

SPECIAL ISSUE PAPER

Interfacing HDF5 with a scalable object-centric storage system on hierarchical storage

Jingqing Mu¹  | Jerome Soumagne¹  | Suren Byna² | Quincey Koziol² | Houjun Tang² | Richard Warren¹

¹The HDF Group, Champaign, Illinois

²Lawrence Berkeley National Laboratory, Berkeley, California

Correspondence

Jingqing Mu, 410 E. University Ave, Champaign, IL 61820.

Email: kmu@hdfgroup.org

Funding information

US Department of Energy, Grant/Award Numbers: DE-AC02-05CH11231, DE-SC0016454

Summary

Object storage technologies that take advantage of multitier storage on HPC systems are emerging. However, to use these technologies at present, applications have to be modified significantly from current I/O libraries. HDF5, a widely used I/O middleware on HPC systems, provides a virtual object layer (VOL) that allows applications to connect to different storage mechanisms transparently without requiring significant code modifications. We recently designed the proactive data containers (PDC) object-centric storage system that provides the capabilities of transparent, asynchronous, and autonomous data movement taking advantage of multiple storage tiers—a decision that has so far been left upon the user on most current systems. To enable PDC's features through HDF5 without modifying application codes, we have developed an HDF5 VOL connector that interfaces with PDC. We present in this article the connector interface and evaluate its performance on Cori, a Cray XC40 supercomputer located at the National Energy Research Scientific Computing Center (NERSC). Our evaluation demonstrates up to an 8× improvement compared with HDF5 that has the most recent optimizations.

KEYWORDS

data handling, HDF5, large-scale systems, memory management, object-centric models

1 | INTRODUCTION

The challenges for scientific applications on upcoming HPC systems, when rapidly moving toward exascale, are known from three directions: extreme parallelism, a deepening heterogeneous memory hierarchy, and massively increasing data by volume and complexity. Current data management and I/O technologies present severe limitations in this regard: the POSIX and MPI I/O standards that are the basis for existing I/O libraries and parallel file systems have fundamental restrictions in the areas of scalable metadata operations, semantics-based data movement, performance tuning, asynchronous operations, and scalable consistency of distributed operations—such that simple and efficient methods of data management that can address these challenges are critical for running scientific applications on future HPC systems.

In particular for I/O libraries, one of the challenges to address the diverse performance characteristics of deep storage hierarchy expected in exascale systems is the capability and efficiency of data movement across storage levels. New architectures are considered to contain multiple layers of storage, such as NVRAM on compute nodes, SSD-based burst buffers shared by compute nodes, disk or SSD-based parallel file systems, and tape-based archival storage. Typically, parallel file systems are unaware of multilevel storage hierarchies and require the use of external middleware to manage the hierarchy. Traditional HPC data management and movement solutions were designed for simpler systems, which are designed to manage each storage layer separately; similarly, scientific data models were designed for a two-tiered storage hierarchy. Another critical deficiency

of traditional file systems is metadata management, where files are managed with a small amount of prescriptive metadata leading to a performance bottleneck from metadata servers.

Object-based storage systems provide semantics that have the potential to reduce the complexity of storage systems, as well as to improve performance. We have developed a user-space object-centric storage and data management system, called *proactive data containers* (PDC).^{1,2} PDC provides scalable distributed metadata management¹ and the capabilities of transparent, asynchronous, and autonomous data movement to multiple storage tiers.² PDC offers a programming interface that applications can use in order to take advantage of the data and metadata management services. A PDC *container* is a container that may reside in a single storage layer (ie, memory, burst buffer, disk) or span across multiple layers. It contains both the metadata and data objects. The PDC system provides an interface for creating, updating, retrieving, and deleting data objects and for managing metadata on those objects. It moves us away from existing file-oriented methods and instead allows us to explore novel object-oriented data management methods in an autonomous way.

HDF5 has been widely used in the context of HPC and big data as an I/O middleware capable of supporting extreme scale and complex data structures. HDF5's virtual object layer (VOL) is a storage abstraction layer within the HDF5 library that is designed to target different storage mechanisms while preserving HDF5 object metadata. The VOL design allows applications to connect to different storage mechanisms transparently without significant code modifications. Several HDF5 VOL connectors have been developed, for instance, PLFS,³ Data Elevator,⁴ or more recently DAOS⁵—offering HDF5 applications an easy way to use data storage systems transparently with a significant I/O performance increase over POSIX I/O.

Toward making PDC services available to legacy HDF5 applications with minimal source code changes, we implement and present in this article an HDF5 VOL connector interface to PDC. This article and contributions that it presents address the following objectives:

1. Enabling object-oriented storage through HDF5 APIs and library.
2. Enabling implicit and asynchronous data movement through existing HDF5 APIs with minimal code modification.
3. Implementing data movement strategies using TCP and Cray GNI transports.
4. Evaluating and demonstrating an object-centric HPC storage system for scientific use cases using HDF5 APIs.

This article is organized as follows: we first discuss related work in Section 2, and then introduce the object-based PDC system in Section 3, which enables asynchronous data movement. In Section 4, we focus on the HDF5 VOL and provide details of our PDC VOL connector implementation. In Section 5, we provide experimental results evaluating new methodologies in HDF5 utilizing I/O patterns representative of science applications on HPC systems. We conclude our work in Section 6.

2 | RELATED WORK

POSIX I/O⁶ is known for describing the file access API, data model, and data consistency semantics. Current parallel file systems, such as PVFS,^{7,8} Lustre,⁹ GPFS,¹⁰ and NFS¹¹ were all designed to comply with the POSIX I/O standard. As the original POSIX I/O design was not intended for highly concurrent programming models, which are common in HPC systems,¹² that design has now become a performance bottleneck. With an increasing number of memory storage layers and complexity of storage system interactions, the issue of I/O performance is getting imperative and significantly hinders the overall performance of applications.^{13,14} Research efforts have been made to relax the POSIX semantics and alleviate the I/O bottleneck from high-level libraries (eg, HDF5,¹⁵ netCDF,¹⁶ ADIOS¹⁷), I/O middleware (eg, MPI-IO,¹⁸ TAPIOCA¹⁹), to I/O forwarding layers.²⁰ All of these provide an array-based data model to organize the data and define data access semantics. However, deep memory and storage hierarchy introduced into modern supercomputer systems further increase POSIX I/O limitations.

Object-based storage has been proposed²¹⁻²⁴ to overcome the limitations of current parallel file systems, which has long been considered a potential solution for managing rich metadata in scalable environments. It describes an abstract data container that consists of many byte-streams (or objects), each with related attributes. RADOS,²⁵ Amazon S3,²⁶ and OpenStack Swift²⁷ have been developed for managing data as objects and storing them in a flat namespace. DAOS²⁸ is an object-based file system solution currently under development, which provides asynchronous data movement and manages objects in a hierarchical storage with multiple layers, whose scalability is still under evaluation and the features are in development.⁵ Furthermore, efforts to implement object-based storage is attempted on individual layers separately, that is, on disks, in NVRAM, and in memory. SSDUP²⁹ proposed to redirect data writes to burst buffer when it detects random accesses for potential high latency if writing to HDDs. The data in burst buffer are later flushed to HDDs when the size is over half of the burst buffer capacity.

Data Elevator⁴ provides automatic caching and data movement across multiple levels of storage hierarchies. It uses shared burst buffer as a caching layer before writing data to file system. UniviStor³⁰ integrates hierarchical and distributed storage devices into a unified view of memory distributed on compute nodes and storage layers in heterogeneous HPC storage and achieves better performance than Data Elevator.⁴ Both Data Elevator and UniviStor retain the supporting file format of the data management software that they use. Data Elevator transfers data asynchronously to the final destination in the background, while UniviStor waits for the data to be written to persistent storage. Toward object-centric

hierarchical storage system, PDC takes advantage of deep storage hierarchies and provides efficient strategies to support autonomous and asynchronous data management by the PDC service, as well as targeting deep storage hierarchies.²

3 | PDC SYSTEM ARCHITECTURE

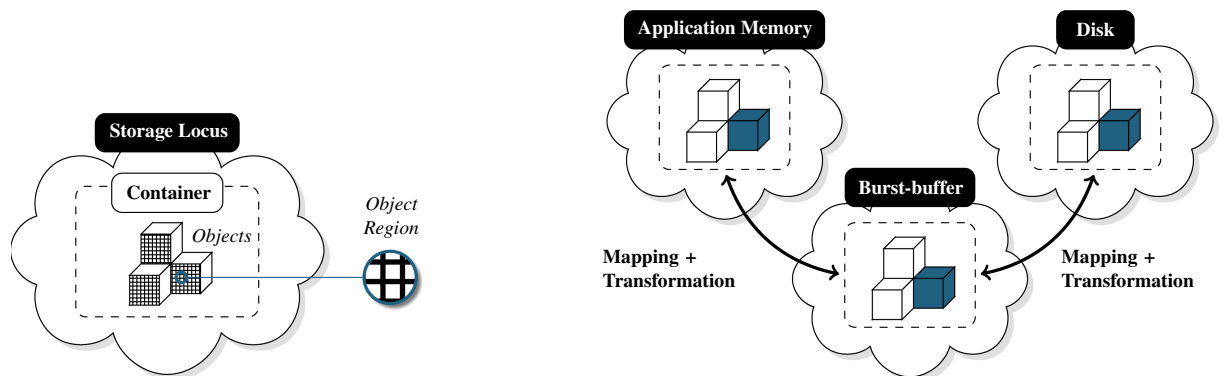
In this section, we provide an overview of the *proactive data containers* system and focus on the capabilities and semantics that it provides. For full details of the PDC system implementation, please refer to our previous publications.^{1,2}

3.1 | PDC data model and semantics

As shown in Figure 1A, PDC organizes data as a set of *objects* within a *Container*. Object is a generic term to describe byte streams in an abstract manner. Parts of objects are described using the term *Regions*, where the actual data as well as the metadata associated with it is stored. Region is the basic and fine-grained unit for data movement operations in PDC—these operations are further described in more details. In addition, all of the previously mentioned entities include *Properties*, regarded as metadata, which contain the descriptive information that is set by the user, or generated by PDC. The properties contain prescriptive metadata such as the data type and dimensions that describe an object and provenance such as user and application information. These objects are managed by PDC services and can be placed at any level of the storage system, and hence, containers, which consist of scientific data, are abstracted within the entire storage stack. This approach spreads the data over different locations within a storage level, which we reference as *storage locus*. As shown in Figure 1B, this representation is also augmented by two types of operations: *object mapping* operations between a memory buffer and an abstract PDC object and *transformation* operations while data are being moved from one storage locus to the other.

We introduced the concept of *object mapping* in PDC in previous work.² As opposed to explicit read and write semantics, object mapping makes data movement operations *implicit* to the user by defining a memory to storage relationship with a PDC entity. Similar to POSIX `mmap()` semantics, where a file can be mapped to a region of memory, PDC's mapping operations allow PDC object regions to be accessed just like an array in a program. All the user needs to do is to create a mapping between a region within the memory of an application and a region within a global PDC object. Once a mapping is established, data movement can occur to keep updates globally visible. However, as opposed to standard POSIX `mmap()` operations, concurrent access can and is expected to occur. Therefore, PDC applications are required to use an explicit lock operation on the object before modifying its associated memory and to release that lock when the modification is done. Unlocking a region allows the PDC system to start data movement and to globally propagate the modified data.

Data movement and I/O in PDC is realized *asynchronously*. Once the data have been transferred to a storage locus, further transfers to deeper levels of the storage hierarchy can be realized by PDC without the need for an application to wait for their completion. This capability provides the opportunity for applications to overlap computation with I/O operations and we can also make the safe assumption that data that are written to deeper storage tiers will always fit into that tier, memory representing the lower level, and disk representing the higher level. It is worth noting though that application's buffers, which are mapped, can only be reused and modified once a lock is reacquired, hence when the transfer to the



(A) High-level representation of the PDC data model with container, object, and region relationship.

(B) Containers and objects within them can be mapped to one another temporarily or permanently, and have transformations occur during I/O.

FIGURE 1 PDC representation. Abstracted data can reside at any level of the storage hierarchy

first level of storage hierarchy has completed. Figure 2 illustrates that mechanism and shows the data flow after a region unlock request has been initiated. If the region is mapped, the data will be first moved from the client's memory to the data server, and once it is safely transferred, the region lock can be released. PDC, in the meantime though, can carry on moving that data to other storage tiers as needed.

3.2 | User-space client-server model

To execute these operations and manage data, PDC uses a client-server model. Designing a client-server middleware for HPC can be a difficult process, both in terms of ease of deployment alongside the user's application and in terms of system resource management. PDC services, though, are designed to run in user-space as an additional service process with minimal disruption to the application. In our client-server architecture, PDC servers are responsible for executing both metadata and data management operations. We have currently implemented two different modes for users to deploy the PDC servers in user-space:

1. *shared mode*, where the server processes run on the compute nodes alongside the client processes and share CPU and memory resources, as shown in Figure 3A;
2. *dedicated mode*, where all server processes are placed on dedicated nodes that are separate from the nodes where the client processes are running, as shown in Figure 3B.

In the first case, the PDC system can take advantage of shared-memory for efficient data movement between node-local clients and servers, while in the second case, the PDC system must make use of the native interconnect for high-speed transfers. Users can start any number of PDC servers suitable for the application workload. In shared mode, users are expected to only reserve one core per compute node to run a PDC server while the rest of the cores may be used to run the application processes. In dedicated mode, servers and clients are all allocated to separate nodes, therefore the number of servers used for PDC tasks directly depends on the user workload and the number of nodes that are available on the system.

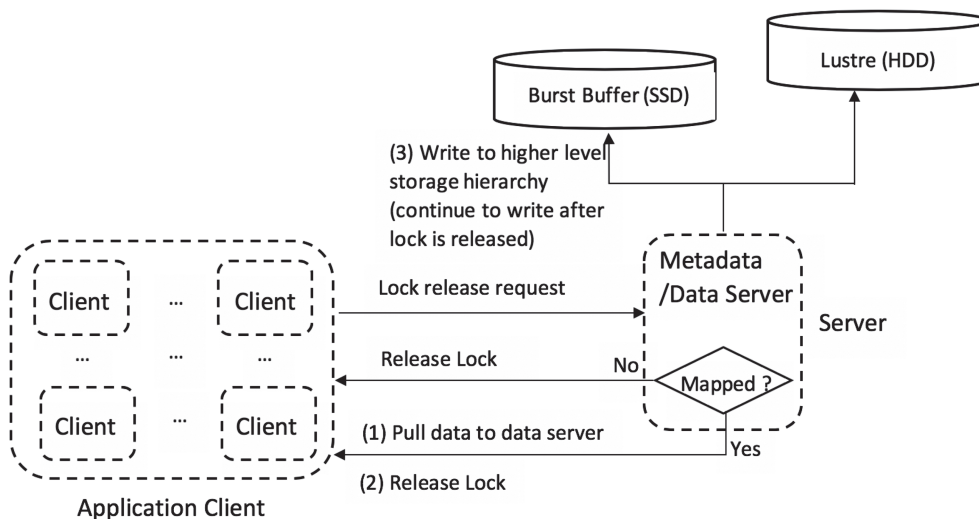
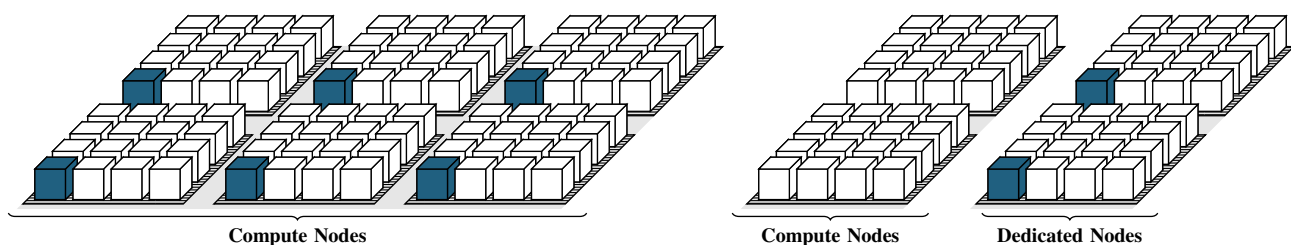


FIGURE 2 Data flow through PDC on region unlock requests



(A) Shared server modes: servers and clients are located on the same node.

(B) Dedicated server modes: servers are on separate nodes.

FIGURE 3 PDC service deployment modes

FIGURE 4 PDC within virtual object layer. All of the HDF5 I/O related calls are routed to the corresponding VOL connector

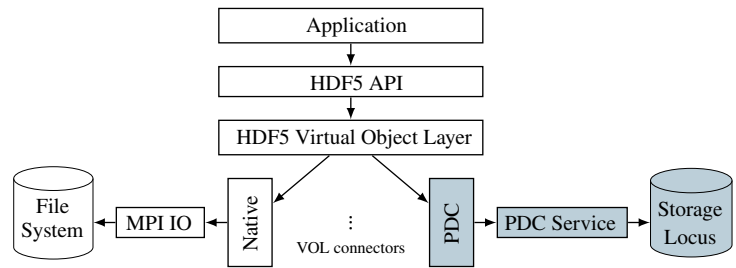
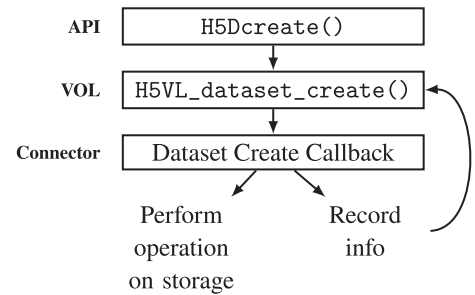


FIGURE 5 Dataset create call within the VOL. *Terminal* connectors interface with storage systems while *pass-through* connectors may record info on the fly and reenter the VOL



4 | CONNECTING PDC TO HDF5

We present in this section the HDF5 VOL connector interface to PDC and its use by applications through the HDF5 API.

4.1 | HDF5 VOL

HDF5 is a well-established I/O middleware package, used by a large number of HPC, scientific, and industrial applications. HDF5 provides them with file portability, reliability, and performance when storing their data. By default, the HDF5 library uses its *native* file format when storing data and makes use of MPI-IO to perform parallel I/O, as shown in Figure 4. While this has been a good choice for many years, it also carries on the burden of POSIX I/O semantics and limits that are inherent to the existing native file format, which defines an HDF5 file as a file structure that is contiguously mapped to a file system. For instance, the native file format has a well-known limitation of requiring collective creation of new HDF5 objects, such that the file metadata is ensured to be coherent between processes.

With the emergence of file and storage systems that do not strictly comply with the POSIX I/O standard, new file formats and ways of performing I/O can be defined with additional degrees of freedom. To provide that capability and give developers the ability to store the data in the form of their choice—while preserving the metadata that is attached to the HDF5 objects—the HDF5 library defines a *virtual object layer* (VOL), which will be released in the upcoming 1.12 version of the library. The VOL effectively allows developers to redefine the HDF5 I/O API calls (ie, related to operations on files, groups, datasets, attributes, etc) by seamlessly rerouting them to the corresponding VOL *connector* backend, which can in turn translate these calls into the operations that it desires to perform.

In the case of PDC and as shown in Figure 4, those operations translate into PDC calls, which in turn interact with the PDC runtime service and PDC storage backends. One of the main advantages of the PDC runtime is that it transparently and automatically provides new capabilities to HDF5 such as asynchronous I/O without requiring any major code change for the application user. One of the difficulties, however, is potentially for the VOL connector developer as the semantics that the underlying layer provides may not always be a direct match with the ones that the HDF5 VOL requires.

Figure 5 shows how the data are accessed in the file through a VOL connector callback. The VOL intercepts all HDF5 API calls that interact with files and reroutes those calls to the associated VOL callback for the requested API function. For example, a call to `H5Dcreate()` would be implemented within the HDF5 library as Figure 5, where the connector is responsible for the actual I/O operation to the storage system. Similarly, other operations in the HDF5 library follow the same execution pattern.

In order for an application to use an HDF5 VOL connector, one must first register the connector to HDF5 by calling the `H5VLregister_connector()` function. This effectively registers and initializes the VOL connector, which if successfully initialized will return a unique VOL connector ID. That ID can then be passed to the `H5Pset_vol()` routine, which notifies the file access property list of the VOL

connector it needs to use when creating or opening the file (since different VOL connectors could be initialized and used within the same application). For convenience, the library also defines environment variables that allow this information to be set by a user and avoids code modification to the application. Once the VOL connector information is set, the application can carry on with regular HDF5 function use.

4.2 | HDF5 PDC VOL connector implementation

Our PDC VOL connector currently only implements a subset of the HDF5 APIs and in this article we focus on file and dataset operations. HDF5 files can be easily mapped to PDC containers, while HDF5 datasets are naturally mapped to PDC objects and PDC regions are similar in essence to HDF5 selections. File create, open, and close operations are a direct match to PDC container operations, and therefore no particular implementation challenges were faced. Dataset operations, however, differ from PDC's API as the model chosen for PDC is to make data movement implicit, whereas HDF5 chooses to make data movement explicit by providing explicit read and write semantics. We therefore walk through the details of our implementation in this section.

As presented in the pseudocode below, dataset create and open operations map to `PDCobj_create()` and `PDCobj_open()`.

```
static void *
H5VL_pdc_dataset_create(void *obj, ..., const char *name, hid_t dcpl_id, hid_t dapl_id,
                        hid_t dxpl_id, ...)
{
    ...
    /* Create a new object */
    dset->obj.obj_id = PDCobj_create(o->file->cont_id, name, obj_prop);
    ...
}
```

```
static void *
H5VL_pdc_dataset_open(void *obj, ..., const char *name, hid_t dapl_id, hid_t dxpl_id,
                      ...)
{
    ...
    /* Open an existing object */
    dset->obj.obj_id = PDCobj_open(name, pdc_id);
    ...
}
```

Dataset write and read operations, however, require some extra handling as PDC does not provide explicit read and write semantics. Therefore, in these connectors callbacks, `PDCbuf_obj_map()` is used to first map the region in memory to the PDC object. A pair of `lock()` and `unlock()` calls are then also used, `unlock()` triggering asynchronous data movement. Finally, the region is unmapped using the `PDCbuf_obj_unmap()` call. Once `H5Dwrite()` returns after the `unlock` call, the data region has been transferred from the application buffer to the PDC data server, and the second step of data movement to further storage level can be taken care of by PDC, allowing further computation to be overlapped.

Write implementation is illustrated in the pseudocode below:

```
static herr_t
H5VL_pdc_dataset_write(void *_dset, hid_t mem_type_id, hid_t mem_space_id, hid_t
                      file_space_id, hid_t dxpl_id, const void *buf, ...)
{
    ...
    /* HDF5 selection to PDC region translation */
    ...
    PDCbuf_obj_map((void *)buf, mem_type, mem_reg, dset->obj.obj_id, obj_reg);
    ...
    PDCreg_lock(dset->obj.obj_id, obj_reg, WRITE, NOBLOCK);
    ...
    PDCreg_unlock(dset->obj.obj_id, obj_reg, WRITE);
    ...
    PDCbuf_obj_unmap((void *)buf, mem_reg, dset->obj.obj_id, obj_reg);
    ...
}
```

Reads are implemented similarly to writes but because PDC differentiates read locks from write locks (as read locks require less constraints than write locks), we pass the `READ` flag to both lock and unlock calls. This is illustrated in the pseudocode below:

```

static herr_t
H5VL_pdc_dataset_read(void *_dset, hid_t mem_type_id, hid_t mem_space_id, hid_t
    file_space_id, hid_t dxpl_id, void *buf, ...)
{
    ...
    /* HDF5 selection to PDC region translation */
    ...
    PDCbuf_obj_map((void *)buf, mem_type, mem_reg, dset->obj.obj_id, obj_reg);
    ...
    PDCreg_lock(dset->obj.obj_id, obj_reg, READ, NOBLOCK);
    ...
    PDCreg_unlock(dset->obj.obj_id, obj_reg, READ);
    ...
    PDCbuf_obj_unmap((void *)buf, mem_reg, dset->obj.obj_id, obj_reg);
    ...
}

```

The dataset is then closed using `PDCobj_close()` as illustrated below:

```

static herr_t
H5VL_pdc_dataset_close(void *_dset, hid_t dxpl_id, ...)
{
    ...
    /* Close the object */
    PDCobj_close(dset->obj.obj_id);
    ...
}

```

4.3 | Application usage example with the HDF5 PDC VOL

We provide a detailed example of how an application I/O kernel can be modified to support the HDF5 PDC VOL connector in Figure 6 and mark the two extra steps of function calls that are required in order to make use of the PDC VOL connector. As previously mentioned, HDF5 also provides a way of specifying that information through environment variables and allowing for no code changes. For an application to use the HDF5 PDC VOL connector for reading a dataset, it simply needs to follow the write example, using `H5Dopen()` instead of `H5Dcreate()`, and then call `H5Dread()` instead of `H5Dwrite()`. All of the function mapping details from HDF5 to the underlying PDC APIs are hidden by the HDF5 VOL design and are once again transparent to the user. The underlying VOL connector internally initiates map, lock, lock release, and unmap operations to that PDC object to trigger data movement and enables asynchronous data movement, transparently allowing for data movement to be overlapped by the following application computation step.

As mentioned at the end of Section 4.1, the library also defines environment variables, which allow for no code modification to the application. To enable this feature, the user only needs to set two environmental variables: `HDF5_VOL_CONNECTOR` and `HDF5_PLUGIN_PATH`, such that the corresponding VOL can be picked up during application execution. In the case of PDC, we set environment variables as follows:

```

export HDF5_VOL_CONNECTOR=pdc
export HDF5_PLUGIN_PATH=/path/to/vol_lib

```

5 | EXPERIMENTAL EVALUATION

We evaluate in this section the performance of the PDC VOL connector and also demonstrate the impact of the number of servers on the performance when PDC is deployed in dedicated mode. We compare the performance of writing multiple time steps using the PDC VOL connector with native HDF5. Finally, we compare the read performance of the PDC VOL connector with that of native HDF5.

5.1 | Experiment setup

To evaluate the performance of the PDC VOL, we ran the experiments with different configurations. We installed PDC on the Cori supercomputer at the National Energy Research Scientific Computing Center (NERSC), which is a Cray XC40 supercomputer with 1630 Intel Xeon Haswell nodes.


```

hid_t pdc_vol_id, file_id, fapl_id;
H5VL_pdc_info_t pdc_vol; // to pass VOL connector info
...
/* Register PDC VOL */
pdc_vol_id = H5VLregister_connector(&H5VL_pdc_g, H5P_DEFAULT); // extra step to use PDC

/* Create a new file access property */
fapl_id = H5Pcreate(H5P_FILE_ACCESS);

/* Set the VOL */
H5Pset_vol(fapl_id, pdc_vol_id, &pdc_vol); // extra step to use PDC VOL

/* Create file */
file_id = H5Fcreate(argv[1], H5F_ACC_TRUNC, H5P_DEFAULT, fapl_id);
...
/* Close property */
H5Pclose(fapl_id);

/* Create dataset */
dset_id1 = H5Dcreate(file_id, "x", H5T_NATIVE_FLOAT, filespace, H5P_DEFAULT, H5P_DEFAULT,
H5P_DEFAULT);

/* Write the data */
H5Dwrite(dset_id1, H5T_NATIVE_FLOAT, memspace, filespace, fapl_id, x);

/* Close dataset */
H5Dclose(dset_id1);

/* Close file */
H5Fclose(file_id);

```

FIGURE 6 Application usage example of the PDC VOL connector. There are three lines of code (including the creation of the file access property list) added to the original application code to use the PDC VOL.

Each node consists of 32 cores and 128 GB of memory. The supporting storage system, Lustre, has 248 object storage targets (OSTs) and is shared by all users.

5.1.1 | Deployment

We ran the experiments using both shared and dedicated deployment modes. With a shared server and client configuration (shared mode), we have one PDC server on each node, which utilizes one core, leaving the remaining 31 cores for user application execution. In dedicated mode, the PDC servers and user's application are on separate nodes. PDC servers in this configuration have only one server per node that provides both metadata server and data server services. In both configuration cases, we have relied on the Mercury³¹ RPC library, an HPC-optimized C library for Remote Procedure Calls, as the communication mechanism. In our experiments, we configure Mercury with two communication protocols using the libfabric plugin³² over TCP and Cray GNI. Note that in the latter case, the PDC server is configured to make use of Cray Dynamic RDMA Credentials (DRC)³³ to allow the user's applications and PDC server to share credentials and communicate together through GNI. GNI job runs are therefore currently a little more complex in terms of deployment. To use Cray GNI on Cori, the PDC server/service has to first acquire a credential and wait for the client application to start. The client application then first contacts the server over TCP so that the job can be granted access and use the DRC token. The generated DRC credential is later passed down to Mercury and used by both server and client sides for execution. Once the communication is established, the server and client can proceed and resume their normal execution. Using GNI currently requires the server and client to be in separate sessions but to start at the same time. This is achieved through the `srun pack-group` option. The job script for that run is attached in Appendix Figures A1 to A3.

5.1.2 | Applications and methodology

We used a plasma-physics application's I/O kernel, called VPIC-IO to evaluate the PDC system's performance. VPIC-IO is extracted from VPIC,³⁴ a code developed for simulating several plasma physics phenomena, including magnetic reconnection in space weather. In VPIC-IO, each MPI process writes a region of 8M (8×2^{20}) particles and each particle has eight properties. Each region is represented with a 1-D array with a size of 8M on one process. VPIC data structures use 1-D arrays for representing each property and each property is retreated as an object in our design. We

also evaluate the read performance by the BD-CATS I/O kernel,³⁵ which is extracted from a parallel clustering algorithm, used for analyzing the data produced by particle simulations. It reads data generated by VPIC or VPIC-IO using the same I/O trace as the BD-CATS implementation of the DBSCAN algorithm. In this kernel, data related to the particles are read among all of the MPI processes in a load-balanced distribution. The original kernels use HDF5 for performing I/O and are highly tuned using MPI-IO and Lustre optimizations.^{36,37} In this article, we simply reuse those I/O kernels to make them use the PDC VOL connector instead of going through native HDF5 and MPI I/O. The total data size being accounted for goes from 248 GB for 992 processes on 32 nodes, to 3968 GB for 15 872 processes on 512 nodes.

5.2 | HDF5 write performance comparison

We compare the `H5Dwrite()` performance of VPIC-IO in Figure 7 using the following methodologies: HDF5 collective I/O, HDF5 independent I/O, HDF5 Data Elevator VOL, PDC VOL in shared server mode, PDC VOL in dedicated server mode using TCP protocol and PDC VOL in dedicated server mode using Cray GNI. The `H5Dwrite()` function using the PDC VOL connector involves the times to map the memory buffer to a remote object and to lock and then release the lock on an object without waiting for data to be flushed to disk, while the native `H5Dwrite()` function is issuing MPI I/O calls directly to the file system. For all the evaluations presented in this section, we run the experiments at least 10 times and report best numbers. Since there are eight properties in VPIC-IO, the `H5Dwrite()` function is called eight times. The time is measured by adding an MPI barrier before the first call and another after the last call to the `H5Dwrite()` function. The total write time for all the eight properties is collected and used for evaluation in this section. The x-axis shows the number of client processes with the number of PDC servers (in brackets, in these plots as well as in the remaining plots in this section, unless specified otherwise). The native HDF5 I/O performance (collective and independent) was observed on Cori at our time of experiment, which could vary depending on the system software stack installed and system load.

The performance of the PDC VOL connector in shared server mode is 1.7× to 4.9× faster compared with independent native HDF5 method, with an average of 3.3×; is 2.9× to 4.2× faster compared with collective native HDF5 method, with an average of 3.5×; and is 1.6× to 2.7× faster compared with Data Elevator, with an average of 2.2×. In dedicated server mode, with additional nodes utilized as servers, the PDC VOL connector achieves 3.1× to 6.7× better performance compared with collective HDF5 I/O, with an average of 4.7×; achieves 3.8× to 4.8× better performance compared with independent native HDF5, with an average of 4.4×; and achieves 2.4× to 3.3× better performance compared with Data Elevator, with an average of 2.8×. To further improve the performance, we evaluate the performance with Cray GNI. It allows 4.6× to 15.6× speedup compared with collective native HDF5, with an average of 8.3×; allows 6.1× to 7.6× speedup compared with independent native HDF5, with an average of 6.6×; and allows 3.7× to 5.9× speedup compared with Data Elevator, with an average of 4.8×. With native HDF5, collective or independent, the data are written directly to the lustre file system. With Data Elevator, the data are first staged on SSD-based burst buffer, and then asynchronously transferred to the file system in the background. With all the other three PDC VOL asynchronous methods, the data are moved to PDC servers when `H5Dwrite()` returns, allowing for further data movement to the file system to be overlapped with following computation. Figure 7 shows the actual wait time in the `H5Dwrite()` call for the user and not the total time for data movement down to the file system.

5.3 | Varying servers in dedicated mode

In the previous experiments in dedicated mode, we had a configuration of one server per node and the number of server processes was the same as the number of client nodes. In this section, we evaluate the impact of the number of servers on the performance of `H5Dwrite()`. The experiments

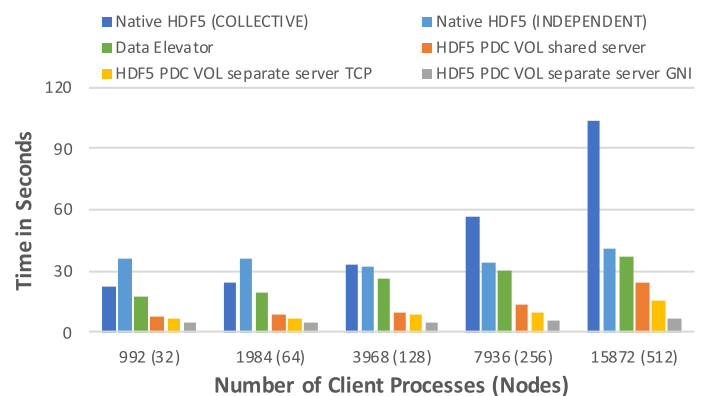


FIGURE 7 `H5Dwrite()` performance using different methodologies for one time step. Total data size goes from 248 GB for 992 processes to 3968 GB for 15 872 processes. The number of PDC servers for each configuration is equal to the number of compute nodes, indicated in parentheses

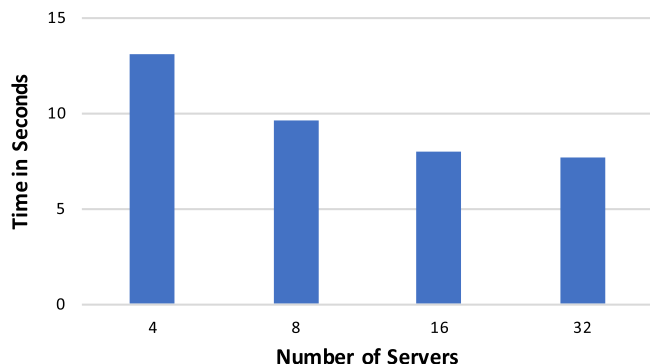


FIGURE 8 `H5Dwrite()` performance using dedicated server mode with varying the number of servers. Total data size is 248 GB for 992 processes

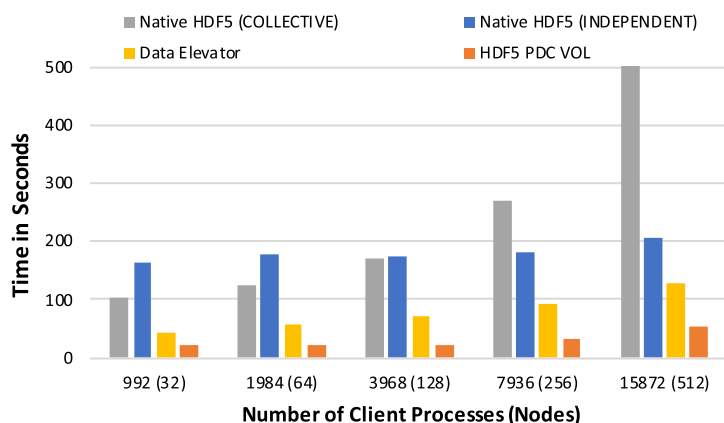


FIGURE 9 Total time including `H5Dwrite()` for five time steps and `H5Dclose()`

are run with 992 client processes on 32 nodes and involve varying numbers of additional nodes, from 4 to 32, for 32 servers in total. The total size of data to be written is 248 GB. Figure 8 shows one time step of `H5Dwrite()` time, increasing when fewer servers are available. More servers would naturally provide more bandwidth and achieve the best performance depending on the system resources availability, though fewer servers are still able to provide a reasonable performance.

5.4 | Multiple time steps

We mentioned in Section 5.2 that data movement is still happening after the `H5Dwrite()` call returns, unless the application user chooses to wait for the data to be flushed to disk. We experimented with five successive time steps of I/O using `H5Dwrite()` calls for each dataset to highlight that behavior and show the results in Figure 9. In this case, all data will be flushed between time steps. We use the PDC VOL connector in shared server mode for this experiment and observed 4.8 \times to 9.5 \times speedup compared with native HDF5 collective I/O, 3.8 \times to 8.3 \times speedup compared with native HDF5-independent I/O, and 2.1 \times to 3.0 \times speedup compared with HDF5 Data Elevator VOL, with an average speedup of 7.3 \times , 6.6 \times and, 2.6 \times , respectively.

5.5 | Total execution time for VPIC-IO

To reflect the total time consumed by the VPIC-IO application, we measured the time from the first HDF5 file create operation until file close. In this experiment we only covered one time step for each property within VPIC-IO. The more time steps the application executes, the more performance benefits it gains by utilizing the PDC VOL connector. Figure 10 shows the real total execution time of VPIC-IO covering just one time step. We can see that the PDC VOL in shared server mode is 1.5 \times to 3.1 \times faster compared with native HDF5 collective I/O, 1.5 \times to 2.2 \times faster than HDF5 independent I/O, and 1.4 \times to 2.7 \times faster compared with Data Elevator. For the best case, using a PDC separate server and GNI, it achieves 1.4 \times to 5.4 \times , 2.0 \times to 2.6 \times , and 1.0 \times to 1.9 \times performance speedup compared with HDF5 collective, independent I/O, and Data Elevator, respectively. This time reflects when the user chooses to wait for the data to be flushed to disk and then exits from the application. If the user chooses not to wait for the data to be flushed but lets the server handling the transfer, as shown in Figure 11, the performance becomes 2.1 \times to 3.8 \times and 1.8 \times to 3.2 \times better compared with HDF5 collective and independent I/O methods, respectively.

FIGURE 10 Total elapsed time for the execution of VPIC-IO

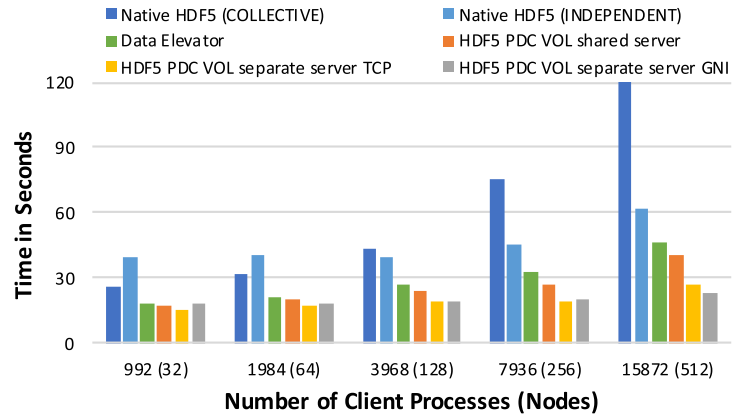


FIGURE 11 Total elapsed time for the execution of VPIC-IO if the user chooses not to wait for data to be flushed to disk

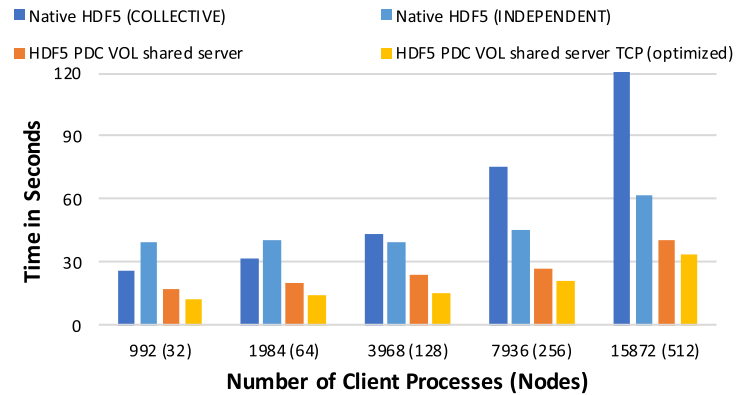
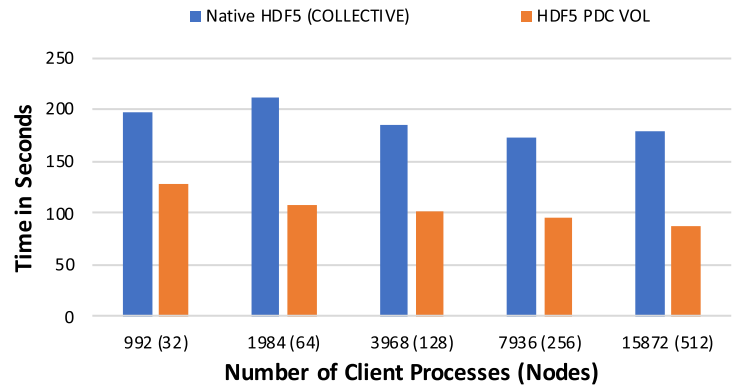


FIGURE 12 `H5Dwrite()` strong scaling performance using native HDF5 and HDF5 PDC VOL for H5BOSS benchmark writing 250 000 datasets from 992 processes to 15 872 processes. The number of PDC servers for each configuration is equal to the number of compute nodes, indicated in parentheses



5.6 | HDF5 write performance with H5BOSS

In this section, we evaluated the performance of the H5BOSS benchmark,³⁸ which writes 250 000 datasets out of 25 000 000 datasets in total. H5BOSS enables parallel processing and searching of millions of objects in Baryon Oscillation Spectroscopic Survey (BOSS) based on a HDF5-based python package. Compared with large file size of VPIC-IO, each of the H5BOSS I/O transaction size becomes relatively small, ranging from several hundreds of kilobytes to a few megabytes. As shown in Figure 12, we use the PDC VOL connector in shared server mode for this experiment and observed 1.5× to 2.0× faster (in terms of execution time) compared with native HDF5 collective I/O (native HDF5 collective I/O achieves a better performance than independent I/O in this case), with an average of 1.8×.

5.7 | HDF5 read performance comparison

BD-CATS-IO is a read I/O kernel, which reads data produced by VPIC-IO in a load balanced way. In Figure 13, we show the performance of reading a single time step of data that were written in previous VPIC experiments, calling instead the `H5Dread()` function. The PDC VOL connector in shared

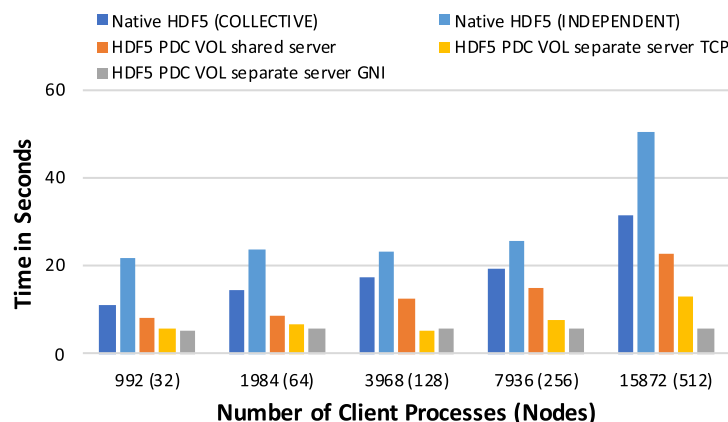


FIGURE 13 `H5Dread()` performance using different methodologies for one time step. Total data size goes from 248 GB for 992 processes to 3968 GB for 15 872 processes. The number of servers for each setup is the number in parentheses

server mode achieves 1.3× to 1.7× better performance compared with native HDF5 collective I/O and 1.4× to 2.8× better performance compared with native HDF5-independent I/O. The PDC VOL connector in dedicated server mode is 1.8× to 3.4× faster compared with native HDF5 collective I/O and is 2.7× to 4.6× faster compared with native HDF5-independent I/O. With Cray GNI, the PDC VOL in dedicated server mode executed 2.1× to 5.3× faster compared with native HDF5 collective I/O and 4.1× to 5.8× faster compared with HDF5-independent I/O. The read performance speedup is not as large as for the write due to the fact that `H5Dread()` requires the data to be fetched by the PDC data server from the backend file system and then transferred back to the application.

6 | CONCLUSION AND FUTURE WORK

In this article, we presented how to take advantage of PDC through HDF5 by developing an HDF5 VOL connector, which enables implicit and asynchronous data movement to different storage tiers with minimal to zero code modification, and is able to be deployed in different scenarios using native network fabric transports such as Cray GNI on modern supercomputers. We also evaluated and demonstrated that interface at scale, which showed a significant performance gain over native HDF5, as file system accesses are no longer issued directly by the application, but are instead handled by the PDC service.

When mapping HDF5 to PDC, one apparent limitation of HDF5 that transpired is its current inability to provide to the application a way of directly exposing the user's memory, as PDC is able to do through map operations. We will in future work study how that type of semantic could be brought into HDF5, allowing users to establish a direct mapping of the application memory to the storage, which similarly to `mmap()` operations can allow more efficient transfers and paging to take place and be handled by the PDC system, such that users do not have to explicitly direct of the amount data that need to be written at a given time. This naturally requires applications to adapt their code in order to reflect this type of change and would hence be more intrusive than the current solution presented in the article.

ACKNOWLEDGMENTS

This work is supported by the Director, Office of Science, Office of Advanced Scientific Computing Research, of the US Department of Energy under Contract No. DE-AC02-05CH11231 and DE-SC0016454. This research used resources of the National Energy Research Scientific Computing Center (NERSC).

ORCID

Jingqing Mu  <https://orcid.org/0000-0001-6758-2576>

Jerome Soumagne  <https://orcid.org/0000-0002-5480-1669>

REFERENCES

1. Tang H, Byna S, Tessier F, et al. Toward scalable and asynchronous object-centric data management for HPC. Paper presented at: Proceedings of the IEEE/ACM International Symposium on Cluster, Cloud and Internet Computing CCGRID; 2018.
2. Mu J, Soumagne J, Tang H, Byna S, Koziol Q, Warren R. A transparent server-managed object storage system for HPC. Paper presented at: Proceedings of the 2018 IEEE International Conference on Cluster Computing (CLUSTER); 2018:477-481.
3. Mehta K, Bent J, Torres A, Grider G, Gabriel E. A plugin for HDF5 using PLFS for improved I/O performance and semantic analysis. Paper presented at: Proceedings of the 2012 SC Companion: High Performance Computing, Networking Storage and Analysis; 2012:746-752.
4. Dong B, Byna S, Wu K, et al. Data elevator: low-contention data movement in hierarchical storage system. Paper presented at: Proceedings of the 2016 IEEE 23rd International Conference on High Performance Computing (HiPC); 2016:152-161.
5. Breitenfeld MS, Fortner N, Henderson J, et al. DAOS for extreme-scale systems in scientific applications. *CoRR*. 2017; abs/1712.00423.

6. Walli SR. The POSIX family of standards. *Standard View*. 1995;3(1):11-17.
7. Carns P, Ligon W III, Ross R, Thakur R. PVFS: a parallel virtual file system for linux clusters. *Linux J*. 2000;317-327.
8. Moore M, Ligon, B., Marshall, M et al. OrangeFS: advancing PVFS. *FAST poster session*; 2011.
9. Braam PJ. The Lustre storage architecture. *White Paper*. 2004. <https://arxiv.org/pdf/1903.01955.pdf>. Accessed March 5 2019.
10. Schmuck FB, Haskin RL. GPFS: a shared-disk file system for large computing clusters. *FAST*. 2002;2:231-244.
11. Microsystems S. *NFS: Network File System Protocol Specification RFC 1094*. Sun Microsystem; 1989. <https://doi.org/10.17487/RFC1094>.
12. Barker KJ, Davis K, Hoisie A et al. Entering the petaflop era: the architecture and performance of roadrunner. Paper presented at: Proceedings of the 2008 ACM/IEEE Conference on Supercomputing (SC'08); 2008:1-11. IEEE.
13. Jung M, Choi W, Shalf J, Kandemir MT. Triple-A: a Non-SSD based autonomic all-flash array for high performance storage systems. *SIGPLAN Not*. 2014;49:441-454.
14. Schürmann F, Delalandre F, Kumbhar PS. et.al. Rebasng I/O for scientific computing: leveraging storage class memory in an ibm bluegene/q supercomputer. Paper presented at: Proceedings of the International Supercomputing Conference; ISC; 2014: 331-347; Springer, Cham.
15. Folk M, Heber G, Koziol Q, Pourmal E, Robinson D. An overview of the HDF5 technology suite and its applications. Paper presented at: Proceedings of the EDBT/ICDT 2011 Workshop on Array Databases; 2011:36-47.
16. Li J, Liao KW, Choudhary A, et al. Parallel netCDF: a high-performance scientific I/O Interface. Paper presented at: Proceedings of the 2003 ACM/IEEE Conference on Supercomputing; 2003:39.
17. Liu Q, Logan J, Tian Y, et al. Hello ADIOS: the challenges and lessons of developing leadership class I/O frameworks. *Concurr Comput Pract Exp*. 2014;26(7):1453-1473.
18. Thakur R, Gropp W, Lusk E. On implementing MPI-IO portably and with high performance. Paper presented at: Proceedings of the 6th Workshop on I/O in Parallel and Distributed Systems ICPADS; 1999:23-32.
19. Tessier F, Vishwanath V, Jeannot E. TAPIOCA: an I/O library for optimized topology-aware data aggregation on large-scale supercomputers. Paper presented at: Proceedings of the 2017 IEEE International Conference on Cluster Computing (CLUSTER); 2017:70-80; IEEE.
20. Welton B, Kimpe D, Cope J, Patrick CM, Iskra K, Ross R. Improving I/O forwarding throughput with data compression. Paper presented at: Proceedings of the 2011 IEEE International Conference on Cluster Computing CLUSTER; 2011:438-445; IEEE.
21. Brown AL, Morrison R. A generic persistent object store. *Softw Eng J*. 1992;7(2):161-168. <https://doi.org/10.1049/sej.1992.0017>.
22. Moss JEB. Design of the mnome persistent object store. *ACM Trans Inf Syst*. 1990;8(2):103-139. <https://doi.org/10.1145/96105.96109>.
23. Cockshot WP, Atkinson MP, Chisholm KJ, Bailey PJ, Morrison R. Persistent object management system. *Softw Pract Exp*. 1984;14(1):49-71. <https://doi.org/10.1002/spe.4380140106>.
24. Gibson GA, Nagle DF, Amiri K, et al. A cost-effective, high-bandwidth storage architecture. *SIGPLAN Not*. 1998;33:92-103.
25. Weil SA, Leung AW, Brandt SA, Maltzahn C. RADOS: a scalable, reliable storage service for petabyte-scale storage clusters. Paper presented at: Proceedings of the 2nd International Workshop on Petascale Data Storage: Held in Conjunction with Supercomputing'07 PDSW; 2007:35-44.
26. Amazon. Amazon Web Services. <http://s3.amazonaws.com>. Accessed May 6 2019.
27. Arnold J. *OpenStack Swift: Using, Administering, and Developing for Swift Object Storage*. Sebastopol, CA: O'Reilly Media Inc; 2014.
28. Lofstead J, Jimenez I, Maltzahn C, Koziol Q, Bent J, Barton E. DAOS and friends: a proposal for an exascale storage system. *Supercomputing*. 2016;50(12):1-50.
29. Shi X, Li M, Liu W, Jin H, Yu C, Chen Y. SSDUP: a traffic-aware ssd burst buffer for HPC systems. Paper presented at: Proceedings of the International Conference on Supercomputing ICS; 2017.
30. Wang T, Byna S, Dong B, Tang H. UniviStor: integrated hierarchical and distributed storage for HPC. Paper presented at: Proceedings of the 2018 IEEE International Conference on Cluster Computing (CLUSTER); 2018.
31. Soumagne J, Kimpe D, Zounmevo J, et al. Mercury: enabling remote procedure call for high-performance computing. Paper presented at: Proceedings of the 2013 IEEE International Conference on Cluster Computing (CLUSTER); 2013:1-8.
32. OFIWG. Libfabric. <https://ofiwg.github.io/libfabric/>. Accessed April 20 2019.
33. Shimek J, Swaro J. Dynamic RDMA Credentials. Paper presented at: Proceedings of the Cray User Group (CUG) Meeting; 2016.
34. Bowers KJ, Albright BJ, Yin L, Bergen B, TJT K. Ultrahigh performance three-dimensional electromagnetic relativistic kinetic plasma simulation. *Phys Plasmas*. 2008;15(5):055703. <https://doi.org/10.1063/1.2840133>.
35. Patwary MMA, Byna S, Satish, N. R. BD-CATS: big data clustering at trillion particle scale. Paper presented at: Proceedings of the SC'15 International Conference for High Performance Computing, Networking, Storage and Analysis; 2015.
36. Byna S, Chou J, Rübel O, et al. Parallel I/O, analysis, and visualization of a trillion particle simulation. Paper presented at: Proceedings of the SC'12 International Conference on High Performance Computing, Networking, Storage and Analysis; 2012:59:1-59.
37. Behzad B, Byna S, Wild SM, Prabhat M, Snir M. Improving parallel I/O autotuning with performance modeling. Paper presented at: Proceedings of the 23rd International Symposium on High-Performance Parallel and Distributed Computing HPDC; 2014:253-256.
38. Liu J, Bard D, Koziol Q, Bailey S. Searching for millions of objects in the BOSS spectroscopic survey data with H5Boss. 2017 *New York Scientific Data Summit (NYSDS)*. New York, NY: IEEE; 2017:1-9.

How to cite this article: Mu J, Soumagne J, Byna S, Koziol Q, Tang H, Warren R. Interfacing HDF5 with a scalable object-centric storage system on hierarchical storage. *Concurrency Computat Pract Exper*. 2020;32:e5715. <https://doi.org/10.1002/cpe.5715>

APPENDIX

In Figures A1, A2, and A3, we provide sample Slurm job scripts for running the PDC VOL using Cray GNI on the Cori supercomputer at NERSC. The first script (Figure A1) is used to launch a test application with the PDC VOL using Cray GNI. It has four `srun` commands: Dynamic RDMA credentials (DRC) server; DRC client, PDC server, and the HDF5 application. In Figure A2, we show the script (`drc_server.sh`) to obtain DRC server tokens and in Figure A3, we show the script for running the client and request access to the server using TCP

```
srun --pack-group=0 -n 1 drc_server.sh &
sleep 5
srun --pack-group=1 -n 1 drc_client.sh
export PDC_DRC_KEY='cat $SCRATCH/drc.txt'
....
srun --pack-group=0 -n 32 pdc_server.exe &
sleep 5
srun --pack-group=1 -n 992 h5_pdc_write cc
```

FIGURE A1 A sample Slurm job script used to launch a test application with PDC VOL connector using Cray GNI

```
procIdx=${SLURM_PROCID}
if [ $procIdx -eq 0 ]; then
/mercury/build/bin/hg_test_drc_auth -c ofi -p tcp -H ipogif0 -L -a
fi
```

FIGURE A2 A sample job server script (`drc_server.sh`), used to obtain dynamic RDMA credentials

```
procIdx=${SLURM_PROCID}
if [ $procIdx -eq 0 ]; then
/mercury/build/bin/hg_test_drc_auth -c ofi -p tcp -H ipogif0 -a
fi
```

FIGURE A3 A sample job client script (`drc_client.sh`), which is used to request access to the server