on disk without the disk I/O subsystem becoming the limiting factor. With sequential storage access speeds on the same order of magnitude as main memory performance, this allows out-of- core computation to become practical due to disk storage being several orders of magnitude cheaper per GB than main memory storage.

2 Problem Statement

Computational Chemistry packages such as the General Atomic and Molecular Electronic Structure System (GAMESS) [?, ?] can generate very large amounts of intermediate data that can either be stored or recomputed as needed. Some of these algorithms have been designed with both direct (recomputation of integrals) and conventional (integrals computed once, stored and reused) modes. Conventional methods require far fewer computations, but can require multigigabyte data files. In the past conventional methods were used quite frequently, but increases in CPU speed have outpaced increases in IO bandwidth significantly.

In addition, the industry trend to large numbers of light weight nodes with often no local secondary storage has led to more emphasis on direct methods. However, some of the most accurate methods are not currently available in a parallel implementation and often are the most resource intensive types of runs. These calculations are usually substantially I/O bound. Thus, to support these types of calculations we need to be able to provide secondary storage at the highest possible bandwidth to nodes within our cluster.

From a management perspective we want to be able to centralize the storage and deliver it precisely to the nodes that need it. Centralized storage allows flexibility to the system scheduler and makes it easier to manage the failure prone disk subsystems. In order to make this practical we need to be able to deliver bandwidth to secondary storage that matches or exceeds that available from typically locally attached storage systems on a cluster. We feel that a target of 150MB/second is reasonable. In order to achieve this target we plan to utilize multiple high-speed storage InfiniBand attached storage servers with a virtual file system to make them appear as one storage target to the compute node clients. With this setup hope to support I/O rates that can keep processor utilization at greater than 90% and achieve sustained sequential disk read rates of 50% of the aggregate raw disk spindle read performance.

3 Configuration

3.1 Hardware Configuration

Our hardware configuration consists of a set of high-performance, low-cost Opteron based storage

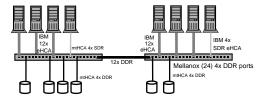


Figure 1. Hardware Configuration

servers and a set of IBM OpenPower Power5 based servers as the clients utilizing InfiniBand for the interconnect as shown in Figure 1. The storage component is provided by six dual AMD Opteron based servers. Each server is configured with 4GB of RAM, two PCI-X based Areca 8-port SATA RAID controllers and 16 Seagate 250GB SATA disks configured in two hardware RAID5 arrays that are then combined into a single RAID0 array in software. The network for the storage servers is provided by a Mellanox 8x PCI-express based 4X Double Data Rate (DDR) InfiniBand providing a peak data payload of 16Gbps. The clients are a set of ten IBM quad-Power5 OpenPower 720 servers configured with 8 or 16 GB of RAM and a GX processor bus attached IBM 12X Single Data Rate (SDR) Galaxy InfiniBand adapter. In addition, two clients were configured with 4X SDR PCI-X based Mellanox InfiniBand adapters. This setup allows us to test clients running at peak data payloads of 8 Gbps, 16 Gbps, and 24 Gbps and allows varying the server count from 1 to six and the client count from 1 to 10. The nodes and servers are distributed between two 24 4X DDR port Mellanox switches. The switches are connected to each other via a 12X DDR link providing up to 48Gbps of switch to switch bandwidth. For the purposes of this work this should provide an effectively flat network without switch imposed bandwidth limitations.

3.2 Data and Storage Layout

Utilizing PVFS2 the data storage on the Opteron based servers was combined into sets of 1, 2, 4 and 6 servers. In each case data is automatically striped across the number of servers in the virtual file system. The stripe size was also examined.

4 Software

Our software core is based on the latest development version of PVFS2 running on AMD64 and PowerPC64 based Debian Linux. To quantify our performance we have used a new version of NetPIPE and the GAMESS quantum chemistry application. To further validate these results we have checked them against the output of the vmstat linux command.

4.1 PVFS2

The core piece of software for our test setup is the Parallel Virtual File System 2 (PVFS2) [?]. PVFS2 is designed to stripe data storage across multiple independent servers while providing the appearance of a single unified file system to the clients. It has been demonstrated to have quite good scalability in the past with large numbers of clients and servers. However, the performance to an individual client has been less of a focus, in part because the network connection speeds have imposed an upper limit of about 100MB/sec on a network file system. With the arrival of InfiniBand in the past few years the network capability has been increased from 1-2 Gbps to 8 and now 16-24 Gbps. Achieving that level of performance requires the direct use of the native verbs interface. Thus, we have been assisting in the development of an OpenIB [?] native message layer for PVFS2. While an initial version of this work was released in the version 1.5 release of PVFS2 work is still ongoing and significant performance improvements were made during the course of this effort. PVFS2 provides up to three different application interfaces. The simplest is provided through a kernel extension and client process that provides a normal file system interface to applications. The second option is through the use of MPI-IO and the third is through direct calls to the PVFS2 user land library. The normal file system interface is obviously the simplest for applications to use, but potentially imposes additional overhead and limits the ability of the application to tune various IO related parameters. The native libpyfs2 usage is the most invasive, but allows the most tuning. Thus, for our tests we have chosen to test a native libpyfs2 implementation and a normal file system implementation.

4.2 NetPIPE

The Network Protocol Independent Performance Evaluator, or NetPIPE, is a utility originally designed to produce a more extensive analysis of interconnect related communication performance at the MPI and TCP layers [?]. It has since been developed to include support for other message-passing libraries, in addition to native software layers such as OpenIB. It is organized into a central program which provides an identical testing environment for each supported implementation over modularized interfaces.

In most cases, NetPIPE measures the interconnect or message-passing latencies and bandwidths via a synchronized, 2-sided communications system. The system utilizes a ping-pong style measurement across different ranges of message sizes, while each message size is repeated over a number of iterations and a worst-case analysis of the system is produced [?].

A NetPIPE [?] module has been recently developed specifically for testing the performance characteristics of file systems and the PVFS2 software stack. This includes methods by which the throughput of the PVFS2 file system implemented over the OpenIB interface can be analyzed through both native calls and the use of the VFS layer, as well its performance when accessing local disk via unix read() and write() APIs.

The PVFS2-NetPIPE module was developed specifically to address the need to examine the performance characteristics of a parallel file system implemented via a system of high-performance data servers capable of saturating the interconnect with raw disk I/O. Since the primary goal of the implementation is to provide each client with a maximum sustained bandwidth to secondary storage rather than an optimal aggregate bandwidth over many clients, the benchmark was designed to demonstrate the ability of a single client to achieve large sustained data throughput directly to disks.

The NetPIPE module was developed against the current NetPIPE-3.6 and PVFS2-1.5.1 releases, and utilizes PVFS2 function calls from libpvfs to directly interface with a remote PVFS2 file system. These internal calls reside in the PVFS2 layer, above the BMI interface, which allows the NetPIPE module to be used with any BMI-supported interconnect. The connection to the PVFS2 file system is done in the Setup

and Init functions of the program, also implemented in these functions are various PVFS2-side file checks necessary to ensure files are not overwritten. The Send and Recv functions provide the data transfer mechanisms to and from the data servers.

During a write test, a new file is created and opened remotely on the PVFS2 data servers for writing. For a read test, an existing file created by a previous Net-PIPE write test (or any other sufficiently sized existing file) is opened for read-only tests. Once the file is opened, NetPIPE will manage local buffers and memory segments aligned to the current message size, and read or write the contents from or to the open files. When using cache invalidate mode, after each send or receive, file pointers are moved forward to the next 4k aligned file offset, causing a seek, and reducing the possibility that data is being cached on the servers. This does little when the message sizes are very small, however typical tests produce files greatly in excess of system memory size. After the tests have completed, all files that were opened via the PVFS2 interface are closed and finalized via the internal PVFS2 PVFS_finalize function.

The NetPIPE disk I/O module is substantially similiar to the PVFS2 client module, and the file is opened in the NetPIPE module Setup function. During read and write tests, the only substantial difference between the PVFS and Disk modules is that the return value of the unix read function must be checked to ensure that all the data has been read in case of an interrupted system call.

4.3 GAMESS

As mentioned previously GAMESS is a large quantum chemistry package that encompasses many different algorithms and code paths. For the purposes of this work we have chosen to utilize one of the more basic algorithms, the Hartree- Fock energy calculation. This type of calculation is the base computation required for many of the higher level methods in quantum chemistry and is thus very important and widely used. GAMESS offers both a direct and conventional version of the algorithm. In the conventional algorithm the primary disk usage is for the storage of the set of four center two electron integrals. These integrals are precomputed and written to disk once and read in a

number of times while molecular wavefunction is optimized in an iterative fashion. This algorithm tends to be quite IO bound in both the read and write phases, though because the data is read many more times than it is written the read performance is more important. GAMESS includes timing information for many of the individual steps, but due to the overlap of I/O and computation the rates computed from these timings represent a lower bound on the I/O performance.

While GAMESS normally uses the standard FOR-TRAN I/O interfaces we also wanted to be able to test the native PVFS2 system. Thus, a shim layer was developed that passes I/O from GAMESS to the native PVFS2 interfaces. The shim layer also allows us to tune the I/O buffer size passed to the file system or PVFS2 library as well the PVFS2 stripe size.

For our test runs two sets of input were chosen. The first is the insecticide rotenone, $C_{23}O_6H_{22}$, with a 6-31G* basis set or 479 atomic orbitals (AOs). This run results in a 17.1 GB integral file. This file size will fit in the file cache when run with 6 data servers and thus allows us to evaluate the performance of the network and PVFS2 software separately from the disk performance. The second test case uses the anticancer drug taxol, $C_{47}O_{14}N_1H_{51}$, with a 6-31G* basis set or 1032 AOs. This produces a 120GB file that is clearly several times our aggregate file cache size. Using these test cases we hope to validate the NetPIPE results by relating real application performance to specific points on the NetPIPE curves.

5 Solution

Our solution to the problem combines InfiniBand network attached storage, the PVFS2 file system, and the GAMESS computational chemistry application. By building a system with these components with many more disks than processors we can support I/O rates to a single node exceeding 10 Gigabits. Delivering this I/O rate in turn allows for building a system that can support conventional, secondary storage for scratch space, types of algorithms without locally attached disks and their accompanying problems.

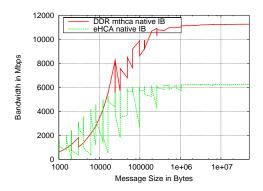


Figure 2. Native OpenIB verbs performance

6 Experiments

We began our experiments by establishing the base level of performance for the network and the disk subsystems on the Opteron based storage servers. This was done using NetPIPE to measure the native OpenIB verbs performance on the IBM eHCAs and on the Opteron based storage servers. Those results are plotted in Figure 2. The curve for the Mellanox DDR NIC in the Opteron systems appears much as we would expect with a respectable peak performance of around 11 Gbps. However, the curve for the IBM eHCA is much less impressive. At small message sizes on regular block sizes the performance is quite good as a result of its excellent message latency, but the peak performance is only slightly over 6 Gbps. After discussing this with the developers the reason is that this adapter is really designed for many simultaneous message streams. As such it has multiple DMA engines on the card, but a single message stream can use only a single DMA engine which is limited to about 6 Gbps. While this result appears quite bad, it turns out not to be a serious issue for PVFS2 when we are normally talking to multiple storage servers. Since PVFS2 must use at least one message stream per storage server the multiple DMA engines in the eHCA engine are used in parallel to good effect as we shall see later.

Next we examined the local disk performance on our Opteron storage servers. We found a peak performance of about 3.5 Gbps (435MB/sec) per server. This is somewhat less than we had hoped for, but still a good rate in aggregate as the 6 servers would be able to provide 2.6 GB/sec of total I/O bandwidth.

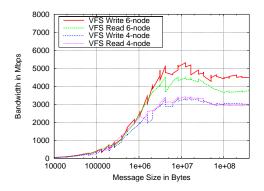


Figure 3. PVFS2 VFS performance

The first PVFS2 tests examined the VFS performance. These results are plotted in Figure 3 showing the performance for 4 and 6 storage server configurations. These tests allow the servers to utilize file cache so they represent a test of the network and software performance. For validation the smaller GAMESS computation averaged 181 MB/s over the total run time of the job on 1 CPU and 336 MB/s on 4 CPUs. If we subtract off the CPU time and recompute the I/O rate we get of 447 MB/sec (3.5 Gbps) on one CPU and 473 MB/sec (3.75 Gbps) on four CPU run, both within a single node. This gives fairly good agreement with the peak NetPIPE numbers.

Figure 4 illustrates the same runs using the native PVFS2 interfaces. These show a doubling of peak performance, but not until a fairly large message size. Using GAMESS as a comparison is more difficult in this case because we can not easily separate out the CPU time used by the computation and the CPU time used for polling on I/O. However, we can compare total wall times which went from 1344 seconds over VFS to 850 seconds native. The minimum performance is 288 MB/sec (2.3 Gbps), but in comparison with Net-PIPE and previous GAMESS results we estimate that it is closer to the NetPIPE peak performance around 1.1 GB/sec (9 Gbps).

Figure 5 shows NetPIPE running in cache invalidate mode. Cache invalidate mode causes seeks inside the file at every iteration of a send/recv call. This is shown most effectively when we begin seeking outside of the data that remains in system cache on any given data node, which occurs beyond 2MBs for the number of iterations done in this test. This type of test shows how

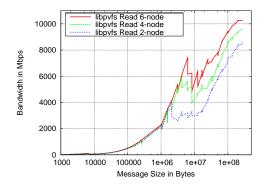


Figure 4. PVFS2 native performance

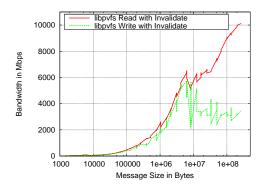


Figure 5. NetPIPE in cache invalidate mode

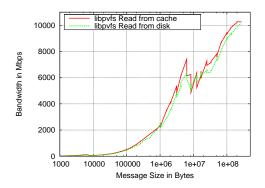


Figure 6. Comparison of read performance with and without cache invalidate

a GAMESS job would access storage while reading an integral of potentially several hundred gigabytes.

The lines in Figure 5 help show the effects of caching in various places between the client and the disk subsystem. An important thing to note is that as the message size increases, the read bandwidth is still increasing and has yet to level off. This may imply that we have not reached a bandwidth limitation imposed by the disk subsystem for reads. We were unable to examine this due to limitations with registering memory on the eHCAs. We do however know that we have yet to reach the limits measured for the disk subsystems.

As the test show, the effects of caching tail off when the test file size becomes large. The test file used to produce this graph was over 400GB which helped eliminate caching effects beyond 50-75MB message sizes. The write line shows our expected bandwidth directly to disk of about 3.5Gbit/sec with little to no caching effects due to the fact that the test file is many times the aggregate file cache size of the storage servers.

Figure 6 compares the measured NetPIPE read performance with and without cache invalidate. Thus, this is a comparison of PVFS2 performance from storage server file cache versus all the way to the disk. Clearly the performance from disk is quite close to the performance from cache with the addition of a little higher latency. This indicates that our peak performance is being limited more by the network and PVFS2 protocols than by the performance of the disk subsystems.

Figure 7 illustrates the sustained aggregate read performance as reported by vmstat on the storage servers. As such it is more reliable than the results for the smaller test case. Tests are included for runs with the VFS interface for both 1 and 4 processes per node and for the native interface with 1 process per node. Unfortunately bugs in the eHCA driver prevented us from running more than one native process per node. This restriction does not effect the VFS layer because PVFS2 currently uses a single client process to service the client VFS layer. Thus multiple processes per node can be used to good effect with the VFS layer. The first item of note is that the single CPU read performance is about 220 MB/sec (1.76 Gbps) through the VFS layer. Through the native interface we achieve 375 MB/sec (3 Gbps), which is almost as fast as 4 CPUs running

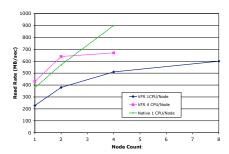


Figure 7. Scaling of GAMESS taxol Computation

through through the VFS interface. The scaling for the VFS interface seems to top out lower than we expect. However, the native interface appears to scale much better and achieves 900 MB/sec (7.2 Gbps) with 1 process on each of four nodes.

6.1 Conclusions

We have demonstrated through various tests that PVFS2 running over the OpenIB stack can be used to deliver high-bandwidth to secondary storage to an individual node. The VFS layer delivers 220 MB/sec and the native layer delivers 375 MB/sec from disk to a single process. These are quite respectable I/O rates and are well in excess of our 150 MB/sec target. However, we recognise that there is still room for improvement. Tests on files that fit within the file cache of the storage servers indicate peak VFS performance of 550 MB/sec and peak native performance of over 1 GB/sec both to a single process.

We believe that there are several opportunites for further performance improvements. These include multithreading the VFS client and server processes, the use of read ahead in the server, and the use of AIO in the server.

7 Judging Criteria

The focus of our effort is on obtaining large I/O bandwidths to a single node and/or a single process.

Thus, our objective for overall performance is limited to demonstrating exceptional per node I/O bandwidth in a framework that can be used on a distributed memory cluster to deliver that I/O bandwidth to the specific processes that require it. To that end we have demonstrated real application read rates from disk at 375 MB/sec for a single process. This is significantly larger than what is available today on most any cluster system.

We also demonstrated peak read rates in excess of 1.25 GB/sec (10 Gbps) with NetPIPE (a single process). This result is close to 50% of the aggregate local disk performance measured on each server, which seems quite good for a distributed file system.

PVFS2 and other parallel file systems have demonstrated good performance and good scaling when dealing with simultaneous access by large numbers of nodes to network attached storage. However, until 10 Gbps and faster networks recently arrived in the cluster arena they haven't been an option that could compete with the bandwidth of locally attached disk. The implementation of a native PVFS2 layer for OpenIB has changed that equation substantially. As mentioned previously we have demonstrated PVFS2 over OpenIB performance that exceeds not only the bandwidth of all but the fastest disk subsystems, but also exceeds the bandwidth available on 10 Gigabit Ethernet or 4X single data rate InfiniBand.

Utilizing InfiniBand attached storage with PVFS2 can be an extremely effective way to deliver large amounts of very fast scratch storage to any node in a cluster. This flexibility along with the significant management advantages of centrally located disk subsystems over locally attached disks on each node make it an attractive solution for secondary storage on distrbuted memory systems.

8 Proposal Deviations

The testing followed the proposal fairly closely with the exception that we were not able to test the native mode interface of PVFS2 as thoughly as we had planned due to problems with the drivers for our eHCAs. These problems limited us to running only one native mode process per node, which did give good results, but we feel that we could still see significantly higher per node bandwidth with multiple CPUs

active at once.

9 Price-Performance Analysis

Our test setup can be divided into three parts, the storage servers, the network, and the compute clients. The storage servers are probably the most relavent part to examine in terms of price/performance. Our servers are about one year old and are made up of dual processor AMD Opteron CPUs, two Areca PCI-X RAID controllers and 16 250 GB SATA HDs. These units cost approximately \$8000 each. Today the same money would likely buy two dual-core CPUs and PCI-Express RAID controllers.

The Interconnect is provided for by Mellanox DDR PCI-Express InfiniBand NICs and Mellanox unmanaged 24 port DDR switches. While these parts were purchased prior to full production we feel that a cost of about \$1000 per node is accurate.

The IBM Power5 servers were purchased for approximately \$18,000 each. The processor bus attached 12X InfiniBand adapter is to the best of our knowledge not yet publically available. However, we have been advised of a price of approximately \$2,500 each.

10 Appendices

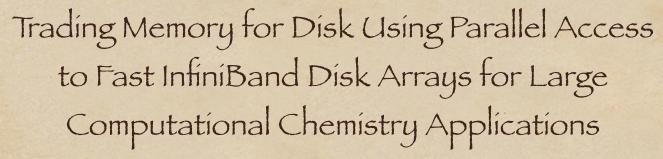
10.1 Acknowledgements

This work was performed under the auspices of the U.S. Department of Energy under contract W-7405-Eng-82 at Ames Laboratory operated by the Iowa State University of Science and Technology. Funding was provided by the Mathematical, Information and Computational Science division of the Office of Advanced Scientific Computing Research. We would also like to thank Brad Benton and Chet Mehta at IBM for early access and support with the 12x IBM eHCA, and Mellanox for early access to DDR-capable InfiniBand HCA's and switches.

Finally, this work would not have been possible without the efforts of the entire PVFS2 development team. In particular, Pete Wyckoff at OSC did much of work on the PVFS2 port to OpenIB upon which this work is based.

10.2 Slides For Finalist Presentation

The following pages include the slides as we expect to present them in the public storage challenge finalists presentation.



Kyle Schochenmaier, Troy Benjegerdes and Brett M. Bode Scalable Computing Laboratory Ames Laboratory, U.S. DOE Iowa State University



Iowa State University

Outline

- Problem Statement
- Hardware Configuration
- Software Configuration
- Software Tools
- Results
- Conclusions

Problem Statement

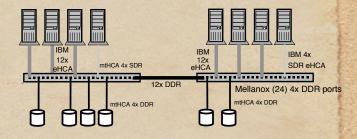
- Our primary application, the GAMESS quantum chemistry app. has many code paths that are quite I/O bound writing and especially reading large temporary storage.
- Most HPC systems are moving towards providing minimal locally attached secondary storage with a corresponding meager I/O bandwidth.
 - This is particularly troubling as the number of CPUs per node is increased.

Proposed Solution

- Network interconnects have reached the point where they can potentially deliver access to secondary storage faster than locally attached storage subsystems.
- This also requires scalable client/server software capable of delivering very high bandwidth to a single node while simultaneously scaling to large numbers of clients.
 - We have chosen to use PVFS2 on Linux clients and servers interconnected by InfiniBand.

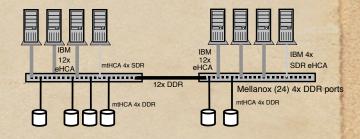
Hardware Configuration

- Six storage servers
 - dual AMD Opteron processors
 - 4 GB RAM
 - 2 Areca PCI-X SATA RAID controllers
 - 16 250 GB Seagate SATA HDs
 - Mellanox 4X DDR PCI-Express InfiniBand adapter (16 Gbps)



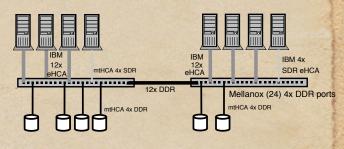
Hardware Configuration

- Eight compute clients
 - quad Power5 processors
 - 8 or 16 GB RAM
 - IBM 12X GX processor bus attached (eHCA) InfiniBand adapter (24 Gbps)
- Interconnect
 - 2 Mellanox 24 port (4X SDR/ DDR) switches
 - Connected together with a 12X DDR link (48 Gbps max data payload)



Software Configuration

- AMD64 version of Debian Linux on storage servers
- PPC64 version of Debian Linux on IBM power5 clients
- PVFS2 running on OpenIB verbs natively. Version 1.5.1 ++ (from latest development tree)



NetPIPE

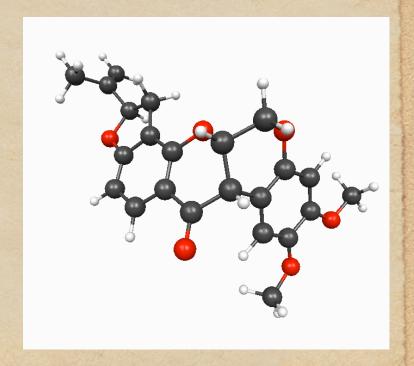
- Tool for measuring network bandwidth versus message size.
- New modules to test I/O bandwidth
 - Can be set to allow testing of file system cache (reread the same data over and over)
 - Can also stride through a file to obtain performance numbers all the way to disk.

GAMESS

- Our motivating application
- Large (750k lines) FORTRAN application
- Has many different algorithms including both direct (~diskless) and conventional (potentially very large temporary files).
- MPI version, but not normally used.
- Used the common Hartree-Fock energy calculation for our tests.

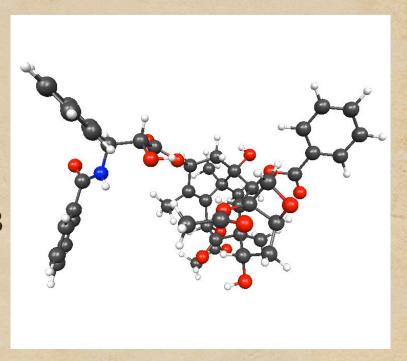
Small Test - Rotenone

- 479 AOs, 104
 occupied MOs
- Produces 16.2
 GB scratch file.



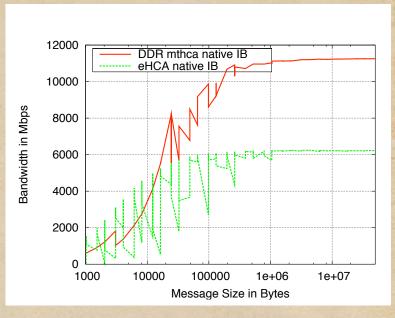
Large Test - Taxol

- 1032 AOs, 226
 occupied MOs
- Produces 120 GB scratch file.



Base Network Performance

- Performance for the storage servers exceeds 11 Gbps
- IBM eHCA performance is a disappointing 6.2 Gbps.
 - eHCA has 6 DMA engines
 - eHCA can parallelize multiple streams with the multiple DMA engine

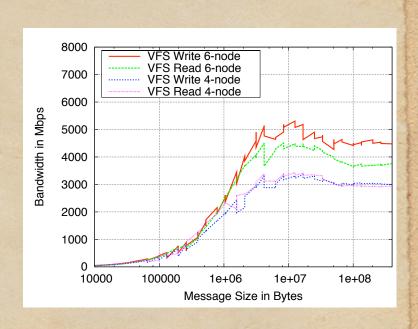


Base Disk Performance

• Directly testing of the I/O performance on the Opteron storage servers indicated a peak read performance of 435 MB/sec. measured using NetPIPE (a single stream). Much higher bandwidth can be obtained with Linux AIO approaching 600 MB/sec.

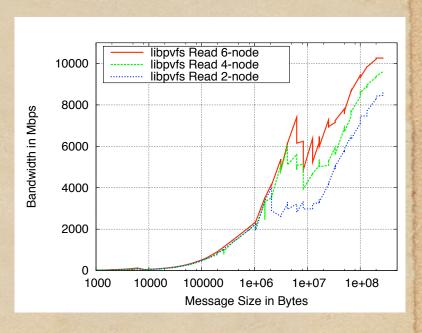
VFS Results from Cache

- Peak read
 performance of
 greater than 500
 MB/sec
- GAMESS tests on small test case show similar peak numbers.



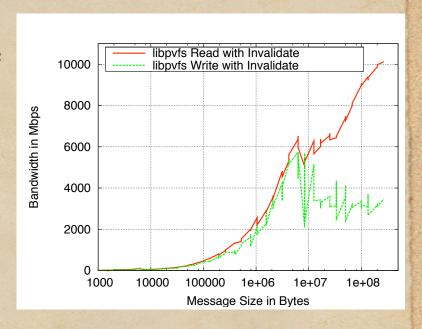
Native Results from Cache

- Peak read
 performance of
 greater than 1 GB/
 sec
- GAMESS performance is also nearly double the VFS result.



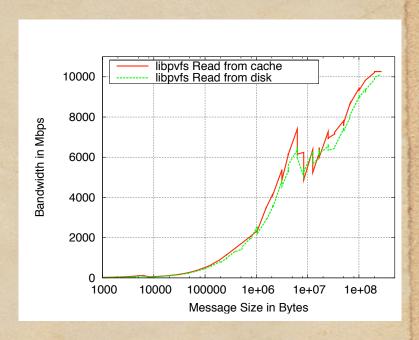
Results from Disk

- Write performance all the way to disk is significantly reduced.
- Read performance seems barely effected.



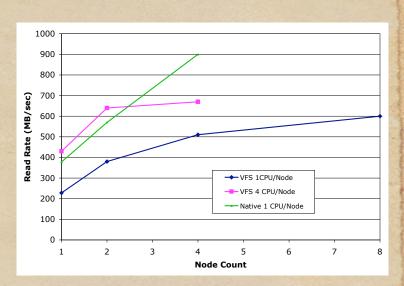
Results from Disk

Indeed read
 performance is
 only slightly
 reduced with some
 additional latency.



GAMESS Large Run Performance

- Single process
 performance is 228
 MB/sec over VFS,
 375 MB/sec native.
- 4 processes on 4 nodes gives 900 MB/sec in native mode!



Conclusions

- PVFS2 over OpenIB can be used to deliver I/O to a single node and a single process at rates that significantly exceed the performance of locally attached disk subsystems typically used in clusters.
- This setup offers the possibility of using inexpensive storage servers to provide very fast I/O to high-end compute servers.

Acknowledgments

- Funding:
 - U.S. Department of Energy
 - IBM
- Brad Benton and Chet Mehta at IBM
- Pete Wyckoff at OSC
- Rob Latham, Sam Lang and Rob Ross on the PVFS2 development team.

