

# April: A Run-Time Library for Tape-Resident Data

Gokhan Memik, Mahmut T. Kandemir, Alok Choudhary, and Valerie E. Taylor

Center for Parallel and Distributed Computing

Northwestern University, Evanston, IL 60202

{memik,mtk,choudhar,taylor}@ece.nwu.edu

tel: +1 847 467-2299

fax: +1 847 491-4455

## Abstract

*Over the last decade, processors have made enormous gains in speed. But increase in the speed of the secondary and tertiary storage devices could not cope with these gains. The result is that the secondary and tertiary storage access times dominate execution time of data intensive computations. Therefore, in scientific computations, efficient data access functionality for data stored in secondary and tertiary storage is a must. In this paper, we give an overview of APRIL, a parallel runtime library that can be used in applications that process tape-resident data. We present user interface and underlying optimization strategy. We also discuss performance improvements provided by the library on the High Performance Storage System (HPSS). The preliminary results reveal that the optimizations can improve response times by up to 97.2%.*

## 1 Introduction

We address the problem of managing the movement of very large data sets between different levels of a hierarchical storage system. It is now widely acknowledged that the data set sizes manipulated by scientific codes are getting larger as programmers have access to faster processors and larger main memories. The data sets whose sizes exceed main memories should be stored in secondary and tertiary storages. Although the prices for secondary storage devices are decreasing, tertiary storage devices are becoming increasingly attractive especially for applications that require vast amount of storage capacity which cannot be satisfied by secondary storage devices and for applications which cannot afford the cost or system complexity of a large number of disk drives. There has been a considerable amount of work in addressing the flow of data to and from secondary storage devices (e.g., magnetic disks) [1, 2, 3, 4, 5, 6, 7, 8, 9]. There has also been a significant amount of work on the management of large scale data in a storage hierarchy involving tertiary storage devices (e.g., tapes devices) [10, 11, 12, 13, 14]. Striping has been studied to improve the response time of tertiary storage devices [15, 16].

The Department of Energy's ASCI plan draws an outline of the expected storage requirements for large-scale computational challenges. According to this plan, a large scientific application today is producing 3-30 terabytes of simulation datasets for a run, requiring 3 petabytes of archive capacity. These sizes are expected to grow to 100-1000 terabytes per run and to 100 petabytes of archive capacity in the year 2004. Even with the assumptions

of aggressive improvements in the evolution of storage devices, the data accesses in these applications will take a significant proportion of the overall execution time [17]. On top of this, aggregate data sizes may require the employment of tertiary storage devices. Many of these applications do not demand the entire datasets to be accessed at a given time. So, having means to access portions of the tape-resident datasets efficiently may decrease the time spent in data accesses significantly.

In this paper, we present APRIL, a parallel run-time library, that can be used to facilitate the explicit control of data flow for tape-resident data. Our library can be used by application programmers as well as optimizing compilers that manipulate large scale data. The objective is to allow programmers to access data located on tape via a convenient interface expressed in terms of arrays and array portions (regions) rather than files and offsets. In this sense the library can be considered as a natural extension of state-of-the-art run-time libraries that manipulate disk-resident datasets (e.g., [2, 18]). The library implements a data storage model on tapes that enables users to access portions of multi-dimensional data in a fast and simple way. In order to eliminate most of the latency in accessing tape-resident data, we employ a *sub-filing strategy* in which a large multi-dimensional tape-resident *global* array is stored not as a single file but as a number of smaller *sub-files*, whose existence is transparent to the programmer. The main advantage of doing this is that the data requests for relatively small portions of the global array can be satisfied without transferring the entire global array from tape to disk as is customary in many hierarchical storage management systems. In addition to read/write access routines, the library also supports pre-staging and migration capabilities which can prove very useful in environments where the data access patterns are predictable and the amount of disk space is limited.

The main contributions of this paper are as follows:

- The presentation of a high-level parallel I/O library for tape-resident data. APRIL library provides a simple interface to the tape-resident data, which relieves the programmers from orchestrating I/O from tertiary storage devices such as robotic tapes and optical disks.
- The description of the implementation of the library using HPSS [19] and MPI-IO [3]. We show that it is both simple and elegant to build an I/O library for tape-resident data on top of these two state-of-the-art systems. In this paper, however, we focus on a single processor performance.
- The presentation of preliminary performance numbers using representative array regions and sub-file sizes. The results demonstrate that the library is quite effective in exploiting the secondary storage – tertiary storage hierarchy without undue programmer effort.

Section 2 gives an overview of the APRIL library. Section 3 describes sub-filing and its use in the library. Section 4 briefly explains the implementation and Section 5 presents the user interface. Section 6 gives preliminary experimental results and Section 7 concludes the paper with a summary and an outline of future work.

## 2 Library Overview

The library provides routines to efficiently perform I/O required in sequential and parallel applications. It can be used for both in-core and out-of-core applications. It uses a high-level interface which can be used by application programmers and compilers. For example, an application programmer can specify what section of an array she wants to read in terms of lower and upper bounds in each dimension, and the library will fetch it in an efficient manner, first from tape to disk and then from disk to main memory. It provides a portable interface on top of HPSS [19] and MPI-IO [3]. It can also be used by an optimizing compiler that targets programs whose data sets require transfers between secondary storage and tertiary storage. It might even be possible to employ the library within a database management system for multi-dimensional data.

At the heart of the library is an optimization technique called *sub-filing*, which is explained in greater detail in the next section. It also uses collective I/O using a two-phase method, data pre-staging, pre-fetching, and data migration. The main advantage of sub-filing is that it provides low-overhead random access image for the tape-resident data. Sub-filing is invisible to the user and helps to efficiently manage the storage hierarchy which can consist of a tape sub-system, a disk sub-system and a main memory. The main advantage of the collective I/O, on the other hand, is that it results in high-granularity data transfers between processors and disks, and it also makes use of the higher bandwidth of the processor inter-connection network.

In general, a processor has to wait while a requested tape-resident data set is being read from tape. The time taken by the program can be reduced if the computation and tape I/O can be overlapped somehow. The pre-staging achieves this by bringing the required data ahead of the time it will be used. It issues asynchronous read calls to the tape sub-system, which help to overlap the reading of the next data portion with the computation being performed on the current data set. The data pre-fetching is similar except that it overlaps the disk I/O time with the computation time.

## 3 Sub-filing

Each *global tape-resident* array is divided into *chunks*, each of which is stored in a *separate sub-file* on tape. The chunks are of equal sizes in most cases. Figure 1 shows a two-dimensional global array divided into 64 chunks. Each chunk is assigned a unique *chunk coordinate*  $(x_1, x_2)$ , the first (upper-leftmost) chunk having  $(0,0)$  as its coordinate. For the sake of ensuing discussion we assume that the sub-files corresponding to the chunks are stored in row-major as depicted in the figure by horizontal arrows.

A typical access pattern is shown in Figure 2. In this access a small two-dimensional portion of the global array is requested. In receiving such a request, the library performs three important tasks:

- Determining the sub-files that collectively contain the requested portion,
- Transferring the sub-files that are *not* already on disk from tape to disk, and
- Extracting the required data items (array elements) from the relevant sub-files from

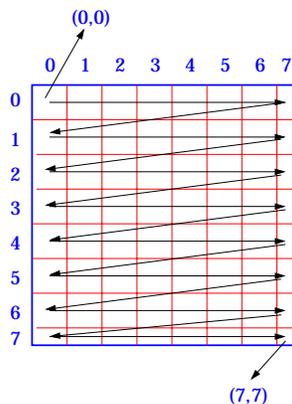


Figure 1: A global tape-resident array divided into 64 chunks.

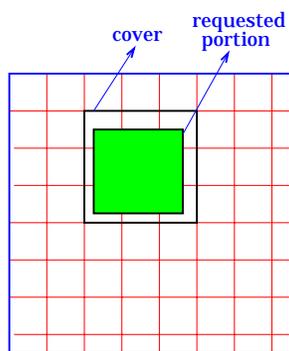


Figure 2: An access pattern (shaded portion) and its cover.

disk and copying the requested portion to a buffer in memory provided by the user call.

In the first step, the set of sub-files that collectively contain the requested portion is called *cover*. In Figure 2, the cover contains the sub-files (1,2), (1,3), (1,4), (2,2), (2,3), (2,4), (3,2), (3,3), and (3,4). Assuming for now that all of these sub-files are currently residing on tape, in the second step, the library brings these sub-files to disk. In the third step, the required portion is extracted from each sub-file and returned to the user buffer. Note that the last step involves some computational overhead incurred for each sub-file. Instead, had we used just one file per global array this computational overhead would be incurred only once. Therefore, the performance gain obtained by dividing the global array into sub-files should be carefully weighed against the extra computational overhead incurred in extracting the requested portions from each sub-file. Our preliminary experiments show that this computational overhead is not too much.

Parallel reads by multiple processors pose additional problems. Consider now Figure 3(a) where four processors are requesting four different sub-columns of a region. The underlying cover contains 28 sub-files. After bringing these sub-files from tape to disk, we have a problem of reading the required sub-portions (sub-columns) for each processor. As stated by del Rosario et al. [20], *collective I/O* is a technique in which processors perform I/O on

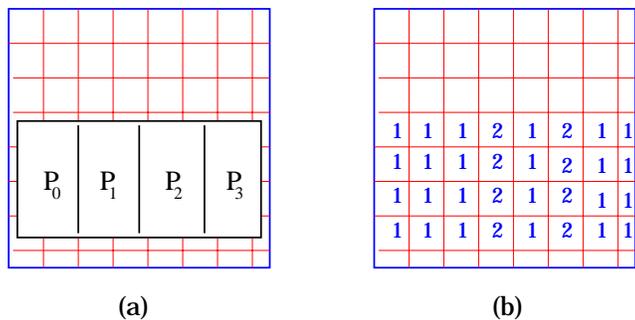


Figure 3: (a) An access pattern involving four processors. (b) The global array with each sub-file marked with the number of processors that share it.

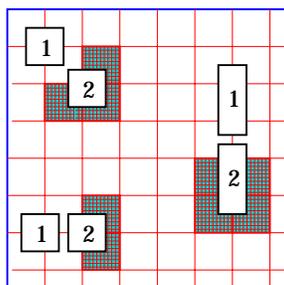


Figure 4: Successive array accesses.

behalf of each other in order to reduce the time spent in disk-I/O at the expense of some extra communication. Two-phase I/O is a specific implementation of collective I/O, which uses the information available about the access and storage patterns. It employs two phases. In the first phase, the processors access the data in a *layout conformant* way (to exploit spatial locality on disk as much as possible) and in the second phase they re-distribute the data in memory among themselves such that the desired access pattern is obtained. While it is quite straightforward how to use collective I/O when we have a single file, in our multiple file case it is not clear how to utilize it. One simple approach might be to use collective I/O for each sub-file on disk. In our example, that would mean calling a collective I/O routine 28 times. A better alternative might be to read the data from disk to memory in two steps. In the first step, the processors that have exclusive access to some sub-files perform these independent accesses. In the second step, for each of the remaining sub-files, we can perform collective I/O using only the processors that request some data from the sub-file in question. Considering Figure 3(b), this collective I/O scheme corresponds to first reading the sub-files marked '1' and then collectively (using two processors) reading the sub-files marked '2'. We plan to implement this last collective I/O strategy in the future.

During successive reads from the same global file it might happen that the same sub-file can be required by two different reads. Assuming that the sub-file in question still resides on the disk after the first read, it is unnecessary to read it again from tape in the second read. In such a case only the other (additional) sub-files required by the current access are read from tape. The situation is shown in Figure 4 for three scenarios. In each case, the

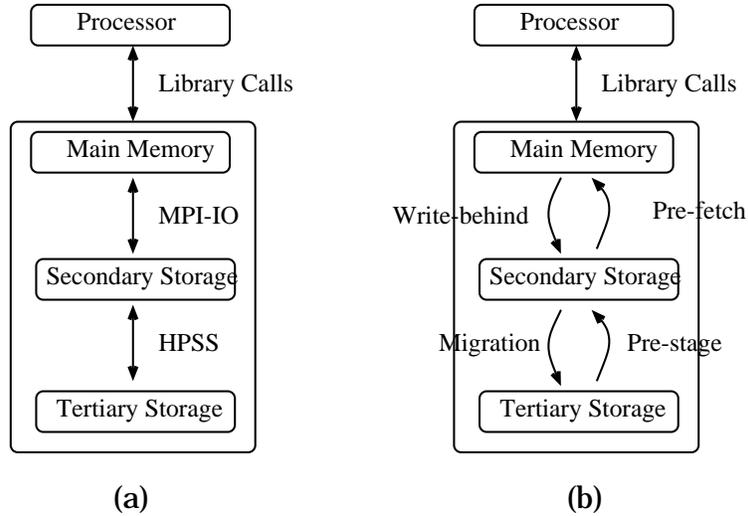


Figure 5: (a) Library architecture. (b) Pre-fetching, pre-staging, and migration.

first portion read is marked ‘1’ and the second portion read is marked ‘2’. The shaded parts around the second portions correspond to additional sub-files that are needed for the second read. In other words, the library effectively uses the storage hierarchy.

## 4 Implementation

We are implementing the APRIL library on top of HPSS [19] and MPI-IO [3]. The connections between different components are shown in Figure 5(a). In a read call, the sub-files are first read from tape to disk using HPSS and then from disk to main memory using MPI-IO. As mentioned earlier, we employ collective I/O between disk and memory. In a write call the direction of data-flow is reversed. Figure 5(b) shows the corresponding storage levels for each action described in the following sections.

To store the information about the file and the chunks, we are currently using the Postgres95 database [21]. When a new file is created, the user may enter the necessary information about the chunks. The detailed information about the file creation is given in Section 5. Then this meta-data is stored in the database for later usage. When a user opens a previously created file, the corresponding meta-data about the file and the chunks are read from the database and cached in the memory. Then the following accesses uses this meta-data. The database is informed about the changes when the file is closed. In other words, the database is accessed only in file open and file close.

## 5 User Interface

The routines in the library can be divided into four major groups based on their functionality – Initialization/Finalization Routines, File Manipulation Routines, Array Access Routines, and Stage/Migration Routines. Table 1 lists some of the basic routines in the library and their functionalities. All the routines listed here are the low level instructions, which should explicitly be called by the user. We are currently adding high level routines to our library,

Table 1: Some of the library routines.

<b>Initialization/Finalization Routines</b>	
Routine	Functionality
T_INITIALIZE	Initializes the library structures
T_FINALIZE	Finalizes the library structures
<b>File Manipulation Routines</b>	
Routine	Functionality
T_OPEN	Opens a tape-resident global file for read/write
T_CLOSE	Closes a tape-resident global file
T_REMOVE	Removes both the sub-files of the file and the corresponding info.
<b>Array Access Routines</b>	
Routine	Functionality
T_READ_SECTION	Reads a rectilinear section
T_WRITE_SECTION	Writes a rectilinear section
<b>Stage/Migration Routines</b>	
Routine	Functionality
T_STAGE_SECTION	Stages a rectilinear file section from tape to disk
T_STAGE_WAIT	Waits for a Stage to complete
T_PREFETCH_SECTION	Pre-fetches a rectilinear file section from tape (or disk) to memory
T_PREFETCH_WAIT	Waits for a Pre-fetch to complete
T_MIGRATE_SECTION	Migrates a rectilinear file section from disk to tape

which will call the low level routines implicitly.

**Initialization/Finalization Routines:** These routines are used to initialize the library buffers and meta-data structures and finalize them when all the work is done. The routine to initialize the system has the format

```
int T_INITIALIZE ().
```

This routine initializes the connections to the HPSS and the database. It returns a positive number upon successful completion. Similarly, `int T_FINALIZE ()` closes the above mentioned connections.

**File Manipulation Routines:** These routines are used for creating files, opening existing files, closing open files and removing all the chunks and the information related to a global file. `T_OPEN` is used for creating new files and for opening existing files. It returns a file handle for later referral to the file. The synopsis of `T_OPEN` is as follows:

```
T_FILE T_OPEN (char *filename, char *mode, T_INFO *tapeinfo).
```

‘Filename’ stands for the name of the file to be opened. ‘Mode’ indicates whether the file is opened for read, write, or read/write. ‘Tapeinfo’ is the structure used for entering the necessary information about the file and chunks. It has fields for the elementsize, number of dimensions, the size of each dimension of the chunk and the size of each dimension of the global file.

**Array Access Routines:** These routines handle the movement of data to and from the tape

```

int main(int argc, char **argv)
{
    T_FILE exfile;
    T_INFO exinfo;
    int start[2];
    int end[2];

    /* Initialize the library */
    T_INITIALIZE();

    /* Open the file for read.  The exinfo will be filled by the library.
    For creating the file (i.e. if the file is opened for the first time),
    information about the file should be supplied to T_OPEN via exinfo.*/
    exfile = T_OPEN ("file_1","r", &exinfo);

    start[0] = 0;
    start[1] = 0;
    end[0] = 24000;
    end[1] = 80;

    /* Perform the operation */
    T_READ_SECTION (&exfile, &buf, starts, ends);

    /* Close the file */
    T_CLOSE (&exfile);

    T_FINALIZE();
}

```

Figure 6: An example code for reading from a two dimensional file.

subsystem. An arbitrary rectilinear portion of a tape-resident array can be read or written using these access routines. Let us focus now on T\_READ\_SECTION. The signature of this routine is

```

int T_READ_SECTION (T_FILE *fd, void *buffer, int *start_coordinate,
int *end_coordinate)

```

‘fd’ is the file descriptor returned by T\_OPEN. ‘start\_coordinate’ and ‘end\_coordinate’ are arrays that hold the boundary coordinates for the section to be read. There are as many elements as the dimensionality of the associated tape-resident global array. This command reads the corresponding elements and stores them in ‘buffer’. As discussed earlier, what actually happens here is that the relevant sub-files are read from tape to disk (if they are not on the disk already), and the required sections are read from these sub-files on disk and forwarded to the corresponding positions in the buffer in main memory. An example code for T\_READ\_SECTION is given in Figure 6. In this example, a 24000×80 portion of the file is read to the buffer. The syntax for the T\_WRITE\_SECTION routine is almost the same except that the direction of the transfer is reversed.

**Stage/Migration Routines:** These routines are used to stage and migrate the data between

the tapes and the disk sub-system. The command

```
int T_STAGE_SECTION (T_FILE *fd, int *start_coordinate, int *end_coordinate)
```

immediately returns and starts a data staging operation in the background from tape to disk. It returns an integer to the application which can be interpreted as a descriptor for the associated pre-stage operation. Note that what is actually performed here is to bring the relevant *sub-files* from tape to disk. Note also that there is no ‘buffer’ parameter in the signature. The routine

```
int T_STAGE_WAIT (int pre-stage_descriptor)
```

can be used to wait for a previously initiated pre-stage operation to complete.

```
int T_PREFETCH_SECTION (T_FILE *fd, void *buffer, int *start_coordinate,  
int *end_coordinate)
```

is used to start a pre-fetch operation from disk to memory. The parameters are the same as for T\_READ\_SECTION. It returns an integer which can be used as a pre-fetch descriptor in a later T\_PREFETCH\_WAIT call.

```
int T_MIGRATE_SECTION (T_FILE *fd, int *start_coordinate, int *end_coordinate)
```

starts to migrate the relevant sub-files (i.e., those corresponding to the section described in the signature) from disk to tape. It should be used with care as these sub-files may contain portions of data that will be requested by a later library call.

## 6 Experiments

The experiments are performed using the HPPS at the San Diego Supercomputing Center (SDSC). We have used the low level routines of the SDSC Storage Resource Broker (SRB) to access the HPSS files. SRB is a client-server middleware that provides a uniform interface for connecting to heterogeneous data resources over a network and accessing replicated data sets [22].

We experimented with different *access patterns* in order to evaluate the benefits of the library. Table 2 gives the start and end coordinates (on a two dimensional global array) as well as the number elements read/written for each access pattern (A through H). Note that the coordinate (0,0) corresponds to the upper-left corner of the array. In each case, the accessed array consists of  $50000 \times 50000$  floating point elements (10 GB total data). We used two different sub-file (chunk) sizes: small ( $1000 \times 1000$  elements) and large ( $2000 \times 2000$  elements).

Table 3 shows the performance results obtained. For each operation (read or write) we give the response times (in *seconds*) for a naive access strategy and the gains obtained against it using our library which employs sub-filing. The naive strategy reads/writes the required portion from/to the array directly, i.e., it does not use sub-filing and the entire  $50000 \times 50000$  array is stored as a single large file. For the sub-filing cases we show the *percentage reduction* in response time of the naive scheme. For example, in access pattern A, the sub-filing with small chunk size improved (reduced) the response time for the read operation by 85.2%. Figures 7 and 8 show the results obtained in graphical form. Note that the y-axes on the figures are *logarithmically scaled*.

Table 2: Access patterns.

Access Pattern	Pattern Information		
	Start Coordinate	End Coordinate	Total floating points
A	(0,0)	(1000,1000)	$1 * 10^6$
B	(0,0)	(4000,1000)	$4 * 10^6$
C	(0,0)	(24000,1000)	$24 * 10^6$
D	(5000,5000)	(6000,6000)	$1 * 10^6$
E	(0,0)	(50000,80)	$4 * 10^6$
F	(0,0)	(80,50000)	$4 * 10^6$
G	(0,0)	(1000,4000)	$4 * 10^6$
H	(6000,6000)	(8000,8000)	$4 * 10^6$

Table 3: Execution times and percentage gains.

Acc. Ptr.	Write Operations			Read Operations		
	Times w/o chunking	Small Chunk Gain (%)	Large Chunk Gain (%)	Times w/o chunking	Small Chunk Gain (%)	Large Chunk Gain (%)
A	2774.0	96.1	94.5	784.7	85.2	77.1
B	2805.9	83.8	84.9	810.1	43.2	55.6
C	2960.3	8.8	37.9	793.3	-240.5	-172.4
D	3321.2	96.7	95.4	798.4	84.1	79.7
E	151.7	-3525.1	-2437.6	165.2	-3229.3	-2623.9
F	138723.3	96.0	97.2	39214.1	85.9	88.5
G	11096.3	95.9	96.4	3242.9	88.3	88.6
H	5095.2	91.2	96.5	1612.9	76.6	89.9

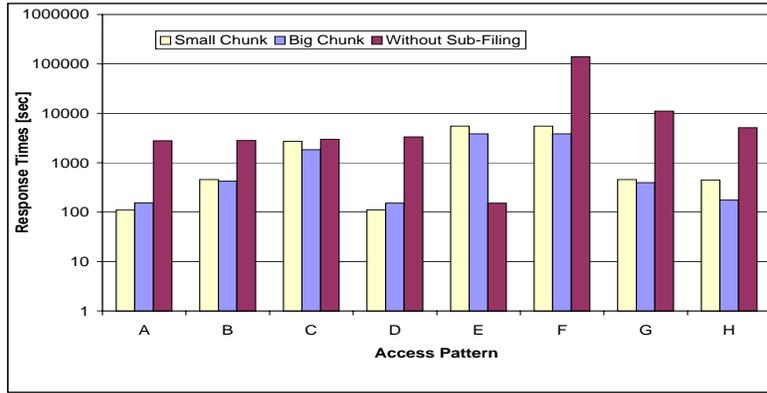


Figure 7: Execution times for **write** operations.

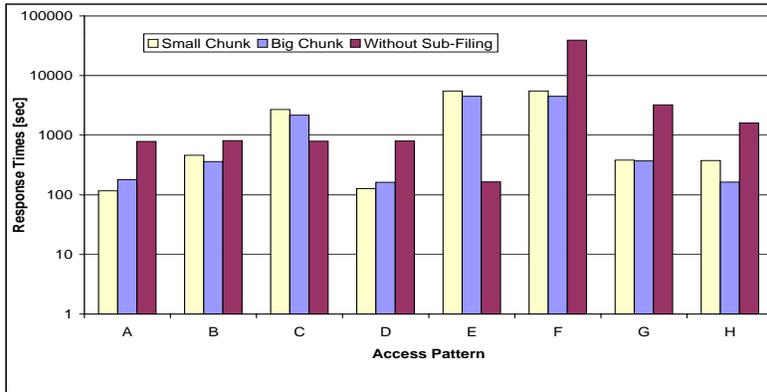


Figure 8: Execution times for **read** operations.

In the patterns A and D, where a 4 MB square chunk is accessed on the left corner and around the middle, respectively, the small chunk size outperforms the large chunk size as the latter accesses extra data elements that do not belong to the required portion. In the pattern H, on the other hand, increasing the chunk size reduces the number of I/O calls which in turn results in the best response time. In B and G, 16 MB of data are accessed in orthogonal directions. In G, since we access a sub-column portion of a row-major array, we need to issue 4000 I/O calls in the naive case. In B, the naive strategy issues only 1000 I/O calls to access the same volume of data. Consequently, the impact of sub-filing is more pronounced in G. By comparing the response times of A, B, C, and E, we note that the response times are dominated by the number of I/O calls (in the naive version) and of chunks (in the sub-filed versions—that also corresponds to I/O calls—) rather than by the volume of data accessed. Finally, in the pattern F (whose response time in the naive case was calculated using interpolation from A and G), the sub-filing strategy has the best performance of

all and brings a 97.2% improvement in write calls.

In access pattern E, however, the naive strategy outperforms the sub-filing. The sub-filing strategy has two drawbacks for this access pattern. First, the naive strategy completes the whole access with a single I/O call, whereas the sub-filing strategy requires 50 calls to different sub-files to satisfy the access. Secondly, a high percentage of data read by the sub-filing is not used to satisfy the request. As a result of these two drawbacks, the naive strategy performs better than the sub-filing for this access pattern. However, we show in Section 6.1 that by choosing an appropriate sub-file size, the sub-filing strategy can perform as good as the naive strategy even for the access pattern E. Note that, HPSS allows the users to access portions of the data residing in tape. The response time of the naive solution for the access pattern E will increase dramatically for the tape architectures, where the granularity of access is a file, because the whole file should be brought to the disk from tape.

An important aspect of our library is its handling of random I/O accesses. When we compare the times for the access patterns A and D, we see that there is a 20% increase in the response time of the naive strategy for the write operation. On the other hand, the times for sub-filed versions remains the same.

Overall, the sub-filing strategy performs very well compared to the naive strategy which performs individual accesses to a large file, except for the cases where the access pattern and the storage pattern of the array match exactly. For large chunk size the average improvement for writing is 93.48%, and for reading it is 73.48%. These data show that, in average our library brings substantial amount of improvement over the naive strategy. The next section shows even in the case where the access and storage pattern match exactly, a suitable chunk shape allows our scheme to match the response time of the naive strategy.

## 6.1 Adaptive Chunk Size

The preliminary results show that our library can bring substantial amount of improvement over the naive case. In the access pattern E, however, the naive strategy performs better. In this section, we experiment with a different chunk size to explore the possibility of matching the performance of the naive scheme in this access pattern.

In the new experiments, the chunk size is set to  $50000 \times 80$  floating points, which is similar to the access pattern E. The other parameters remain as in Section 6. The response time for write operation drops to 148.9 seconds, which is 1.85% better than the naive scheme. For read operation, the response time is 166.2 seconds, which is 0.61% worse than the naive scheme. These results show that our library can perform as good as the naive scheme even in cases, where the access pattern and storage pattern exactly match.

## 7 Conclusions and Future Work

We presented a portable interface to the tape-resident data. The interface makes it easier for the user to specify the tape I/O required in sequential and parallel applications. The

experience gained during its design and development will, hopefully, also help in reaching a set of standard routines for accessing the tape-resident data. We are in the process of implementing the library. We completed the read, write, pre-stage, and pre-fetch routines and made some initial experiments with them. We are currently implementing different migration routines and collective I/O strategies and will later experiment with I/O-intensive applications that manipulate tape-resident data.

## Acknowledgements

We would like to thank NPACI for providing the hardware and software components used in this study. This work is supported by Department of Energy under the Accelerated Strategic Computing Initiative (ASCI) Academic Strategic Alliance Program (ASAP), under subcontract no. (W-7405-ENG-48) from Lawrence Livermore National Laboratories. We would also like to thank Regan Moore and Micheal Wan of SDSC for their help with the SRB system.

## References

- [1] P. Cao, E. Felten, and K. Li. Application-controlled file caching policies. In *Proc. the 1994 Summer USENIX Technical Conference*, pages 171–182, June 1994.
- [2] A. Choudhary, R. Bordawekar, M. Harry, R. Krishnaiyer, R. Ponnusamy, T. Singh, and R. Thakur. PASSION: parallel and scalable software for input-output. *NPAC Technical Report SCCS-636*, 1994.
- [3] P. Corbett, D. Fietelson, S. Fineberg, Y. Hsu, B. Nitzberg, J. Prost, M. Snir, B. Traversat, and P. Wong. Overview of the MPI-IO parallel I/O interface, In *Proc. Third Workshop on I/O in Paral. and Distr. Sys., IPPS'95*, Santa Barbara, CA, April 1995.
- [4] J. del Rosario and A. Choudhary. High performance I/O for parallel computers: problems and prospects. *IEEE Computer*, March 1994.
- [5] C. S. Ellis and D. Kotz. Prefetching in file systems for MIMD multiprocessors. In *Proc. of the 1989 Intl. Conf. on Paral. Proc.*, pages 1:306–314, St. Charles, IL, August 1989.
- [6] J. F. Karpovich, A. S. Grimshaw, and J. C. French. Extensible file systems (ELFS): An object-oriented approach to high performance file I/O. In *Proc. the Ninth Annual Conference on Object-Oriented Prog. Sys., Lang., and Appl.*, pp. 191–204, Oct 1994.
- [7] D. Kotz. Multiprocessor file system interfaces. In *Proc. the Second Intl. Conf. on Paral. and Distr. Info. Sys.*, pages 194–201. IEEE Computer Society Press, 1993.
- [8] R. H. Patterson, G. A. Gibson, and M. Satyanarayanan. A status report on research in transparent informed prefetching. *ACM Operating Systems Review*, V 27(2), pp 21–34, April 1993.
- [9] R. Bordawekar, A. Choudhary, K. Kennedy, C. Koelbel, and M. Paleczny. A model and compilation strategy for out-of-core data parallel programs. *Proceedings of the ACM Symp. on Prin. and Prac. of Paral. Prg.*, pages 1–10, July 1995.
- [10] S. Sarawagi: Execution reordering for tertiary memory access. *Data Engineering Bulletin* 20(3): 46–54, 1997.
- [11] B. Hillyer and A. Silberschatz. Random I/O scheduling in Online Tertiary Storage Systems. In *Proc. of the 1996 ACM SIGMOD Conference*, pages 195–204, 1996.
- [12] T. Johnson and E. L. Miller. Performance measurements of Tertiary Storage Devices. In *Proc. of 24rd International Conference on Very Large Data Bases*, pages 50–61, August 1998.

- [13] T. Johnson and E. L. Miller. Benchmarking Tape System Performance. In *Proc. of the 15th IEEE Mass Storage Systems Symposium*, 1996.
- [14] B. Kobler, J. Berbert, P. Caulk and P. C. Hariharan. Architecture and Design of Storage and Data Management for the NASA Earth Observing System Data and Information System (EOSDIS). In *Proc. of the 14th IEEE Mass Storage Systems Symposium*, pages 65–78, 1995.
- [15] A. L. Drapeau and R. H. Katz. Analysis of Striped Tape Systems. In *Proc. of the 12th IEEE Mass Storage Symposium*, Monterey, CA, March 1993.
- [16] L. Golubchik, R. R. Muntz, and R. W. Watson. Analysis of Striping Techniques in Robotic Storage Libraries. In *Proc. of the 14th IEEE Mass Storage Sys. Symp.*, pages 225–238, 1995.
- [17] P. H. Smith and J. van Rosendale. Data and visualization corridors. *Technical Report CACR-164*, CACR, Caltech, Sept 1998.
- [18] S. Toledo and F. G. Gustavson. The design and implementation of SOLAR, a portable library for scalable out-of-core linear algebra computations, In *Proc. Workshop on I/O in Paral. and Distr. Sys.*, May 1996.
- [19] R. A. Coyne, H. Hulen, and R. Watson. The high performance storage system. In *Proc. Supercomputing 93*, Portland, OR, November 1993.
- [20] J. del Rosario, R. Bordawekar, and A. Choudhary. Improved parallel I/O via a two-phase run-time access strategy. In *Proc. the 1993 IPPS Workshop on Input/Output in Paral. Comp. Sys.*, April 1993.
- [21] A. Yu, and J. Chen. *The POSTGRES95 User Manual*. Dept. of EECS, University of California at Berkeley, July 1995.
- [22] S. Baru. Storage Resource Broker (SRB) Reference Manual. Enabling Technologies Group, San Diego Supercomputer Center, La Jolla, CA, August 1997.