

Dynamic File Striping and Data Layout Transformation on Parallel System with Fluctuating I/O Workload

Seung Woo Son* Saba Sehrish[†] Wei-keng Liao* Ron Oldfield[‡] Alok Choudhary*
*Northwestern University [†]Fermi National Accelerator Laboratory [‡]Sandia National Laboratories

Abstract—As the number of compute cores on modern parallel machines increases to more than hundreds of thousands, scalable and consistent I/O performance is becoming hard to obtain due to fluctuating file system performance. This fluctuation is often caused by rebuilding RAID disk from hardware failures or concurrent jobs competing for I/O. We present a mechanism that stripes across a dynamically-selected subset of I/O servers with the lightest workload to achieve the best I/O bandwidth available from the system. We implement this mechanism into an I/O software layer that enables memory-to-file data layout transformation and allows transparent file partitioning. File partitioning is a technique that divides data among a set of files and manages file access, making data appear as a single file to users. Experimental results on NERSC’s Hopper indicate that our approach effectively isolates I/O variation on shared systems and improves overall I/O performance significantly.

Keywords—Collective I/O; Parallel NetCDF; File partitioning;

I. INTRODUCTION

Scientists and engineers are increasingly using highly parallel machines in order to run their large, often data-intensive applications, such as thermonuclear reactions, combustion, climate modeling, and so on [10], [28], [29], [30]. Scalable parallel I/O libraries are one of the key components to scaling those applications [6], [18]. The I/O requirements of such applications can be staggering, ranging from terabytes to petabytes, and managing such massive data sets presents a significant bottleneck [8], [19].

There are many approaches proposed to coordinate I/O requests from multiple processes, and collective I/O in MPI-IO [27] has been widely used to allow collaboration among participating processes and rearrange their I/O requests to achieve high performance. There have been many optimizations to improve collective I/O performance [7], [33], [15], [31], [22], [26], [23], [21], but even with these improvements, collective I/O operations in large-scale are facing new challenges on modern parallel machines. As the size of parallel machines grows, various access contentions can significantly degrade the I/O performance, such as communication network contention because of the high ratio of application processes to file servers, and file locking contention among processes in a single job because of the shared-file access.

Furthermore, despite the use of state-of-the-art techniques described above, significant challenges still exist in achieving scalable yet consistent I/O performance. The file servers often exhibit unbalanced I/O load from various applications sharing the storage resources, resulting in fluctuating file system performance [5], [24], [34]. In petascale systems at scale, the amount of I/O throughput available to any particular job can fluctuate to a large extent based on the behaviors of other

running jobs accessing the shared file system. Another source for this kind of fluctuation is a RAID rebuild from a hardware failure. Since the performance of collective I/O is determined by the slowest participating process, it is important to ensure no process remarkably lags behind.

The study presented in this paper supports the view of conventional collective I/O, yet provides more scalable I/O performance in the presence of fluctuating file server performance. We make the following main contributions:

- We demonstrate that I/O performance could suffer from fluctuating file system behavior because of contention on shared I/O resources.
- We propose a dynamic bandwidth monitoring to probe the file servers and isolate the impact of accessing slower I/O servers by excluding them from being used for file striping.
- We propose a transparent file partitioning and data layout transformation mechanism that divides the data into a set of files, each of which is mapped and stored onto a single I/O node.

We have implemented the proposed scheme into a high-level I/O library, parallel netCDF [20], as a prototype. Our experimental evaluations on NERSC’s Hopper [3] using several benchmarks running up to 8,192 processes have shown significant I/O performance improvements. We show that our approach effectively isolates the impact of accessing slower I/O nodes and reduces write I/O time significantly with less variation. Since the partition is done at high-level I/O library layer (PnetCDF), each file partition is also a self-describing file. Maintaining portable data representation is important because it provides seamless access to data structures, and layouts across all I/O software layers. Also, the richer information available at high-level I/O library made much flexible partitioning like per-array partitioning or use of different dimension for partitioning. Our evaluations with real I/O applications demonstrate that our transparent file partitioning brings significant I/O performance improvement in both write and read while maintaining comparative number of partitioned files. We also show that our one-to-one mapping between file partitions and I/O nodes could maximize the benefit of prefetching, thereby increasing read performance significantly.

The remainder of this paper is organized as follows. The next section extends the discussion on our motivation. The design of our approach and our enhancement to PnetCDF to implement our idea are described in Section III. Section IV presents our experimental evaluation results. We discuss related work in Section V. Finally, Section VI summarizes the paper and discusses future work.

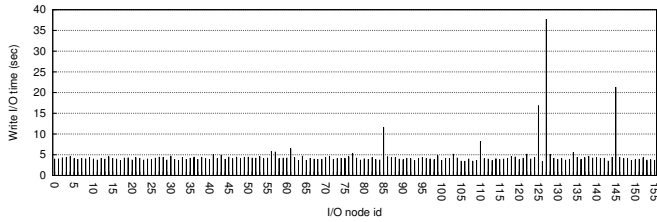


Fig. 1: The write I/O time distribution among all I/O servers (nodes) while the amount of bytes written to each I/O node remain the same.

II. BACKGROUND

Collective I/O is an optimization in many MPI-IO implementations that improves the I/O performance to shared files. In ROMIO, an implementation of MPI I/O functions adopted by many MPI implementations, the choice of aggregators depends on the file systems. For most file systems, one MPI process per compute node is picked to serve as an aggregator. In the systems containing multi-core CPUs in each node, this strategy avoids the intra-node resource contention that could be caused by two or more processors making I/O calls concurrently. For the Lustre file system, the current implementation of ROMIO picks the number of aggregators equal to the file striping count (or `striping_factor`). This design produces an one-to-one mapping between the aggregators and the file servers in order to eliminate the possible lock conflicts on the servers [21], [35]. The striping count of a file is the number of I/O servers, or Object Storage Targets (OSTs) for Lustre, where a file is stored. Like all parallel file systems, files are striped into fixed-length blocks, and they are stored in the OSTs in a round-robin fashion.

While collective I/O often offers huge improvements for I/O performance on shared files, recent studies revealed that it continues to face significant challenges at scale [36], [24], [34], [5] for several reasons. First, as demonstrated by several prior studies, global synchronization cost and lock contention among aggregators accessing the shared file within the assigned file domain during collective I/O operations pose a limit to the I/O performance. Similar observations have been made in recent studies [36], [21], but the problem will only exacerbate as the number of processes increases to thousands and more. More importantly (especially in accessing shared storage systems), there are higher levels of variability in I/O performance in petascale machines. This variability is hard to avoid because of the different ways applications access “shared” file systems. For example, multiple applications running simultaneously on the petascale machine use the file system at the same time. Another example of such a case occurs when analysis code is trying to read the data stored in the shared storage while simulation code is writing their output data. This I/O variability is a big barrier to achieve scalable collective I/O operations because I/O performance is tied to the slowest storage nodes. In other words, even if most storage nodes perform relatively fast, the overall collective I/O time is determined by the slowest nodes.

To quantify our hypothesis, we wrote a small program where each process opens a file striped on a single I/O node and writes 1GB of data on it. We have collected the write I/O time observed at each I/O server. Details of our experimental

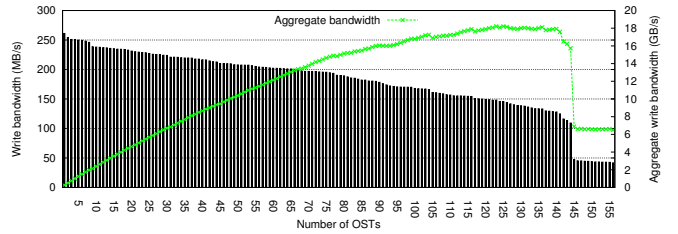


Fig. 2: Distribution of write bandwidth observed for all 156 OSTs when 16MB of dummy data is written to each OST, and the maximum achievable aggregate bandwidth. Each OST’s measured bandwidth is sorted in descending order.

setup is given in Section IV. Figure 1 shows that, although the amounts of bytes written to each I/O node are equal, a couple of I/O nodes exhibit an excessively high write I/O time relative to others. The slowest I/O node is in fact almost 9x slower than most other I/O nodes. Such an imbalance is a significant barrier to achieve scalable I/O performance if a file is simply striped across over any of those slower I/O nodes.

III. DESIGN OF FILE PARTITIONING LAYER AT PNETCDF

A. Runtime Storage Nodes Selection

Our mechanism to isolate the impact of accessing imbalanced I/O nodes is to use a runtime bandwidth probing to identify each I/O server’s load before the file striping layout is determined. The goal of this step is twofold. First, we would like to monitor each OST’s current bandwidth availability. Because the I/O pattern in HPC systems is typically bursty, we probe each I/O server’s bandwidth by writing a small dummy dataset just before writing actual file. Second, once the behaviors are identified, we would like to select the list of OSTs that can be used for storing each partitioned file.

We consider two criteria when designing our runtime probing module. First, the impact of probing should be minimized as it will not be part of the actual I/O. Second, the sampled bandwidth should reflect the temporal behavior of each I/O node. Combining these two, we determine each I/O node’s bandwidth by writing 1% of dummy data out of the total data I/O requests to each I/O server. We use POSIX I/O with the `O_DIRECT` flag because our probing module writes relatively small files, so they could sit on clients’ buffer cache unless we explicitly bypass them. We also have to make sure the sampling data size we use is large enough to fill the RPC buffer size; otherwise the data will sit on the client and will not be transferred to the I/O server.

The bar graphs in Figure 2 show the distribution of measured write bandwidth across all 156 I/O nodes (OSTs) available on NERSC’s Hopper, sorted in descending order. These bars indicate that certain OSTs exhibit relatively slower bandwidth than the others. We again attribute this to an inherent imbalance when accessing shared storage, as extensively discussed in recent studies [24], [34], [5]. Given these observed I/O bandwidths, we use the following algorithm to select the I/O nodes for file striping. Assuming B_i to be the sorted bandwidth observed for each I/O server, i , we denote the aggregate I/O bandwidth, \mathbb{B}_i , using i OSTs by $\mathbb{B}_i = i \times B_i$, where $1 \leq i \leq 156$. The linepoints in Figure 2 show the estimated maximum achievable aggregate I/O bandwidth based on this formula. As we can see, the

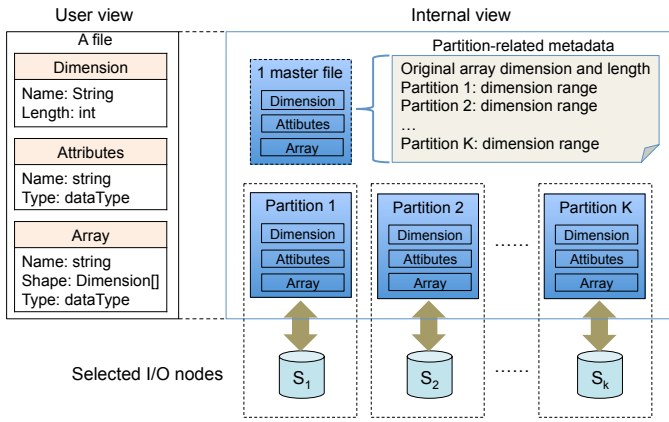


Fig. 3: Overview of our file partitioning mechanism. In our approach, each array is internally divided into K file partitions, each of which is stored in a single I/O node, that is, there is “1-to-1 mapping” between each partition and I/O node. All files (both master and partitioned files) are in self-describing file format.

aggregate bandwidth gradually increases as more I/O nodes are added, but eventually saturates and then declines because the aggregate bandwidth is confined to the slowest node. To select the maximum number of I/O nodes that provide us the best achievable bandwidth, we calculate the derivative of \mathbb{B}_i , which represents the slope of \mathbb{B}_i at each value of i . Since our goal here is to maximize the number of I/O nodes, we select i when \mathbb{B}'_i is negative and is less than a certain threshold, δ . The threshold value is basically meant for capturing the degree of slowness in the aggregate bandwidth when a certain probed bandwidth is added. In our implementation, we used the δ value of -15%. Using the results shown in Figure 2, our algorithm excludes 12 OSTs with less than 50MB/s for striping partitioned files. The aggregate bandwidth was estimated to peak when the first 128 OSTs were added, but the significant bandwidth drop occurs when 145th OST is added. We note that if all probed bandwidth values are similar to each other, our algorithm will end up selecting most of the available OSTs.

We note that the probing happens only once at the file create time, once for each new file creation. The overhead of probing is less than 1% of the total I/O time. This cost is included in the timings reported in this paper. We also note that there exist chances where two jobs running probing simultaneously. We however speculate that such chances are slim because probing happens only once, at creating a file. If two probes did occur concurrently, the bandwidths obtained should be halved. However, this also means two jobs are most likely competing for the file system. Hence, lower I/O performances are expected for both jobs. An ideal solution should be at system-level to monitor individual server workload, but ours is a user-level solution that makes best use of available information to try produce maximum achievable performance.

B. Mapping Arrays to File Partitions

When the subset of OSTs that are more efficient than others has been isolated, the ideal solution would be to stripe files across those OSTs. However, on Lustre, users have no way to stripe files across a set of specific OSTs. This leads to our partitioning solution where each partitioned

file is stored on only one OST (i.e., `stripe_count` of 1). Lustre does allow users to select the starting OST for a file. Figure 3 gives an overview of our file partitioning scheme. The basic concept of our scheme is that, from an application’s perspective, partitioning is transparent; that is, all processes open and access a single file throughout program execution. Then, our partitioning mechanism internally splits application processes into set of subprocesses, each of which creates its own file partition collectively. The file partition created by each subprocess group is accessed *solely* by that group.

We perform partitioning when array definition is finished and the data in memory is ready for write/read. We choose this time because each array’s shape (number of dimensions, length of each dimension, and datatype of each element) is finalized at this point. The header information is also written at the end of file partition. In order to convey the users’ intention of their file partitioning policy, we use the MPI hint mechanism.

The default partitioning policy is along the most significant dimension. For example, an array of Z - Y - X dimension, each with the same length will be partitioned along the dimension Z . There are however certain applications that prevent applying the default policy. For instance, in the S3D I/O application, the dataset called u is a 4D array with the most significant dimension has length 3. Such a small dimension length limits the number of file partitions, preventing the application from exploiting potential benefits of partitioning in larger partition counts. In this case, we partition the arrays along the second most significant dimension.

The details of the file creation are as follows. It first obtains users’ intention of partitioning through `MPI_Info_get()`. We store acquired information as a metadata in both master and the partitioned file’s header information. If no hints were provided regarding file partitions, the normal procedure will be executed; it creates a single file without partitions. Otherwise, it splits the communicator because each process is divided into a subprocess group. The split processes then collectively create their own file partition using a dataset function provided in the high-level I/O library, for instance, `ncmpi_create` in PnetCDF. After creating a partitioned file, our algorithm traverses each defined array in the original definition and determines which dimension ID it needs to use for partitioning. It calculates a new dimension length for each partition. Note that only the partitioning dimension will be affected; all the remaining dimensions will have the same length as the original. Once a new dimension length is determined, we define a new dimension for the partitioned file and create an array with the new dimension lists. If the array is partitioned, we update the original array definition in the master file with a scalar value. In other words, the master file does not have a physical space allocated for the partitioned array as the actual data will be stored in the partitioned files. We repeat these procedures until all arrays in the original file are processed.

Figure 4 shows a typical example of PnetCDF code that includes the sequence of dimension and array (variable) definition followed by the code to write data on it. In PnetCDF, all processes in the communicator must make an explicit call (`ncmpi_enddef`) at the end of the define mode in order to verify that the values passed in by all processes match. From our design viewpoint, this is the time when all shapes of arrays are known, therefore, our array partitioning is internally

```

MPI_Info_set (info, "nc_partitioning_enabled", "true");
ncmpi_create(comm, ..., info, &ncid);
...
/* dimension definition */
ncmpi_def_dim(ncid, "z", 100L, &cube_dim[0]);
ncmpi_def_dim(ncid, "y", 100L, &cube_dim[1]);
ncmpi_def_dim(ncid, "x", 100L, &cube_dim[2]);
...
/* variable (array) definition */
ncmpi_def_var(ncid, "cube", NC_INT, 3, cube_dim, &cube_id);
...
ncmpi_enddef();
...
/* perform I/O */
ncmpi_put_vara_all(ncid, cube_id, start[], count[], buf,
                  bufcount, MPI_INT);
...

```

Fig. 4: A PnetCDF example code that creates a file with partitioning enabled set to true. The number of file partitions is determined through the profiling mechanism explained in Section III-A. This example creates a variable named “cube” of Z-Y-X dimension, each with 100 length. From an application writer’s viewpoint, it only requires adding a hint to specify the intent of partitioning to store a variable.

<pre> netcdf test { dimensions: z = 100; y = 100; x = 100; variables: double cube (z, y, x); // global attributes: data: cube = ; } </pre> <p>(a)</p>	<pre> netcdf test { dimensions: z = 100; y = 100; x = 100; variables: double cube; cube: num_partitions = 2; cube: ndims_orig = 3; // global attributes: :partition 0: "test.0"; :partition 1: "test.1"; data: cube = 0; } </pre> <p>(b)</p>
<pre> netcdf test.0 { dimensions: z.cube = 50; y.cube = 100; x.cube = 100; variables: double cube(z.cube, y.cube, x.cube); cube: range(z) = 0,49; // global attributes: :partition_index = 0; data: cube = ; } </pre> <p>(c)</p>	<pre> netcdf test.1 { dimensions: z.cube = 50; y.cube = 100; x.cube = 100; variables: double cube(z.cube, y.cube, x.cube); cube: range(z) = 50,99; // global attributes: :partition_index = 1; data: cube = ; } </pre> <p>(d)</p>

Fig. 5: NetCDF file header information by `ncmpidump` when the file is divided into 2. (a) Original NetCDF file (i.e., non-partitioned case). (b) The master NetCDF file after partition. Note that the data section is 0, meaning empty. (c) First partitioned NetCDF file. (d) Second partitioned NetCDF file.

executed at the end of this call.

The NetCDF header information for this example of partitioning case is given in Figure 5. Note that, after partitioning, both master and partitioned files have more additional attributes than the original file. For instance, the master file (Figure 5(b)) has global attributes that indicate the file name for each partitioned file, the number of partitions, and the original dimension size for a variable, “cube”. The partitioned file header, on the other hand, has attributes for describing the range of partitioned dimension as well as the partition index.

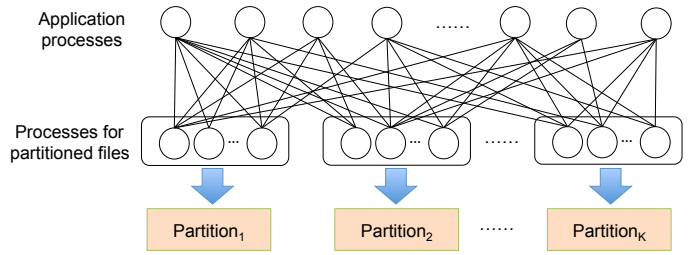


Fig. 6: Any I/O requests from applications belonging to a file partition owned by other processes need to communicate among processes before going to I/O servers. The data exchange between application and partitioned file processes can be all-to-all personalized. The local-to-global partitioning information is kept in each file’s metadata.

C. Memory-to-File Layout Transformation

Once a file is partitioned, we need to provide a transparent way to access those partitioned files. Note that, from an application’s viewpoint, all I/O accesses still go through the master file as it has sufficient information about how each array is partitioned in each file. In other words, there is no change in user’s I/O routines. We also note that reading datasets already stored in partitioned files can be performed transparently using the same metadata retrieval process.

Figure 6 shows an overview of the memory-to-file layout transformation mechanism. The transformation mechanism to partitioned files is mainly composed of two steps: i) calculating each process’s requests to partitioned files and exchanging it among all processes; ii) exchanging requests among processes in each split communicator and issues I/O requests using I/O calls (either synchronous or asynchronous ones). We note that, since our partitioning is done at the higher-level I/O library layer, all user’s array partitioning is represented as start, count, and stride offset list for each dimension.

In the first step, each process calculates the list of start and count offsets to each file partition, dividing the data in memory among the processes who own the partitions. This is done by (logically) dividing the start and count offset, denoted as `my_req[]`, into file partitions, each of which can be directly accessed by the processes within a sub-communicator. In our implementation, we do not restrict the number of such delegate processes in each subprocess groups. In fact, any process can be a delegate so that we do not make load imbalance at an application layer by selecting limited number of delegates because non-delegate processes do not read/write files directly. This phase requires one `MPI_Allreduce()` among all processes.

The second step is based on everyone’s `my_req`, and calculates what requests of other processes lie in this process’s file partitions. `others_req[i].{start,count}` indicates how many noncontiguous requests of process *i* accessing this process’s file partition. All these incur an `MPI_Alltoall` and many `isend/irecv/wait_all`. This step ensures delegates collect the request information from all other processes.

Then each process sends requests to the appropriate remote delegate. Only delegates may have multiple I/O requests. Non-delegate processes will not participate in this loop, but will call to the data exchange routine if they have certain requests to delegates. Delegate processes iterate until they receive requests

from all other processes, and issue a non-blocking I/O. Each iteration goes through all `others_req[*]` and continues until all requests are processed. We ensure they are all processed by calling `wait_all()` at the I/O library layer.

We illustrate how I/O requests to the partitioned files are processed using the example code in Figure 4. Let us assume there are 4 processes to access this array and the number of file partitions is 2. Each I/O request is composed mainly of start offset, count and stride for each dimension. Since our example dataset is 3 dimensional, we have `start[3]`, `count[3]`, and `stride[3]`. For illustrative purposes, let us assume that stride count is 1, meaning all array elements are accessed contiguously. Given this, one partition (50 by 100 by 100) is owned by P_0 and P_1 whereas the other partition (50 by 100 by 100) is owned by P_2 and P_3 . Assuming a block-block access pattern and user’s file partition, we calculate each process’s request to each file partition. For instance, P_0 ’s original request, denoted as `start{0,0,0}` and `count{100,50,50}`, is now divided into two portions: a portion belonging to its own file partition (denoted as `start{0,0,0}` and `count{50,50,50}`) and the other (denoted as `start{50,0,0}` and `count{50,50,50}`) to be sent to the remote process that owns that file partition. Once all this information is obtained, all processes now exchange information (using `alltoall`) in order to figure out which process has a portion of the data not belonging to its own partition. Afterwards, all processes know which sub-I/Os they need to handle by themselves. The code then communicates the corresponding buffers and issues all those received I/O requests using PnetCDF’s nonblocking I/O calls. The I/O to partitioned files returns when all the issued nonblocking I/O calls are completed.

IV. EXPERIMENTAL EVALUATIONS

All our experiments are performed on the Cray XE6 machine, Hopper, at NERSC. Hopper has a peak performance of 1.28 Petaflops/sec, 153,216 processors cores for running scientific applications, 212 TB of memory, and 2 Petabytes of online disk storage. The Hopper system has two locally attached high-performance scratch disk spaces, `/scratch` and `/scratch2`, each of 1 PB capacity. They both have the same configuration: 26 OSSs (Object Storage Servers), each of which hosts 6 OSTs (Object Storage Target), making a total of 156 OSTs. The parallel file system deployed in Hopper is Lustre [2] mounted as both scratch disk spaces. When a file is created in `/scratch`, it is striped across two OSTs by default. Lustre provides users with a tunable striping configuration for a directory and files; both directory and files have the same striping configuration. In our experiment, we use all available OSTs for striping and 1 MB as default stripe sizes.

We implemented our proposed approach into the parallel netCDF 1.3.1. Our feature added approximately 1,500 lines of new code to PnetCDF. Our implementation is configured to link with Cray’s `xt-mpich2` version 5.6.0. We used a separate ROMIO module described in [21] as a standalone library, which is then linked with the native MPI library. Our previous experience indicates this optimized ROMIO is about 30% faster than the system’s default one. In other words, our base collective I/O performance is already optimized for our evaluation platform. All applications including benchmarks and our modified PnetCDF are compiled using PGI compiler version 12.9.0 with the “-fast” compilation flag.

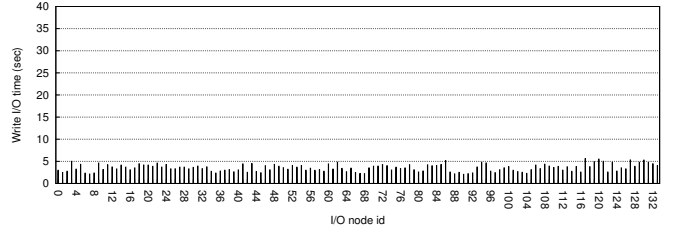


Fig. 7: Balanced write I/O time observed when only subset of I/O nodes that were detected through our dynamic bandwidth probing.

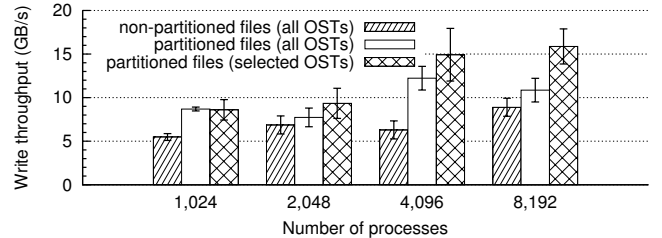


Fig. 8: Write throughput results for `coll_perf`.

We evaluate our file partitioning scheme against the base scheme, where all arrays are stored in a normal file (non-partitioned) striped across *all* available OSTs. Because users have no way to stripe files across a set of specific OSTs on Lustre, we cannot provide the case of non-partitioned files with selected OST in the experiments. To show the effectiveness of our dynamic bandwidth probing, we ran two schemes of our partitioning cases: striped over all OSTs and striped over selected OSTs. While there are several other techniques whose goals are similar to ours like PLFS [4] and ADIOS [25], we do not compare our approach against them because fair comparison is hard to make; they do not preserve canonical order of the original dataset whereas all partitioned files in our approach are stored in portable NetCDF format. Furthermore, both PLFS and ADIOS lack the ability to map the file partition and underlying I/O nodes selectively.

Our evaluation is conducted using up to 8,192 processes because our solution is designed to deal with I/O systems with fluctuating performance due to multiple jobs competing for shared I/O resource. If we ran bigger jobs that use most of the available compute nodes, the opportunity of seeing such I/O competition shall decrease. For large-scale runs, we anticipate all available OSTs be selected to serve the I/O.

A. Collective I/O Performance Benchmark

Before presenting our evaluation with the collective I/O performance benchmark, we first show how our approach effectively isolates slower I/O nodes. In order to do this, we wrote a small test case that writes 1GB of data to an individual I/O node selected by our dynamic probing module. Figure 7 shows the write I/O time, collected using the TAU profiling tool [32], observed at each I/O node that was selected by our sampling module. In this example, 134 out of 156 OSTs were selected for writing. As compared with Figure 1, it clearly demonstrates more balanced write I/O time across all selected OSTs.

To understand the performance of our approach against the

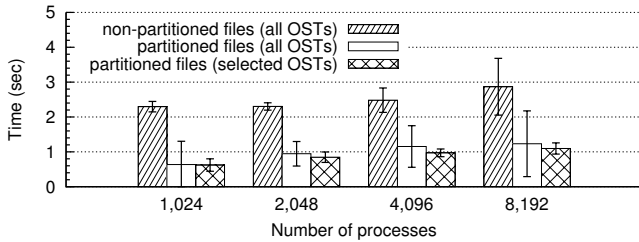


Fig. 9: The average write I/O time for coll_perf with error bars. Regardless of file is partitioned or not, using all OSTs shows much “higher” deviation.

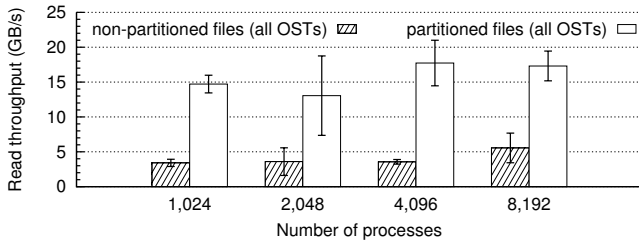


Fig. 10: Read throughput results for coll_perf.

base case, we ran a collective I/O test program, coll_perf, originally from ROMIO test suite, that writes and reads the three-dimensional arrays, all in a block-partitioned manner. We made it write/read four 3D variables. The data partitioning is done by assigning a number of processes to each Cartesian dimension. In our experiments, we set the subarray size in each process to $128 \times 128 \times 128$ of 4-byte integers, corresponding to 8MB. All data is written to a single file for the base case (non-partition). For our partitioning case, all four variables are partitioned along the most significant dimension.

Figure 8 shows the write throughput of coll_perf with and without our file partitioning schemes. The results indicate that writing data into a single file does not scale with larger number of processes; the write throughput actually went up and down when the number of processes are increased. On the other hand, our partitioning schemes improve the write throughput significantly by 12%–94% when used with all 156 OSTs and 36%–137% when used with selected OSTs, respectively.

To understand the performance improvement obtained by our approach, we have collected the performance breakdown of coll_perf during collective I/O using the TAU profiling tool [32]. Figure 9 shows that the time spent in POSIX write() time by each aggregator process gradually increases as the number of processes increase. This is because the amount of data written increases with larger number of processes. An important observation we made here is that writing to partitioned files using either all OSTs or selected OSTs reduces the write I/O time significantly, about 70% on average. This indicates that writing to partitioned files clearly lessens the contention on the file server. Another important insight from this graph is the high variations on the write I/O time when all OSTs are used, and the variations increased with larger process counts. The partitioned files with selected OSTs show low deviation from average mainly because relatively slower OSTs were eliminated before the time of writing.

In our next experiments, we would like to understand

how the read from partitioned files behaves. To do this, we perform the same weak scalability tests on the read case, each case collectively reads the entire files in a block-block partitioned manner. Since partitioning on selected OSTs does not have fixed the number of OSTs per run, we evaluate only reading from all OSTs for a fair comparison. To ensure data is read from the storage nodes, all caches are flushed before each run. Figure 10 shows that the non-partitioned file case is not scalable while our partitioning scheme shows much higher performance improvement than the write case. Also, the observed read throughput is about 30% lower than that of the write throughput. Our TAU profiling result indicates a notable increase in read I/O time; reading from the normal (i.e., non-partitioned) is about 6x slower than reading from partitioned files. We attribute this to the pretty aggressive readahead mechanism used in Lustre file system. In the case of reading from non-partitioned files on all OSTs and given the default stripe size of 1MB, the majority of prefetched data by an aggregator is irrelevant parts of the data, thus slowing down the overall performance. In our partitioned file case, the readahead mechanism is entirely reading from a single OST, so the benefit of readahead is maximized.

Our partitioning approach introduces additional communication during memory-to-file layout transformation time: MPI_Isend, MPI_Irecv(), MPI_Alltoall(), and MPI_wait(). In order to quantify this overhead, we have measured time spent on those additional communication costs using TAU. The results indicate that the coordination overhead incurred by the additional communication is negligible; the extra communication overhead accounts for less than 1% of the collective I/O operations. The time spent on the all-to-all communication is small because, during that phase, we only exchange each process’s requests to each file partition. The buffer exchange phase also does not incur much overhead because only participating process pairs exchange small amount of buffer. Since our algorithm selects the delegation process in other subprocess groups in a balanced manner, the pairwise communication is also mostly balanced.

B. FLASH I/O Benchmark

The FLASH I/O benchmark [38], [19] is the I/O kernel of a block-structured adaptive mesh hydrodynamics code that solves the compressible Euler equations on a block structured adaptive mesh and incorporates the necessary physics to describe the environment, including the equation of state, reaction network, and diffusion [10]. We use a FLASH I/O format where all mesh variables (including density, pressure and temperature) are written to the same dataset (variable) in the output file. Both checkpointing and plot files are written in this file format. In case of checkpoint files, among 24 variables defined in FLASH, only 10 variables correspond to those mesh variables, each of which is a four-dimensional (4D) array of double-precision typed data. All unknown variables are defined as a 5D array, the first dimension being the number of unknown variables. Since this dimension length is only 10, we partition these unknown variables along the second most significant dimension. We partition only unknown variables in our approach; All other variables are stored in the master file without partitioning. The plot files have three mesh variables and we again applied partitioning for the unknown variables.

Figure 11 shows the I/O bandwidth of FLASH for the

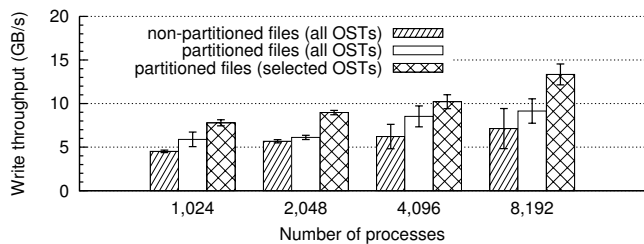


Fig. 11: FLASH I/O write throughput.

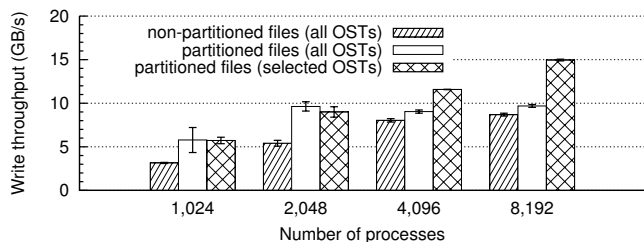


Fig. 12: S3D I/O write throughput.

non-partitioned case and our two approaches. Using a non-partitioned file did not scale well even with increased process counts. The maximum I/O bandwidth observed with 8,192 processes is about 8 GB/s. This is significantly below the maximum I/O bandwidth on Hopper. The partitioned files with all OSTs slightly outperform the non-partitioned file case, by 28% on average, but there is higher variation with larger process counts. Overall, the partitioned files with selected OSTs can achieve about 70% I/O bandwidth improvement than the non-partitioned case.

C. S3D I/O Benchmark

The S3D application [29] simulates turbulent combustion using direct numerical simulation of a comprehensive Navier-Stokes flow. The domain is decomposed among processes in 3D. All processes periodically participate in writing out a restart file. This file can be used both as a mechanism to resume computation and as an input for visualization and post-simulation analysis. We used $50 \times 50 \times 50$ fixed subarrays.

The checkpoint files consist of four global arrays: two 3-dimensional, temp (z, y, x) and pressure (z, y, x) in double precision, and two 4-dimensional arrays (double yspecies (nsc, z, y, x) and double u (three, z, y, x)). Since the length of the most significant dimension in 4D variables are relatively small, 3 and 11 for *three* and *nsc* respectively, we partition these variables along z-dimension, that is, the second most significant dimension.

Figure 12 shows the I/O bandwidth of S3D for all three cases we evaluated. We have observed that the non-partitioned file case is marginally scalable. The partitioned files using all OSTs can achieve higher performance improvement than the non-partitioned file case up to 2,048 processes, but only marginal improvement beyond that point. The partitioned files case with selected OSTs consistently outperforms than the non-partitioned file case, by 60% on average.

V. RELATED WORK

PLFS [4] introduced a virtual layer that remaps an application’s preferred data layout into one optimized for the

underlying parallel file system. Like PLFS, Yu et al. [37] also use a library approach to reduce contention from concurrent access at runtime. However, the split files are merged at close time, preventing later accesses from leveraging the benefits of partitioned files. It also requires application modification. Yu and Vetter proposed an augmented collective I/O, called ParColl, with file area partitioning and I/O aggregator distribution [36]. PIDX [17], [16] is a parallelization of IDX data format, and uses a novel aggregation technique to improve its scalability. Dickens and Logan [8] proposed an approach, called Y-Lib, to collective I/O in Lustre, which improves performance by reducing contention among processes participating in collective operations. SIONlib [9] provides a transparent mapping of a large number of task-local files onto a small number of files, but it again requires internal metadata handling and block alignment, and is required to use a set of their new APIs.

Our earlier study by Gao et al. [13] is similar to our approach, but it requires user intervention of how each subfile is partitioned using a set of new APIs. Also, it only allows partitioning along the most significant dimensions of an array, and does not support record variables. In our new design and implementation, we remove these restrictions to enable any further layout transformation between memory and partitioned files. All these data transformations would require all-to-all personalized communications among application and subfile processes, which does not occur in the subfiling. Further, unlike the subfiling, our approach gives more flexibility by allowing application writers to specify per-variable partitioning. A similar idea of subfiling is also provided in the ADIOS BP file format [25]. However, ADIOS has limited flexibility in selecting how the data is stored across subfiles, and also it does not store arrays in canonical order. Fu et al. [12], [11] proposed an application-level two-phase I/O, called reduced-blocking I/O (rbIO), and demonstrated that rbIO performs better than the n to n approach. rbIO is similar to our approach in that it reduces conflicts using the partitioned files and application 2-phase I/O. However, the partition in rbIO is done by the application writers, and the coordination does not cross the partitioned process group. Kendall et al. also used an application-level 2-phase I/O in order to organize I/O requests to multiple-file dataset [14]. Their optimization, however, is targeted mainly for visualization workloads, and application writers manually provide the list of starts and sizes of a block that each process needs to read or write.

Many recent studies have identified that staggering file servers are one of the main reasons of inconsistent I/O performance in large petascale and beyond systems [24], [34], [5]. [34] characterizes the I/O bottlenecks in supercomputers, and it demonstrates that slower I/O servers limit the aggregate and striping bandwidth and reduce the parallelism. Also, due to locking protocols, lower bandwidths are observed while writing to a shared file. In [24], it is shown that the I/O load variation on I/O servers leads to performance degradation, and adaptive I/O methods are proposed using a grouping approach to balance the workload; i.e., for a group of writer processes, assign a sub-coordinator to each group, and assign a coordinator for all the sub-coordinators. In a recent study on Hopper [5], it is shown that once the I/O stragglers are isolated from the I/O, and using one file for all processes, the performance can be significantly improved. Our approach

does take the slower I/O servers into account and dynamically isolates these servers from the collective I/O operation. Using one partition per file server can potentially achieve better performance by minimizing file system locking contention.

VI. CONCLUSION AND FUTURE WORK

This paper has proposed a transparent file partitioning mechanism to provide scalable collective I/O performance while keeping a conventional view of large multi-dimensional arrays to a user. We use a dynamic bandwidth probing to detect slower I/O nodes and isolate the impact of these slower I/O nodes. Our implementation is incorporated into PnetCDF, a high-level I/O library, and we evaluate its performance using a set of I/O benchmarks on NERSC's Hopper. Our experimental results demonstrate that our partitioning scheme consistently improves the performance of collective I/O significantly by reducing write I/O time with less variation. We also show that storing each partition onto a single I/O node could maximize the effect of the read-ahead mechanism, resulting in significantly improved read I/O performance.

We will continue to evaluate our approach on other platforms like Intrepid, IBM Blue Gene/P, at Argonne National Laboratory [1], and other high-level I/O libraries. Future research will focus on investigating how the data exchange mechanism we proposed in this paper can be applied on more general layout transformation techniques like transposing array dimensions.

ACKNOWLEDGMENT

This work is supported in part by the following grants: NSF awards CCF-0833131, CNS-0830927, IIS-0905205, CCF-0938000, CCF-1029166, and OCI-1144061; DOE awards DE-FG02-08ER25848, DE-SC0001283, DE-SC0005309, DESC0005340, and DESC0007456; AFOSR award FA9550-12-1-0458. This research used resources of the National Energy Research Scientific Computing Center, which is supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC02-05CH11231.

REFERENCES

- [1] "Argonne Leadership Computing Facility," <http://www.alcf.anl.gov/intrepid/>.
- [2] "Lustre File System," <http://www.lustre.org>.
- [3] "National Energy Research Scientific Computing Center," <http://www.nersc.gov/users/computational-systems/hopper/>.
- [4] J. Bent *et al.*, "PLFS: A Checkpoint Filesystem for Parallel Applications," in *SC*, 2009.
- [5] S. Byna *et al.*, "Trillion Particles, 120,000 Cores, and 350 TBs: Lessons Learned from a Hero I/O Run on Hopper," in *CUG*, 2013.
- [6] P. Carns *et al.*, "24/7 Characterization of Petascale I/O Workloads," in *IASDS*, 2009.
- [7] J. M. del Rosario *et al.*, "Improved Parallel I/O via a Two-phase Runtime Access Strategy," in *Proceedings of Workshop on Input/Output in Parallel Computer Systems*, 1993, pp. 56–70.
- [8] P. M. Dickens and J. Logan, "Y-lib: A User Level Library to Increase the Performance of MPI-IO in a Lustre File System Environment," in *HPDC*, 2009, pp. 31–38.
- [9] W. Frings *et al.*, "Scalable Massively Parallel I/O to Task-Local Files," in *SC*, 2009, pp. 17:1–17:11.
- [10] B. Fryxell *et al.*, "FLASH: An Adaptive Mesh Hydrodynamics Code for Modeling Astrophysical Thermonuclear Flashes," *The Astrophysical Journal Supplement Series*, vol. 131, no. 1, p. 273, 2000.
- [11] J. Fu *et al.*, "Scalable Parallel I/O Alternatives for Massively Parallel Partitioned Solver Systems," in *LSP*, 2010.
- [12] —, "Parallel I/O Performance for Application-Level Checkpointing on the Blue Gene/P System," in *IASDS*, 2011, pp. 465–473.
- [13] K. Gao *et al.*, "Using Subfiling to Improve Programming Flexibility and Performance of Parallel Shared-file I/O," in *ICPP*, 2009, pp. 470–477.
- [14] W. Kendall *et al.*, "Visualization Viewpoint: Towards a General I/O Layer for Parallel Visualization Applications," *IEEE Computer Graphics and Applications*, vol. 31, no. 6, pp. 6–10, 2011.
- [15] D. Kotz, "Disk-directed I/O for MIMD multiprocessors," *ACM Trans. Comput. Syst.*, vol. 15, no. 1, pp. 41–74, Feb. 1997.
- [16] S. Kumar *et al.*, "Efficient Data Restructuring and Aggregation for IO Acceleration in PIDX," in *SC*, 2012.
- [17] —, "PIDX: Efficient Parallel I/O for Multi-resolution Multi-dimensional Scientific Datasets," in *Cluster*, 2011, pp. 103–111.
- [18] S. Lang *et al.*, "I/O Performance Challenges at Leadership Scale," in *SC*, 2009, pp. 40:1–40:12.
- [19] R. Latham *et al.*, "A Case Study for Scientific I/O: Improving the FLASH Astrophysics Code," *Computational Science & Discovery*, vol. 5, no. 1, 2012.
- [20] J. Li *et al.*, "Parallel netCDF: A High-Performance Scientific I/O Interface," in *SC*, 2003.
- [21] W.-k. Liao and A. Choudhary, "Dynamically Adapting File Domain Partitioning Methods for Collective I/O Based on Underlying Parallel File System Locking Protocols," in *SC*, 2008.
- [22] W.-k. Liao *et al.*, "Scalable Design and Implementations for MPI Parallel Overlapping I/O," *Parallel and Distributed Systems, IEEE Transactions on*, vol. 17, no. 11, pp. 1264–1276, nov. 2006.
- [23] —, "Collective Caching: Application-aware Client-side File Caching," in *HPDC*, 2005, pp. 81–90.
- [24] J. Lofstead *et al.*, "Managing Variability in the IO Performance of Petascale Storage Systems," in *SC*, ser. SC '10, 2010, pp. 1–12.
- [25] J. F. Lofstead *et al.*, "Flexible IO and Integration for Scientific Codes through the Adaptable IO system (ADIOS)," in *CLADE*, 2008, pp. 15–24.
- [26] X. Ma *et al.*, "Improving MPI-IO Output Performance with Active Buffering Plus Threads," in *IPDPS*, 2003.
- [27] Message Passing Interface Forum, "MPI-2: Extensions to the Message Passing Interface," <http://www.mpi-forum.org/docs/docs.html>.
- [28] D. Randall *et al.*, "Breaking the Cloud Parameterization Deadlock," *Bull. Amer. Meteor. Soc.*, vol. 84, pp. 1547–1564, 2003.
- [29] R. Sankaran *et al.*, "Direct Numerical Simulations of Turbulent Lean Premixed Combustion," *Journal of Physics: Conference Series*, vol. 46, no. 1, p. 38, 2006.
- [30] K. Schuchardt *et al.*, "IO Strategies and Data Services for Petascale Data Sets from a Global Cloud Resolving Model," *Journal of Physics: Conference Series*, vol. 78, 2007.
- [31] K. E. Seamons *et al.*, "Server-directed Collective I/O in Panda," in *SC*, 1995.
- [32] S. S. Shende and A. D. Malony, "The TAU Parallel Performance System," *Int. J. High Perform. Comput. Appl.*, vol. 20, no. 2, pp. 287–311, May 2006.
- [33] R. Thakur and A. Choudhary, "An Extended Two-phase Method for Accessing Sections of Out-of-core Arrays," *Sci. Program.*, vol. 5, no. 4, pp. 301–317, Dec. 1996.
- [34] B. Xie *et al.*, "Characterizing Output Bottlenecks in a Supercomputer," in *SC*, 2012, pp. 8:1–8:11.
- [35] L. Ying, "Lustre ADIO Collective Write Driver – White Paper," Sun and ORNL, Tech. Rep., 2008.
- [36] W. Yu and J. Vetter, "ParColl: Partitioned Collective I/O on the Cray XT," in *ICPP*, 2008, pp. 562–569.
- [37] W. Yu *et al.*, "Exploiting Lustre File Joining for Effective Collective IO," in *CCGrid*, 2007, pp. 267–274.
- [38] M. Zingale, "FLASH I/O Benchmark Routine - Parallel HDF 5," http://www.ucoick.org/~zingale/flash_benchmark_io/, March 2001.