

# Efficient Structured Data Access in Parallel File Systems

Avery Ching Alok Choudhary Wei-keng Liao  
*Electrical and Computer Engineering Department*  
*Northwestern University*  
*Evanston, IL 60208*  
{aching, choudhar, wkliao}@ece.nwu.edu

Robert Ross William Gropp  
*Mathematics and Computer Science Division*  
*Argonne National Laboratory*  
*Argonne, IL 60439*  
{rross, gropp}@mcs.anl.gov

## Abstract

*Parallel scientific applications store and retrieve very large, structured datasets. Directly supporting these structured accesses is an important step in providing high-performance I/O solutions for these applications. High-level interfaces such as HDF5 and Parallel netCDF provide convenient APIs for accessing structured datasets, and the MPI-IO interface also supports efficient access to structured data. However, parallel file systems do not traditionally support such access.*

*In this work we present an implementation of structured data access support in the context of the Parallel Virtual File System (PVFS). We call this support “datatype I/O” because of its similarity to MPI datatypes. This support is built by using a reusable datatype-processing component from the MPICH2 MPI implementation. We describe how this component is leveraged to efficiently process structured data representations resulting from MPI-IO operations. We quantitatively assess the solution using three test applications. We also point to further optimizations in the processing path that could be leveraged for even more efficient operation.*

## 1. Introduction

Scientific applications have begun to rely heavily on high-level I/O APIs such as HDF5 [7] and parallel netCDF [9] for their storage needs. These APIs allow scientists to describe their data in terms that are meaningful to them, as structured, typed data, and to store and retrieve this data in a manner that is portable across all the platforms they might find useful. Because scientists have a richer language with which to describe their data, I/O for an application as a whole can be described in terms of the datatypes and organizations that the scientist is really using, rather than posing I/O operations in

terms of independent reads or writes of bytes on many processors.

These APIs also allow I/O experts to embed the knowledge of how to efficiently access storage resources in a library that many applications can use. The result is a big win for both groups. Implementors of high-level I/O libraries in turn use MPI-IO as their interface to storage resources. This lower-level interface maps higher-level accesses to file system operations and provides a collection of key optimizations. MPI-IO also understands structured data access, providing high-level I/O API programmers the ability to describe noncontiguous accesses as single units, just as the scientist did, and to interface to underlying resources through a portable API.

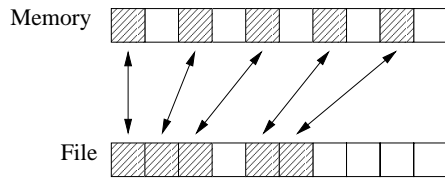
Today’s parallel file systems do not, for the most part, support structured or even noncontiguous accesses. Instead they support the POSIX interface, allowing for only contiguous regions to be accessed and modified. This approach severely limits the ability of the MPI-IO layer to succinctly and efficiently perform the accesses that have been described by these higher layers.

A significant step in the direction of efficient structured data access is the *list I/O* interface [15], implemented in the Parallel Virtual File System (PVFS) and supported under MPI-IO [4, 3]. This new interface, when well supported by the parallel file system, allows structured accesses to be described and serves as a solid building block for an MPI-IO implementation. However, because it does not retain any information on the regularity of access, such as stride information, the representation of structured accesses can be very large. Building, transmitting, and processing this representation can significantly limit performance when accesses consist of many small regions [17].

In this work we investigate the next step in efficient support for structured data access. The approach, *datatype I/O*, provides a mechanism for mapping MPI datatypes (passed to MPI-IO routines) into a type representation understood by the file system. This new representation maintains the concise descriptions possible with MPI type constructors

---

### Multiple I/O Example



**Figure 1. Example POSIX I/O call. Using traditional POSIX interfaces for this access pattern cost five I/O calls, one per contiguous region.**

---

such as `MPI_Type_vector`. This representation is passed over the network to I/O servers, which may process this directly, avoiding the overhead of building lists of I/O regions at the MPI-IO layer, passing these lists over the network as part of the file system request, or processing these lists during I/O.

In Section 2 we explain the existing approaches for performing noncontiguous access, including list I/O. In Section 3 we describe our prototype implementation of datatype I/O, the component on which it is built, how datatype I/O is integrated into the parallel file system, and some limitations of the implementation. In Section 4 we examine the performance of our prototype using three benchmarks: a tile reading application, a 3D block decomposition I/O kernel, and a simulation of the FLASH I/O checkpoint process. In Section 5 we discuss related efforts and future directions for this work.

## 2. Current Noncontiguous I/O Approaches

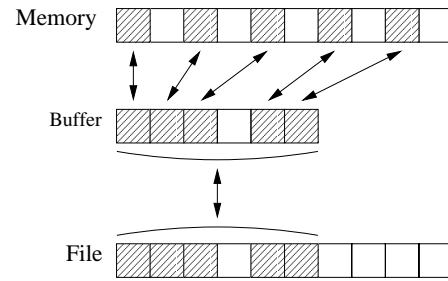
*Noncontiguous* data is simply data that resides in different areas, with gaps between them. Noncontiguous I/O refers to an I/O operation in which data in memory, in file, or in both is noncontiguous. Several approaches have been implemented for supporting noncontiguous I/O access. The first, and most naive, is the approach we call *POSIX I/O*.

### 2.1. POSIX I/O

Most parallel file systems implement the POSIX I/O interface [8]. This interface provides the capability to perform contiguous data access only. To support noncontiguous access with POSIX I/O, one must break the noncontiguous I/O into a sequence of contiguous I/O operations. This approach to noncontiguous I/O access requires significant overhead in the number of I/O requests that must be processed by the underlying file system. As shown in Figure 1, even simple noncontiguous access patterns can re-

---

### Data Sieving I/O Example



**Figure 2. Example data sieving I/O call. By first reading a large contiguous file region into a buffer, data movement is subsequently performed between memory and the buffer.**

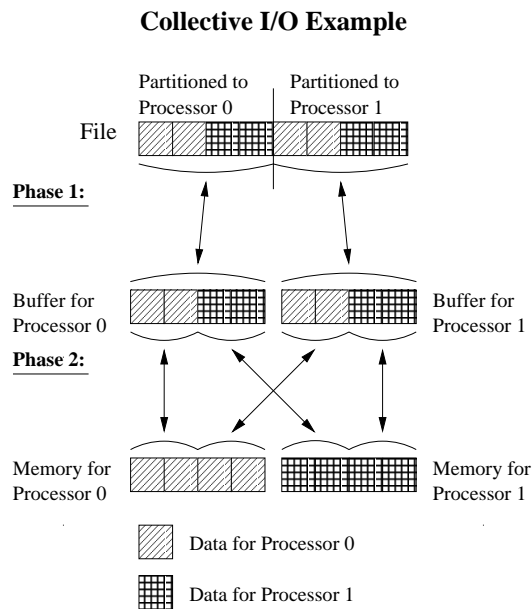
---

sult in numerous contiguous I/O operations. Because operations in parallel file systems often require data movement over a network, latency for I/O operations can be high. For this reason performing many small I/O operations to service a noncontiguous access is very inefficient. Fortunately for users of these file systems, two important optimizations have been devised for more efficiently performing noncontiguous I/O using only POSIX I/O calls: data sieving and two-phase I/O.

### 2.2. Data Sieving I/O

To address the problem of excessive overhead resulting from using POSIX I/O for contiguous I/O access, *data sieving* reduces the number of I/O operations [13]. When using data sieving, a large region encompassing all the data in the file is accessed with a minimum number of POSIX I/O operations. For read operations, a large contiguous data region containing desired data is first read into a temporary buffer, and the desired data is then extracted into the user's buffer. For write operations, a read-modify-write sequence is performed. A large contiguous region is read into a temporary buffer, new data is placed into the appropriate positions in this buffer, and the buffer is then written back to storage. In order to ensure consistency during concurrent operations, a lock must be held on the region to be modified during the read-modify-write process.

This approach is efficient when the desired noncontiguous regions exhibit good spatial locality (i.e., are close together). When data is more dispersed, however, data sieving accesses a great deal of additional data in order to perform the operation and, at some point, becomes less effective than simply using a sequence of POSIX I/O calls.



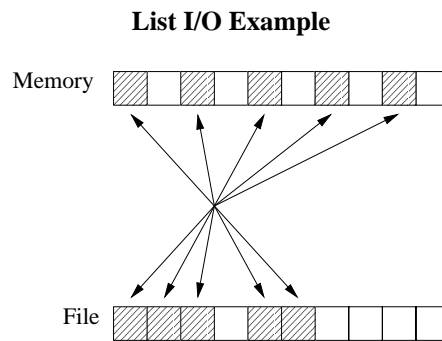
**Figure 3. Example collective I/O call. Interleaved file access patterns can be effectively accessed in larger file I/O operations with the two-phase I/O method.**

### 2.3. Two-Phase I/O

When used to their fullest, interfaces such as MPI-IO give a great deal of information about how the application as a whole is accessing storage. One example of this is the collective I/O calls that are part of the MPI-IO API. By making collective I/O calls, applications tell the MPI-IO library not only that each process is performing I/O but also that these I/O operations are part of a larger whole. This information provides additional opportunities for optimization over application processes performing independent operations.

The *two-phase* I/O optimization, developed by Thakur et al., is one example of a collective I/O optimization [12]. The two-phase method builds on POSIX I/O operations and data sieving. The two-phase method identifies a subset of the processes that will actually perform I/O; these processes are known as *aggregators*. Each aggregator is responsible for I/O to a specific portion of the file; the implementation in ROMIO calculates these regions dynamically based on the size and location of the accesses in the collective operation.

Read operations using the two-phase method are performed as shown in Figure 3. First, aggregators read a contiguous region containing desired data from storage and put this data in a temporary buffer. Next, data is redistributed from these temporary buffers to the final destination pro-



**Figure 4. Example list I/O call. Only a single I/O request is necessary to handle this non-contiguous access because of more descriptive I/O requests.**

cesses. Write operations are performed in a similar manner. First, data is gathered from all processes into temporary buffers on aggregators. Next, this data is written back to storage using POSIX I/O operations. An approach similar to data sieving is used to optimize this write back to storage in the case where there are still gaps in the data. Alternatively, other noncontiguous access methods, such those described in upcoming sections, can be leveraged for further optimization.

Two-phase I/O has a distinct advantage over data sieving alone in that it is significantly more likely to see dense regions of desired data because of combining the regions of many processes. Hence, the reads and writes in two-phase I/O are more efficient than data sieving in many cases. However, two-phase I/O also relies on the MPI implementation providing high-performance data movement. If the MPI implementation is not significantly faster than the aggregate I/O bandwidth in the system, the overhead of the additional data movement in two-phase I/O is likely to prevent it from outperforming the direct access optimizations (data sieving I/O, list I/O, and datatype I/O).

### 2.4. List I/O

The *list I/O* interface is an enhanced parallel file system interface designed to support noncontiguous accesses [15]. List I/O provides an interface capable of describing accesses that are noncontiguous in both memory and file (see prototypes in Figure 5). With this interface an MPI-IO implementation can *flatten* the memory and file datatypes (convert them into lists of contiguous regions) and then describe an MPI-IO operation with a single list I/O call. Given an efficient implementation in the parallel file system, this interface can significantly boost performance. In previous works we discussed the implementation of list I/O in PVFS and

support for list I/O under the ROMIO MPI-IO implementation [4, 3].

The major drawbacks of list I/O are the creation and processing of these large lists and the transmission of these lists from client to server inside the parallel file system layer. Additionally, given that we want to bound the size of I/O requests within the file system, only a fixed number of regions can be described in one request. Thus, while list I/O does significantly reduce the number of I/O operations (in our implementation by a factor of 64), a linear relationship still exists between the number of noncontiguous regions and the number of I/O operations (within the file system layer). Hence, while list I/O is an important addition to the optimizations available under MPI-IO, it does not replace two-phase or data sieving, but rather augments them.

### 3. Datatype I/O

*Datatype I/O* is an effort to address the deficiencies seen in the list I/O interface when faced with accesses that are made up of many small regions, particularly ones that exhibit some degree of regularity. Datatype I/O borrows from the datatype concept that has proven invaluable for both message passing and I/O in MPI applications. The constructors used in MPI datatypes allow for concise descriptions of the regular, noncontiguous data patterns seen in many scientific applications, such as extracting a row from a two-dimensional dataset.

The datatype I/O interface, shown in Figure 6, replaces the lists of I/O regions seen in the list I/O interface with an address, count, and datatype for memory and a displacement, datatype, and offset into the datatype for file. These correspond directly to the address, count, datatype, and offset into the file view passed into an MPI-IO call and the displacement and file view datatype previously defined for the file. The datatype I/O interface is not meant to be used by application programmers; it is an interface specifically for use by I/O library developers. Helper routines are used to convert MPI datatypes into the format used by the datatype I/O functions. A full-featured implementation of datatype I/O would

- maintain a concise datatype representation locally and avoid datatype flattening,
- use this concise datatype representation when describing accesses, and
- service accesses using a system that processes this representation.

Our prototype implementation of datatype I/O was written as an extension to the Parallel Virtual File System. The ROMIO MPI-IO implementation was likewise modified to use datatype I/O calls for PVFS file system operations.

We emphasize while we present this work in the context of MPI-IO and MPI datatypes, nothing precludes our using the same approach to directly describe datatypes from other APIs, such as HDF5 hyperslabs.

#### 3.1. The Parallel Virtual File System and ROMIO MPI-IO Implementations

The Parallel Virtual File System is a parallel file system for commodity Linux clusters [1]. It provides both a cluster-wide consistent name space and user-defined file striping. PVFS is a client-server system consisting of clients, a metadata server, and I/O servers. Clients retrieve a list of the I/O servers that contain the file data from the metadata server at file open time. Subsequent reading or writing is processed directly by the I/O servers without manager interaction.

The approach that PVFS uses for processing requests is detailed in [10]. In short, PVFS builds a data structure called a *job* on each client and server for every client/server pair involved in an I/O operation. This structure points to a list of *accesses*, which are contiguous regions in memory (on a client) or in file (on a server) that must be moved across the network. This is essentially the flattened representation of the datatype being used to move data. While this is not ideal from a processing overhead standpoint, we will retain this representation in this testing; it would be very time consuming to reimplement this component of PVFS.

ROMIO is the MPI-IO implementation developed at Argonne National Laboratory [14]. It builds on the MPI-1 message-passing operations and supports many underlying file systems through the use of an abstract interface for I/O (ADIO). ADIO allows the use of file system specific optimizations such as the list I/O and datatype I/O interfaces described here. Additionally, ROMIO implements the data sieving and two-phase optimizations described in Sections 2.2 and 2.3. It also implements a datatype flattening system that is used to support list I/O for PVFS.

#### 3.2. Datatype I/O Implementation in PVFS and ROMIO

Our datatype I/O prototype builds on the datatype processing component in MPICH2 [11]. Three key characteristics of this implementation make it ideal for reuse in this role:

- Simplified type representation (over MPI datatypes)
- Support for partial processing of datatypes
- Separation of type parsing from action to perform on data

Types are described by combining a concise set of descriptors called *dataloops*. Dataloops can be of five types: contig, vector, blockindexed, indexed, and struct [6]. These five

```

int listio_read(int fd, int mem_list_count, void *mem_offsets[],
               int mem_lengths[], int file_list_count, int file_offsets[],
               int file_lengths[])

int listio_write(int fd, int mem_list_count, void *mem_offsets[],
                 int mem_lengths[], int file_list_count, int file_offsets[],
                 int file_lengths[])

```

Figure 5. List I/O prototypes

```

int dtype_read(int fd, void *mem_addr, int mem_dtype_count, dtype *mem_dtype,
               int file_dtype_disp, int offset_into_dtype, dtype *file_dtype)

int dtype_write(int fd, void *mem_addr, int mem_dtype_count, dtype *mem_dtype,
                int file_dtype_disp, int offset_into_dtype, dtype *file_dtype)

```

Figure 6. Datatype I/O prototypes

### Datatype I/O Example Execution

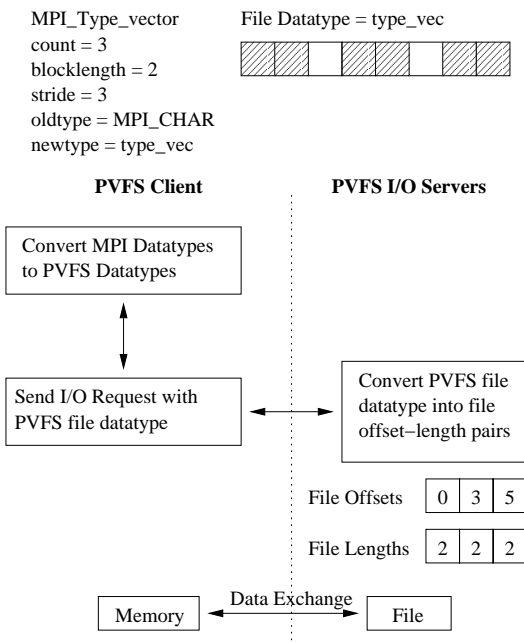


Figure 7. Example datatype I/O call. Since file datatype are broken into file offset-length pairs at the I/O servers, the number of I/O requests is dramatically reduced for regular access patterns.

types capture the maximum amount of regularity possible, keeping the representation concise. At the same time these are sufficient to describe the entire range of MPI datatypes. Simplifying the set of descriptors aids greatly in implementing support for fast datatype processing because it reduces the number of cases that the processing code must handle. The type's extent is retained in this representation (a general concept) while the MPI-specific LB and UB values are eliminated. This simplification has the added benefit of allowing resized type processing with no additional overhead in our representation. We use dataloops as the native representation of types in our PVFS implementation.

The MPICH2 datatype component provides the functions necessary to process dataloop representations. We provide functions to convert MPI datatypes into dataloops and functions that are called during processing to create the offset-length pairs we need to build the PVFS job and access structures. Additionally we provide functionality for shipping dataloops as part of I/O requests. In our prototype, MPI datatypes are converted to dataloops by a recursive process built by using the functions `MPI_Type_get_envelope` and `MPI_Type_get_contents`. By utilizing these MPI functions, we can ensure the portability of our datatype I/O method across different MPI implementations. The resulting dataloop representation is passed into the datatype I/O calls and from there sent to the relevant I/O servers. The dataloops are converted into the job and access structures on servers and clients side to create the traditional PVFS job and access structures. Figure 7 outlines this process. PVFS-specific functions for creating these offset-length pairs are passed to the dataloop processing component. These functions are written to efficiently convert contiguous, vector, and indexed dataloops into offset-length pairs, and they include optimizations to coalesce adjacent regions. The par-

tial processing capabilities of the datatype processing component are used to limit the overhead of storing the intermediate offset-length pairs that are created through dataloop processing.

This is only a partial implementation of the datatype I/O approach, in that a complete approach would avoid the creation of lists of regions on server and client. However, we will show that even without this final capability our prototype exhibits clear performance benefits over the other approaches.

Because the MPI datatypes are converted at every MPI I/O operation into dataloops, we expect there to be slightly higher overhead in the local portion of servicing these operations in comparison to list I/O. On the other hand, because we are concisely describing these types, we expect to see significantly less time spent moving the I/O description across the network.

#### 4. Performance Evaluation

To evaluate the performance of the datatype I/O optimization against the other noncontiguous I/O methods, we ran a series of noncontiguous MPI-IO tests including a three-dimensional block access test, a tile reader benchmark, and the FLASH I/O simulation. For each test we provide a table summarizing the I/O characteristics of the test for each of the tested access methods.

##### 4.1. Benchmark Configuration

Our results were gathered on Chiba City at Argonne National Laboratory [2]. Chiba City has 256 nodes available with dual Pentium III processors, 512 MBytes of RAM, a single 9 GByte Quantum Atlas IV SCSI drive, and a 100 Mbits/sec Intel EtherExpress Pro fast-ethernet card operating in full-duplex mode. Each node is using RedHat 7.3 with kernel 2.4.21-rc1 compiled for SMP use. Our PVFS server configuration for all test cases included 16 I/O servers (one also doubled as a metadata server). PVFS files were created with a 64 KByte strip size (1 MByte stripes across all servers). In the tile reader tests we allocate one process per node because so few nodes are involved. In the other two cases we allocate two processes per node. Our prototype was built using the ROMIO version 1.2.4 and PVFS version 1.5.5. All data sieving and collective operations were conducted with a 4 Mbyte buffer size. Our results are the average of three test runs.

All read benchmarks are conducted with POSIX I/O, data sieving I/O, two-phase collective I/O, list I/O, and datatype I/O. ROMIO can support write operations with data sieving I/O only if file locking is supported by the underlying file system. Since PVFS does not support file locking, we cannot perform data sieving writes on PVFS. We

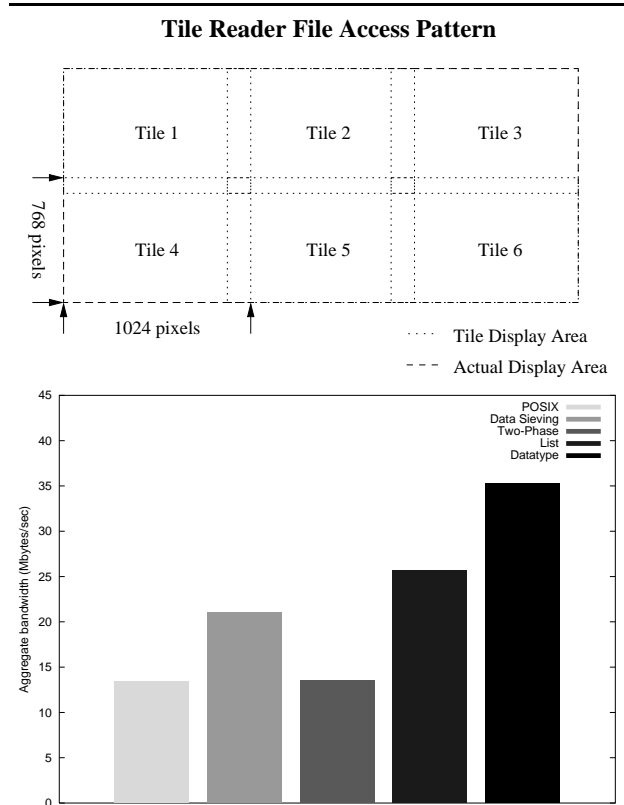


Figure 8. Tile reader access pattern and performance

note, however, that for file systems that do allow file locking, data sieving performance for writes will have worse performance than data sieving reads for the same access pattern, for two reasons. Data sieving writes perform the same data movement from file system into the data sieving buffer at the client and then data movement from data sieving buffer to memory buffer, but they also have to write this data back, thereby doubling the network data transfer of a data sieving read. Also, locking the modified regions can cause serialization of I/O requests for overlapping requests, another serious overhead. On the other hand, the MPI-IO consistency semantics do allow us to perform a read-modify-write during collective I/O. Thus, an approach similar to data sieving is used in the two-phase write case.

##### 4.2. Tile Reader Benchmark

Often the amount of detail in scientific application visualization exceeds the display capabilities of a single desktop monitor. In these cases an array of displays, usually LCDs or projectors, is used to create a more high-resolution display than would otherwise be available. Because of the high

	Desired Data per Client	Data Accessed per Client	# of I/O Ops per Client	Resent Data per Client
POSIX I/O	2.25 MB	2.25 MB	768	—
Data Sieving I/O	2.25 MB	5.56 MB	2	—
Two-Phase I/O	2.25 MB	1.70 MB	1	1.50 MB
List I/O	2.25 MB	2.25 MB	12	—
Datatype I/O	2.25 MB	2.25 MB	1	—

**Table 1. I/O Characteristics of the tile reader benchmark**

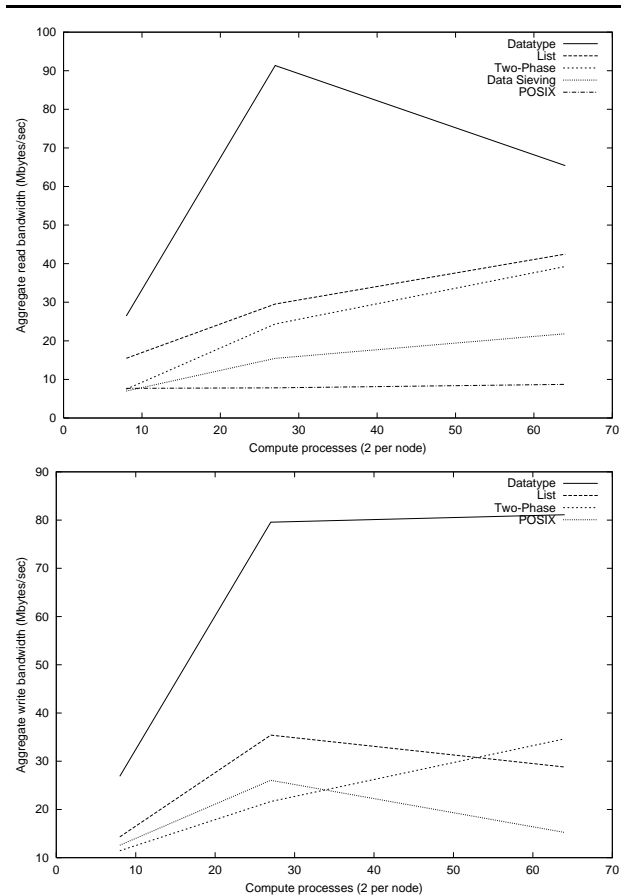
density of these displays, the I/O rates necessary to display data are very high. Typically an array of PCs is used to provide inputs to the individual devices, or *tiles*, that make up the display. The tile reader benchmark is a tool for measuring the rate at which an I/O system can serve data to such a system.

The tile reader benchmark uses a 3 by 2 array of compute nodes (shown in Figure 8) to each display a fraction of the entire frame. By accessing only the tile data for its own display, the compute nodes exhibit a simple noncontiguous file access pattern. The six compute nodes each render a portion of the display with 1024 pixel by 768 pixel resolution and 24-bit color. In order to hide the merging of display edges, there is a 270-pixel horizontal overlap and a 128-pixel vertical overlap between tiles. Each frame is 10.2 MBytes. This data is read into a contiguous buffer in memory. Our tests were conducted with a frame set consisting of 100 frames played back in sequence.

We can see in Figure 8 that datatype I/O is the clear winner in terms of performance for this benchmark, 37% faster than list I/O. The characteristics of the resulting reads using the optimizations tested are shown in Table 1. The datatype I/O result is due to the combination of a single, concise I/O operation and no extra data file data being transferred, and no data passing over the network more than once. In contrast, POSIX I/O requires 768 read operations, data sieving requires more than twice as much data to be read as is desired, two-phase I/O requires resending about 88% of the data read, and list I/O sends a list of 768 offset-length pairs as part of the requests (9 KBytes of total data in I/O requests from each client).

### 4.3. ROMIO Three-Dimensional Block Test

The ROMIO test suite comprises a number of correctness and performance tests. One of these, the `coll.perf.c` test, measures the I/O bandwidth for both reading and writing to a file with a file access pattern of a three-dimensional block-distributed array. The three-dimensional array, shown graphically in Figure 9, has dimensions 600 x 600 x 600 with an element size of an integer (4 bytes on our test platform). Each process reads or writes a single one of these blocks. The memory datatype is contiguous.

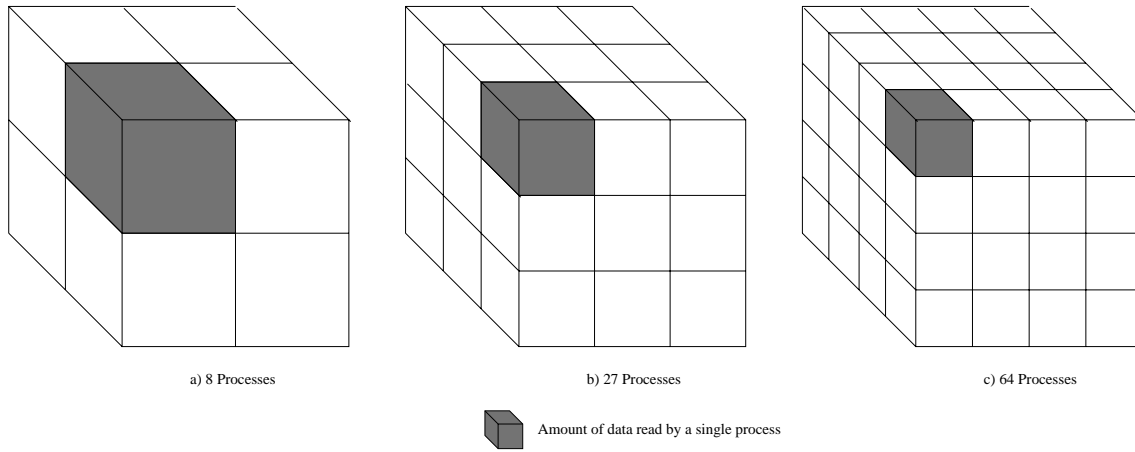


**Figure 10. 3D block read and write performance**

Table 2 characterizes the resulting I/O patterns using our tested optimizations, and Figure 10 shows the results of our tests. Again, datatype I/O is the clear performance winner; peak performance is more than double that of the next-best approach. Of note is the unusual drop in performance in the read case as number of processes increases. We believe that this is due to the increased overhead of offset-length list processing on the server side. Because the servers are the source of data, and clients are operating on a contiguous

	Desired Data per Client	Data Accessed per Client	# of I/O Ops per Client	Resent Data per Client
8 Clients				
POSIX I/O	103 MB	103 MB	90,000	—
Data Sieving I/O	103 MB	412 MB	103	—
Two-Phase I/O	103 MB	103 MB	26	77.2 MB
List I/O	103 MB	103 MB	1408	—
Datatype I/O	103 MB	103 MB	1	—
27 Clients				
POSIX I/O	30.5 MB	30.5 MB	40,000	—
Data Sieving I/O	30.5 MB	274.7 MB	69	—
Two-Phase I/O	30.5 MB	30.5 MB	8	27.1 MB
List I/O	30.5 MB	30.5 MB	626	—
Datatype I/O	30.5 MB	30.5 MB	1	—
64 Clients				
POSIX I/O	12.9 MB	12.9 MB	22,500	—
Data Sieving I/O	12.9 MB	206.0 MB	52	—
Two-Phase I/O	12.9 MB	12.9 MB	4	12.1 MB
List I/O	12.9 MB	12.9 MB	352	—
Datatype I/O	12.9 MB	12.9 MB	1	—

**Table 2. I/O Characteristics of the ROMIO three-dimensional block test**



**Figure 9. Data distribution for 3D block accesses**

ous region of memory, any delays caused by list processing will directly impact performance. On the other hand, in the write case the servers are data sinks. Buffering in the TCP stack helps hide this inefficiency, although it might appear at larger numbers of processes. This overhead is not visible in the list I/O results because the number of I/O operations and the size of the I/O requests obscures this effect. A full-featured datatype I/O implementation that operated directly on the dataloop representation would likely not exhibit this behavior.

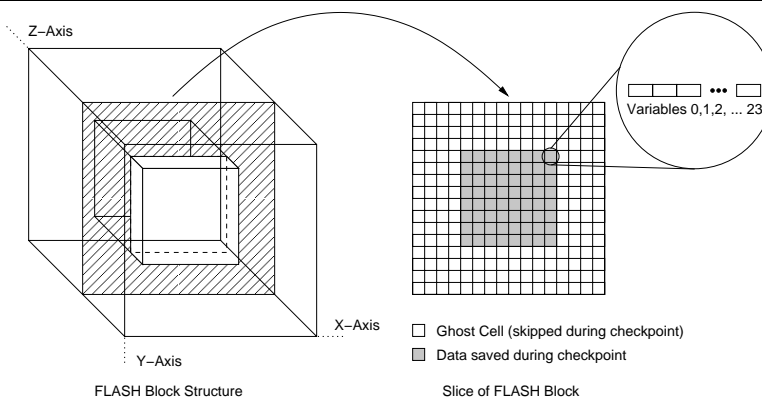
#### 4.4. FLASH I/O Simulation

The FLASH code is an adaptive mesh refinement application that solves fully compressible, reactive hydrodynamic equations, developed mainly for the study of nuclear flashes on neutron stars and white dwarfs [5]. Because the FLASH code has a very long execution time, checkpointing is a necessary component of the application. During checkpointing, blocks of data are reorganized from the in-memory organization, which includes ghost cells, into a new in-file organization. The in-file organization is stored



	Desired Data per Client	Data Accessed per Client	# of I/O Ops per Client	Resent Data per Client
POSIX I/O	7.50 MB	7.50 MB	983,040	—
Data Sieving I/O	—	—	—	—
Two-Phase I/O	7.50 MB	7.50 MB	2	$7.5 \text{ MB} * \frac{n-1}{n}$
List I/O	7.50 MB	7.50 MB	15,360	—
Datatype I/O	7.50 MB	7.50 MB	1	—

**Table 3. I/O Characteristics of the FLASH I/O simulation (n is the # of clients)**

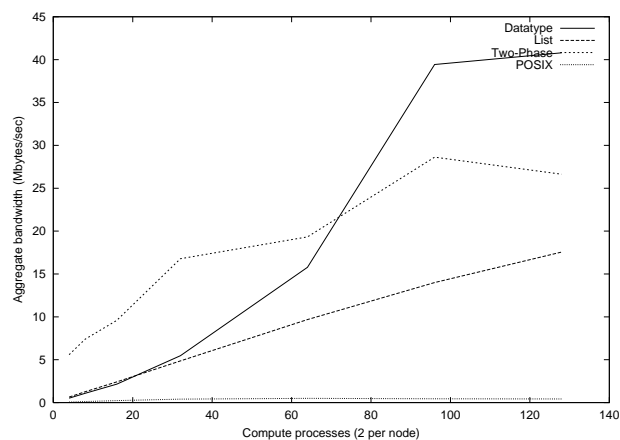


**Figure 11. FLASH memory layout**

by variable for convenience during post processing. As a result the access pattern is noncontiguous in both memory and file.

Figure 11 shows the layout of data in the memory of FLASH processes. Each process holds 80 blocks. Each block is a three-dimensional array of data elements surrounded by guard cells. Each data element consists of 24 variables. When writing data we reorganize the data so that all values for variable 0 are stored first, then variable 1, and so on. Since every processor writes 80 FLASH blocks to file, as we increase the number of clients, the dataset size increases linearly as well. Every processor adds 7 MBytes to the file, so the dataset ranges between 14 MBytes (at 2 clients) to 896 MBytes (at 128 clients).

Table 3 provides the I/O characteristics of the test using the available optimizations. Figure 12 shows the results of these tests. This is the first test in which the memory datatype is noncontiguous; thus it is the first time that the overhead of list processing might affect the clients. We see this in both the list I/O and datatype I/O cases; both underperform at small numbers of clients. As the number of clients increases, the clients are eventually able to feed the servers adequately. At 96 processes datatype I/O performance rises to nearly 40 Mbytes/sec, 37% faster than two-phase. This trend continues at higher numbers of processes. We would expect that a datatype I/O system that op-



**Figure 12. FLASH I/O Performance**

erated directly on the dataloop representation would allow clients to more effectively push data to servers, resulting in improved performance at smaller numbers of clients. List I/O, because of the size and number of I/O requests, is not able to overtake two-phase for the tested numbers of processes.

This case shows that two-phase I/O still has a place as an I/O optimization. Because data is not wasted and I/O ac-

cesses are large, two-phase is able to provide good performance despite moving the majority of the data over the network twice.

## 5. Conclusions and Future Work

Datatype I/O provides the opportunity for extremely efficient processing of structured, independent I/O requests. Our tests show that this approach outperforms both list I/O and data sieving I/O in virtually all situations. Further, it supplants two-phase I/O as the preferred optimization in many cases as well. Datatype I/O in conjunction with the two-phase collective I/O optimization makes a strong MPI-IO optimization suite. We note that in almost every case POSIX I/O alone would result in a nearly unusable system from the performance perspective; these optimizations are a necessary part of scientific parallel I/O.

This prototype does not fully exploit the potential of the datatype I/O approach. We are implementing a more full-featured version of the approach in our second-generation parallel file system, PVFS2. This version will remove the creation of the I/O lists on both client and server, further widening the performance gap between datatype I/O and other optimizations. Datatype caching similar to that seen in some remote memory access implementations [16] could boost the performance of PVFS datatype I/O by further reducing I/O request overhead. Further optimization of the approach can be provided in ROMIO as well. Caching the dataloop representations of types locally would be one way to improve datatype I/O. Leveraging datatype I/O underneath two-phase I/O would boost performance of the collectives further.

## Acknowledgements

This work was funded in part by National Aeronautics and Space Administration grants NAG5-3835 and NGT5-21 and in part by the Mathematical, Information, and Computational Sciences Division subprogram of the Office of Advanced Scientific Computing Research, U.S. Department of Energy, under Contract W-31-109-Eng-38.

## References

- [1] P. H. Carns, W. B. Ligon III, R. B. Ross, and R. Thakur. PVFS: A parallel file system for Linux clusters. In *Proceedings of the 4th Annual Linux Showcase and Conference*, pages 317–327, Atlanta, GA, October 2000. USENIX Association.
- [2] Chiba City, the Argonne scalable cluster. <http://www.mcs.anl.gov/chiba/>.
- [3] A. Ching, A. Choudhary, K. Coloma, W. Liao, R. Ross, and W. Gropp. Noncontiguous I/O accesses through MPI-IO. In *Proceedings of the Third IEEE/ACM International Symposium on Cluster Computing and the Grid (CCGrid2003)*, May 2003.
- [4] A. Ching, A. Choudhary, W. Liao, R. Ross, and W. Gropp. Noncontiguous I/O through PVFS. In *Proceedings of the 2002 IEEE International Conference on Cluster Computing*, September 2002.
- [5] B. Fryxell, K. Olson, P. Ricker, F. X. Timmes, M. Zingale, D. Q. Lamb, P. MacNeice, R. Rosner, and H. Tufo. FLASH: An adaptive mesh hydrodynamics code for modelling astrophysical thermonuclear flashes. *Astrophysical Journal Supplement*, 131:273, 2000.
- [6] W. Gropp, E. Lusk, and D. Swider. Improving the performance of MPI derived datatypes. In A. Skjellum, P. V. Bangalore, and Y. S. Dandass, editors, *Proceedings of the Third MPI Developer's and User's Conference*, pages 25–30. MPI Software Technology Press, 1999.
- [7] HDF5. <http://hdf.ncsa.uiuc.edu/HDF5/>.
- [8] IEEE/ANSI Std. 1003.1. Portable operating system interface (POSIX)—part 1: System application program interface (API) [C language], 1996 edition.
- [9] J. Li, W. Liao, A. Choudhary, R. Ross, R. Thakur, W. Gropp, and R. Latham. Parallel netCDF: A scientific high-performance I/O interface. Technical Report ANL/MCS-P1048-0503, Mathematics and Computer Science Division, Argonne National Laboratory, May 2003.
- [10] W. B. Ligon and R. B. Ross. Implementation and performance of a parallel file system for high performance distributed applications. In *Proceedings of the Fifth IEEE International Symposium on High Performance Distributed Computing*, pages 471–480. IEEE Computer Society Press, August 1996.
- [11] R. Ross, N. Miller, and W. Gropp. Implementing fast and reusable datatype processing. In *Proceedings of the 10th EuroPVM/MPI Conference*, September 2003.
- [12] 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.
- [13] R. Thakur, A. Choudhary, R. Bordawekar, S. More, and S. Kuditipudi. Passion: Optimized I/O for parallel applications. *IEEE Computer*, 29(6):70–78, June 1996.
- [14] R. Thakur, W. Gropp, and E. Lusk. Data sieving and collective I/O in ROMIO. In *Proceedings of the Seventh Symposium on the Frontiers of Massively Parallel Computation*, pages 182–189. IEEE Computer Society Press, February 1999.
- [15] R. Thakur, W. Gropp, and E. Lusk. On implementing MPI-IO portably and with high performance. In *Proceedings of the 6th Workshop on I/O in Parallel and Distributed Systems*, pages 23–32. ACM Press, May 1999.
- [16] J. L. Traff, H. Ritzdorf, and R. Hempel. The implementation of MPI-2 one-sided communication for the NEC SX-5. In *Proceedings of Supercomputing 2000*, November 2000.
- [17] J. Worringer, J. L. Traff, and H. Ritzdorf. Improving generic non-contiguous file access for MPI-IO. In *Proceedings of the 10th EuroPVM/MPI Conference*, September 2003.