

Remote I/O: Fast Access to Distant Storage

Ian Foster, David Kohr, Jr., Rakesh Krishnaiyer, and Jace Mogill
Mathematics and Computer Science Division
Argonne National Laboratory
Argonne, IL 60439

RECEIVED**JUL 26 1999****OSTI**

Abstract

As high-speed networks make it easier to use distributed resources, it becomes increasingly common that applications and their data are not colocated. Users have traditionally addressed this problem by manually staging data to and from remote computers. We argue instead for a new *remote I/O* paradigm in which programs use familiar parallel I/O interfaces to access remote filesystems. In addition to simplifying remote execution, remote I/O can improve performance relative to staging by overlapping computation and data transfer or by reducing communication requirements. However, remote I/O also introduces new technical challenges in the areas of portability, performance, and integration with distributed computing systems. We propose techniques designed to address these challenges and describe a remote I/O library called RIO that we have developed to evaluate the effectiveness of these techniques. RIO addresses issues of portability by adopting the quasi-standard MPI-IO interface and by defining a RIO device and RIO server within the ADIO abstract I/O device architecture. It addresses performance issues by providing traditional I/O optimizations such as asynchronous operations and through implementation techniques such as buffering and message forwarding to offload communication overheads. RIO uses the Nexus communication library to obtain access to configuration and security mechanisms provided by the Globus wide area computing toolkit. Microbenchmarks and application experiments demonstrate that our techniques achieve acceptable performance in most situations and can improve turnaround time relative to staging.

1 Introduction

Improvements in networking and software infrastructure are making it easier for programmers to execute programs at remote sites and to write programs that use resources at multiple locations. One consequence of remote execution is that a program may be geographically separated from the files that it accesses. This separation can significantly increase conceptual and temporal overheads in program development and execution. Ideally, we would like to enable programs to access data in a manner independent of data and program location. In our experience, the key challenges that must be addressed before we can provide this capability are portability (across different networks and filesystems), performance (in potentially high-latency, low-bandwidth, heterogeneous networks), and integration into distributed computing environments.

Historically, the high-performance computing community has achieved remote data access by manually *staging* input data from its home filesystem to the computer where a program is to execute; this process is then reversed for output data. However, this approach is clumsy, prevents overlapping of communication and computation, and can result in excessive data transfer in situations where a program accesses only part of a file. Distributed filesystems [15] also support remote data access, but performance and administrative problems often render them inappropriate for high-performance computing.

DISCLAIMER

This report was prepared as an account of work sponsored by an agency of the United States Government. Neither the United States Government nor any agency thereof, nor any of their employees, make any warranty, express or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or any agency thereof. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or any agency thereof.

DISCLAIMER

Portions of this document may be illegible in electronic image products. Images are produced from the best available original document.

We propose a new approach to remote data access in which programs use *remote I/O libraries* to access files located on remote filesystems in a manner that is independent of physical location. In contrast to distributed filesystems, remote I/O libraries use parallel I/O interfaces and focus on high-performance transfer. We believe that this narrow focus can allow remote I/O libraries to meet requirements for performance, flexibility, and convenience without introducing undue complexity in their implementation.

As part of an investigation of the remote I/O concept, we have designed and implemented a remote I/O library called RIO. RIO achieves portability by adopting the I/O interface defined by MPI-IO [6] and by exploiting features of the ADIO abstract I/O device [21], providing a RIO device that translates ADIO calls into communications to remote RIO servers. Performance issues are addressed by the use of dedicated forwarder nodes, buffering, and support for asynchronous and collective operations. RIO uses the Nexus communication library [10] for client/server communication, hence providing access to configuration and security mechanisms provided by the Globus wide area computing toolkit.

We have performed experiments in a controlled multicomputer environment to evaluate the effectiveness of our techniques. Microbenchmarks demonstrate that RIO can drive networks at close to their peak performance; these experiments also allow us to quantify the benefits of optimizations such as buffering. Application experiments illustrate the feasibility of remote I/O in a representative application. In particular, we demonstrate enhanced performance relative to staging.

The principal contributions of this article are fourfold. Specifically, we

- introduce the concept of remote I/O, explain why it is important, and motivate its requirements;
- discuss networking issues that make remote I/O challenging, and propose library facilities that address these challenges;
- present experimental results that demonstrate the efficacy of our design techniques, and indicate where more work is needed; and
- show how to integrate a remote I/O library with mechanisms that support operation in a distributed environment.

2 The Remote I/O Problem

We first expand upon why remote I/O is important, discuss various networking issues that remote I/O libraries must address, and review other approaches to the remote I/O problem.

2.1 Motivation

Remote computational or data resources may be used because they provide a unique capability (e.g., a supercomputer or database) or simply because they are available (e.g., in a computational grid that uses a load-sharing system such as Condor [14] or LSF to map tasks to idle resources). In either case, filesystems may be geographically separated from computers. This need to access “remote” filesystems arises frequently even within a single site, as it is rare that all disks are crossmounted.

Programmers have traditionally resorted to staging techniques when a program and its data are not colocated. However, there can be significant performance, flexibility, and convenience advantages to having a uniform interface to local and nonlocal filesystems.

Performance. Remote I/O can reduce total wall-clock time by allowing overlapping of data transfer and computation. In a best-case situation, overlapping can reduce execution time by a factor of two. For example, a climate model may require 8 hours to perform a 10-year simulation and produce 29 GB of output data, which at 10 Mb/sec takes 6.4 hours to transfer. If, on the other hand, computation and communication can be overlapped, total turnaround time is reduced from 14.4 to 8 hours.

Flexibility. Remote I/O provides a higher-level specification of I/O operations than does staging and hence permits greater flexibility in terms of how I/O is performed. For example, a program may need to access just selected components of remote data sets. If the identity of those data elements is computed during program execution, a staging approach often transfers more data than is necessary, wasting both disk and network resources. In contrast, a remote I/O library can choose to stage (i.e., prefetch) the entire dataset or transfer only required elements directly to memory, depending on available resources. Hibbard et al. [11] prototyped the latter strategy in the I-WAY networking experiment, fetching data from a remote IBM SP data server only when a user zoomed in on a particular area within a virtual reality browser.

Convenience. Remote I/O allows programs to execute at remote sites without programmer management of data transfer. In contrast, staging can require that the user learn details of remote filesystems, transfer data among potentially complex directory hierarchies, translate data formats, and manage multiple copies of their datasets. Norman et al. [17] report that such issues were a major source of complexity in their distributed simulations of galactic collisions. In a different area, conversations with users reveal that some are uncomfortable leaving sensitive data in local file systems, but are happy to transfer such data over networks to an application, perhaps over a secure network. Remote I/O makes this transfer possible without requiring that data be encrypted prior to writing it to files at a remote site.

2.2 Wide-Area Computing Issues

Remote I/O libraries, like parallel I/O libraries, must orchestrate efficiently the transfer of data from a user application running on multiple processors to a filesystem. Remote I/O is complicated, however, by several issues not encountered in typical parallel computing environments.

Performance Characteristics. A remote I/O library running over a continental network can see a combined roundtrip communication and I/O latency of 100 msec. This is three to four orders of magnitude more than the roundtrip communication time found in a typical parallel computer (tens or hundreds of microseconds) and one order of magnitude more than the typical time for an I/O node to perform a disk access on behalf of a compute node (~10 msec). In addition, the bandwidth offered by the network over which a remote I/O library operates may be (but is not always) significantly lower than the internal communication network of a parallel computer and/or the remote filesystem. The connectivity of the remote I/O network is also often low (e.g., a single fiber). In contrast, parallel computers typically offer many paths from processors to disks.

Heterogeneity and Configuration. The computer, network, and storage systems used by a remote I/O system often include a heterogeneous mixture of hardware, software, and protocols. In such environments, selecting optimal I/O strategies is more difficult than in the relatively homogeneous environments in which parallel I/O libraries typically operate. A related issue is that the quality of service (QoS: e.g., average bit rate, or reliability) offered by the remote I/O network may be extremely variable, in which case we may require specialized techniques to shield an application from

this variability, for example, buffering, or creating local copies to avoid loss of data over unreliable links. Alternatively, QoS may be controllable by a remote I/O library or user application, in which case accurate estimates of QoS requirements can improve both overall application performance and resource utilization.

Naming and Security. Remote I/O systems often connect computers and filesystems located in different administrative domains. This causes difficulties for both naming and security. We require a global name space for files, but different sites will use different filesystem structures. While distributed filesystems such as AFS create a global structure, we may not have that luxury. Resource Locators (URLs) represent an alternative approach. In addition, authentication, authorization, and privacy all become problematic issues in a distributed environment.

2.3 Approaches to Remote Data Access

Past approaches to the remote data access problem fall into three general categories: distributed filesystems, parallel filesystems, and remote execution systems.

Traditional distributed filesystems (NFS [18] to some extent, and AFS [15] and DFS to a greater extent) provide a convenient interface for remote I/O: a uniform file name space is provided, and files are accessed with conventional read and write statements. However, these systems typically do not achieve good performance for high-performance computing workloads: they were designed primarily for a different class of users (e.g., software developers). For example, NFS bandwidth over an Ethernet LAN may be 1-3 Mb/sec, but an optimized communication library can achieve close to 10 Mb/sec. The lack of explicit interfaces for collective I/O also hinders performance optimization. In addition, they introduce significant implementation complexity and administrative overhead, which tend to hinder their widespread deployment. Web-based distributed file systems [1, 24] reduce implementation and administration costs but do not improve performance. Data servers such as DPSS [22] and MARS [4] use networked disk servers to provide high-speed streaming access to distributed data, but do not support access from parallel programs.

Parallel filesystems (e.g., [16, 7]) and I/O libraries (e.g., [3, 6, 19]) address performance issues by defining I/O interfaces that allow identification and optimization of collective I/O operations, by incorporating specialized buffering techniques, by supporting asynchronous operations, and by incorporating techniques (e.g., disk-directed [13], server-directed [19], and two-phase [20] I/O) for transferring data efficiently from compute nodes to disks. However, these systems are not designed to address the complex configurations, unique performance tradeoffs, and security problems that arise in wide area environments.

Remote execution systems (Condor [14] is one example) redirect Unix filesystem calls to a home filesystem, hence enabling location-independent execution of tasks scheduled to remote computers. However, these systems do not support parallel I/O interfaces or access to parallel filesystems.

In summary, what is lacking is an approach that provides the high-performance characteristics of parallel I/O libraries while addressing the unique requirements of networked environments. This is the goal of our remote I/O work.

3 The RIO Remote I/O Library

To support our investigations of remote I/O, we have developed a remote I/O library called RIO. In this section, we describe how RIO addresses issues of portability, performance, and integration with wide area computing environments.

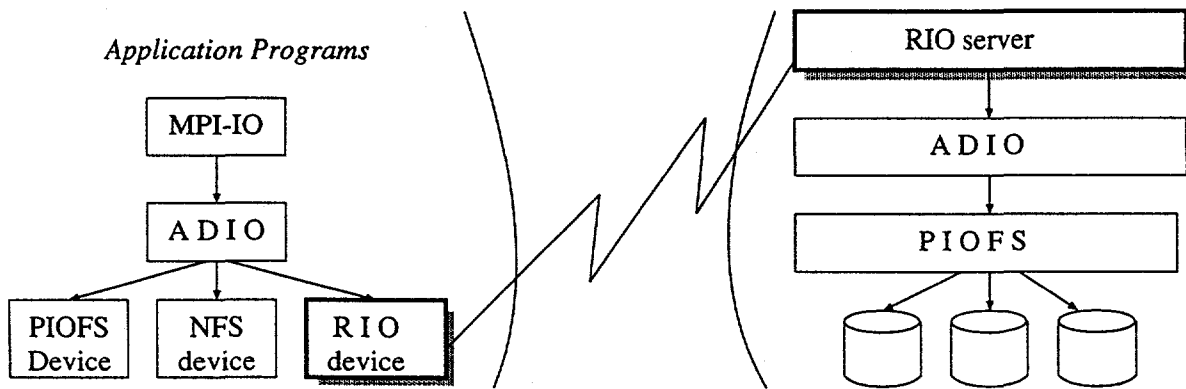


Figure 1: RIO architecture, showing how RIO layers below ADIO in the client and above ADIO in the server

3.1 Portability

An I/O library is most useful if it supports both a wide range of application I/O patterns and multiple filesystems. It is not our goal to innovate in the area of I/O interfaces, and so we adopt the quasi-standard MPI-IO [6]. This interface incorporates support for collective operations, asynchronous operations, and other I/O abstractions that have been found useful for high-performance parallel I/O. Whether the requirements of remote I/O can motivate modifications or extensions to the MPI-IO interface remains to be seen, but our initial approach is to use MPI-IO unchanged.

Portability is a challenging problem in a remote I/O library because there may be no commonality in architecture between the computer on which an application runs and the potentially many remote filesystems that the program accesses. We address the portability problem by exploiting features of the ADIO implementation of MPI-IO [21]. ADIO adopts a modular design in an attempt to maximize code reuse across filesystems. High-level I/O libraries (in our case, MPI-IO) invoke services provided by a set of ADIO “devices,” each providing low-level support for a particular I/O system (e.g., Unix, Intel PFS, IBM PIOFS).

As illustrated in Figure 1, RIO exploits the ADIO framework in two ways. On the client side, we provide a RIO device that implements ADIO calls as interactions with remote RIO servers. The servers themselves also use ADIO calls, in this case to access the remote filesystem in a system-independent fashion. This approach of simultaneously layering *below* ADIO (on the client side) and *above* ADIO (on the server side) greatly reduces implementation costs. On the client side, we need not implement all of MPI-IO nor be concerned with remote filesystem details. Instead, we can focus our attention on a small number of portable low-level functions. On the server side, we can operate on any system supported by ADIO.

3.2 Performance

A remote I/O library can use various strategies to transfer data between client and server. Research in parallel I/O has identified collective operations, non-blocking operations, and buffering as important techniques for maximizing performance on parallel filesystems. However, the characteristics of networked systems listed in Section 2.2 lead to different tradeoffs.

Our RIO prototype uses the Nexus communication library [10] for client-server communications; Nexus, MPI, or potentially other communication mechanisms may be used within an application.

When opening a file, a designated client process first attempts to connect to a server gateway process. The client and server then exchange information about file type and file access patterns, and the server issues an ADIO open call to open the relevant file(s). The client and server then establish the communication structure to be used for subsequent read and write operations. Finally, both client and server establish local data structures representing the open file; on the client side, a "file descriptor" is returned, encoding a reference to the client-side data structure.

Let P_C denote the number of processes executing at the client and P_S the number at the server. Following an open call, each client process can read and write at a distinct location in the file, either independently or as part of a collective operation involving multiple client processes. In a simple RIO implementation, each client process keeps track of its own location within the file, and implements a read or write operation as a remote procedure call (i.e., a round-trip communication) to a server process. A round trip is required even for write calls, in order to provide a return code.

An analysis of the various inefficiencies inherent in this simple approach allows us to introduce some of the optimizations used in RIO.

Forwarder Nodes. Client and server processes communicate directly. A disadvantage of this strategy is that a single process may have to use two communication methods: e.g., on the IBM SP, a vendor-supplied MPI library and TCP/IP. This simultaneous use can introduce significant overheads due to the need to manage two communication interfaces [8] or may be disallowed entirely if network interfaces can be accessed only from dedicated service nodes. Hence, we introduce *forwarder nodes* (analogous to the dedicated I/O nodes used in some I/O systems), to which each client process forwards communications destined for the server, and which handles communications from the server to client processes. These forwarder nodes must use both MPI and TCP, but are dedicated and hence can be optimized for this purpose. The forwarder nodes can also be used to throttle traffic to avoid network saturation [23]. In our current work, the client and server each use a single forwarder. However, multiple forwarders can be advantageous if there are multiple network interfaces or if compression, message digest, or encryption techniques are to be applied to data.

Exploitation of Collective Operations. Each client process communicates independently with the server, even when engaged in a collective operation. Hence, a single client-side collective call requires P_C messages and results in P_C independent I/O operations at the server. Both the multiple communications and multiple I/O operations can be inefficient in some situations. Multiple communications can be avoided by collecting the communications performed by the P_C clients (e.g., at the forwarder) and transferring them to the client in a single message. Multiple server I/O operations can be avoided by tagging client messages to indicate when they refer to collective calls, and then invoking a collective I/O operation at the server. The latter strategy is straightforward if $P_S = P_C$, since the server can issue open, read, and write operations identical to those performed by the client. The situation remains straightforward if P_C is an integer multiple of P_S , or vice versa, as the calls issued by the client are easily mapped to server processes. In other situations, it can be hard to translate a client-side collective operation into an efficient collective operation at the server. Our RIO prototype assumes that $P_S = P_C$.

Reduction of Round-Trip Messages. The round trip performed for each read and write operation can take 100 msec or more in a wide area environment, significantly more than an I/O operation. RIO seeks to reduce these costs by incorporating support for asynchronous I/O operations. Asynchronous operations allow several I/O operations to be outstanding at once, hence enabling pipelining of I/O operations in the network and I/O system, and overlapping of computation and I/O in the application. Another approach that we have yet to evaluate is to reduce the number of communication operations

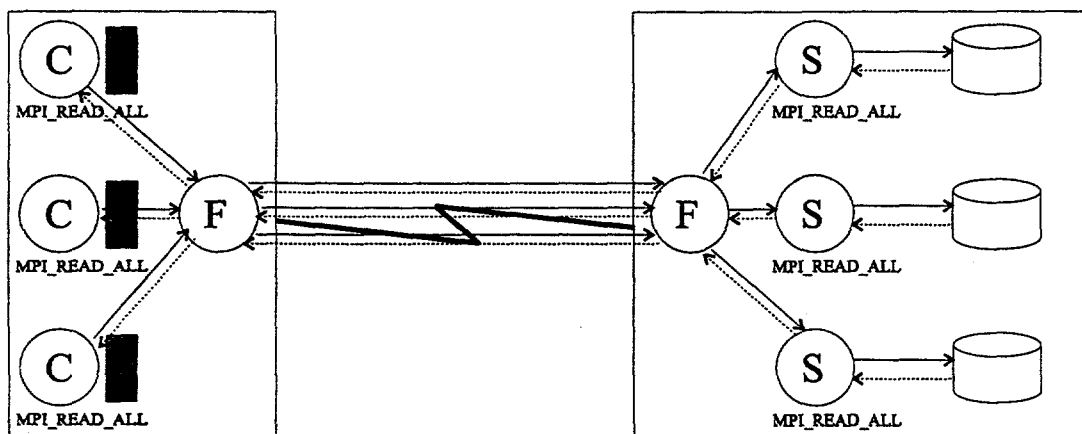


Figure 2: RIO's optimized I/O strategy, showing the client (C), forwarder (F), and server (S) processes, and the communications performed following a collective read operation.

by client-side buffering, either independently by each client process, collectively by multiple client processes, or by the forwarder.

Figure 2 outlines the structure that results when these various optimizations are introduced. The figure shows the client-side buffers (here associated with client processes), the forwarder processes, and the translation of a client-side collective I/O call `MPI_READ_ALL` into a collective I/O call at the server.

3.3 Integration

Because RIO is designed to execute in a wide area environment, its implementation must address issues of naming, configuration, and security. In the following, we explain how these issues can be addressed by using mechanisms provided by the Globus distributed computing toolkit [9].

Naming. RIO uses a URL-like notation to provide a uniform name space for files. A file is opened with a call of the form

```
MPI_Open(..., "x-rio://host-name:port-num/pathname", ...)
```

where the *host-name* and *port-num* identify a RIO server and *pathname* identifies a file managed by that server. In the future, we may substitute Uniform Resource Names (URNs) for URLs, to permit location-independent naming of cachable or replicated resources such as databases.

Configuration. RIO permits the use of Globus configuration mechanisms. For example, when establishing a Nexus connection between client and server, RIO can use the Globus Metacomputing Directory Service (MDS) to determine which networks are available, their current load, and access mechanisms. RIO also can interact with Globus schedulers to reserve capacity on networks that support quality of service negotiation.

Security. A remote I/O system may be required to verify a user's identity (authentication), to determine whether and how a user is able to access a file (authorization), and to ensure the integrity

and privacy of data transferred over public networks. We design RIO to incorporate the solutions to these problems provided by Globus.

The current Globus system supports a global "Globus id" but requires that a user have an account at a site before it can use that site's resources. Globus provides a cryptographically secure mapping from Globus id to local ids, hence allowing a user to authenticate once (to Globus) and subsequently access resources at any Globus site where the user has an account. These mechanisms can easily be adapted for use by RIO. Authentication is performed by Nexus when a RIO client connects to a RIO server. If authentication succeeds, the local user id of the Globus user is also established, and hence the file access rights of the Globus user at that site are determined. Once authentication is in place, Globus/Nexus mechanisms can be used to apply digital signatures for message integrity and/or encryption for privacy. If desired, these mechanisms can be applied only when communicating over networks defined to be insecure.

In the longer term, we expect Globus (and hence RIO) to eliminate the requirement that a user have a local account at every site. Access control lists are one approach to authorization in this regime. Cryptographically signed "use condition certificates" [12] represent another promising approach.

4 Experimental Studies

We report on experiments designed to determine the basic performance characteristics of RIO and to provide a preliminary evaluation of RIO's utility for applications. These experiments comprise a series of microbenchmarks similar to those used traditionally for evaluation of I/O library performance, plus a single application.

4.1 Experimental Platform

In selecting an experimental platform, we must trade off our interest in exploring true remote I/O against the need for a controlled environment in which the impacts of different performance issues can be easily measured. These considerations motivate us to define a testbed comprising two partitions of the same IBM SP multicomputer. Within each partition, communication can occur via vendor-supplied MPI, while TCP/IP is used between partitions. The client runs in one partition and the server in the other. Because of our use of forwarder nodes, this simple configuration has performance characteristics quite similar to two IBM SPs connected by a high-bandwidth local or metropolitan area network. While intrapartition communication peaks at over 30 MB/sec with latencies of around 50 μ sec, interpartition communication peaks at 8 MB/sec with latencies of around 2000 μ sec.

All experiments were performed on the IBM SP2 at the Cornell Theory Center and used IBM PIOFS version 1.2 as the "remote" filesystem. All nodes used in our experiments were SP thin nodes (roughly equivalent to RS/6000 Model 390, with at least 128 MB memory) running AIX 4.1. PIOFS distributes files across multiple PIOFS servers [2]. At Cornell, there are eight such servers. Each file consists of a set of cells, and each cell is stored on a particular server node. The default number of cells is the number of PIOFS servers; if the number is greater, cells are striped across servers in a round-robin fashion. A file is divided into basic striping units (BSUs), which are assigned to cells in a round-robin fashion. The default BSU size is 32 KB. In some situations, tuning of these various parameters can significantly affect performance. We used default values in all experiments.

PIOFS performance is sensitive to the size of the data being read and written. For small (< 8 KB) accesses, access time is about 4 to 5 msec, presumably because of the round-trip communications between the node performing I/O and the PIOFS server nodes. High performance can be achieved for large read and write sizes.

4.2 Microbenchmark Results

Our microbenchmarks are designed to reveal how RIO read and write bandwidths vary as functions of P_C (note that $P_S = P_C$ in all experiments) and read or write size. Each microbenchmark uses blocking operations to transfer contiguous data from/to a single shared file.

Figure 3 shows the transfer rates (totals summed over P_C processes) that we measured for $P_C = P_S = 1, 2, \text{ and } 4$, and for different access sizes. We give results for RIO and for PIOFS only; our PIOFS results match those of other researchers. We also include in the graph horizontal lines representing the bandwidth measured with two simple ping-pong programs. The line labeled "Client/server forwarding" was obtained with a program that bounces large messages between a client process and a server process, via the intervening forwarders. Hence, it approximates the best data transfer rate that can be obtained for synchronous operations between a single client and a single server in our architecture. The line labeled "TCP peak" was obtained with a simpler program that uses TCP to bounce large messages back and forth between two processes. This line approximates the best data transfer rate that can be achieved with Nexus and TCP on the IBM SP. The first number (4.25 MB/sec) is less than the second (8.6 MB/sec) because a round-trip client/server communication involves six messages, as compared to just two in the latter case.

Examining Figure 3, we see that sustained bandwidth generally increases with P_C , but peaks at around 4.2 MB/sec for the larger transfer sizes, when we saturate the forwarder-to-forwarder connection. We are able to exceed 4.25 MB/sec slightly in some situations because of pipelining of communications. Other preliminary experiments show modest (10-20 percent) improvements in RIO performance when nonblocking operations are used, because pipelining is enhanced.

These results show that RIO is able to drive the network connecting client and server at close to its peak bandwidth, at least for large messages. We see also that in our experimental configuration, the principal obstacle to improved performance is the capacity of this network. Faster networks and improved forwarder structures are two possible approaches to improving performance.

4.3 Application Results

We use the BTIO benchmark from the NAS I/O benchmark suite [5], specifically, the program `BTIO-simple-mpiio`. This benchmark simulates the I/O required by a pseudo-timestepping flow solver. It implements an approximate factorization algorithm with the requirement that after every k iterations, the three-dimensional solution vector (of size N^3) is written to a disk file (no reads are performed). A total of I iterations of the algorithm are performed. The application code is in Fortran and uses the MPI-IO interface to write output data to a single file.

In our experiments, we fix $k = 5$ and $I = 200$ and consider problem sizes $N = 32$ and $N = 64$; total data written in these two cases is 52 MB and 420 MB, respectively. We define the elapsed time as the wall clock execution time for the application, and the application sustained I/O transfer rate as (total amount of I/O performed)/(elapsed time). The elapsed time includes both the time for computation and I/O inside the application. Note that this I/O transfer rate is different from that measured in the microbenchmarks, where no significant computation is performed. We use four application processors in all experiments, for a total of 10 processors, with 2 forwarding nodes and 4 server nodes. We have observed similar behavior with different numbers of application processors, and hence for brevity we do not report those results.

`BTIO-simple-mpiio` performs many small writes in an irregular pattern and hence performs poorly on PIOFS, due to the high PIOFS overhead associated with small writes. Hence, we produced a modified version of the benchmark that redistributes the output data before the solution vector is written to disk. In the redistribution code, each node essentially collects data from other nodes into

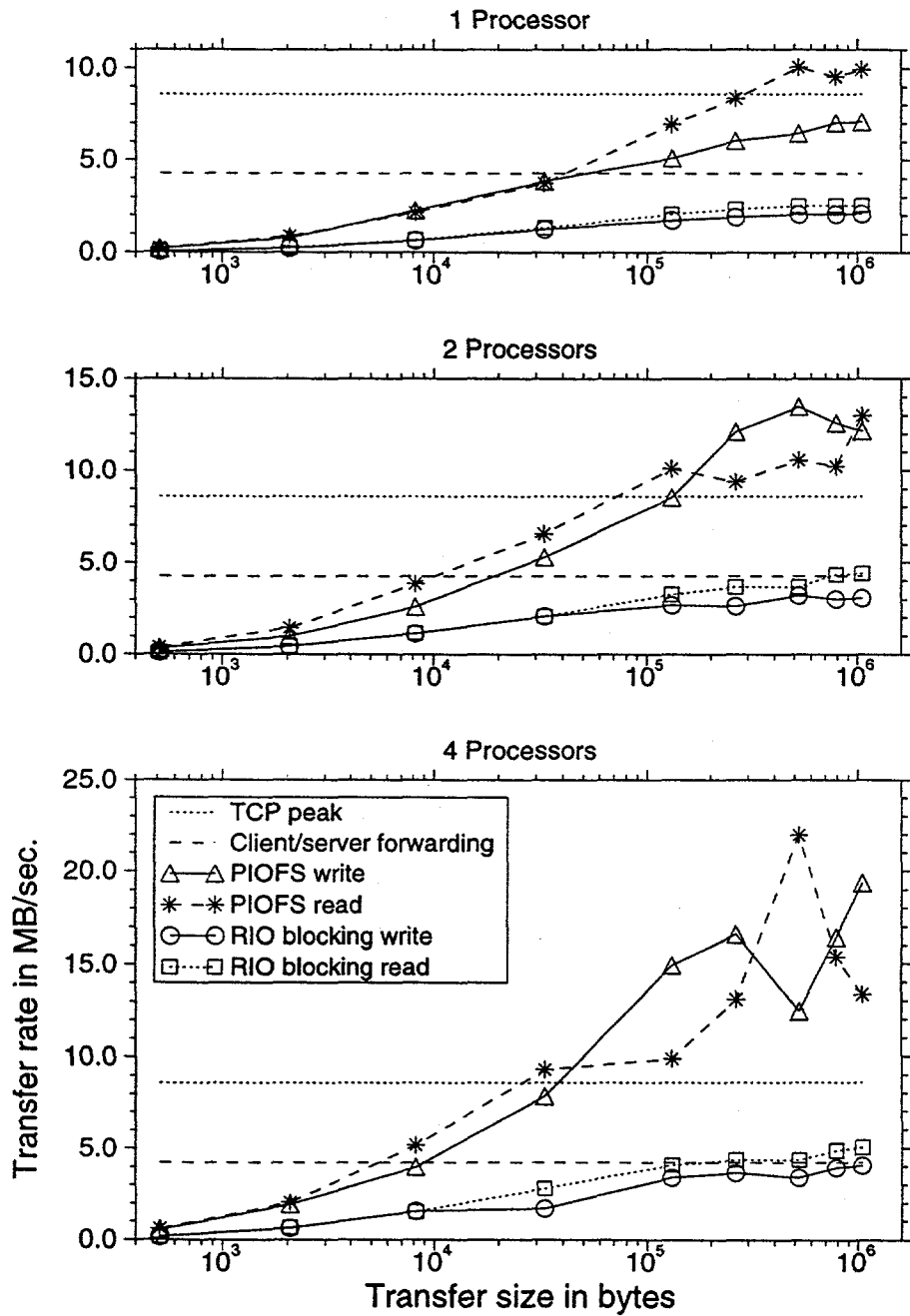


Figure 3: Microbenchmark performance results

Table 1: Execution times for the original and optimized BTIO application (secs)

Version	N	Local PIOFS	PIOFS+ftp	RIO (blocking)	RIO (nonblock)
Original	32	285.28	310.46	457.39	314.54
Original	64	1711.48	1912.94	2278.97	1886.00
Optimized	32	178.68	203.86	186.07	180.99
Optimized	64	1446.05	1647.51	1524.20	1389.35

Table 2: Application sustained I/O transfer rates for the original and optimized BTIO (MB/sec). Note that these rates are computed based on the total elapsed times, not just the I/O times

Version	N	Local PIOFS	PIOFS+ftp	RIO (blocking)	RIO (nonblock)
Original	32	0.1838	0.1689	0.1146	0.1667
Original	64	0.2451	0.2193	0.1840	0.2224
Optimized	32	0.2934	0.2572	0.2818	0.2897
Optimized	64	0.2900	0.2546	0.2752	0.3019

a temporary contiguous write buffer; each processor then performs a single write operation at each dump.

We measure elapsed time for four configurations: when using PIOFS directly (i.e., without using RIO); when using RIO (blocking calls) to transfer data from the application to the RIO server, which then makes the PIOFS calls; when using RIO (nonblocking calls); and when data is first written to PIOFS directly, without RIO, and then transferred to a user filesystem with ftp. The latter configuration corresponds to the use of staging. In the nonblocking version of the original code, multiple (up to 64) I/O operations may be outstanding. In the nonblocking RIO version of the optimized code, we issue an asynchronous write for each dump and then proceed with computation, waiting for completion only at the start of the next dump. This is possible because writes are performed from the write buffer.

Tables 1 and 2 show the elapsed times and application sustained transfer rates measured for both the original and optimized versions of the program. The optimized version of BTIO performs better than the original, due to the reduced number of write operations. We see that nonblocking calls significantly affect performance in all cases. This is because in the absence of nonblocking calls, the round-trip message exchange between application and server is a significant source of overhead. When nonblocking operations are used, performance improves either because multiple I/O operations are pipelined (in the original code) or because I/O and round-trip overheads are overlapped with computation (in the optimized code). As a result, throughput is close to what we get when accessing PIOFS directly: in fact, because PIOFS does not support nonblocking operations, RIO performs better than local PIOFS in some cases. The maximum application sustained transfer rate is 0.29 MB/sec for local PIOFS and 0.30 MB/sec with RIO. If we consider only I/O time (and hence compute burst I/O rates), we obtain I/O rates of 2-4 MB/sec.

Finally, we see that the total execution time when using RIO is, in most cases, less than the total turnaround time when staging is used (PIOFS+ftp). For the optimized code with $N = 64$, RIO is 19 percent faster. This result is due to the overlapping of computation and data transfer achieved by RIO and illustrates how remote I/O can improve application performance as well as providing a more convenient interface.

5 Conclusions

We have argued for the importance of remote I/O as a tool for high-performance, low-overhead distributed computing. Remote I/O libraries allow programs to use familiar parallel I/O interfaces to access data contained in remote filesystems. In principle, they can improve performance, enhance flexibility, and reduce complexity in applications that must access nonlocal data. We have identified some of the challenges that must be overcome before these benefits can be realized; these include high latencies, low bandwidths, complex configurations, and security. We have also described a prototype remote I/O library called RIO that incorporates solutions to some of these problems. Performance experiments in a controlled multicomputer environment show that RIO introduces little overhead and can achieve improved turnaround time compared to remote execution combined with staging.

The work presented here is just a first step toward a truly usable remote I/O facility for high-performance computing applications. Our next step will be to deploy the RIO prototype in a wide area computing testbed. Our first target is the sites connected by the ESnet and CAIRN networks, in particular Argonne, Berkeley, and USC/ISI. This environment will enable us to evaluate our techniques more realistically and will also support experiments with network quality of service reservation. We also plan detailed comparisons with distributed filesystems for a range of scientific applications.

Acknowledgments

We thank David Lifka, Rajeev Thakur, and Steven Tuecke for their invaluable assistance with this work, and the Cornell Theory Center and Argonne National Laboratory's Center for Computational Science and Technology for access to their IBM SP systems.

References

- [1] A. D. Alexandrov, M. Ibel, K. E. Schausser, and C. J. Scheiman. Extending the operating system at the user level: The UFO global file system. In *1997 Annual Technical Conference on UNIX and Advanced Computing Systems (USENIX'97)*, January 1997.
- [2] F. Bassow. *IBM AIX Parallel I/O File System: Installation, Administration, and Use*. IBM, Kingston, N.Y., May 1995. Document Number SH34-6065-00.
- [3] R. Bennett, K. Bryant, A. Sussman, R. Das, and J. Saltz. Jovian: A framework for optimizing parallel I/O. In *Proceedings of the 1994 Scalable Parallel Libraries Conference*, pages 10–20. IEEE Computer Society Press, October 1994.
- [4] M. Buddhikot, G. Parulkar, and J. Cox. Design of a large scale multimedia storage server. In *Proc. INET '94*, 1994.
- [5] R. Carter, B. Ciotti, S. Fineberg, and B. Nitzberg. NHT-1 I/O benchmarks. Report RND-92-016, NAS, NASA Ames Research Center, Moffett Field, CA, Nov 1992.
- [6] P. Corbett, D. Feitelson, Y. Hsu, J.-P. Prost, M. Snir, S. Fineberg, B. Nitzberg, B. Traversat, and P. Wong. MPI-IO: A parallel file I/O interface for MPI. Technical Report NAS-95-002, NAS, NASA Ames Research Center, Moffett Field, CA, January 1995. Version 0.3.
- [7] P. Corbett, D. Feitelson, J.-P. Prost, and S. Baylor. Overview of the Vesta parallel file system. In *IPPS '93 Workshop on Input/Output in Parallel Computer Systems*, pages 1–16, 1993.

- [8] I. Foster, J. Geisler, C. Kesselman, and S. Tuecke. Managing multiple communication methods in high-performance networked computing systems. *Journal of Parallel and Distributed Computing*, 40:35-48, 1997.
- [9] I. Foster and C. Kesselman. Globus: A metacomputing infrastructure toolkit. *International Journal of Supercomputer Applications*, 1997. To appear.
- [10] I. Foster, C. Kesselman, and S. Tuecke. The Nexus approach to integrating multithreading and communication. *Journal of Parallel and Distributed Computing*, 37:70-82, 1996.
- [11] W. Hibbard, J. Anderson, I. Foster, B. Paul, C. Schafer, and M. Tyree. Exploring coupled atmosphere-ocean models using Vis5D. *International Journal of Supercomputer Applications*, 10(2):211-222, 1996.
- [12] W. Johnston and C. Larsen. A use-condition centered approach to authenticated global capabilities: Security architectures for large-scale distributed collaboratory environments. Technical Report 3885, LBNL, 1996.
- [13] D. Kotz. Disk-directed I/O for an out-of-core computation. In *Proc. 4th IEEE Symp. on High Performance Distributed Computing*, pages 159-166, August 1995.
- [14] M. Litzkow, M. Livney, and M. Mutka. Condor - a hunter of idle workstations. In *Proc. 8th Intl Conf. on Distributed Computing Systems*, pages 104-111, 1988.
- [15] J.H. Morris et al. Andrew: A distributed personal computing environment. *Communications of the ACM*, 29(3), 1986.
- [16] N. Nieuwejaar and D. Kotz. Performance of the Galley file system. In *Proceedings of the Fourth Annual workshop on I/O in Parallel and Distributed Systems*, Philadelphia, PA, May 1996.
- [17] M. Norman, P. Beckman, G. Bryan, J. Dubinski, D. Gannon, L. Hernquist, K. Keahey, J. Ostriker, J. Shalf, J. Welling, and S. Yang. Galaxies collide on the I-WAY: An example of heterogeneous wide-area collaborative supercomputing. *International Journal of Supercomputer Applications*, 10(2):131-140, 1996.
- [18] R. Sandberg, D. Goldberg, S. Kleiman, D. Walsh, and B. Lyon. Design and implementation of the Sun network filesystem. In *Proc. Summer USENIX*, pages 119-130, June 1985.
- [19] K. E. Seamons, Y. Chen, P. Jones, J. Jozwiak, and M. Winslett. Server directed collective I/O in Panda. In *Proceedings of Supercomputing '95*, San Diego, California, December 1995.
- [20] R. Thakur and A. Choudhary. An extended two-phase method for accessing sections of out-of-core arrays. *Scientific Programming*, 5(4):301-317, Winter 1996.
- [21] R. Thakur, W. Gropp, and E. Lusk. An abstract-device interface for implementing portable parallel-I/O interfaces. In *Proceedings of The 6th Symposium on the Frontiers of Massively Parallel Computation*, October 1996.
- [22] B. Tierney, W. Johnston, L. Chen, H. Herzog, G. Hoo, G. Jin, and J. Lee. Distributed parallel data storage systems: A scalable approach to high speed image servers. In *Proc. ACM Multimedia 94*. ACM Press, 1994.
- [23] B. Tierney, W. Johnston, J. Lee, and G. Hoo. Performance analysis in high-speed wide area IP over ATM networks: Top-to-bottom end-to-end monitoring. Technical report, LBNL, 1996.

- [24] A. Vahdat, P. Eastham, and T. Anderson. WebFS: A global cache coherent filesystem. Technical report, Department of Computer Science, UC Berkeley, 1996.