# Improving the performance of file-based coupling in multi-component HPC workflows

Tatiana V. Martsinkevich[1,a]    Balazs Gerofi[1,b]    Wei-keng Liao[2,c]    Yutaka Ishikawa[1,d]

Alok Choudhary[2,e]

**Abstract:** Efficient coupling of components in multi-component HPC workflows is a common issue especially if the components are loosely coupled via files. A number of solutions have been proposed, mainly in the form of a coupling software. However, such software usually requires the programmer to use its API, which would often mean rewriting the I/O kernel code in case the original program used some other I/O library. We propose a coupling framework called the DTF that replaces file I/O in components with direct sending of data transparently for the user. The DTF works for applications which use the PnetCDF library for file I/O. In this work we present the evaluation of the latest version of the DTF with improved scalability and faster data transfer.

## 1. Introduction

A common way of tackling a complex scientific problem is to use a multi-component workflow approach where each component performs a certain task and all components work together to achieve a common goal. An example of such components can be different physics models, in-situ data analytics or visualization components and so on. Tuning High Performance Computing (HPC) systems so that they could efficiently run such multi-component workflows is a topic of a number of studies [4], [11].

Often a multi-component workflow consists of independently developed pieces of software that are stand-alone programs by themselves. In this case, it is not rare to couple the components via files rather than use a special coupling software: One component outputs its result to a file and the file becomes the input for the next component. However, as the amount of data to pass between the components grows, the I/O quickly becomes a bottleneck in such applications. This is particularly the issue in the case when components are loosely coupled through files.

An example real-world application that experiences this problem is SCALE-LETKF [18]. SCALE-LETKF is a real-time severe weather prediction application that combines weather simulation with assimilation of weather radar observations. It consists of two components — SCALE and LETKF — that are developed independently. SCALE is a numerical weather prediction application based on the ensemble simulation; LETKF performs data assimilation of real-world observation data together with simula-

tion results produced by SCALE. In one iteration SCALE passes the results of its computations to LETKF via files written using the Parallel NetCDF API [2]. The file output of LETKF, in its turn, becomes the input for SCALE in the next iteration.

A particular feature of SCALE-LETKF is that it has a strict timeliness requirement. The target execution scenario is to assimilate the observations arriving at an interval of 30 seconds. However, this requirement is hard to fulfill once the amount of file I/O grows too big for a fine model resolution. One way to overcome this would be to switch from file I/O to some coupling software that would allow the components to exchange the data directly. However, this would require rewriting the I/O kernels in both components using the API of the coupler as such a software usually requires. This can be a daunting task for a software as large and complex as SCALE-LETKF.

We have previously proposed a framework called Data Transfer Framework (DTF) [19] intended to be used in multi-component workflows. The framework allows to silently bypass file I/O and to send the data between the components directly over network. The DTF assumes that the workflow components use PnetCDF API for file I/O and it uses the Message Passing Interface (MPI) [17] to transfer the data. Applications like SCALE-LETKF can benefit from DTF because it allows the developers of the application to easily switch from file I/O to direct data transfer without having to rewrite the I/O code.

In the current work we (i) present an evolved version of the DTF in which we solved some scalability issues encountered in the previous version; (ii) we show the results of the performance evaluation of the latest version of DTF measured using a benchmark program.

The rest of this paper is organized as follows. In Section 2 we revisit the design of our data transfer framework and discuss its implementation and what have changed since the previous ver-

---

1    RIKEN AICS, Japan
2    Northwestern University, Chicago, USA
a)    tatiana.mar@riken.jp
b)    bgerofi@riken.jp
c)    wkliao@eecs.northwestern.edu
d)    yutaka.ishikawa@riken.jp
e)    choudhar@eecs.northwestern.edu

sion in Section 3. The results of the performance evaluation are presented in Section 4. In Section 5 we overview existing solutions proposed to facilitate the data movement between the components in multi-component workflows. Finally, we conclude with Section 6.

## 2. Data Transfer Framework

In this section we briefly review the design of the DTF.

We note that from now on we will call the component that writes the file and the component that reads it as the writer and the reader components, respectively.

### 2.1 Parallel NetCDF Semantics

The DTF can be used in workflows in which the components use the PnetCDF library for file I/O. PnetCDF is the parallel implementation of the NetCDF [20] library that allows processes to describe, write to a file and read from it array-based scientific data. Because PnetCDF is built on top of the parallel MPI-IO, the I/O can be performed on a shared file.

Before performing I/O, the user must first define the structure of the file, that is, define variables, their attributes, variable dimensions and dimension lengths. Once the structure of the file is defined, the user may call PnetCDF's API to read or write variables. In a typical PnetCDF call, the user must specify the file id and variable id, which were assigned by PnetCDF during the definition phase, specify the start coordinate and block size in each dimension for multi-dimensional variables, and pass the input or output user buffer.

### 2.2 General Overview

DTF aims to provide users of multi-component workflows with a tool that would allow them to quickly switch from file I/O to direct data transfer without needing to cardinally change the source code of the components.

First, the user must provide a simple configuration file that describes the file dependency in the workflow (example in Figure 4). It only needs to list the files that create a direct dependency between two components, i.e. if the components are coupled through this file. The DTF intercepts PnetCDF calls in the program and, if the file for which the call was made is listed in the configuration file as subject to the data transfer, the DTF handles the call accordingly. Otherwise, PnetCDF call is executed normally.

In order to transfer the data from one component to another, we treat every PnetCDF read or write call as an *I/O request*. The data transfer is performed via what we call *the I/O request matching*. First, designated processes, called *I/O request matchers*[*1], collect all read and write requests for a given file. Then, each matcher finds out who holds the requested piece of data by matching each read request against one or several write requests. Finally, the matcher instructs the processes who have the data to send it to the corresponding process who requested it. All the inter-process communication happens using MPI.

---

[*1] We note that in [19] the processes were called *masters*. We have changed the naming in order to avoid association with the master-slave programming paradigm.
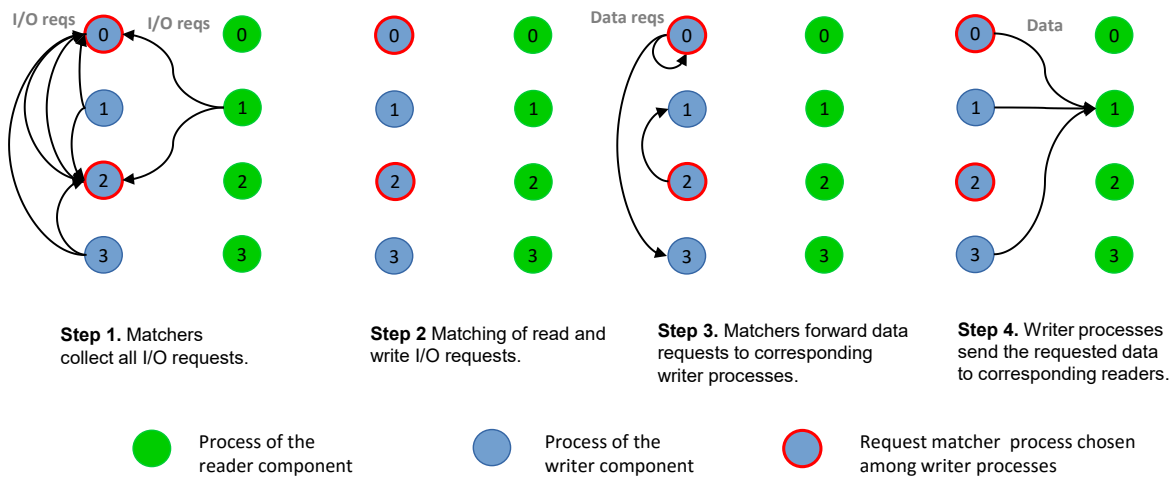
The I/O patterns of the component that writes to the file and the component that reads from it may be drastically different, however, dynamic I/O request matching makes DTF flexible and allows it to handle any kind of I/O patterns transparently for the user.

### 2.3 I/O Request Matching

When the writer component creates a file, matchers that will be handling the I/O request matching are chosen among its processes. The number of matchers can be configured by the user or else a default value is set.

When a process calls a read or write PnetCDF function for a file intended for data transfer, the DTF intercepts this call and internally creates an I/O request that stores the call metadata such as the PnetCDF variable id, data type, pointer to the user buffer and, in case the process writes a part of the variable, the corner coordinate of the sub-block and the length of the block in each dimension.

The request matching process can be divided in four steps (Figure 1). First, all the processes of the reader and writer component send all their I/O requests posted so far to corresponding matching processes (Step 1). Then, a matching process takes the next read I/O request and, based on the corner coordinate `start` of the requested array block and the block size `count`, searches for matching write requests (Step 2). The I/O pattern of the reader and writer components are not necessarily identical, therefore, one read request may be matched with several write requests, each of them - for a sub-block of the requested array block. Once a match is found, the matcher sends a message to the writer process holding the requested portion of data and asks it to send this data to the corresponding reader process (Step 3). Finally, when a writer process receives a data request from the matcher, it finds the requested data in the memory, copies it to the send buffer along with the metadata and sends it to the reader (Step 4). When the reader receives the message, it parses the metadata and unpacks the data to the user buffer.

For better performance, the requests are distributed among the matching processes and each matcher is in charge of matching requests for a particular sub-block of a multi-dimensional variable. The size of the sub-block is determined by dividing the length of the variable in the lowest (zeroth) dimension by the number of matching processes. If there is a request that overlaps blocks handled by different matchers, such a request will be split into several requests for sub-blocks, and each matcher will match the corresponding part. There is a trade-off in this approach: On one hand, the matching happens in a distributed fashion, on the other hand, if there are too many matchers the request may end up being split too many times resulting in more communication between readers and writers. Therefore, it is recommended to do some test runs of the workflow with different number of matchers to find a reasonable configuration for DTF.

## 3. Implementation

The Data Transfer Framework is implemented as a library providing API to user programs. To let the DTF intercept PnetCDF calls, we also modified the PnetCDF-1.7.0 library. The modifica-

Step 1. Matchers collect all I/O requests.

Step 2 Matching of read and write I/O requests.

Step 3. Matchers forward data requests to corresponding writer processes.

Step 4. Writer processes send the requested data to corresponding readers.

Process of the reader component

Process of the writer component

Request matcher process chosen among writer processes

Fig. 1: I/O request matching. Request matchers are marked with a red shape outline. For simplicity, only one reader process is showed to have read I/O requests.

tions were relatively small and consisted of around 100 lines of code.

The MPI library is used to transfer the data. To establish the communication between processes in the reader and writer components, we use the standard MPI API for creating an inter-communicator during the DTF initialization stage in both components. This implies that the coupled components must run concurrently.

Current version of the DTF implements a synchronous data transfer, meaning that all the processes in the two components stop their computations to perform the data transfer and resume only when all the data requested by the reader has been received. Generally, it is preferable to transfer the data to the reader as soon as it becomes available on the writer's side so that the reader could proceed with computations. However, because the I/O patterns of the two components may differ significantly, it is hard to automatically determine when it is safe to start the matching process. Therefore, we require that the user signals to the DTF when to perform a request matching for a given file by explicitly invoking a special API function in both components.

Overall, to enable the data transfer, the user needs to modify the source code of all the components of the workflow by adding API to initialize, finalize the DTF, as well as explicitly invoke the data transfer. However, we believe that these modifications are rather minimal compared to what traditional coupling software usually requires.

### 3.1 Handling of I/O Requests

Depending on the scale of the execution and the I/O pattern, matching processes sometimes may have to handle thousands of I/O requests. Using a suitable data structure to arrange the requests in such a way that matching read and write requests can be found fast is important.

Unless the variable is a scalar, an I/O request is issued for a multi-dimensional block of data. In the previous version of the DTF we used a flat representation of array coordinates and
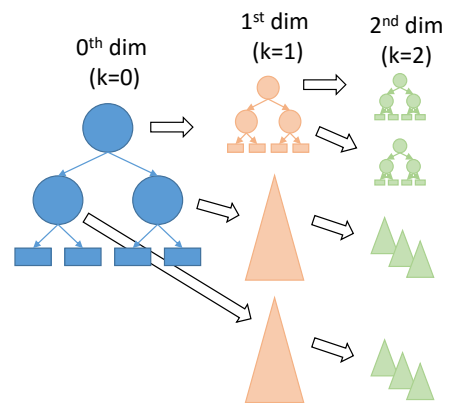


Fig. 2: An example layout of a k-d tree to arrange sub-blocks of a 3-dimensional variable.

described a multi-dimensional array block as a set of contiguous one-dimensional arrays. These one-dimensional blocks were then distributed between request matchers based on their start coordinate. Matchers arranged the blocks in a binary tree and request matching consisted of finding overlapping one-dimensional arrays. However, this approach proved to be infeasible as we started to increase the scale of the experiments: The number one-dimensional arrays grew so significantly that time to handle them became unacceptable.

In the current version of the DTF we abandoned the approach of flat array representation and represent k-dimensional data blocks as a set of k intervals. We use an augmented k-dimensional interval tree [16] to arrange these blocks in such a way that would allow us to find a block (write I/O request) that overlaps with a quired block (read I/O request) in a reasonable amount of time. Figure 2 shows an example layout of a tree that stores write requests for a 3-dimensional variable. A tree on each level (k=0,1,2) arranges intervals of the variable sub-blocks for which a write request was issued in the corresponding dimension. Each node of the tree links to the tree in the k+1 dimension. We

note that only write I/O requests are stored as an interval tree.

Read I/O requests are stored as a linked list sorted by the rank of the reader. Every time new requests metadata arrives, the matcher updates the request database and goes through the read I/O requests one by one and tries to match as many of them for a given rank as possible.

### 3.2 User API

The three main API functions provided by the DTF are the following:

- `dtf_init(config_file, component_name)` - initializes the DTF. The user must specify the path to the DTF configuration file and state the name of the current component which should match one of the component names in the configuration file;
- `dtf_finalize()` - finalizes the DTF;
- `dtf_transfer(filename)` - invokes the data transfer for file *filename*;

All the API functions are collective: `dtf_init()` and `dtf_finalize()` must be invoked by all processes in both components, while `dtf_transfer()` must be invoked only by processes that share the file.

During the initialization, based on the DTF configuration file, each component finds out all other components with whom it has an I/O dependency and establishes a separate MPI inter-communicator for every such dependency. All the further inter-component communication happens via this inter-communicator.

A `dtf_transfer()` call should be added after corresponding PnetCDF read/write calls in the source code of both, reader and writer components. The call will not complete until the reader receives the data for all the read I/O requests posted before `dtf_transfer()` was invoked, therefore, it is user's responsibility to ensure that the components call the function in the correct place in the code, that is, that the writer does not start matching I/O until all the write calls for the data that will be requested in the current transfer phase have been posted as well. `dtf_transfer()` function can be invoked arbitrary number of times but this number should be the same for both components. We note that, because this function acts like a synchronizer between the reader and writer components, the recommended practice is to invoke it just once after all the I/O calls and before the file is closed.

By default, the DTF does not buffer the user data internally. Therefore, the user should ensure that the content of the user buffer is not modified between the moment the write PnetCDF call was made until the moment the data transfer starts. Otherwise, data buffering can be enabled in the DTF configuration file. In this case, all the data buffered on the writer's side will be deleted when a corresponding transfer function is completed.

### 3.3 Example Program

A simplified example of a writer and reader components is presented Figures 3a and 3b, as well as their common DTF configuration file (Figure 4). To enable the direct data transfer it was enough to add three lines of code to each component — to initialize, finalize the library and to invoke the data transfer — and

```
/* Initialize DTF */
dtf_init(dtf_inifile, "wrt");
/* Create file */
ncmpi_create("restart.nc",...);
<...>
 /* Write some data */
ncmpi_put_vara_float(...);
/* Write some more data */
ncmpi_put_vara_float(...);
/* Perform I/O request matching */
dtf_transfer("restart.nc");
/* Close the file */
ncmpi_close(...);
/* Finalize DTF */
dtf_finalize();
```

(a) Component writing to file

```
/* Initialize DTF */
dtf_init(dtf_inifile, "rdr");
/* Open the file */
ncmpi_open("restart.nc",...);
<...>
/* Read all data at once */
ncmpi_get_vara_float(...);
/* Perform I/O request matching */
dtf_transfer("restart.nc");
/* Close the file */
ncmpi_close(...);
/* Finalize DTF */
dtf_finalize();
```

(b) Component reading from file

Fig. 3: Sample code using the DTF API

provide a simple configuration file.

```
[INFO]
ncomp=2   ! number of components
comp_name="rdr"  ! component name
comp_name="wrt"
ioreq_distrib_mode="range"  !divide by dim length
buffer_data=0
[FILE]
filename="restart.nc"
writer="wrt"  !component that writes to the file
reader="rdr"  !component that reads from the file
iomode="memory"  !enable direct transfer
```

Fig. 4: DTF configuration file

## 4. Evaluation

We demonstrate the performance of DTF using the S3D-IO[*2] benchmark program.

S3D-IO[14] is the I/O kernel of the S3D combustion simulation code developed at Sandia National Laboratory. In the benchmark, a checkpoint file is written at regular intervals. The checkpoint consists of four variables — two three-dimensional and two four-dimensional — representing mass, velocity, pressure, and temperature. All four variables share the lowest three spatial dimensions X, Y and Z which are partitioned among the processes in block fashion. The value of the fourth dimension is fixed.

We imitate a multi-component execution in S3D-IO by running concurrently two instances of the benchmark: Processes of

---

*2 Available at http://cucis.ece.northwestern.edu/projects/PnetCDF/#Benchmarks

the first instance write to a shared file, processes in the second instance read from it. Each test is executed at least eight times and an average value of the measured parameter is computed. To determine the number of matchers necessary to get the best performance for data transfer, we first execute several tests of S3D-IO varying the number of matching processes and use the result in the subsequent tests.

In the tests with the direct data transfer, the I/O time was measured in the following manner. On the reader side, it is the time from the moment the reader calls the data transfer function to the moment all its processes received all the data they had requested. On the writer's side, the I/O time is the time between the data transfer function and the moment the writer receives a notification from the reader indicating that it had got all the requested data. The I/O time also includes the time to register the metadata of the PnetCDF I/O calls and to buffer the data, if this option is enabled. In all the tests the writer component invokes the I/O matching function before the reader and completes the matching only after the reader has finished, therefore, by data transfer time we hereafter assume the I/O time of the writer component unless stated otherwise. The runtime of the workflow is measured from the moment two components create an MPI inter-communicator inside the `dtf_init()` function and the moment it is destroyed in `dtf_finalize()` as these two functions work as a synchronization mechanism between the reader and writer components.

All the experiments were executed on K computer [10]. Each node has an 8-core 2.0 GHz SPARC64 VIIIfx CPU equipped with 16 GB of memory. Nodes are connected by a 6D mesh/torus network called Tofu [1] with 5 GB/s x 2 bandwidth in each link. Compute nodes in K computer have access to a local per-node file system as well as a global shared file system based on Lustre file system.

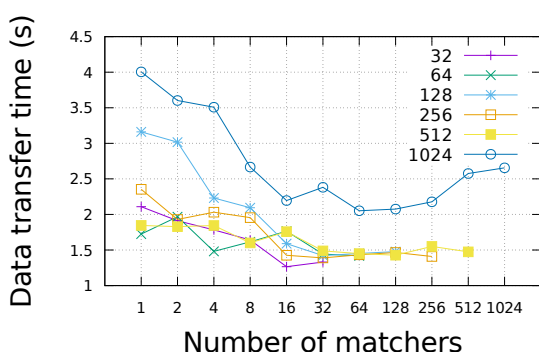### 4.1 Choosing the Number of Matching Processes



Fig. 5: Data transfer time for various test sizes and number of matching processes per component.

To get the best performance, it is recommended that the user chooses the number of matching processes that will perform I/O matching instead of using the default configuration of one matcher per 64 processes. This number is application-dependent. The load on a matching process is determined by the number of read and write I/O requests the process has to match. For example, if all reader and writer processes perform I/O symmetrically

and the size of the variable in the zeroth dimension divides by the number of matchers, the number of I/O requests one matcher will have to match roughly equals the number of I/O requests one process generates multiplied by the number of processes in both components.

Depending on the I/O pattern, increasing the number of matchers does not always decrease the number of I/O requests per matcher, but it generally improves the throughput of data transfer. The reason is that rather than waiting for one matching process to match requests for one block of a multi-dimensional array, multiple processes can match sub-blocks of it in parallel. Consequently, the reader may start receiving the data earlier.

To find an optimal number of matchers, we run tests of different sizes — from 32 processes per component up to 1024 — with a problem size such that each process reads or writes 1 GB of data. In each test we then vary the number of matchers and measure the time to transfer the data. The results in Figure 5 show that increasing the number of matchers up to some point improves the transfer time and then the performance starts decreasing. The reason for this is that an I/O request for a block of data may be split into several requests for sub-blocks between multiple matchers and, if the number of matchers is too big, the request is over-split and it takes more smaller messages to deliver all the data to the reader. Based on this result, for our further tests we use the following setting: for tests with up to 256 processes in one component, each writer process functions as a matcher, for tests with 512 processes per component — four processes in one work-group, i.e. 128 matchers in total. Finally, for tests with 1024 processes per component the work-group size is 16, i.e. there are 64 matchers in total.

### 4.2 Scalability

We first demonstrate how the DTF scales compared to file I/O (PnetCDF) by measuring the read and write bandwidth for weak and strong scaling tests. In this test, processes write to a shared file using non-blocking PnetCDF calls. To measure the I/O bandwidth, we divide the total file size by the respective read or write I/O time. The results for the strong and weak scaling are presented on Figures 6 and 7. The X axis denotes the number of processes in one component. We point out that the Y-axis is logarithmic in these plots. Figures 6a and 7a show the total execution time of the coupled workflow.

In all tests each process executes a PnetCDF read or write function four times — one per variable, i.e. each process generates four I/O requests.

In the strong scaling test, we fix the file size to 256 GB and vary the number of processes in one component. We note that, due to the node memory in K computer limited to 16 GB, the results in Figure 6 start from the test with 32 processes per component. In the weak scaling tests, we fix the size of the data written or read by one process to 256 MB, thus, in the test with one process per component the file size is 256 MB, in the test with 1024 processes — 256 GB.

As we see, DTF significantly outperforms file I/O in all tests. We also notice that the read bandwidth in all tests with the direct data transfer is always higher than the write bandwidth. The rea-

son for this is the timing when the matching starts in the reader component. When the reader opens the file, it first inquires the writer component about the list of the matcher ranks and waits for the reply. The writer component often processes this inquiry when it has already entered the data transfer function, so by the time the reader initiates the data transfer, the writer has already posted all write requests and is ready to do the matching and transfer the data immediately. For this reason, from the reader's point of view, the matching finishes faster than from the writer's point of view, hence, the read bandwidth is higher.

The bandwidth using the data transfer does flatten eventually in the strong scaling test (Figure 6b and 6c), because the size of the data sent by one process decreases and the overhead of doing the request matching and issuing data requests starts to dominate the transfer time. In the weak scaling tests in Figure 7 the data transfer time grows slower as the amount of data to transfer by one process stays the same and the overhead of the I/O request matching is relatively small. Hence, the total I/O bandwidth increases faster than in the strong scalability test.

### 4.3 DTF Performance under I/O Load

Other major factors that impact the data transfer time apart from the number of matching processes are the size of data to transfer and the total number of I/O requests to be matched. To measure the former we perform data transfer for files of various sizes while the number of I/O requests per matcher stays the same. To evaluate the impact of the number of I/O requests, we fix the file size to 256 GB and increase the number of I/O requests a matcher process matches by manipulating how the I/O requests are distributed among matchers. By default, the size of the variable sub-block for which a matcher process matches read and write requests is defined by dividing the variable in the zeroth dimension by the number of matchers. An I/O request is split in the zeroth dimension based on this stripe size and distributed among the matchers in a round robin fashion. In this experiment, we vary the value of the stripe size which effectively changes the number of I/O requests each matcher has to handle.

In both experiments there are 1024 processes per component and there is one matcher per 16 processes. We also note that two out of four variables in S3D-IO have the zero dimension length fixed to 11 and 3, respectively. This is smaller than the number of matchers (64) and results in asymmetrical distribution of work among matchers. For this reason, in the two experiments, on top of the average number of I/O requests per matcher, a small group of matchers has to match approximately 4,000 more I/O requests.

Figure 8 shows the results of the first experiment. The file size was gradually increased from 8 GB to 2 TB. Each matcher process matched on average 576 requests in every test. We measured the time for actual matching of read and write requests — it took only around 2% of the whole data transfer time, thus, we conclude that most of the time was spent on packing and sending the data to reader processes. Thanks to the fast torus network in K computer, sending 2 TB of data over network took less than 3 seconds.

In the second experiment (Table 1) the file size is fixed, i.e. in every test each process transfers the same amount of data. The

Table 1: DTF performance for different number of I/O requests

| Average number of I/O requests per matcher | Data transfer time (s) | Time to match read and write requests (s) |
|---|---|---|
| 576 | 1.799 | 0.041 |
| 1,088 | 1.498 | 0.031 |
| 2,144 | 2.107 | 0.046 |
| 4,224 | 2.061 | 0.045 |
| 8,448 | 2.108 | 0.058 |
| 16,832 | 1.777 | 0.085 |

matching processes handled from 576 to 16,832 I/O requests, plus the additional requests for some matchers due to the imbalance. We expect that in this experiment it is the request matching process that will have the biggest impact on the data transfer time as the number of requests grow. However, according to the Table 1, the actual request matching took on average no more than 2-3% of the data transfer time and only in the test with 16,832 requests per matcher the matching took around 5% of the data transfer.

Moreover, we observe that despite the growing number of I/O requests per matcher, the time to perform the data transfer actually decreases in some cases. One explanation for this could be that, when we decrease the stripe size by which the I/O requests are distributed, one matcher becomes in charge of several smaller sub-blocks located at a distance from each other along the zeroth dimension, rather than just one contiguous big sub-block. And this striping may accidentally align better with the I/O pattern of the program, so the matcher ends up matching requests for the data that was written by it. Then, instead of having to forward a data request to another writer process, the matcher immediately can send the data to the reader.

Overall, we conclude that the DTF shows stable performance under increased load of the amount of data that needs to be transfered as well as the load on the matching processes.

## 5. Related Work

A number of works has addressed the data movement problem, the file I/O bottleneck in particular, in multi-component workflows, and different coupling toolkits have been designed for such workflows [7], especially in Earth sciences [12], [21], [22] applications. Such toolkits often provide not only the data movement feature but also allow to perform various data processing during the coupling phase, such as data interpolation or changing the grid size.

For example, DART [6] is a software layer for asynchronous data streaming, it uses dedicated nodes for I/O offloading and asynchronously transferring the data from compute nodes to I/O nodes, visualization software, coupling software, etc. The ADIOS [15] I/O library is built on top of DART and provides additional data processing functionality. However, both, DART and ADIOS require to use a special API for I/O. In additional, ADIOS uses its own non-standard data format for files.

Other coupling approaches include implementing a virtual shared address space accessible by all the components [5], or using dedicated staging nodes to transfer the data from compute job to post-process analysis software during the runtime [23]. In [3], the authors propose a toolkit utilizing the type-based publisher/-subscriber paradigm to couple HPC applications with their analytics services. The toolkit uses a somewhat similar concept to
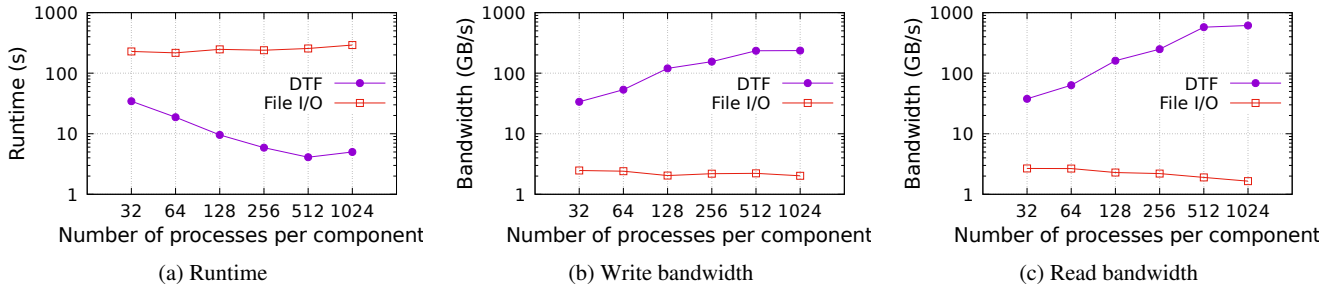
(a) Runtime　　　　　　　(b) Write bandwidth　　　　　　　(c) Read bandwidth

Fig. 6: Strong scaling of S3D-IO. Y-axis is in logarithmic scale in all plots.



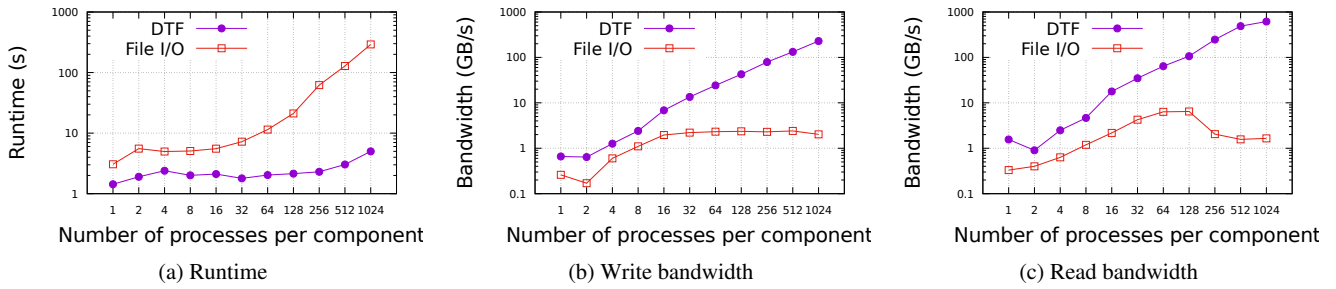(a) Runtime　　　　　　　(b) Write bandwidth　　　　　　　(c) Read bandwidth

Fig. 7: Weak scaling of S3D-IO. Y-axis is in logarithmic scale in all plots.
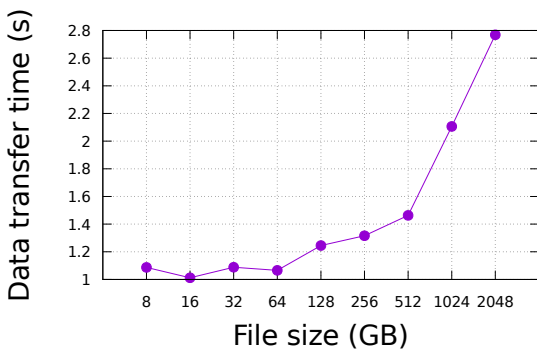


Fig. 8: DTF performance for various file sizes.

inter-component data transferring as proposed in this work, however, they rely on the ADIOS library underneath which the coupling toolkit was built which includes the description of the I/O pattern of the components. Additionally, in our work the matching process is simpler in a way that it takes fewer steps to perform the data transfer.

Providing support to multi-component executions on a system-level is another approach to facilitating the inter-component interaction[8], [9]. Current HPC systems usually do not allow overlapping of resources allocated for one executable file. Thus, each component in a multi-component workflow ends up executing on a separate set of nodes and, consequently, the problem of data movement between the components arises. But, for example, in cloud computing, several virtual machines can run on the same node and communicate with each other via shared memory or virtual networking. It has been previously proposed to use virtualization techniques in HPC as well. For example, in [9], the authors show that such virtualization can be used in an HPC environment to allow more efficient execution of multi-component workflows with minimal costs. However, the virtualization is not yet widely adopted in HPC systems.

The difference of our solution with the aforementioned approaches is that our goal was to provide a simple framework that would allow to switch from file I/O to data transfer with minimal efforts and without having to rewrite the I/O kernels of the workflow components. The closest solution that we are aware of is the I/O Arbitration Framework (FARB) proposed in [13]. However, the framework was implemented for applications using NetCDF I/O library, that is, it assumes the file-per-process I/O pattern and a process-to-process mapping of data movement. Moreover, during the coupling stage in FARB, contents of the whole file were transferred to the other component's processes regardless of whether the process actually required the whole data or not. In our work, we determine at runtime what data needs to be transferred and only send this data.

## 6. Conclusion

Tackling of a large scientific task can be done using a multi-component approach where each component solves a smaller sub-task and exchanges the result with other components. Combined with the computational power of modern HPC systems, this approach can help the scientists to solve sophisticated problems and perform complex data processing. However, efficient exchange of computation results between the components in such workflows is a common issue, especially when components are loosely coupled via files. A number of solutions has been proposed, each of them having its pros and cons.

In this work we proposed a data transfer framework called DTF that can be used as a coupling solution for multi-component workflows that use PnetCDF API for file I/O. The DTF intercepts the PnetCDF calls and bypasses the file system by sending the data directly to the corresponding processes that require it. It automatically detects what data should be sent to which processes in the other component through a process of, what we call, I/O re-

quest matching. The DTF requires minimal modifications of the program and, more importantly, does not require the modification of component's original I/O code that uses PnetCDF.

We demonstrated that the DTF shows stable performance under different conditions and that the latest version scales better compared to the previous version thanks to the improved I/O matching process. In the future we plan to improve the load balancing of the I/O request matching by finding a better way to distribute I/O requests among the matchers. However, the results we obtained so far are promising and should help multi-component workflows to cut on I/O time without having to significantly modify them.

## References

[1] Ajima, Y., Sumimoto, S. and Shimizu, T.: Tofu: A 6D Mesh/Torus Interconnect for Exascale Computers, *Computer*, Vol. 42, No. 11, pp. 36–40 (online), DOI: 10.1109/MC.2009.370 (2009).

[2] Argonne National Laboratory and Northwestern University: Parallel NetCDF (Software), http://cucis.ece.northwestern.edu/projects/PnetCDF/.

[3] Dayal, J., Bratcher, D., Eisenhauer, G., Schwan, K., Wolf, M., Zhang, X., Abbasi, H., Klasky, S. and Podhorszki, N.: Flexpath: Type-Based Publish/Subscribe System for Large-Scale Science Analytics, *2014 14th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing*, pp. 246–255 (2014).

[4] Deelman, E., Peterka, T., Altintas, I., Carothers, C. D., van Dam, K. K., Moreland, K., Parashar, M., Ramakrishnan, L., Taufer, M. and Vetter, J.: The future of scientific workflows, *The International Journal of High Performance Computing Applications*, Vol. 0, No. 0, p. 1094342017704893 (online), DOI: 10.1177/1094342017704893.

[5] Docan, C., Parashar, M. and Klasky, S.: DataSpaces: An Interaction and Coordination Framework for Coupled Simulation Workflows, *Proceedings of the 19th ACM International Symposium on High Performance Distributed Computing*, HPDC '10, ACM, (online), DOI: 10.1145/1851476.1851481 (2010).

[6] Docan, C., Parashar, M. and Klasky, S.: Enabling high-speed asynchronous data extraction and transfer using DART, *Concurrency and Computation: Practice and Experience*, Vol. 22, No. 9, pp. 1181–1204 (online), DOI: 10.1002/cpe.1567 (2010).

[7] Dorier, M., Dreher, M., Peterka, T., Wozniak, J. M., Antoniu, G. and Raffin, B.: Lessons Learned from Building In Situ Coupling Frameworks, *Proceedings of the First Workshop on In Situ Infrastructures for Enabling Extreme-Scale Analysis and Visualization*, New York, NY, USA, ACM, (online), DOI: 10.1145/2828612.2828622 (2015).

[8] Kocoloski, B. and Lange, J.: XEMEM: Efficient Shared Memory for Composed Applications on Multi-OS/R Exascale Systems, *Proceedings of the 24th International Symposium on High-Performance Parallel and Distributed Computing*, HPDC '15, New York, NY, USA, ACM, (online), DOI: 10.1145/2749246.2749274 (2015).

[9] Kocoloski, B., Lange, J., Abbasi, H., Bernholdt, D. E., Jones, T. R., Dayal, J., Evans, N., Lang, M., Lofstead, J., Pedretti, K. and Bridges, P. G.: System-Level Support for Composition of Applications, *Proceedings of the 5th International Workshop on Runtime and Operating Systems for Supercomputers*, ROSS '15, New York, NY, USA, ACM, (online), DOI: 10.1145/2768405.2768412 (2015).

[10] Kurokawa, M.: The K computer: 10 Peta-FLOPS supercomputer, *The 10th International Conference on Optical Internet (COIN2012)* (2012).

[11] LANL, NERSC, S.: *APEX Workflows (white paper)* (2016).

[12] Larson, J., Jacob, R. and Ong, E.: The Model Coupling Toolkit: A New Fortran90 Toolkit for Building Multiphysics Parallel Coupled Models, *The International Journal of High Performance Computing Applications*, Vol. 19, No. 3, pp. 277–292 (online), DOI: 10.1177/1094342005056115 (2005).

[13] Liao, J., Gerofi, B., Lien, G.-Y., Nishizawa, S., Miyoshi, T., Tomita, H. and Ishikawa, Y.: Toward a General I/O Arbitration Framework for NetCDF Based Big Data Processing, *Proceedings of the 22nd International Conference on Euro-Par*, (online), DOI: 10.1007/978-3-319-43659-3_22 (2016).

[14] Liao, W.-k. and Choudhary, A.: Dynamically Adapting File Domain Partitioning Methods for Collective I/O Based on Underlying Parallel File System Locking Protocols, *Proceedings of the 2008 ACM/IEEE Conference on Supercomputing*, SC '08, Piscataway, NJ, USA, IEEE Press, (online), available from ⟨http://dl.acm.org/citation.cfm?id=1413370.1413374⟩ (2008).

[15] Lofstead, J., Zheng, F., Klasky, S. and Schwan, K.: Adaptable, metadata rich IO methods for portable high performance IO, *2009 IEEE International Symposium on Parallel Distributed Processing*, pp. 1–10 (online), DOI: 10.1109/IPDPS.2009.5161052 (2009).

[16] Mehta, D. P. and Sahni, S.: *Handbook Of Data Structures And Applications*, Chapman & Hall/CRC (2004).

[17] Message Passing Interface Forum: MPI: A Message-Passing Interface Standard, Version 3.1, www.mpi-forum.org/docs/ (1995).

[18] Miyoshi, T., Lien, G. Y., Satoh, S., Ushio, T., Bessho, K., Tomita, H., Nishizawa, S., Yoshida, R., Adachi, S. A., Liao, J., Gerofi, B., Ishikawa, Y., Kunii, M., Ruiz, J., Maejima, Y., Otsuka, S., Otsuka, M., Okamoto, K. and Seko, H.: Big Data Assimilation; Toward Post-Petascale Severe Weather Prediction: An Overview and Progress, *Proceedings of the IEEE*, Vol. 104, No. 11 (online), DOI: 10.1109/JPROC.2016.2602560 (2016).

[19] Tatiana, V. M., Wei-Keng, L., Balazs, G., Yutaka, I. and Alok, C.: Improving Multi-component Application Performance by Silently Replacing File I/O with Direct Data Transfer: Preliminary Results, Technical Report IPSJ 2017-HPC-158.

[20] UNIDATA: Network Common Data Form, http://www.unidata.ucar.edu/software/netcdf/.

[21] Valcke, S.: The OASIS3 coupler: a European climate modeling community software, *Geosci. Model Dev.*, Vol. 6 (online), DOI: doi:10.5194/gmd-6-373-2013 (2013).

[22] Valcke, S., Balaji, V., Craig, A., DeLuca, C., Dunlap, R., Ford, R. W., Jacob, R., Larson, J., O'Kuinghttons, R., Riley, G. D. and Vertenstein, M.: Coupling technologies for Earth System Modelling, *Geoscientific Model Development*, Vol. 5, No. 6, pp. 1589–1596 (online), DOI: 10.5194/gmd-5-1589-2012 (2012).

[23] Vishwanath, V., Hereld, M. and Papka, M. E.: Toward simulation-time data analysis and I/O acceleration on leadership-class systems, *2011 IEEE Symposium on Large Data Analysis and Visualization*, (online), DOI: 10.1109/LDAV.2011.6092178 (2011).