

Distributed Parallel Data Storage Systems: A Scalable Approach to High Speed Image Servers

*Brian Tierney (blierney@lbl.gov),
William E. Johnston¹(wejohnston@lbl.gov),
Hanan Herzog, Gary Hoo, Guojun Jin, Jason Lee,
Ling Tony Chen*, Doron Rotem**

*Imaging and Distributed Computing Group and *Data Management Research Group
Lawrence Berkeley Laboratory²
Berkeley, CA 94720*

Abstract

We have designed, built, and analyzed a distributed parallel storage system that will supply image streams fast enough to permit multi-user, “real-time”, video-like applications in a wide-area ATM network-based Internet environment. We have based the implementation on user-level code in order to secure portability; we have characterized the performance bottlenecks arising from operating system and hardware issues, and based on this have optimized our design to make the best use of the available performance. Although at this time we have only operated with a few classes of data, the approach appears to be capable of providing a scalable, high-performance, and economical mechanism to provide a data storage system for several classes of data (including mixed multimedia streams), and for applications (clients) that operate in a high-speed network environment.

1.0 Introduction

In recent years, many technological advances have made possible distributed multimedia servers that will allow bringing “on-line” large amounts of information, including images, audio and video, and hypermedia databases. In-

creasingly, there also are applications that demand high-bandwidth access to this data, either in single user streams (e.g., large image browsing, uncompressible scientific and medical video, and multiple coordinated multimedia streams) or, more commonly, in aggregate for multiple users. Our work focuses on two examples of high-bandwidth, single-user applications. First, the terrain visualization application described below requires 300-400 Mbits/s of data to provide a realistic dynamic visualization. Second, there are applications in the scientific and medical imaging fields where uncompressed video (e.g. from typical laboratory monochrome video cameras that produce 115 Mbits/s data streams) needs to be stored and played back at real-time rates. In these example applications compression is not practical: in the case of terrain visualization, the com-

1. Correspondence should be directed to W. Johnston, Lawrence Berkeley Laboratory, MS: 50B - 2239, Berkeley, CA, 94720. Tel: 510-486-5014, fax: 510-486-6363; or Brian Tierney, Tel: 510-486-7381.

2. This work is jointly supported by ARPA - CSTO, and by the U. S. Dept. of Energy, Energy Research Division, Office of Scientific Computing, under contract DE-AC03-76SF00098 with the University of California. This document is LBL report LBL-35408. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement or recommendation by the United States Government or the University of California. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or the University of California, and shall not be used for advertising or product endorsement purposes. The following terms are acknowledged as trademarks: UNIX (Novell, Inc.), Sun and SPARCStation (Sun Microsystems, Inc.), DEC and Alpha (Digital Equipment Corp.), SGI and Indigo (Silicon Graphics, Inc.).

putational cost of decompression is prohibitive; in the case of medical and scientific images, data loss, coupled with the possible introduction of artifacts during decompression, frequently precludes the use of current compression techniques. (See, for example, [7].)

Although one of the future uses of the system described here is for multimedia digital libraries containing multiple audio and compressed video streams, the primary design goal for this system is to be able to deliver high data rates: initially for uncompressed images, later for other types of data. Based on the performance that we have observed, we believe, but have not yet verified, that the approach described below will also be useful for the video server problem of delivering many compressed streams to many users simultaneously.

Background

Current disk technology delivers about 4 Mbytes/s (32 Mbits/s), a rate that has improved at about 7% each year since 1980 [8], and there is reason to believe that it will be some time before a single disk is capable of delivering streams at the rates needed for the applications mentioned. While RAID [8] and other parallel disk array technologies can deliver higher throughput, they

are still relatively expensive, and do not scale well economically, especially in an environment of multiple network distributed users, where we assume that the sources of data, as well as the multiple users, will be widely distributed. Asynchronous Transfer Mode (ATM) networking technology, due to the architecture of the SONET infrastructure that will underlie large scale ATM networks of the future, will provide the bandwidth that will enable the approach of using ATM network-based distributed, parallel data servers to provide high-speed, scalable storage systems.

The approach described here differs in many ways from RAID, and should not be confused with it. RAID is a particular data strategy used to secure reliable data storage and parallel disk operation. Our approach, while using parallel disks and servers, deliberately imposes no particular layout strategy, and is implemented entirely in software (though the data redundancy idea of RAID might be usefully applied across servers to provide reliability in the face of network problems).

Overview

The Image Server System (ISS) is an implementation of a distributed parallel data storage architecture. It is essentially a “block” server that

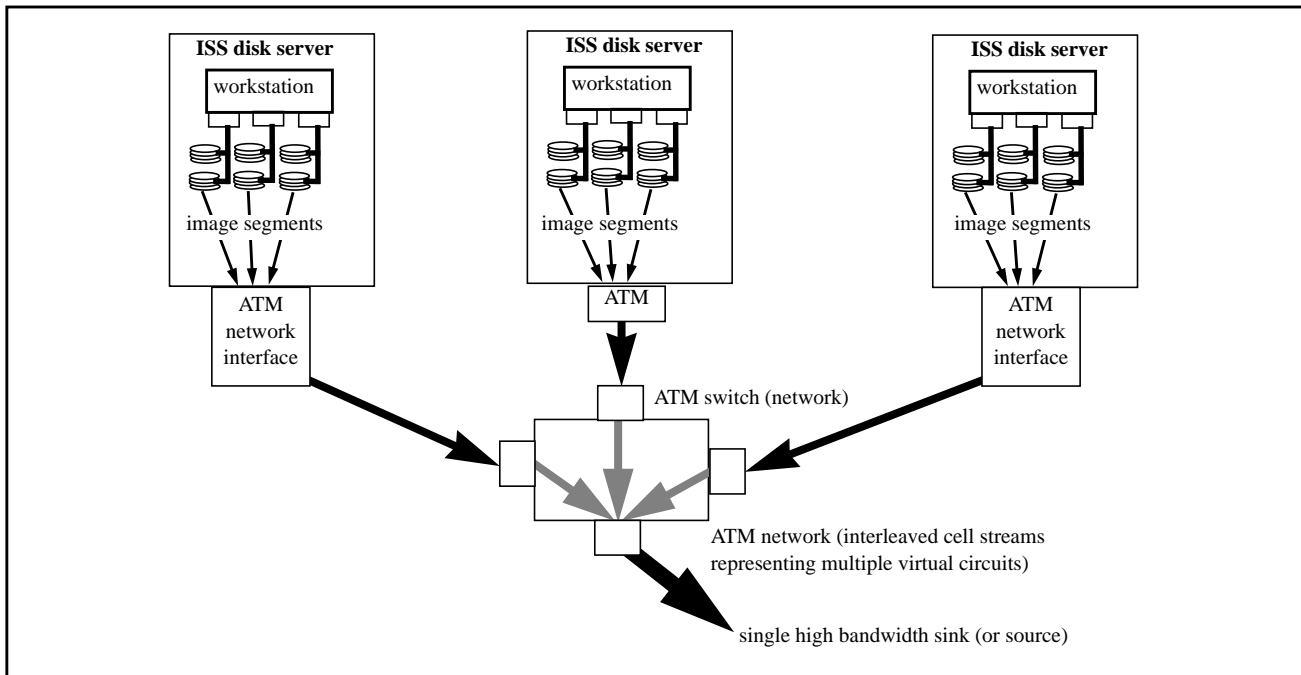


Figure 1: Parallel Data and Server Architecture Approach to the Image Server System

is distributed across a wide area network to supply data to applications located anywhere in the network. See Figure 1: *Parallel Data and Server Architecture Approach to the Image Server System*. There is no inherent organization to the blocks, and in particular, they would never be organized sequentially on a server. The data organization is determined by the application as a function of data type and access patterns, and is implemented during the data load process. The usual goal of the data organization is that data is declustered (dispersed in such a way that as many system elements as possible can operate simultaneously to satisfy a given request) across both disks and servers. This strategy allows a large collection of disks to seek in parallel, and all servers to send the resulting data to the application in parallel, enabling the ISS to perform as a high-speed image server.

The functional design strategy is to provide a high-speed “block” server, where a block is a unit of data request and storage. The ISS essentially provides only one function - it responds to requests for blocks. However, for greater efficiency and increased usability, we have attempted to identify a limited set of functions that extend the core ISS functionality while allowing support for a range of applications. First, the blocks are

“named.” In other words, the view from an application is that of a *logical* block server. Second, block requests are in the form of lists that are taken by the ISS to be in priority order. Therefore the ISS attempts (but does not guarantee) to return the higher priority blocks first. Third, the application interface provides the ability to ascertain certain configuration parameters (e.g., disk server names, performance, disk configuration, etc.) in order to permit parameterization of block placement-strategy algorithms (for example, see [1]). Fourth, the ISS is instrumented to permit monitoring of almost every aspect of its functioning during operation. This monitoring functionality is designed to facilitate performance tuning and network performance research; however, a data layout algorithm might use this facility to determine performance parameters.

At the present state of development and experience, the ISS that we describe here is used primarily as a large, fast “cache”. Reliability with respect to data corruption is provided only by the usual OS and disk mechanisms, and data delivery reliability of the overall system is a function of user-level strategies of data replication. The data of interest (tens to hundreds of GBytes) is typically loaded onto the ISS from archival tertiary storage, or written into the system from live video

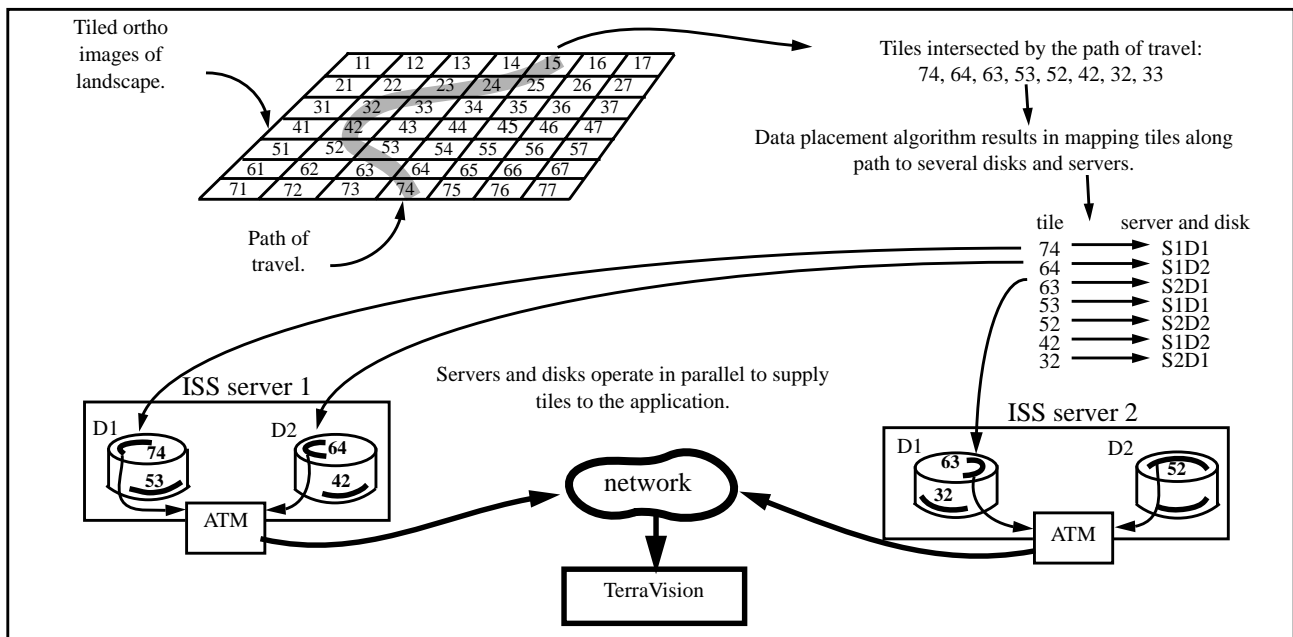


Figure 2: ISS Parallel Data Access Strategy as Illustrated by the TerraVision Application

sources. In the latter case, the data is also archived to bulk storage in real-time.

Client Use

The client-side (application) use of the ISS is provided through a library that handles initialization (for example, an “open” of a data set requires discovering all of the disk servers with which the application will have to communicate), and the basic block request / receive interface. It is the responsibility of the client (or its agent) to maintain information about any higher-level organization of the data blocks, to maintain sufficient local buffering so that “smooth playout” requirements may be met locally, and to run predictor algorithms that will pre-request blocks so that application response time requirements can be met. None of this has to be explicitly visible to the user-level application, but some agent in the client environment must deal with these issues, because the ISS always operates on a best-effort basis: if it did not deliver a requested block in the expected time or order, it was because it was not possible to do so.

Implementation

In our prototype implementations, the typical ISS consists of several (four - five) UNIX workstations (e.g. Sun SPARCStation, DEC Alpha, SGI Indigo, etc.), each with several (four - six) fast-SCSI disks on multiple (two - three) SCSI host adaptors. Each workstation is also equipped with an ATM network interface. An ISS configuration such as this can deliver an aggregated data stream to an application at about 400 Mbits/s (50 Mbytes/s) using these relatively low-cost, “off the shelf” components by exploiting the parallelism provided by approximately five servers, twenty disks, ten SCSI host adaptors, and five network interfaces.

Prototypes of the ISS have been built and operated in the MAGIC³ network testbed. In this paper we describe mainly architecture and approach, as well as optimization strategies. A previous paper [11] describes the major implementation issues, and a paper to be published [12] will describe other ISS applications and ISS performance issues.

2.0 Related Work

There are other research groups working on solving problems related to distributed storage and fast multimedia data retrieval. For example, Ghandeharizadeh, Ramos, et al., at USC are working on declustering methods for multimedia data [2], and Rowe, et al., at UCB are working on a continuous media player based on the MPEG standard [10].

In some respects, the ISS resembles the Zebra network file system, developed by John H. Hartman and John K. Ousterhout at the University of California, Berkeley [3]. Both the ISS and Zebra can separate their data access and management activities across several hosts on a network. Both try to maintain the availability of the system as a whole by building in some redundancy, allowing for the possibility that a disk or host might be unavailable at a critical time. The goal of both is to increase data throughput despite the current limits on both disk and host throughput.

However, the ISS and the Zebra network file system differ in the fundamental nature of the tasks they perform. Zebra is intended to provide traditional file system functionality, ensuring the consistency and correctness of a file system whose contents are changing from moment to moment. The ISS, on the other hand, tries to provide very high-speed, high-throughput access to a relatively static set of data. It is optimized to retrieve data, requiring only minimum overhead to verify data correctness and no overhead to compensate for corrupted data.

3.0 Applications

There are several target applications for the initial implementation of the ISS. These applications fall into two categories: image servers and multimedia / video file servers.

3. MAGIC (Multidimensional Applications and Gigabit Internetwork Consortium) is a gigabit network testbed that was established in June 1992 by the U. S. Government's Advanced Research Projects Agency (ARPA)[9]. MAGIC's charter is to develop a high-speed, wide-area networking testbed that will demonstrate interactive exchange of data at gigabit-per-second rates among multiple distributed servers and clients using a terrain visualization application. More information about MAGIC may be found on the WWW home page at: <http://www.magic.net/>

3.1 Image Server

The initial use of the ISS is to provide data to a terrain visualization application in the MAGIC testbed. This application, known as TerraVision [5], allows a user to navigate through and over a high resolution landscape represented by digital aerial images and elevation models. TerraVision is of interest to the U.S. Army because of its ability to let a commander “see” a battlefield environment. TerraVision is very different from a typical “flight simulator”-like program in that it uses high resolution aerial imagery for the visualization instead of simulated terrain. TerraVision requires large amounts of data, transferred at both bursty and steady rates. The ISS is used to supply image data at hundreds of Mbits/s rates to TerraVision. No data compression is used with this application because the bandwidth requirements are such that real-time decompression is not possible without using special purpose hardware.

In the case of a large-image browsing application like TerraVision, the strategy for using the ISS is straightforward: the image is tiled (broken into smaller, equal-sized pieces), and the tiles are scattered across the disks and servers of the ISS. The order of tiles delivered to the application is determined by the application predicting a “path” through the image (landscape), and requesting the tiles needed to supply a view along the path. The actual delivery order is a function of how quickly a given server can read the tiles from disk and send them over the network. Tiles will be delivered in roughly the requested order, but small variations from the requested order will occur. These variations must be accommodated by buffering, or other strategies, in the client application.

Figure 2: *ISS Parallel Data Access Strategy as Illustrated by the TerraVision Application* shows how image tiles needed by the TerraVision application are declustered across several disks and servers. More detail on this declustering is provided below.

Each ISS server is independently connected to the network, and each supplies an independent data stream into and through the network. These streams are formed into a single network flow by using ATM switches to combine the streams from multiple medium-speed links onto a single high-speed link. This high-speed link is ultimately

connected to a high-speed interface on the visualization platform (client). On the client, data is gathered from buffers and processed into the form needed to produce the user view of the landscape.

This approach could supply data to any sort of large-image browsing application, including applications for displaying large aerial-photo landscapes, satellite images, X-ray images, scanning microscope images, and so forth.

Figure 3: *Use of the ISS for Single High-Bandwidth App.* shows how the network is used to aggregate several medium-speed streams into one high-speed stream for the image browsing application. For the MAGIC TerraVision application,

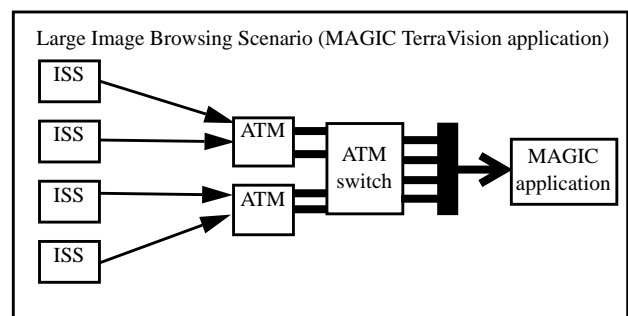


Figure 3: Use of the ISS for Single High-Bandwidth App.

the application host (an SGI Onyx) is using multiple OC-3 (155 Mbit/s) interfaces to achieve the bandwidth requirements necessary. These multiple interfaces will be replaced by a single OC-12 (622 Mbit/s) interface when it becomes available.

In the MAGIC testbed, the ISS has been run in several ATM WAN configurations to drive several different applications, including TerraVision. The configurations include placing ISS servers in Sioux Falls, South Dakota (EROS Data Center), Kansas City, Kansas (Sprint), and Lawrence, Kansas (University of Kansas), and running the TerraVision client at Fort Leavenworth, Kansas (U. S. Army’s Battle Command Battle Lab). The ISS disk server and the TerraVision application are separated by several hundred kilometers, the longest link being about 700 kilometers.

3.2 Video Server

Examples of video server applications include video players, video editors, and multimedia document browsers. A video server might contain several types of stream-like data, including con-

ventional video, compressed video, variable time base video, multimedia hypertext, interactive video, and others. Several users would typically be accessing the same video data at the same time, but would be viewing different streams, and different frames in the same stream. In this case the ISS and the network are effectively being used to “reorder” segments (see Figure 4: *Use of the ISS to Supply many Low-Bandwidth Streams*). This reordering affects many factors in an image

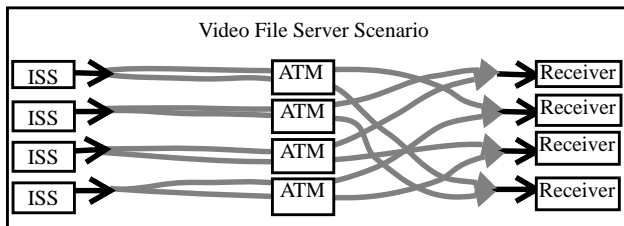


Figure 4: Use of the ISS to Supply many Low-Bandwidth Streams

server system, including the layout of the data on disks. Commercial concerns such as Time Warner and U.S. West are building large-scale commercial video servers such as the Time Warner / Silicon Graphics video server [4]. Because of the relatively low cost and ease of scalability of our approach, it may address a wider scale, as well as a greater diversity, of data organization strategies so as to serve the diverse needs of schools, research institutions, and hospitals for video-image servers in support of various educational and research-oriented digital libraries.

4.0 Design

4.1 Goals

The following are some of our goals in designing the ISS:

- The ISS should be capable of being geographically distributed. In a future environment of large scale, high-speed, mesh-connected national networks, network distributed storage should be capable of providing an uninterrupted stream of data, in much the same way that a power grid is resilient in the face of source failures, and tolerant of peak demands, because of the possibility of multiple sources multiply interconnected.

- The ISS approach should be scalable in all dimensions, including data set size, number of users, number of server sites, and aggregate data delivery speed.
- The ISS should deliver coherent image streams to an application, given that the individual images that make up the stream are scattered (by design) all over the network. In this case, “coherent” means “in the order needed by the application”. No one disk server will ever be capable of delivering the entire stream. The network is the *server*.
- The ISS should be affordable. While something like a HIPPI-based RAID device might be able to provide functionality similar to the ISS, this sort of device is very expensive, is not scalable, and is a single point of failure.

4.2 Approach

A Distributed, Parallel Server

The ISS design is based on the use of multiple low-cost, medium-speed disk servers which use the network to aggregate server output. To achieve high performance we exploit all possible levels of parallelism, including that available at the level of the disks, controllers, processors / memory banks, servers, and the network. Proper data placement strategy is also key to exploiting system parallelism.

At the server level, the approach is that of a collection of disk managers that move requested data from disk to memory cache. Depending on the nature of the data and its organization, the disk managers may have a strategy for moving other closely located and related data from disk to memory. However, in general, we have tried to keep the implementation of data prediction (determining what data will be needed in the near future) separate from the basic data-moving function of the server. Prediction might be done by the application (as it is in TerraVision), or it might be done by a third party that understands the data usage patterns. In any event, the server sees only lists of requested blocks.

As explained in [12], the dominant bottlenecks for this type of application in a typical UNIX workstation are first memory copy speed, and second, network access speed. For these reasons, an important design criterion is to use as few memory copies as possible, and to keep the network interface operating at full bandwidth all the time. Our implementation uses only three copies to get data from disk to network, so maximum server throughput is about $(\text{memory_copy_speed} / 3)$.

Another important aspect of the design is that all components are instrumented for timing and data flow monitoring in order to characterize ISS and network performance. To do this, all communications between ISS components are time-stamped. In the MAGIC testbed, we are using GPS (Global Positioning System) receivers and NTP (Network Time Protocol) [6] to synchronize the clocks of all ISS servers and of the client application in order to accurately measure network throughput and latency.

Data Placement Issues

A limiting factor in handling large data sets is the long delay in managing and accessing subsets of these data sets. Slow I/O rates, rather than processor speed, are chiefly the cause of this delay. One way to address this problem is to use data reorganization techniques based on the application's view of the structure of the data, analysis of data access patterns, and storage device characteristics. By matching the data set organization with the intended use of the data, substantial improvements can be achieved for common patterns of data access[1]. This technique has been applied to large climate-modeling data sets, and we are applying it to TerraVision data stored in the ISS. For image tile data, the placement algorithm declusters tiles so that all disks are evenly accessed by tile requests, but then clusters tiles that are on the same disk based on the tiles' relative nearness to one another in the image. This strategy is a function of both the data structure (tiled images) and the geometry of the access (e.g., paths through the landscape).

The declustering method used for tiles of large images is a lattice-based (i.e., vector-based) declustering scheme, the goal of which is to ensure tiles assigned to the same server are as far

apart as possible on the image plane. This minimizes the chance that the same server will be accessed many times by a single tile request list.

Tiles are distributed among K disks by first determining a pair of integer component vectors which span a parallelogram of area K . Tiles assigned to the same disk are separated by integer multiples of these vectors. Mathematical analysis shows that for common visualization queries this declustering method performs within seven percent of optimal for a wide range of practical multiple disk configurations.

Within a disk, however, it is necessary to cluster the tiles such that tiles near each other in 2-D space are close to each other on disk, thus minimizing disk seek time. The clustering method used here is based on the Hilbert Curve because it has been shown to be the best curve that preserves the 2-D locality of points in a 1-D traversal.

Path Prediction

Path prediction is important to ensure that the ISS is utilized as efficiently as possible. By using a strategy that always requests more tiles than the ISS can actually deliver before the next tile request, we can ensure that no component of the ISS is ever idle. For example, if most of a request list's tiles were on one server, the other servers could still be reading and sending or caching tiles that may be needed in the future, instead of idly waiting. The goal of path prediction is to provide a rational basis for pre-requesting tiles. See [1] for more details on data placement methods.

As a simple example of path prediction, con-

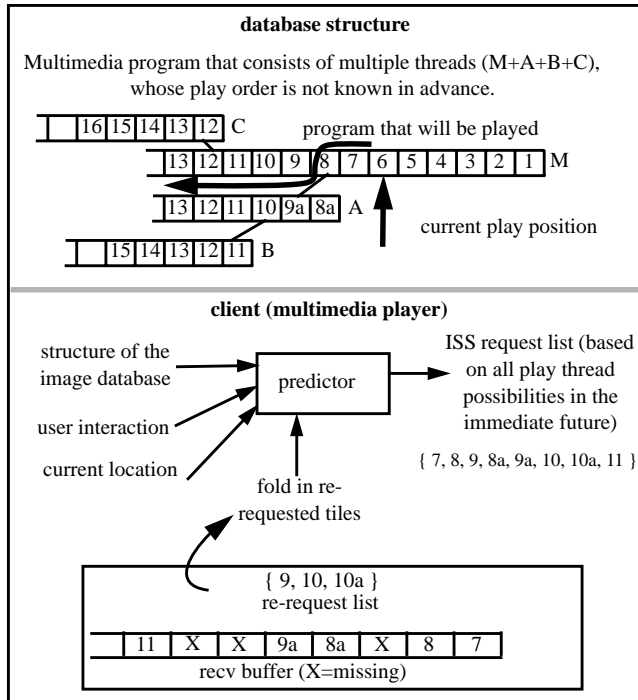


Figure 5: Image Stream Management / Prediction Strategy

sider an interactive video database with a finite number of distinct paths (video clips), and therefore a finite number of possible branch points. (A “branch point” occurs where a user might select one of several possible play clips, see Figure 5: *Image Stream Management / Prediction Strategy*) As a branch point is approached by the player, the predictor (without knowledge of which branch will be taken) will start requesting images (frames) along both branches. These images are cached first at the disk servers, then at the receiving application. As soon as a branch is chosen, the predictor ceases to send requests for images from the other branches. Any “images” (i.e., frames or compressed segments) cached on the ISS, but unsent, are flushed as better predictions fill the cache. This is an example where a relatively independent third party might do the prediction.

The client will keep asking for an image until it shows up, or until it is no longer needed (e.g., in TerraVision, the application may have “passed” the region of landscape that involves the image that was requested, but never received.) Applications will have different strategies to deal with images that do not arrive in time. For exam-

ple, TerraVision keeps a local, low-resolution, data set to fill in for missing tiles.

Prediction is transparent to the ISS, and is manifested only in the order and priority of images in the request list. The prediction algorithm is a function of the client application, and typically runs on the client.

The Significance of ATM Networks

The design of the ISS depends in part on the ability of ATM switches and networks to aggregate multiple data streams from the disk servers into a single high-bandwidth stream to the application. This is feasible because most wide area ATM network aggregate bandwidth upward - that is, the link speeds tend to increase from LANs to WANs, and even within WANs the “backbone” is the highest bandwidth. (This is actually a characteristic of the architecture of the SONET networks that underlie ATM networks.) Aggregation of stream bandwidth occurs at switch output ports. For example, three incoming streams of 50 Mbits/s that are all destined for the same client will aggregate to a 150 Mbit/s stream at the switch output port. The client has data stream connections open to each of the ISS disk servers, and the incoming data from all of these streams typically put data into the same buffer.

5.0 Implementation

In a typical example of ISS operation the application sends requests for data (images, video, sound, etc.) to the name server process which does a lookup to determine the location (server/disk/offset) of the requested data. Requests are sorted on a per-server basis, and the resulting lists are sent to the individual servers. Each server then checks to see if the data is already in its cache, and if not, fetches the data from disk and transfers it to the cache. Once the data is in the cache, it is sent to the requesting application. Figure 6: *ISS Architecture* shows how the components of the ISS are used to handle requests for data blocks.

The disk server handles three image request priority levels:

- high: send first, with an implicit priority given by order within the list.
- medium: send if there is time.

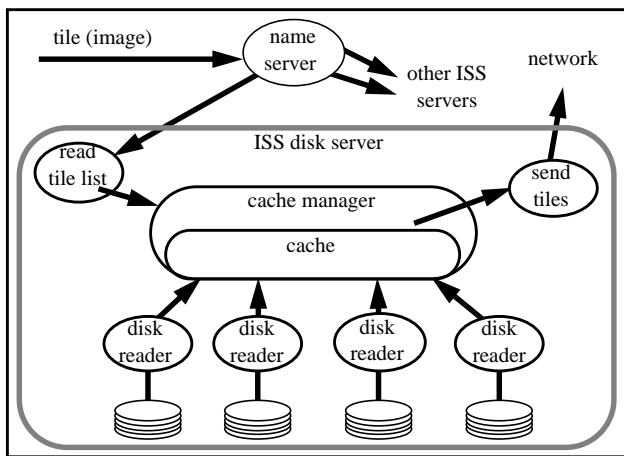


Figure 6: ISS Architecture

- low: fetch into the cache if there is time, but don't send.

The priority of a particular request is set by the requesting application. The application's prediction algorithm can use these priority levels to keep the ISS fully utilized at all times without requesting more data than the application can process. For example, the application could send low priority requests to pull data into the ISS cache, knowing that the ISS would not send the data on to the application until the application was ready. Another example is an application that plays back a movie with a sound track, where audio might be high priority requests, and video medium priority requests.

5.1 Performance Limits

Using a Sun SPARCStation 10-41 with two Fast-SCSI host adaptors and four disks, and reading into memory random 48 Kbyte tiles from all disks simultaneously, we have measured a single server disk-to-memory throughput of 9 Mbytes/s. When we add a process which sends UDP packets to the ATM interface, this reduces the disk-to-memory throughput to 8 Mbytes/s (64 Mbits/s). The network throughput under these conditions is 7.5 Mbytes/s (60 Mbits/s). This number is an upper limit on performance for this platform; it does not include the ISS overhead of buffer management, semaphore locks, and context switching. The SCSI host adaptor and Sbus are not yet saturated, but adding more disks will not help the overall throughput without faster access memory and to the network (e.g., multiple interfaces and multiple independent data paths as are used in

systems like a SPARCServer 1000 or SGI Challenge).

6.0 Current Status

All ISS software is currently tested and running on Sun workstations (SPARCstations and SPARCserver 1000's) running SunOS 4.1.3 and Solaris 2.3, DEC Alpha's running OSF/1, and SGI's running IRIX 5.x. Demonstrations of the ISS with the MAGIC Terrain Visualization application TerraVision have been done using several WAN configurations in the MAGIC testbed [9]. Using enough disks (4-8, depending on the disk and system type), the ISS software has no difficulty saturating current ATM interface cards. We have worked with 100 Mbit and 140Mbit TAXI S-Bus and VME cards from Fore systems, and OC-3 (155 Mbit/s) cards from DEC, and in all cases ISS throughput is only slightly less than $ttcp^4$ speeds.

Table 1 below shows various system $ttcp$ speeds and ISS speeds. The first column is the

TABLE 1.

System	Max ATM LAN $ttcp$	$ttcp$ w/disk read	Max ISS speed
Sun SS10-51	70 Mbits/sec	60 Mbits/sec	55 Mbits/sec
Sun SS1000 (2 processors)	75 Mbits/sec	65 Mbits/sec	60 Mbits/sec
SGI Challenge L (2 processors)	82 Mbits/sec	72 Mbits/sec	65 Mbits/sec
Dec Alpha	127 Mbits/sec	95 Mbits/sec	88 Mbits/sec

maximum $ttcp$ speeds using TCP over a ATM LAN with a large TCP window size. In this case, $ttcp$ just copies data from memory to the network. For the values in the second column, we ran a program that continuously reads from all ISS disks simultaneously with $ttcp$ operation. This gives us a much more realistic value for what network speeds the system is capable of while the ISS is running. The last column is the actual throughput values measured from the ISS. These speeds indicate that the ISS software adds a rela-

4. $ttcp$ is a utility that times the transmission and reception of data between two systems using the UDP or TCP protocols.

tively small overhead in terms of maximum throughput.

6.1 Actual Performance

The current throughput of a single ISS server on a Sun SPARC 10/41 platform is 7.1 Mbytes/s (55 Mbits/s), or 91% of the possible maximum of 7.5 Mbytes/s (60 Mbits/s) derived above. This seems a reasonable result considering the overhead required. We have achieved this speed using a TerraVision-like application simulator which we developed that sends a list of requests for data at a rate of five request lists per second. Five request lists per second does not force the application to predict and buffer too far into the future, but is not so fast that disk read latency is an issue. This application simulator sends request lists that are long enough to ensure that no disk ever is idle. When the ISS receives a request list, all previous requests are discarded. Under these conditions, about one-half of the requests in each request list will never be satisfied (either they will be read into the cache but not written to the network, or they will not be read at all before the next request list arrives).

As an example, a typical TerraVision request list contains fifty tiles. Of these fifty tiles, forty are read into ISS cache, twenty-five are written to the network, and ten are not processed at all. This behavior is reasonable because, as discussed in the section on data path prediction above, the application will keep asking for data until it shows up or is no longer needed. The requesting application will anticipate this behavior, and predict the tiles it needs far enough ahead that “important” tiles are always received by the time they are needed. Tiles are kept in the cache on an LRU basis, and previously requested but unsent tiles will be found in the cache by a subsequent request. The overhead of re-requesting tiles is minimal compared with moving them from disk and sending them over the network.

During ISS operation, the average CPU usage on the disk server platform is 10% user, 60% system, 30% idle, so the CPU is not a bottleneck. With the TerraVision application and 40 Mbyte of disk cache memory on the ISS server, on average 2% of tiles are already in cache relative to any given request. Increasing the cache size will not

increase the throughput, but may improve latency with effective path prediction by the application.

7.0 Future Work

We plan to expand the capabilities of the ISS considerably during the next year or so. These enhancements (and associated investigation of the issues) will include:

- Implementing a multiple data set data layout strategy;
- Implementing a multi-user data layout and access strategy;
- Implementing a capability to write data to the ISS;
- Implementing the ability to monitor the state of all ISS servers and dynamically assign bandwidth of individual servers to avoid overloading the capacity of a given segment of the network (i.e., switches or application host);
- Implementing mechanisms for handling video-like data, including video data placement algorithms and the ability to handle variable size frames (JPEG/MPEG);
- Modifying name server design to accommodate data on server performance and availability and to provide a mechanism to request tiles from the “best” server (fastest or least loaded);
- Investigating the issues involved in dealing with data other than image- or video- like data.

Many of these enhancements will involve extensions to the data placement algorithm and the cache management methods. Also we plan to explore some optimization techniques, including using larger disk reads, and conversion of all buffer and device management processes to threads-based light weight processes.

8.0 References

- [1] Chen L. T. and Rotem D., “Declustering Objects for Visualization”, Proc. of the 19th VLDB (Very Large Database) Conference, 1993.

- [2] Ghandeharizadeh, S. and Ramos, L., "Continuous Retrieval of Multimedia Data Using Parallelism, IEEE Transactions on Knowledge and Data Engineering, Vol 5, No 4, August 1993.
- [3] Hartman, J. H. and Ousterhout, J. K., "Zebra: A Striped Network File System", Proceedings of the USENIX Workshop on File Systems, May 1992.
- [4] Langberg, M., "Silicon Graphics Lands Cable Deal with Time Warner Inc.", San Jose Mercury News, June 8, 1993.
- [5] Leclerc, Y.G. and Lau, S.Q., Jr., "TerraVision: A Terrain Visualization System", SRI International, Technical Note #540, Menlo Park, CA, 1994.
- [6] Mills, D., "Simple Network Time Protocol (SNTP)", RFC 1361, University of Delaware, August 1992.
- [7] Parvin, B., Peng, C., Johnston, W., and Maestre, M., "Tracking of Tubular Objects for Scientific Applications", IEEE Conf. on Computer Vision and Pattern Recognition, June 1994, pp. 295-301
- [8] Patterson, D., Gibson, R., and Katz, R., "The Case for RAID: Redundant Arrays of Inexpensive Disks", Proceedings ACM SIGMOD Conference, Chicago, IL, May, 1988 (pp. 106-113)
- [9] Richer, I and Fuller, B.B, "An Overview of the MAGIC Project," M93B0000173, The MITRE Corp., Bedford, MA, 1 Dec. 1993.
- [10] Rowe L. and Smith B.C, "A Continuous Media Player", Proc. 3rd International Workshop on Network and Operating System Support for Digital Audio and Video, San Diego, CA, Nov. 1992.
- [11] Tierney, B., Johnston, W., Herzog, H., Hoo, G., Jin, G., Lee, J., "System Issues in Implementing High Speed Distributed Parallel Storage Systems", Proceedings of the USENIX Symposium on High Speed Networking, Aug. 1994, LBL-35775.
- [12] Tierney, B., Johnston, W., Chen, L.T., Herzog, H., Hoo, G., Jin, G., Lee, J., "Using High Speed Networks to Enable Distributed Parallel Image Server Systems", Proceedings of Supercomputing '94, Nov. 1994, LBL-35437.