

## Available at www.ComputerScienceWeb.com

Journal of Parallel and Distributed Computing

J. Parallel Distrib. Comput. 63 (2003) 434-447

http://www.elsevier.com/locate/jpdc

# High-performance scientific data management system $\stackrel{\mathcale}{\sim}$

Jaechun No,<sup>a</sup> Rajeev Thakur,<sup>b,\*</sup> and Alok Choudhary<sup>c</sup>

<sup>a</sup> Department of Software Engineering, Sejong University, Seoul, Republic of Korea

<sup>b</sup> Mathematics and Computer Science Division, Argonne National Laboratory, 9700 South Cass Avenue, Argonne (Lemont), IL 60439, USA <sup>c</sup> Department of Electrical and Computer Engineering, Northwestern University, Evanston, IL 60208, USA

Received 25 August 2001; revised 13 January 2003; accepted 28 February 2003

### Abstract

Many scientific applications have large I/O requirements, in terms of both the size of data and the number of files or data sets. Management, storage, efficient access, and analysis of this data present an extremely challenging task. Traditionally, two different solutions have been used for this task: file I/O or databases. File I/O can provide high performance but is tedious to use with large numbers of files and large and complex data sets. Databases can be convenient, flexible, and powerful but do not perform and scale well for parallel supercomputing applications. We have developed a software system, called Scientific Data Manager (SDM), that combines the good features of both file I/O and databases. SDM provides a high-level application programming interface to the user and, internally, uses a parallel file system to store real data (using various I/O optimizations available in MPI-IO) and a database to store application-related metadata. In order to support I/O in irregular applications, SDM makes extensive use of MPI-IO's noncontiguous collective I/O functions. Moreover, SDM uses the concept of a *history file* to optimize the cost of the index distribution using the metadata stored in database. We describe the design and implementation of SDM and present performance results with two regular applications, ASTRO3D and an Euler solver, and with two irregular applications, a CFD code called FUN3D and a Rayleigh–Taylor instability code.

© 2003 Elsevier Science (USA). All rights reserved.

Keywords: Scientific data management; Parallel I/O; MPI-IO; Database; Metadata

### 1. Introduction

Many large-scale scientific experiments and simulations generate very large amounts of data [2,9] (on the order of several hundred giga- to terabytes), spanning thousands of files or data sets. Moreover, such applications have different patterns of accessing data from files. Some applications access data in regular block or cyclic access patterns. Others have more complex irregular access patterns that cannot be detected at compile time. These different ways of accessing data in applications, along with the large number of files and large size of data being generated, make high-performance access,

\*Corresponding author.

management, storage, and analysis of the data an extremely challenging task.

Current techniques for data management are either raw file–I/O interfaces, such as MPI-IO [12,24], or fullfledged databases. File–I/O interfaces provide high performance but are too cumbersome to use with large, complex data sets and large numbers of files. For example, the user must remember file names and the organization of data in a file and must specify the exact location in the file from which the data must be accessed. Databases, on the other hand, provide a convenient, high-level interface and powerful dataretrieval capability, but they do not measure up to the performance requirements of large-scale scientific applications running on supercomputers.

We have developed a software system, called Scientific Data Manager (SDM), that combines the good features of both file I/O and databases [28,29]. SDM provides a high-level, user-friendly interface. Internally, SDM interacts with a database to store application-related metadata and uses MPI-IO to store the real data on a

<sup>&</sup>lt;sup>\*</sup>This work was supported in part by the Mathematical, Information, and Computational Sciences Division subprogram of the Office of Advanced Scientific Computing Research, US Department of Energy, under Contract W-31-109-Eng-38, and in part by a Work-for-Others Subaward No. 751 with the University of Illinois, under NSF Cooperative Agreement #ACI-9619019.

*E-mail address:* thakur@mcs.anl.gov (R. Thakur).

high-performance parallel file system. SDM takes advantage of various I/O optimizations available in MPI-IO, such as collective I/O and noncontiguous requests, in a manner that is transparent to the user.

SDM provides efficient I/O support for irregular applications, while maintaining a high-level unified application programming interface (API) for both regular and irregular applications. In irregular applications, the data accesses make extensive use of arrays, called indirection arrays [8] or map arrays [12], in which each value of the array denotes the corresponding data position in memory or in the file. The data distribution in irregular applications can be done either by using compiler directives with the support of runtime preprocessing [14,15] or by using a runtime library [8]. Most of the previous work in the area of unstructured-grid applications focuses mainly on computation and communication in such applications, not on I/O. SDM takes advantage of MPI-IO's support for noncontiguous data accesses and, therefore, can efficiently handle the reading and writing of data in an irregular mesh, as well as the distribution of index values.

The rest of this paper is organized as follows. In Section 2, we discuss our goals in developing SDM. In Section 3, we present the design and implementation of SDM. Performance results on the SGI Origin2000 and on the IBM SP at Argonne National Laboratory are presented in Section 4. We discuss related work in Section 5 and conclude in Section 6.

### 2. Design objectives

Our main objectives in developing SDM were to provide high-performance parallel I/O, to provide a high-level API, to support a convenient data-retrieval capability, and to optimize the execution time of both regular and irregular applications.

• *High-performance I/O*. To achieve high-performance I/O, we decided to use a parallel file-I/O system to store real data and use MPI-IO to access this data. MPI-IO, the I/O interface defined as part of the MPI-2 standard [12,24], is rapidly emerging as the standard, portable API for I/O in parallel applications. High-performance implementations of MPI-IO, both vendor and public-domain implementations, are available for most platforms [10,19,32,33,42]. MPI-IO is specifically designed to enable the optimizations that are critical for high-performance parallel I/O. Examples of these optimizations include collective I/O, the ability to access noncontiguous data sets, and the ability to pass hints to the implementation about access patterns, file-striping parameters, and so forth.

- *High-level API*. Our goal was to provide a high-level unified API for any kind of application (regular or irregular) while encapsulating the details of either MPI-IO or a database. The user can specify the data with a high-level description, together with annotations, and use a similar API for data retrieval. SDM internally translates the user's request into appropriate MPI-IO calls, including creating MPI-derived datatypes for noncontiguous data [41]. SDM also interacts with the database when necessary, by using embedded SQL functions.
- Convenient data-retrieval capability. SDM allows the user to specify names and other attributes to be associated with a data set. SDM internally selects a file name into which the data will be stored; the mapping between data sets and file names is stored in the database. The user can retrieve a data set by specifying a unique set of attributes for the desired data. SDM also allows the user to query the metadata for a stored data set and then select data with specific attributes.
- Optimization of irregular applications. In irregular applications, the cost of an index distribution is usually expensive, in terms of communication and computation volumes. In SDM, after the index values are partitioned among processes, the local index subsets of all processes are asynchronously written to a *history file*, and the associated metadata is stored in the database. When the same index distribution is needed in subsequent runs, the index values are read from the history file by using the metadata stored in the database, and thereby the user can avoid repeating the communication and computation for the same index distribution. SDM also uses MPI datatypes and collective I/O functions to optimize I/O for irregular access patterns.

### 3. Implementation

With the help of sample regular and irregular problems, we describe the design and implementation of SDM.

### 3.1. Regular applications

We present a sample regular problem and describe the metadata storage in the database, the SDM API, and the organization of data in files.

### 3.1.1. Problem description and SDM API

ASTRO3D is a three-dimensional astrophysics application developed at the University of Chicago. For simplicity of explanation, we consider the two-dimensional version of this three-dimensional application. (The performance results presented in this paper are for

```
SDM_initialize(ASTRO3D):
                                                             SDM_initialize(ASTRO3D):
A = SDM_make_datalist(3, \{a_0, a_1, a_2\});
                                                             A = SDM_make_datalist(3, \{a_0, a_1, a_2\});
A[0].data_type = FLOAT;
                                                             initialize(&date);
A[0].access_pattern[0] = BLOCK;
                                                             date.year = 2000;
A[0].access_pattern[1] = BLOCK;
                                                             date.month = 10;
                                                             date.day = 10;
SDM_associate_attributes(3, &A[0]);
handle = SDM_set_attributes(A, 3);
                                                             handle = SDM_select_attributes(A, 3);
SDM_subarray(handle,3,0,starts,subsizes,NULL);
                                                             SDM_subarray(handle,3,0,starts,subsizes,NULL);
SDM_data_view(handle,3,0,NULL,NULL);
                                                             SDM_data_view(handle,3,0,NULL,NULL);
For (i=1; i < maxStep; i++) {
                                                             For (i=1; i < maxStep; i++) {
     . . . . . .
     Do computation;
                                                                   SDM_read(handle, a_0, i, a_0buf);
                                                                   SDM_read(handle, a_1, i, a_1buf);
     SDM write(handle, a_0, i, a_0buf);
                                                                   SDM read(handle, a_2, i, a_2buf);
     SDM_write(handle, a_1, i, a_1buf);
     SDM_write(handle, a_2, i, a_2buf);
                                                                   Do computation;
                                                                   .....
                                                             }
}
SDM_finalize(handle, 3);
                                                             SDM_finalize(handle, 3);
                      (a)
                                                                                           (b)
```

Fig. 1. SDM API to perform (a) write and (b) read operations for data group A in ASTRO3D.

the full three-dimensional version.) In this application, data are stored in several arrays that are block distributed in each dimension. At various time steps, several of these arrays are written to files for data analysis, restart, and visualization. Six floating-point arrays are written for data analysis and another six for restart; seven character arrays are written for visualization. The frequencies of the writes can be varied.

We use the term *data set* to refer to each array being written and *data group* to refer to all the arrays written at a time step for a particular purpose such as data analysis. For simplicity of explanation, let us assume that three arrays are written for data analysis, another three for restart, and four for visualization. (Note that all arrays—six, six, and seven—were used in the performance experiments reported in this paper.) Let us further assume that the data analysis and restart dumps are performed every six time steps and the visualization dumps are performed every four time steps. Let  $a_0, a_1, a_2$  be the three data sets for data analysis and  $A = (a_0, a_1, a_2)$  be the data group for data analysis. Similarly, we have  $B = (b_0, b_1, b_2, b_3)$  for visualization and  $C = (c_0, c_1, c_2)$  for restart.

Fig. 1 shows how the SDM API is used to perform (a) write and (b) read operations for data analysis (data group A) in a two-dimensional version of ASTRO3D.

### 3.1.2. Implementation details

SDM provides a high-level API and stores application-related metadata in a database. For regular applications where the data access pattern can be predicted before runtime, SDM creates three database tables: run\_table, access\_pattern\_table, and execution\_ table (see Fig. 2). These tables are made for each application. Each time an application writes data sets, SDM enters the problem size, dimension, current date, and a unique identification number (runid) to the run\_table. The access\_pattern\_table includes the properties of each data set, such as datatype, storage order, data access pattern, and global size. SDM uses this information to make appropriate MPI-IO calls to access the real data. The execution\_table stores a globally determined file offset denoting the starting offset in the file of each data set.

In SDM, users can specify groups of data sets by assigning properties to the first data set in a group and by propagating them to the other data sets belonging to the same group. The main reason for making groups of data sets is that SDM can then use different ways of organizing data in files, with different performance implications. For example, each data set can be written in a separate file, or the data sets of a group can be written to a single file. For ASTRO3D, we created the three data groups,  $A = a_0, a_1, a_2, B = b_0, b_1, b_2, b_3$ , and  $C = c_0, c_1, c_2$ , for data analysis, visualization, and restart, respectively. The properties assigned to each data set are stored in the database by calling SDM\_set\_attributes.

In case of read operations, data from a specific run can be retrieved by specifying attributes of the data, such as the date of the run. Also, the properties of the

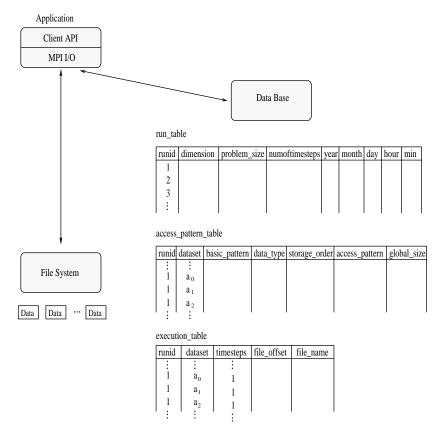


Fig. 2. Database tables used for solving regular applications in SDM.

data sets need not be specified because SDM retrieves this information from the database.

The main SDM functions for writing and reading data are SDM\_write and SDM\_read. Before calling these functions, the user must provide the information necessary for SDM to perform I/O, such as the starting points and sizes of the subarray in each dimension in the case of block distribution, or the size of process grids and distribution arguments in each dimension in the case of cyclic distribution. Also, SDM\_data\_view must be called to perform the data mapping between memory and file.

In order to perform I/O, the handle of a group, position of a data set within the handle (group), current time step, and pointer to the application buffer are passed to the SDM I/O function. Note that the user does not have to provide file names. SDM generates the file name and records the name in the database. SDM calls MPI-IO's collective I/O functions to perform I/O efficiently and in parallel from all processes.

### 3.1.3. File organization

SDM supports three different ways of organizing data in files. In level 1, each data set generated at each time step is written to a separate file, as shown in Fig. 3. This file organization is simple, but it incurs the cost of a file open and close at each time step, which on some file systems can be quite high, as we shall see in the performance results. For a large number of data sets and time steps, this method can be expensive because of the large number of file opens.

In level 2, each data set (within a group) is written to a separate file, but different iterations of the same data set are appended to the same file, as illustrated in Fig. 4. This method results in a smaller number of files and smaller file-open costs. The offset in the file where data is appended is stored in the execution\_table.

In level 3, all iterations of all data sets belonging to a group are stored in a single file, as shown in Fig. 5. As in level 2, the file offset for each data set is stored in the execution\_table by process 0 in the SDM\_write function. If a file system has high open and close costs, SDM can generate a very small number of files by choosing the level-3 file organization. On the other hand, if an application produces a large number of data sets of large size, level 3 would result in very large files, which may affect performance.

We study the performance implications of the three file-organization levels in the performance section.

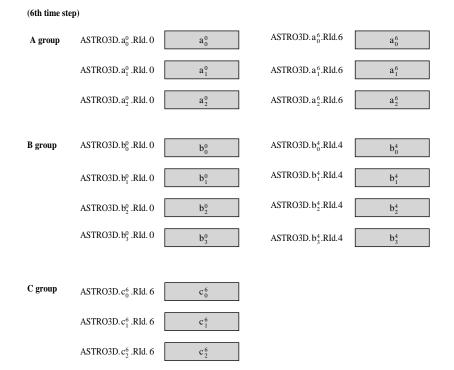


Fig. 3. Level-1 file organization at 6th time step in ASTRO3D. The superscript on a data set denotes the time step at which the data set has been written to a file, RId denotes the current identification number (runid), and each shadowed box (along with the name beside it) shows the SDM-generated file for storing the corresponding data set.

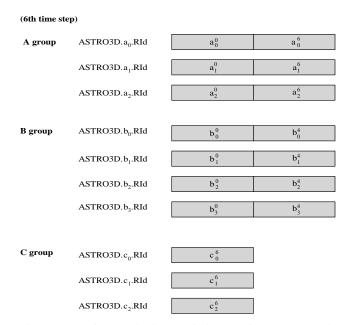


Fig. 4. Level-2 file organization at 6th time step in ASTRO3D. The superscript on a data set denotes the time step at which the data set has been written to a file, RId denotes the current identification number (runid), and each shadowed box (along with the name beside it) shows the SDM-generated file for storing the corresponding data set.

### 3.2. Irregular applications

We now describe the SDM architecture for supporting irregular applications.

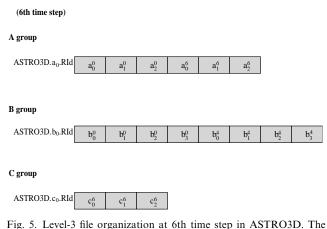


Fig. 5. Level-3 file organization at 6th time step in ASTRO3D. The superscript on a data set denotes the time step at which the data set has been written to a file, RId denotes the current identification number (runid), and each shadowed box (along with the name beside it) shows the SDM-generated file for storing the corresponding data set.

### 3.2.1. Problem description and SDM API

Fig. 6 shows a typical irregular application that sweeps over the edges of an irregular mesh. In this problem, edge1 and edge2 are two arrays representing nodes connected by an edge, and arrays x and y are the actual data associated with each edge and node, respectively. The partitioned arrays of edge1, edge2, x, and y contain a single level of "ghost data" beyond the boundaries to minimize remote accesses. After the

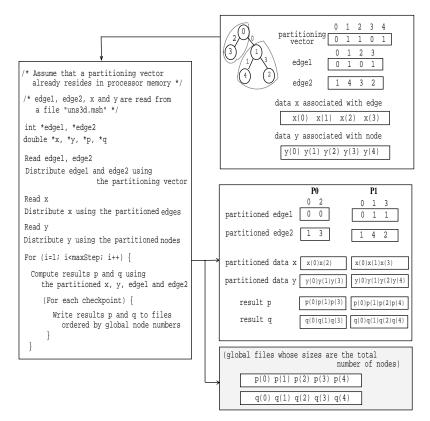


Fig. 6. A sample irregular problem.

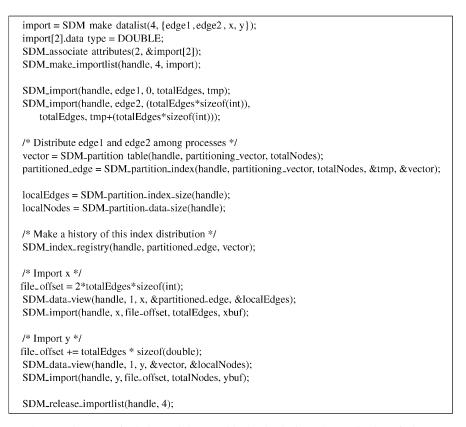


Fig. 7. Using SDM for index and data partitioning in the irregular application of Fig. 6.

```
SDM_initialize(ApplicationName);
result = SDM_make_datalist(2, {p, q});
result[0].data_type = DOUBLE;
SDM_associate_attributes(2, &result[0]);
handle = SDM_set_attributes(result, 2);
......
/* Partition edge1, edge2, x and y among processes (Figure 7) */
.....
SDM_data_view(handle, 2, 0, &vector, &localNodes);
For (i=1; i < maxStep; i++) {
.....
Do computation and produce results p and q;
......
For (each checkpoint) {
SDM_write(handle, p, i, pbuf);
SDM_write(handle, q, i, qbuf);
}
}
SDM_finalize(handle, 2);
```

Fig. 8. Using SDM for writing results in the irregular application.

computation is completed, the results p and q are written to a file in the order of global node numbers.

Fig. 7 shows how the data and index partitioning can be specified in SDM. We use the term *import* to distinguish it from a *read* operation. A read operation reads the data created in SDM, whereas an import operation reads the data created outside of SDM. Fig. 8 shows the use of SDM for writing the results of the computation.

### 3.2.2. Implementation details

Besides providing a high-level API for convenient data retrieval and I/O, SDM provides a capability to partition the index and data arrays being used in irregular applications. To use the SDM partitioning scheme, users must provide a *partitioning vector* in which each value of the vector denotes a processor rank where the node should be assigned. The partitioning vector should be replicated among processes. The current implementation of SDM requires that the partition vector is not currently supported. For the application illustrated in Fig. 6, we used the partitioning vector generated from the graph-partitioning tool Metis [20,34].

The SDM partitioning phase generates *map arrays* that specify the mapping of each element of the local array to the global array. The map arrays are used in I/O operations. This partitioning phase is optional; users can instead provide their own map array, while bypassing this phase.

For irregular applications, in order to partition index and data arrays and to perform I/O, SDM creates six database tables: run\_table, access\_pattern\_table, execution\_table, import\_table, index\_table, and index\_history\_table. The first three tables are the same as described in Section 3.1.2 for regular applications. In order to optimize the cost of index distributions, SDM provides a history file to store the local index subsets of all processes. The associated metadata is stored in the database. When the same index distribution is needed later, each process can retrieve the partitioned index values from the history file by using the metadata stored in the database. This saves the cost of creating the partitions each time.

The disadvantage of the history file is that it cannot be used if the problem runs on a different number of processes from that when the file was created, because the edges and nodes being assigned to each process dynamically change among different numbers of processes. One efficient use of the history file is to create it in advance for the various numbers of processes of interest. As long as the user runs the application with any of those numbers of processes, an appropriate history can be chosen to reduce communication and computation costs.

Fig. 7 describes the steps in SDM to partition the indexes and data. The four arrays, edge1, edge2, x, and y, are imported by creating a data group. The function SDM\_import first accesses the index\_table in the database to see whether a history file exists with this problem size. If so, the metadata, such as each process's partitioned index size and the name of the history file, is retrieved from the index\_table and index\_history\_table. Otherwise, the desired data are imported to the application. Since edge1 and edge2 are being imported in a contiguous way, there is no need to specify data mapping between the file and processor memory.

In SDM\_partition\_table, the global partitioning vector (partitioning\_vector in Fig. 7) is converted to the local vector (vector in Fig. 7) to determine which node should be assigned to which process.

If there is a history file for this problem size, SDM\_partition\_index reads the already-partitioned edge1 and edge2 from the history file and converts them to localized edges by using the partitioning vector. This approach avoids the communication cost to exchange each process's edges and the computation cost to choose the edges to be assigned. If there is no history file, the edges in each process are distributed by reading all the data in parallel and by performing a ring-oriented communication.

For storing the partitioned edges and nodes, including the ghost ones, a certain amount of memory space is initially allocated to each process. When the entire memory space is occupied by the partitioned data, it is automatically doubled by adjusting the memory size. This prevents the system from looking through the entire data in two steps, one step to decide the size of memory space and the other step to actually store the data in the memory space. After the edges and nodes are distributed, the edges in each process are moved to the next process in the virtual ring. As a result, two map

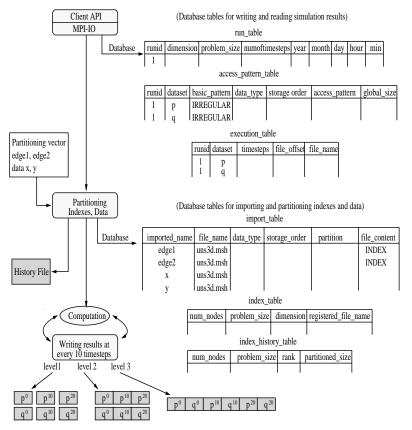


Fig. 9. SDM execution flow to solve the example in Fig. 6.

arrays, partitioned\_edge and vector, are generated that are used to distribute the physical data associated with each edge and node, respectively.

SDM\_index\_registry creates a history file containing the index distribution if no history file had been created earlier for the index distribution. The information to partition index arrays, such as the partition size for each process, is stored in the database tables index\_table and index\_history\_table. Also, the partitioned index values (partitioned\_edge in our example) are written asynchronously to a history file to be retrieved in subsequent runs requiring the same edge distribution.

Fig. 8 shows the steps to write two data sets, p and q, after completing the computations at each checkpoint. Before writing p and q, the data mapping is defined in the SDM\_data\_view by using the map array (vector) associated with the node partition. The data sets are written to files according to one of the file organizations described in Section 3.1.3. All the I/O is performed by using MPI-IO's collective I/O functions and derived datatypes to describe noncontiguous access patterns.

Fig. 9 depicts the metadata storage in the database and the organization of data in files in SDM for the example in Fig. 6.

### 4. Performance results

We used the IBM SP and SGI Origin2000 at Argonne National Laboratory for all our performance tests. At the time the experiments were run, these machines were configured as follows. The IBM SP had 80 compute nodes and 4 I/O nodes. Each I/O node controlled four SSA disks, each of 9 Gbyte capacity. The parallel file system on the machine was IBM's PIOFS [18]. The SGI Origin2000 had 128 processors and 10 Fibre Channel controllers connected to a total of 110 disks, each of 9 Gbyte capacity. The file system on the Origin2000 was SGI's XFS [16,38].

### 4.1. Performance results for regular applications

For performance evaluation of regular applications, we used an optimization called direct I/O on XFS. When certain alignment restrictions are met, the user can choose the direct-I/O option, in which the file system moves data directly between the user's buffer and the storage device, bypassing the file-system cache. Direct I/O thus eliminates an extra memory copy into the cache and can perform well if the I/O size is large and the machine has a high-bandwidth I/O system. Direct I/O can be used from an MPI-IO program: the ROMIO implementation of MPI-IO that we used supports direct I/O [43]. We present performance results with both direct I/O and regular (buffered) I/O.

We used two regular application templates, AS-TRO3D and a three-dimensional Euler solver, in our performance experiments. For ASTRO3D, we used a problem size of  $256 \times 256 \times 256$ . We ran the program for one time step and performed the data analysis, restart, and visualization dumps at that time step. This resulted in a total of around 880 Mbytes of data.

The second application is a three-dimensional Euler solver for the problem of three-dimensional transonic flow about an M6 wing [13]. This application is a mesh-structured code that writes the physical values and residual of each node at certain iterations. The structure of these values is a distributed global vector, and each value has five components (density, energy, and three coordinates of momentum). In addition, the application writes the physical coordinates and pressure at each mesh point. In our experiments, we ran the code for 50 iterations and wrote data at every 5 iterations. The problem size was  $194 \times 34 \times 34$ .

### 4.1.1. Cost of database access

SDM uses TCP/IP to connect to the database servers. We performed our experiments with two different databases, MySQL [25] and PostgreSQL [31]. Fig. 10 shows the database-access cost in the SDM write operation on the Origin2000. The connection to and disconnection from the database server occur once in SDM\_initialize and SDM\_finalize, respectively. In SDM\_set\_attributes, process 0 accesses the run\_table and access\_pattern\_table to store attributes. In the write operation, process 0 stores the file offset in the execution\_table. Access to the execution\_table occurred 19 times in ASTRO3D versus 60 times in the Euler solver. As can be seen in Fig. 10, the database access cost using either of the database servers was less than 0.6 s. This cost, however, will change according to the number of I/O operations occurring in the applications.

We observed that MySQL performs better than PostgreSQL. Therefore, we used only MySQL for the rest of the performance experiments.

### 4.1.2. Results for ASTRO3D

Fig. 11 shows the write and read bandwidths for ASTRO3D on the IBM SP using 32 processors for the three levels of file organization. Since we ran only one iteration of the program, levels 1 and 2 resulted in the same file organization. Level 3 achieved much higher bandwidth because only three different files were created, and, therefore, only three file opens occurred. The high cost of file opens on the PIOFS file system [40] resulted in lower performance for levels 1 and 2, where 19 separate files were created. The impact of file-open time can indeed be quite large.

Figs. 12 and 13 show the write and read bandwidths for ASTRO3D on the SGI using 16 processors. We measured performance for both direct I/O and buffered I/O. For writing data, direct I/O performed better than buffered I/O. There are two reasons for this. First, with buffered I/O, XFS serializes concurrent writes to the same file, whereas with direct I/O, concurrent writes are allowed to proceed in parallel. Second, direct I/O eliminates a copy into the file-system cache. For reading data, buffered I/O performed better. Again, there are two reasons for this. One reason is that XFS does not serialize buffered reads; therefore, direct reads do not have any extra advantage in the area of parallelism. The second reason is that XFS performs a read-ahead (prefetch) in the case of buffered reads, but not in case

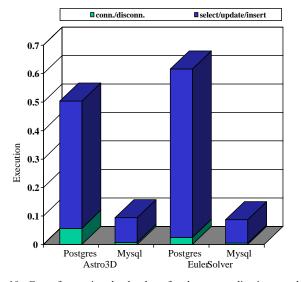


Fig. 10. Cost of accessing the database for the two applications on the SGI Origin2000.

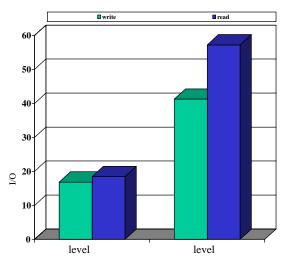


Fig. 11. I/O bandwidth for ASTRO3D on the IBM SP.

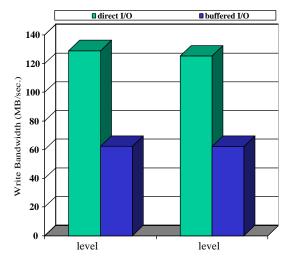


Fig. 12. Write bandwidth for ASTRO3D on the SGI Origin2000.

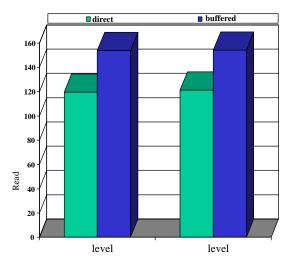


Fig. 13. Read bandwidth for ASTRO3D on the SGI Origin2000.

of direct reads. The read-ahead policy works well for this application, and buffered reads therefore perform better. Since the cost of file opens is small on XFS, the three levels of file organization performed nearly the same.

### 4.1.3. Results for the Euler solver

Fig. 14 shows the write and read bandwidths for the Euler solver on the IBM SP using 32 processors. The total data written was around 240 Mbytes. In level 3, only two files were generated, one for writing the coordinates and pressure at each mesh node and the other for writing the physical values and residual at each node. In level 2, six vectors (that is, the three coordinates, pressure, physical values of each node, and nodal residual) were written separately, resulting in a total of six files. In level 1, the six vectors generated every five iterations were written separately, resulting in

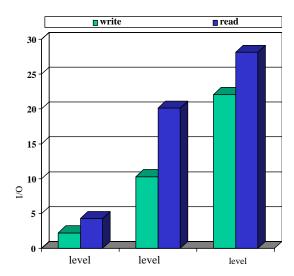


Fig. 14. I/O bandwidth for the Euler solver on the IBM SP.

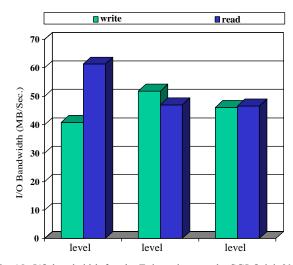


Fig. 15. I/O bandwidth for the Euler solver on the SGI Origin2000.

a total of 60 files. As Fig. 14 shows, level 3 performed the best because of the high open cost on PIOFS. In level 1, the file-open cost took around 80% of the total execution time; in level 2, it took around 30%; and in level 3, it took around 20% of the total execution time.

Fig. 15 shows the write and read bandwidths for the Euler solver using 16 processors on the SGI. For this application, we used only buffered I/O. We could not use direct I/O because the memory allocation for distributed vectors was done inside the numerical library (PETSc [30]) that the application uses, and thus we could not align the buffers to the cache line as required for direct I/O. For the write operation, levels 2 and 3 performed slightly better than level 1. For the read operation, however, level 1 performed the best. The reason is that the read-ahead policy of XFS for buffered reads operates on a per-file basis and therefore works to

the application's advantage when it has a greater number of files.

### 4.2. Performance results for irregular applications

We obtained performance results for irregular applications on the SGI Origin2000 at Argonne National Laboratory. The first application template that we benchmarked was a tetrahedral vertex-centered unstructured grid code called FUN3D developed by Anderson of the NASA Langley Research Center [1]. This application uses a partitioning vector generated from METIS to partition the nodes and edges in a mesh. The mesh we used had about 18 million edges and 2 million nodes. At the initial stage, the application imports edges, four data arrays associated with edges, and another four data arrays associated with nodes. The total imported data size was about 807 Mbytes. As a result of computations, the application wrote about 21 Mbytes of four data sets each and 105 Mbytes of a single data set. Using 64 processors, we iterated the application template for two time steps; at each time step, five data sets were written to files.

The second application template that we ran was a Rayleigh–Taylor instability application [11] that is motivated by a joint project between the University of Chicago and Argonne to study thermonuclear flashes on astrophysical objects. Whenever the simulation time reaches a certain point, the application writes two data sets: a single node data set associated with vertices in a mesh, and a triangle data set associated with triangles on tetrahedral faces. In the application template, we wrote about 36 Mbytes of the node data set at each time step. Since we iterated the template five times, the total data size written was approximately 550 Mbytes.

### 4.2.1. Results for FUN3D

Fig. 16 shows the bandwidth to import and partition 18 million edges, four data sets each of 144 Mbytes of data associated with edges, and another four data sets each of 21 Mbytes of data associated with nodes. The original version of the application—without using SDM—performs all the I/O operations by a single process (process 0), which then broadcasts data to other processes. SDM performs I/O in parallel from all processes using MPI-IO. The bar labeled index distri in Fig. 16 shows the communication and computation costs to partition the edges after importing them to the application. The bar labeled import shows the cost of reading the edges and eight data arrays.

The original application reads the edges in two steps: one step to determine the amount of memory to store the partitioned edges and the other step to actually read the edges. SDM, however, extends the allocated memory dynamically as needed (using C function realloc) and is therefore able to read the partitioned edges in a single step. This contributes to the reduced cost of index distri when using SDM. When partitioning the edges with a history file, the cost of index distri is nothing but reading the history file of the edges in a contiguous way, including the database cost to access the metadata. Since the history file contains the already partitioned edges, there is no need to import the edges; hence, the read cost in import is reduced.

Fig. 17 shows the I/O bandwidth for writing and then reading back the data generated from the application using 64 processors. The total data size was approximately 379 Mbytes. In level 1, each data array is written to separate files, resulting in the creation of 10 different files. Each time the data array is written to files, level 1

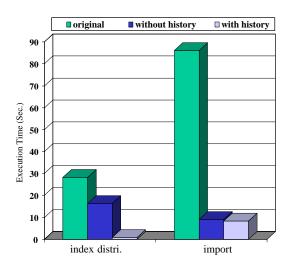


Fig. 16. Execution time for partitioning indices and data in FUN3D on SGI Origin2000.

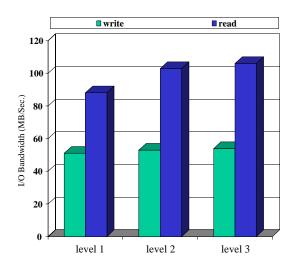


Fig. 17. I/O bandwidth for reading and writing data in FUN3D on SGI Origin2000.

requires the cost for opening a file and defining an MPI-IO *file view* to access the data from the portion of the file pointed by the global file offset. In level 2, however, each data array generated at each time step is appended in five files, generating five file-open and file-view costs. This reduced number of files improves the I/O performance slightly. In level 3, only two files are generated, resulting in the best I/O performance among the three file organizations. The difference between three file organizations is not significant because the file-open cost is small on the SGI Origin2000.

### 4.2.2. Results of RT application

Fig. 18 shows the I/O bandwidth in the RT application for writing approximately 550 Mbytes of data. In the original application, the write operation is performed sequentially. In other words, after seeking to the starting position in a file, processes write their local portion of data one by one. When we ported the application to SDM, the I/O performance increased significantly because of the parallel I/O optimizations of MPI-IO.

In SDM, we wrote the node data set according to the global node number of the partitioned nodes and wrote the triangle data set contiguously. Since two data sets are written to files separately, SDM supports two different ways of file organization: level 1 and level 2/3 (levels 2 and 3 are identical in this case). When the number of processors was increased from 32 to 64, but the total data size remained constant, the I/O performance declined. With 32 processors, the data size being written at each time step was about 1 Mbyte for the node data set and 2 Mbytes for the triangle data set. When the number of processors was increased to 64, the size of each I/O operation became smaller, reducing performance.

# To all a constraints and the second s

Fig. 18. I/O bandwidth for RT on SGI Origin2000.

### 5. Related work

Related work in this area falls into the categories of I/O libraries and parallel file systems for high-performance I/O and other libraries that aim to provide datamanagement capabilities.

Several efforts have sought to optimize I/O in parallel file systems and runtime libraries [4,6,7,17,21,23,27, 35,39]. SDM uses MPI-IO and parallel file systems to access real data, and therefore SDM is a consumer of the results of such research efforts as they become available through MPI-IO implementations.

Other efforts provide data-management capability. For example, HDF [44] and NetCDF [26] are popular libraries for data management in scientific applications, and they provide a self-describing data format. Storage Resource Broker (SRB) [3] provides uniform interface to access various storage systems, such as file systems, Unitree, HPSS and database objects. Shoshani et al. [36,37] describe an architecture to store large volumes of scientific data on tertiary storage systems in a way that the cost of data retrieval is minimized. The Active Data Repository [22] and DataCutter [5] optimize storage, retrieval, and processing of very large multidimensional data sets.

The main differences between our work and the above efforts are that we strive for performance in addition to flexibility and that we aim to be general purpose, not tied to a particular application or particular class of applications. The technique we use of separating data and metadata and storing the data in a parallel file system and metadata in a database could be used to implement other libraries, such as HDF and NetCDF, and we plan to investigate that in the future.

### 6. Conclusions

We have presented the design and implementation of an environment for high-performance scientific data management, called Scientific Data Manager (SDM), that is built on top of MPI-IO and also interacts with a database for storing metadata. SDM provides a simple, high-level interface and performs all necessary I/O optimizations transparently to the user. We also experimented with different ways of organizing data in files, called level 1-level 3. In general, when file-open cost on a particular file system is high, level 3 performs well because it minimizes the number of files created. If the file-open cost is small, the performance of the three levels depends on how the number and size of files affect performance on the particular file system. An appropriate file-organization policy can thereby be chosen for a particular file system.

On the XFS file system, we found that the file-open cost was so small that it did not significantly affect I/O

performance. Instead, our experiment focused on the use of direct I/O and buffered I/O in the ASTRO3D template. For writing data, we found that direct I/O performed much better than buffered I/O by avoiding the overhead of copying the data into the XFS buffer cache and also because XFS allows direct writes to proceed concurrently. For reading data, however, buffered I/O performed better because of its read-ahead policy.

We have also presented the SDM API and architecture for I/O in irregular applications. Besides providing an easy-to-use user interface for managing large data sets, SDM uses the concept of a history file to optimize the cost of the index distribution. We studied the performance of SDM using two irregular applications: FUN3D and RT. When we used SDM in both applications, there was a significant improvement in I/O performance compared with that of the original applications. Also, we observed that using a history file for the index distribution helped to reduce the computation and communication costs.

In the future, we plan to use SDM with more applications and evaluate both the usability and performance. We plan to refine the SDM API so that it can be used on a wide range of applications. We plan to study and improve the scalability of the system, particularly the scalability of the metadata. We also plan to investigate how our implementation approach could be used to implement HDF and NetCDF and whether such an approach would improve performance and flexibility of those libraries.

### Acknowledgments

We thank Lois Curfman McInnes, Lori Freitag, and Dinesh Kaushik at Argonne National Laboratory for their invaluable help in providing the applications used for the performance evaluation of SDM.

### References

- W.K. Anderson, W.D. Gropp, D.K. Kaushik, D.E. Keyes, B.F. Smith, Achieving high sustained performance in an unstructured mesh CFD application, in: Proceedings of SC99: High Performance Networking and Computing, Portland, OR, November 1999.
- [2] Applications Working Group of the Scalable I/O Initiative, Preliminary Survey of I/O Intensive Applications. http:// www.cacr.caltech.edu/SI0/pubs/SI0\_apps.ps, 1994.
- [3] C. Baru, R. Moore, A. Rajasekar, M. Wan, The SDSC storage resource broker, in: Proceedings of CASCON '98, Toronto, Canada, December 1998.
- [4] R. Bennett, K. Bryant, A. Sussman, R. Das, J. Saltz, Jovian: a framework for optimizing parallel I/O, in: Proceedings of the Scalable Parallel Libraries Conference, IEEE Computer Society Press, Silver Spring, MD, October 1994, pp. 10–20.

- [5] M.D. Beynon, R. Ferreira, T. Kurc, A. Sussman, J. Saltz, DataCutter: Middleware for filtering very large scientific datasets on archival storage systems, in: Proceedings of the Eighth Goddard Conference on Mass Storage Systems and Technologies, College Park, MD, March 2000.
- [6] R. Bordawekar, J.M. del Rosario, A. Choudhary, Design and evaluation of primitives for parallel I/O, in: Proceedings of Supercomputing '93, Portland, OR, November 1993, pp. 452–461.
- [7] P.F. Corbett, D.G. Feitelson, The vesta parallel file system, ACM Trans. Comput. Systems 14 (3) (1996) 225–264.
- [8] R. Das, M. Uysal, J. Saltz, Yuan-Shin Hwang, Communication optimizations for irregular scientific computations on distributed memory architectures, J. Parallel Distributed Comput. 22 (3) (1994) 462–479.
- [9] J.M. del Rosario, A. Choudhary, High performance I/O for parallel computers: problems and prospects, Computer 27 (3) (1994) 59–68.
- [10] S.A. Fineberg, P. Wong, B. Nitzberg, C. Kuszmaul, PMPIO—a portable implementation of MPI-IO, in: Proceedings of the Sixth Symposium on the Frontiers of Massively Parallel Computation, October 1996, IEEE Computer Society Press, Silver Spring, MD, 1996, pp. 188–195.
- [11] L. Freitag, M. Jones, P. Plassmann, The scalability of mesh improvement algorithms, IMA Vol. Math. Appl. 105 (1998) 185–212.
- [12] W. Gropp, E. Lusk, R. Thakur, Using MPI-2: Advanced Features of the Message-Passing Interface, MIT Press, Cambridge, MA, 1999.
- [13] W.D. Gropp, D.E. Keyes, L.C. McInnes, M.D. Tidriri, Globalized Newton–Krylov–Schwarz algorithms and software for parallel implicit CFD, Technical Report ANL/MCS-P788-0100, 1998.
- [14] R.V. Hanxleden, K. Kennedy, J. Saltz, Value-based distributions and alignments in Fortran D (Special Issue on Compiling and Run-Time Issues for Distributed Address Space Machines), J. Programming Languages 2 (3) (1994) 259–282.
- [15] High Performance Fortran Forum, High Performance Fortran Language Specification, Version 1.0, Technical Report Version 1.0, Rice University, Houston, TX, May 1993.
- [16] M. Holton R. Das, XFS: a next generation journalled 64-Bit filesystem with guaranteed rate I/O, Technical Report, SGI, Inc., 1994.
- [17] J. Huber, C.L. Elford, D.A. Reed, A.A. Chien, D.S. Blumenthal, PPFS: a high performance portable parallel file system, in: Proceedings of the Ninth ACM International Conference on Supercomputing, July 1995, ACM Press, New York, 1995, pp. 385–394.
- [18] IBM Corp, IBM AIX parallel I/O file system: installation, administration, and use, Document Number SH34-6065-01, August 1995.
- [19] T. Jones, R. Mark, J. Martin, J. May, E. Pierce, L. Stanberry, An MPI-IO Interface to HPSS, in: Proceedings of the Fifth NASA Goddard Conference on Mass Storage Systems, College Park, MD, September 1996, pp. I:37–I:50.
- [20] G. Karypis, V. Kumar, A fast and high quality multilevel scheme for partitioning irregular graphs, J. Sci. Comput. 20 (1) (1998) 359–392.
- [21] D. Kotz, Disk-directed I/O for MIMD multiprocessors, ACM Trans. Comput. Systems 15 (1) (1997) 41–74.
- [22] T. Kurc, Chialin Chang, R. Ferreira, A. Sussman, J. Saltz, Querying very large multi-dimensional datasets in ADR, in: Proceedings of SC99: High Performance Networking and Computing, Portland, OR, November 1999.
- [23] T.M. Madhyastha, D.A. Reed, Intelligent, adaptive file system policy selection, in: Proceedings of the Sixth Symposium on the

Frontiers of Massively Parallel Computation, IEEE Computer Society Press, Silver Spring, MD, October 1996, pp. 172–179.

- [24] Message Passing Interface Forum, MPI-2: Extensions to the Message-Passing Interface, July 1997. http://www.mpi-forum. org/docs/docs.html.
- [25] MySQL Reference Manual. http://www.mysql.com, 1999, Version 3.23.10-alpha.
- [26] NetCDF. http://www.unidata.ucar.edu/packages/netcdf.
- [27] N. Nieuwejaar, D. Kotz, The Galley parallel file system, Parallel Comput. 23 (4) (1997) 447–476.
- [28] J. No, R. Thakur, A. Choudhary, Integrating parallel file I/O and database support for high-performance scientific data management, in: Proceedings of SC2000: High Performance Networking and Computing, Dallas, TX, November 2000.
- [29] J. No, R. Thakur, D. Kaushik, L. Freitag, A. Choudhary, A scientific data management system for irregular applications, in: Proceedings of IPDPS'01, San Francisco, CA, April 2001.
- [30] PETSc 2.0 for MPI. http://www.mcs.anl.gov/petsc.
- [31] Postgres Global Development Group, PostgreSQL User's Guide, 1996.
- [32] J.-P. Prost, MPI-IO/PIOFS. World-Wide Web page at http:// www.research.ibm.com/people/p/prost/sections/ mpiio.html, 1996.
- [33] D. Sanders, Y. Park, M. Brodowicz, Implementation and performance of MPI-IO file access using MPI datatypes, Technical Report UH-CS-96-12, University of Houston, November 1996.
- [34] K. Schloegel, G. Karypis, V. Kumar, Graph Partitioning for High Performance Scientific Simulations, in: J. Dongarra et al. (Eds.), Sourcebook of Parallel Computing, Morgan Kaufmann, 2002, pp. 491–541.
- [35] K.E. Seamons, Y. Chen, P. Jones, J. Jozwiak, M. Winslett, Server-directed collective I/O in panda, in: Proceedings of Supercomputing '95, December 1995, ACM Press, New York, 1995.

- [36] A. Shoshani, L.M. Bernardo, H. Nordberg, D. Rotem, A. Sim, Storage management for high energy physics applications, in: Proceedings of Computing in High Energy Physics (CHEP '98), Chicago, IL, 1998.
- [37] A. Shoshani, L.M. Bernardo, H. Nordberg, D. Rotem, A. Sim, Multidimensional indexing and query coordination for tertiary storage management: in proceedings of SSDBM'99, Cleveland, OH, July 1999, pp. 214–225.
- [38] A. Sweeney, D. Doucette, Wei Hu, C. Anderson, M. Nishimoto, G. Peck, Scalability in the XFS file system, in: Proceedings of USENIX 1996 Annual Technical Conference, San Diego, CA, January 1996.
- [39] R. Thakur, A. Choudhary, An extended two-phase method for accessing sections of out-of-core arrays, Sci. Programming 5 (4) (1996) 301–317.
- [40] R. Thakur, W. Gropp, E. Lusk, An experimental evaluation of the parallel I/O systems of the IBM SP and intel paragon using a production application, in: Proceedings of the Third International Conference of the Austrian Center for Parallel Computation (ACPC) with Special Emphasis on Parallel Databases and Parallel I/O, Lecture Notes in Computer Science, Vol. 1127, Springer, Berlin, September 1996, pp. 24–35.
- [41] R. Thakur, W. Gropp, E. Lusk, A case for using MPI's derived datatypes to improve I/O performance, in: Proceedings of SC98: High Performance Networking and Computing, Orlando, FL, November 1998.
- [42] R. Thakur, W. Gropp, E. Lusk, On implementing MPI-IO portably and with high performance, in: Proceedings of the Sixth Workshop on I/O in Parallel and Distributed Systems, ACM Press, New York, May 1999, pp. 23–32.
- [43] R. Thakur, E. Lusk, W. Gropp, Users guide for ROMIO: a highperformance, portable MPI-IO implementation, Technical Report ANL/MCS-TM-234, Mathematics and Computer Science Division, Argonne National Laboratory, Revised December 1999.
- [44] The NCSA HDF home page. http://hdf.ncsa.uiuc.edu.