

Distributed Environments for Large Data-Objects: The Use of Public ATM Networks for Health Care Imaging Information Systems¹

William Johnston, Jin Guojun, Jason Lee, Brian Tierney, Mary Thompson

Imaging and Distributed Computing Group
Lawrence Berkeley National Laboratory
University of California
Berkeley, CA, 94720

Abstract

We are exploring the feasibility of using highly distributed architectures for all aspects of collecting, storing, analyzing, and otherwise manipulating and making generally available, large data-objects. These objects - typically the result of a single operational cycle of an instrument, and of sizes from tens of MBytes to tens of Gbytes - are the staple of modern analytical systems. It is further the case that many of the instrumentation systems that generate such data-objects are used by a diverse and geographically distributed community: Examples from the scientific community include physics and nuclear science high energy particle accelerators and detector systems, large electron microscopes, ultra-high brilliance X-ray sources, etc. There are correspondingly complex instrumentation systems in the health care community that generate large data-objects.

Apart from the basic issues of data collection, transport, and processing, the project is exploring the general issues of using the network to provide flexible data storage strategies and location independent access to analysis systems.

In this paper, we describe some of the general aspects of such systems, and specifically describe the use of a shared (public) IP over ATM network to facilitate collection, storage, analysis, distribution, and delivery of several kinds of health care imaging data. The data includes still images and video sequences of coronary angiograms, most of which are "large data-objects". The data is collected by an instrumentation system in a hospital in one city, it is then sent over a high speed metropolitan area network and stored in a distributed storage systems at other locations, and subsequently made available to physicians at still other locations.

The health care large data-object project is characterizing and addressing various issues, including the ability of the ATM network to deliver the data end-to-end given certain operational constraints, the capabilities needed to ensure the security of the data, the ability of state-of-the-art computing systems and software that are interconnected by ATM networks to perform at the levels required for routine clinical use, etc.

Within this scenario, we discuss our experiences with the use and behavior of TCP/IP - based data transfer over metropolitan and wide area ATM networks; the impact of the relatively high data-rate use on other users of the network, and the integration issues for a relatively transparent security architecture to protect the confidentiality of the medical data in public networks.

1. This work is supported by the U. S. Dept. of Energy, Energy Research Division, Mathematical, Information, and Computational Sciences office (<http://www.er.doe.gov/production/octr/mics>), under contract DE-AC03-76SF00098 with the University of California, and by ARPA ITO (<http://ftp.arpa.mil/ResearchAreas.html>). Author's address: 50B-2239, Lawrence Berkeley National Laboratory, Berkeley, CA 94720. Tel: +1-510-486-5014, fax: +1-510-486-6363, wejohnston@lbl.gov, <http://www-itg.lbl.gov/~johnston>. This document is report LBNL-nnnn.

Contents

1.0 Introduction	3
2.0 The Health Care Information System Application	3
3.0 System Architecture	8
3.1 Data Collection	8
3.2 Distributed Storage: The Distributed-Parallel Storage System	8
3.2.1 DPSS Architecture	9
3.2.2 DPSS and the TerraVision Application	11
3.3 Data Management	14
3.3.1 The Health Care Application	15
3.4 Performance	16
3.4.1 Performance Issues Related to ATM	16
3.4.2 ATM Network to Workstation Performance	16
3.4.3 ATM Network Performance	18
3.5 Security Architectures for Large-Scale Remote Environments	20
3.5.1 The Problem and Objectives	20
3.5.2 Background	21
3.5.3 The Security Model	21
3.5.4 Implementation Issues	23
3.5.5 Relationship to Other Work	24
3.5.6 Application of the Security Architecture	25
4.0 Conclusions	25
5.0 Acknowledgments	25
6.0 Notes and References	26

1.0 Introduction

The advent of shared, widely available high speed networks is providing the potential for new approaches to the collection, storage, and analysis of large data-objects. Health care information, especially high volume image data used for diagnostic purposes - e.g. X-ray CT, MRI, and cardio-angiography - that are increasingly collected at tertiary (centralized) facilities, may now be routinely stored and used at locations other than the point of collection. The importance of distributed storage is that a hospital (or any other instrumentation environment) may not be the best environment in which to maintain a large-scale digital storage system, and an affordable, easily available, high bandwidth network can provide location independence for such storage. The importance of remote, end-user access is that the health care professionals at the referring facility (frequently remote from the tertiary imaging facility) will have ready access to not only the image analyst's reports, but the original image data itself.

This general scenario extends to many fields other than health care. In particular, the same basic infrastructure is required for remote access to large-scale scientific and analytical instruments, both for data handling and for direct, remote-user operation. See [Johnston95b].

In this paper we describe and illustrate a set of concepts that are contributing to a generalized, distributed information infrastructure, especially as it concerns the types of large data-objects generated in the scientific and medical environments. We will describe the general issues, architecture, and some system components that are currently in use to support distributed large data-objects, and will use a specific health care information system as an example.

The basic elements of the overall architecture include:

- data collection and the instrument-network interface
- on-line storage that is distributed throughout the network (for both performance and reliability)
- processing elements - also distributed throughout the network - for various sorts of data analysis
- data management that provides for the automatic cataloguing (metadata generation) of the data being stored
- data access interfaces, including application-data interface
- tertiary storage ("mass storage") management
- user access to all relevant aspects (application, data, metadata, data management)
- transparent security that provides access control for all of the systems components based on the resource-owner's policy

These elements all need to be provided with flexible, location-independent interfaces so that they can be freely moved around the network as required for operational or other logistical convenience.

2.0 The Health Care Information System Application

An example of a medical application that uses a distributed large data-object architecture is a system that provides for collection, storage, cataloguing, and playback of video-angiography images in a metropolitan area ATM network.

Cardio-angiography² is used to monitor and restore coronary blood flow, and though clinically effective, the required imaging systems and associated facilities

are expensive. To minimize the cost of such procedures, health care providers are beginning to concentrate these services in a few high-volume tertiary care centers. Patients are typically referred to these centers by cardiologists operating at clinics or other hospitals; the centers then must communicate the results back to the local cardiologists as soon as possible after the procedure. The advantages of providing specialized services at distant tertiary centers are significantly reduced if the medical information obtained during the procedure is not delivered rapidly and accurately to the referring physician at the patient's home facility. The delivery systems currently used to transfer patient information between facilities include interoffice mail, U.S. Mail, fax machine, telephone, and courier. Often these systems are inadequate and potentially could introduce delays in patient care. (See [Kaiser-CalREN].)

Using a shared, metropolitan area, ATM network, and a high-speed distributed data handling system, video sequences and still image are collected from the video-angiography imaging systems, stored, and accessed by a remote user. The image data are sent through the network to storage and analysis systems, as well as directly to the users at clinic sites. Thus, data can be stored and catalogued for later use, data can be delivered live from the imaging device to remote clinics in real-time, or these data flows can all be done simultaneously. Whether the storage servers are local or distributed around the network is entirely a function of the optimal logistics: There are arguments in regional health care information systems for centralized storage facilities away from the hospital environment, even though the architecture is that of a distributed system. (See [Johnston93].)

Figure 1 illustrates the configuration of the health care application. This application is a joint project of Lawrence Berkeley National Laboratory, Kaiser Permanente, Philips Research (Palo Alto), and the Pacific Bell CalREN program. (See [Kaiser-CalREN].)

The network aspects of this health care project are closely related to the BAGNet CalREN project (see [BAGnet-1] and [Wiltzius]). BAGNet was the first, and largest, of the CalREN ATM projects, and it established the basic TCP/IP over ATM infrastructure used in this and several other CalREN projects. The physical infrastructure is illustrated in Figure 2, and the logical and network infrastructure are described in [Johnston95c].

Figure 3 illustrates the health care application data flow, and once again, though we are describing a specific application, this data flow architecture is quite general, and we are using an almost identical architecture in several other projects involving the collection, cataloguing, storage, and access of data from several types of (unrelated) scientific experiment scenarios.

The angiography data is being collected directly from a Philips scanner by a computer system in the San Francisco Kaiser Cardiac Catheterization Laboratory. This computing system, in turn, is attached to a Pacific Bell, OC-3, metropolitan area ATM network. When the data collection for a patient is complete (about once every 20-40 minutes) between 10 and 50 Gigabytes of digital video data is sent across the network to LBNL (in Berkeley) and stored first on an experimental, high speed, network-based multimedia storage system (the DPSS) and then is archived to a mass storage system.

2. Cardio-angiography imaging involves a two plane, X-ray video imaging system that produces from several to tens of minutes of digital video sequences for each patient study for each patient session. The digital video is organized as tens of data-objects, each of which are of the order of 100 MBytes.

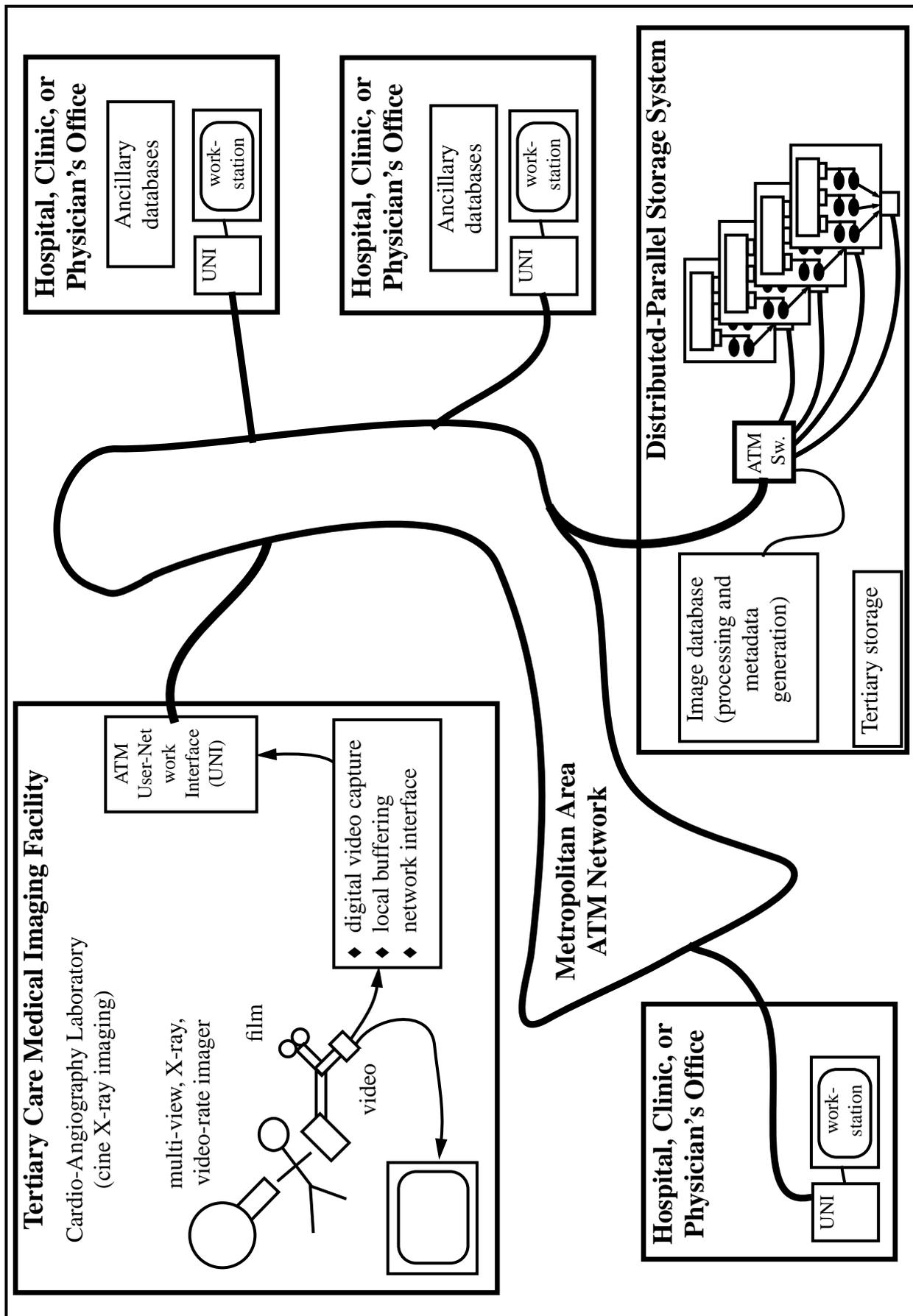


Figure 1 Distributed Large Data-Object Medical Imaging Application

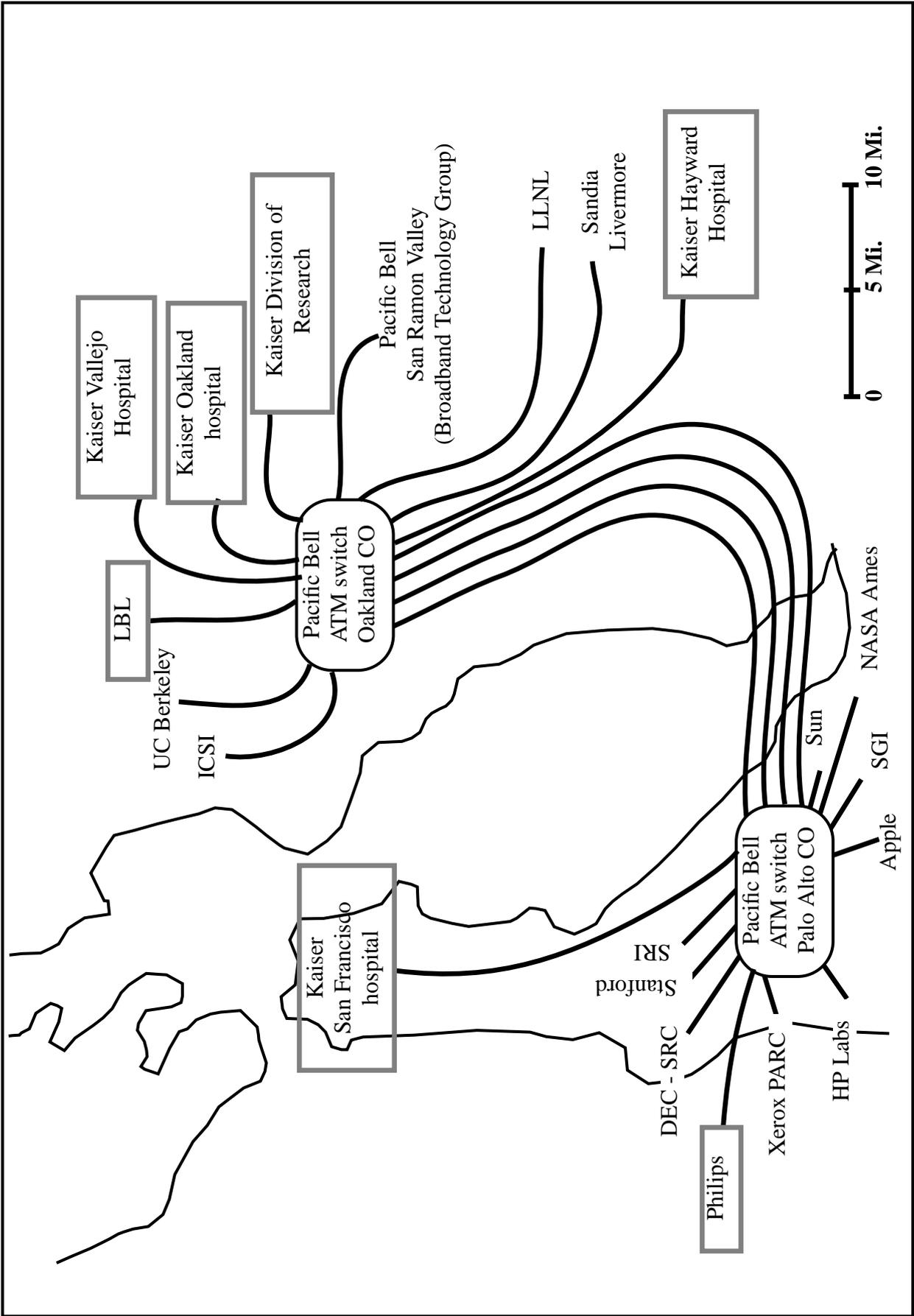


Figure 2 The Pacific Bell, S. F. Bay Area, Metropolitan Area ATM Network: BAGNet and Kaiser CalREN Sites

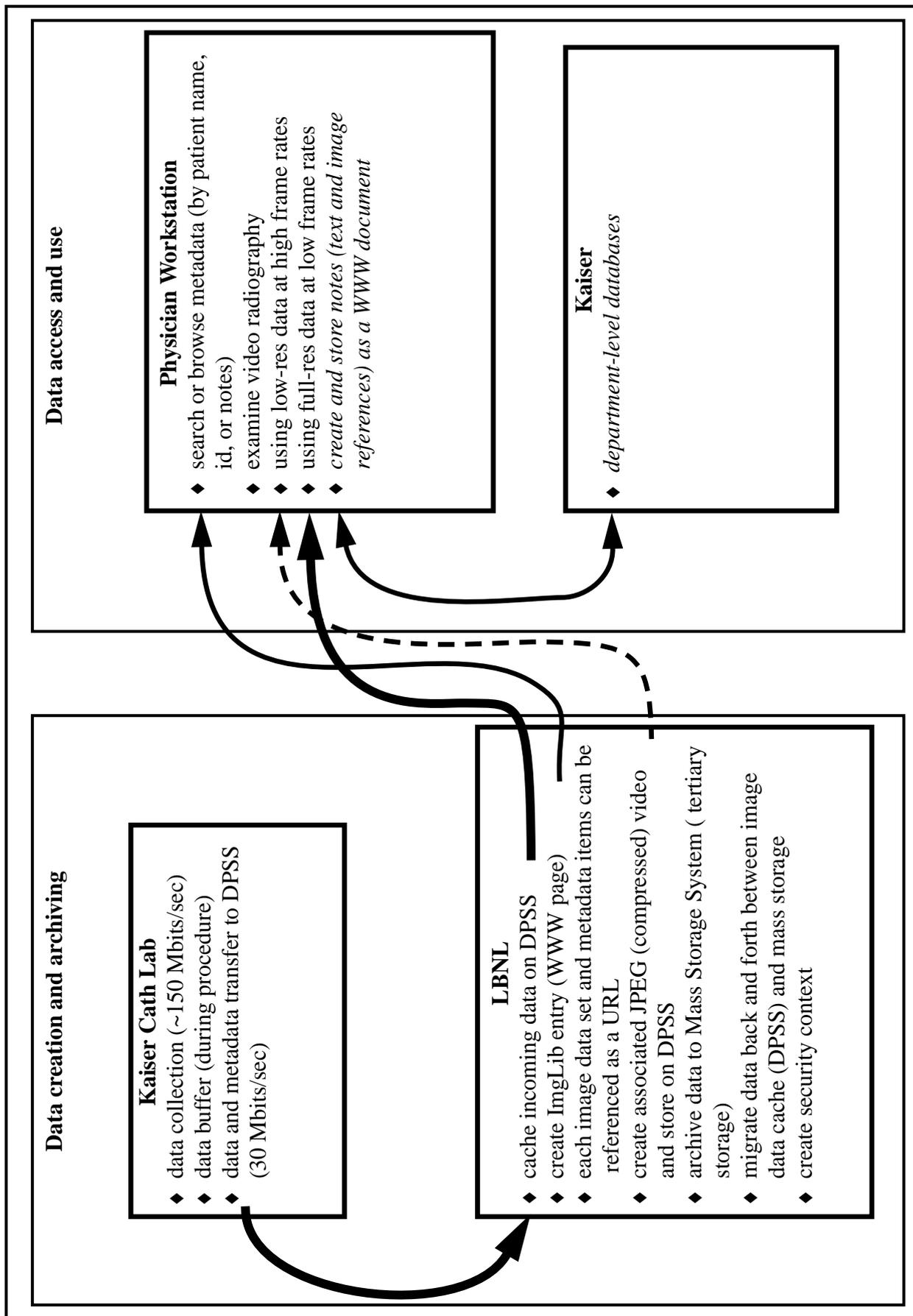


Figure 3 Kaiser Project Data Flow (A typical data flow for distributed large data-object systems)

The data is catalogued in a World Wide Web based image database system (“ImgLib”). Use of the data involves access through ImgLib to locate the data sets of interest, and then to invoke a viewing application.

Department-level Web-based patient databases can refer directly to the data in ImgLib without duplicating the data, or being concerned about tertiary storage management (which is handled by ImgLib).

3.0 System Architecture

3.1 Data Collection

The front-end of most of the distributed large data-object environments that we are working with are instrumentation systems. Examples include particle accelerator detectors and Earth environment monitoring satellites. These sources generate essentially continuous data streams, but ones that have “natural” boundaries that define “objects”: Video sources (like the cardio-angiography studies where data is collected in bursts; large single images, etc. We make this “natural boundary” characterization to distinguish these data types from continuous-media data (such as video and audio multimedia streams): Even though the average data rates may be comparable, or even higher in the large data-object environment, the temporal synchronization requirements and steady-state nature of multimedia streams give them a somewhat different character.

For many of the instrumentation systems that we are interested in, one of the primary issues is getting the data out of the instruments and on to the network. One of the circumstances that has led to both the interest in, and practicality of, the architecture being described here is that general purpose workstations now have memory and I/O bus structures that are fast enough to acquire, structure, and send to the network, significant bandwidth data streams. This is an important capability because it means that the only special hardware that is required to bring instruments on-line is the interface between the workstation and the data source, and even this interface may be able to be done in “software” using off-the-shelf DSP-based I/O boards.

The frontend workstation acquires the raw data, formats it as “objects” by adding or using metadata from the experiment environment, and then sends the objects into the distributed environment. The data collection workstation frequently also serves as a buffer so that brief interpretations or slow-downs in the network do not result in losing data. This frontend architecture is illustrated in Figure 4.

3.2 Distributed Storage: The Distributed-Parallel Storage System

Widely distributed network storage systems are an essential component of a network-based large data-object environment.

The idea of such systems is that by distributing the components of a storage system throughout the network, that capacity, reliability, performance, and security are all increased. Capacity increases in conjunction with a widely deployed, generalized security infrastructure that can support dynamic brokering for resources. (See [Johnston96a].) Reliability increases because storage systems that can be dynamically configured from components that have as little as possible in common (e.g. location) provide the resilience that comes from independence. Performance is increased by the combined characteristics of parallel operation of many sub components, and the independent data paths provided by a large network infrastructure. Security

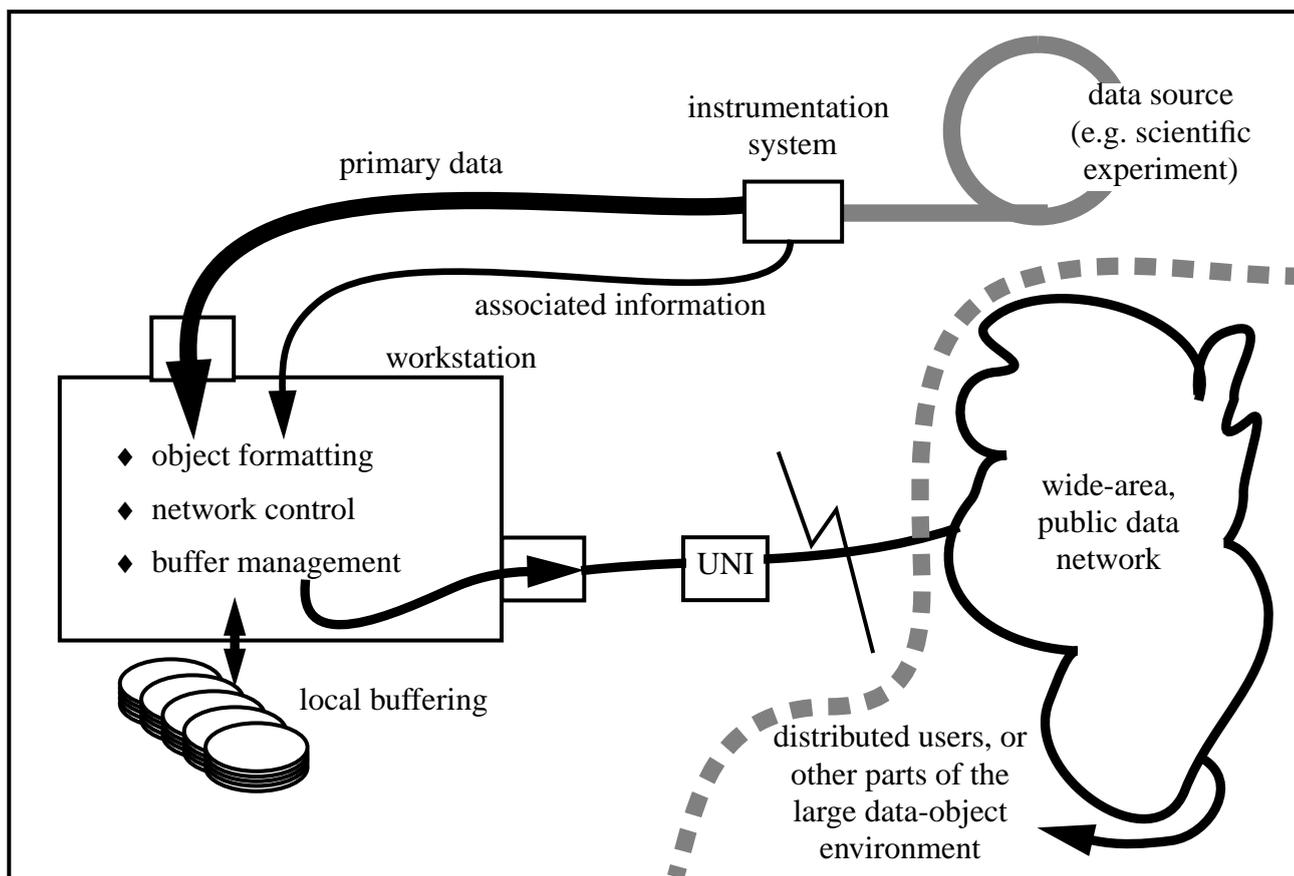


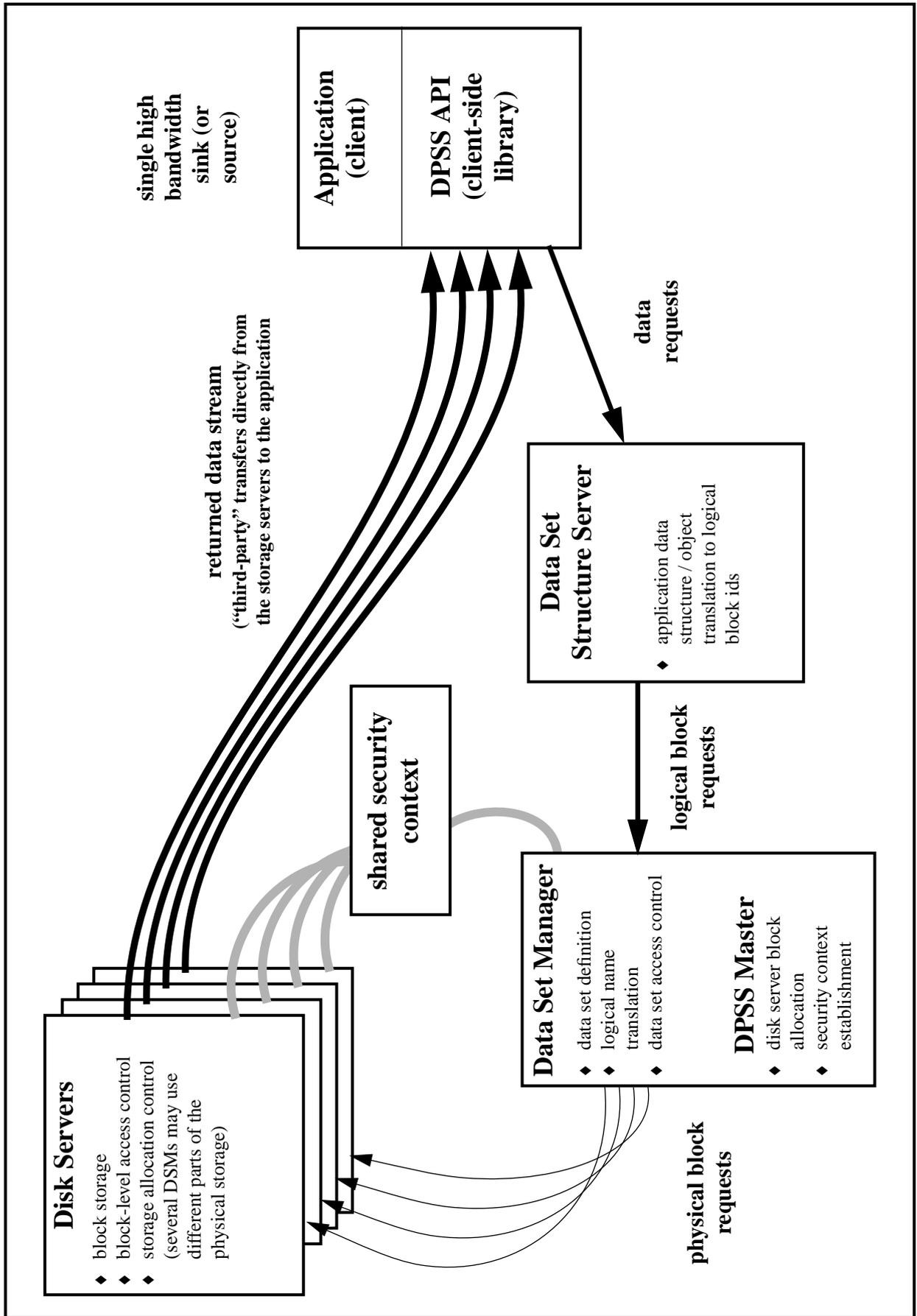
Figure 4 Architecture of the Frontend of a Distributed, Large Data-Object Environment

is also potentially increased by having many independent components each of which has local and independent enforcement mechanisms that can limit the scope of a security breach.

The Distributed-Parallel Storage System (DPSS, aka ISS) is an experimental system in which we are developing, implementing, and testing these ideas. In most configurations, it is used as a network-striped disk array designed to supply and consume high speed data streams to and from other processes in the network. (See [DPSS].)

3.2.1 DPSS Architecture

The DPSS is essentially a “logical block” server whose functional components are distributed across a wide area network. (See Figure 5 illustrating the DPSS architecture.) The DPSS uses parallel operation of distributed servers to supply image streams fast enough to enable various multi-user, “real-time”, virtual reality-like applications in an Internet / ATM environment. 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 DPSS to perform as a high-speed image server.



Distributed-Parallel Storage System Architecture

Figure 5

The implementation is based on the use of multiple low-cost, medium-speed disk servers which use the network to aggregate multiple server outputs for high performance applications. To achieve high performance all types of parallelism are exploited, including those available at the level of the disks, controllers, processors / memory banks, servers, and the network (See Figure 6 illustrating the DPSS implementation).

At the application level, the DPSS is a persistent cache of named objects, and at the storage level it is a logical block server. Although not strictly part of the DPSS architecture, the system is usually used with an application agent called a “data set structure server”. This component provides an object-like encapsulation of complex user-level data structures so that the application does not have to retain this information for each different data set. The function and interface of this structure server are left to the application domain, but one simple example is for video data. In this case the structure server allows applications to request data by “frame” number and “type” (e.g. in the video-angiography application, there are two images - 90 degree offset views - associated with every frame). The structure server converts the application requests into logical block requests. These logical block requests are then sent to the DPSS “data set manager”. The DSM maintains data set definitions, and is responsible for mapping the logical block requests to physical block requests. The DPSS “master” deals with interactions with the servers to determine available storage (a server is an independent entity and may deal with several DPSS masters) and to establish the “security context”.

The security model for the DPSS involves accommodating several different resource owners: The context established between the Master and the disk/storage server reflects an agreement between the owner of the physical resources (disks) and an agent that is providing storage to a user community. This context enforces the disk usage agreement. The context established between the DSM and the storage servers reflects the use-conditions imposed by the data “owner”, and provides for ensuring access control that enforces those use-conditions.

The overall data flow involves “third-party” transfers from the storage servers directly to the data-consuming application. (A model used by most high performance storage systems.) Thus, the application requests data, these requests are eventually translated to physical block addresses (server name, disk number, and disk block), and the servers deliver data directly to the application. (The servers and the application host are made known to each other during the data set “open” process, and sufficient information from the original request is carried along so that the application can identify the data provided by the servers.)

The DPSS is not an inherently reliable tertiary storage system, though the potential exists for this capability by adding tertiary backing storage on the storage servers.

3.2.2 DPSS and the TerraVision Application

Figure 7 illustrates the configuration and operation of the DPSS in the ARPA-funded MAGIC high speed ATM testbed [MAGIC] where the DPSS is being developed. The primary characteristics of this environment are the various networks (MAGIC backbone, Spartan [Spartan], and NTON [NTON]), the TerraVision terrain visualization application [TerraVision], and the model of a data curator (EDC). All components of the DPSS are geographically dispersed via the wide area, high speed, network infrastructure.

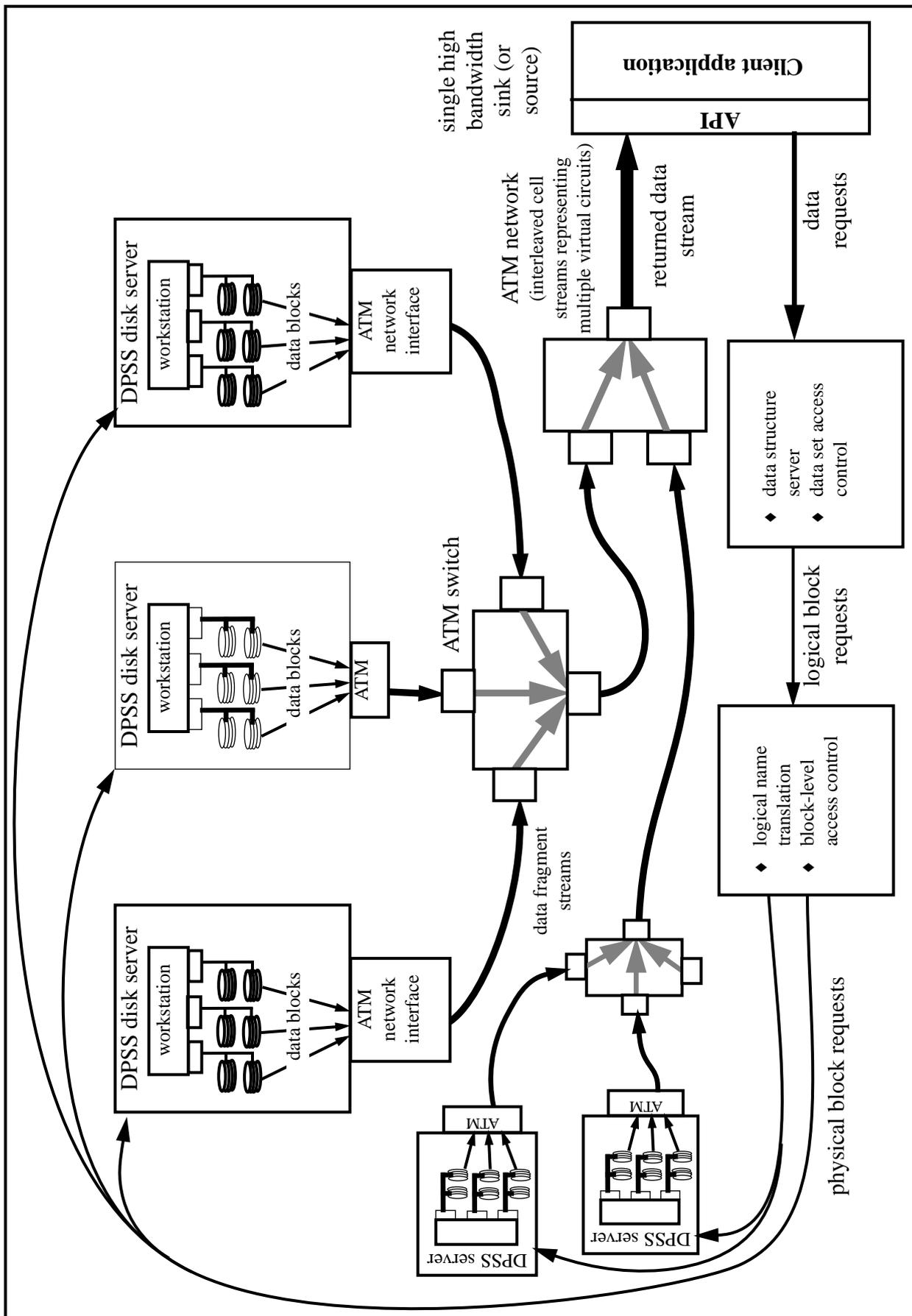


Figure 6 Distributed-Parallel Storage System Implementation

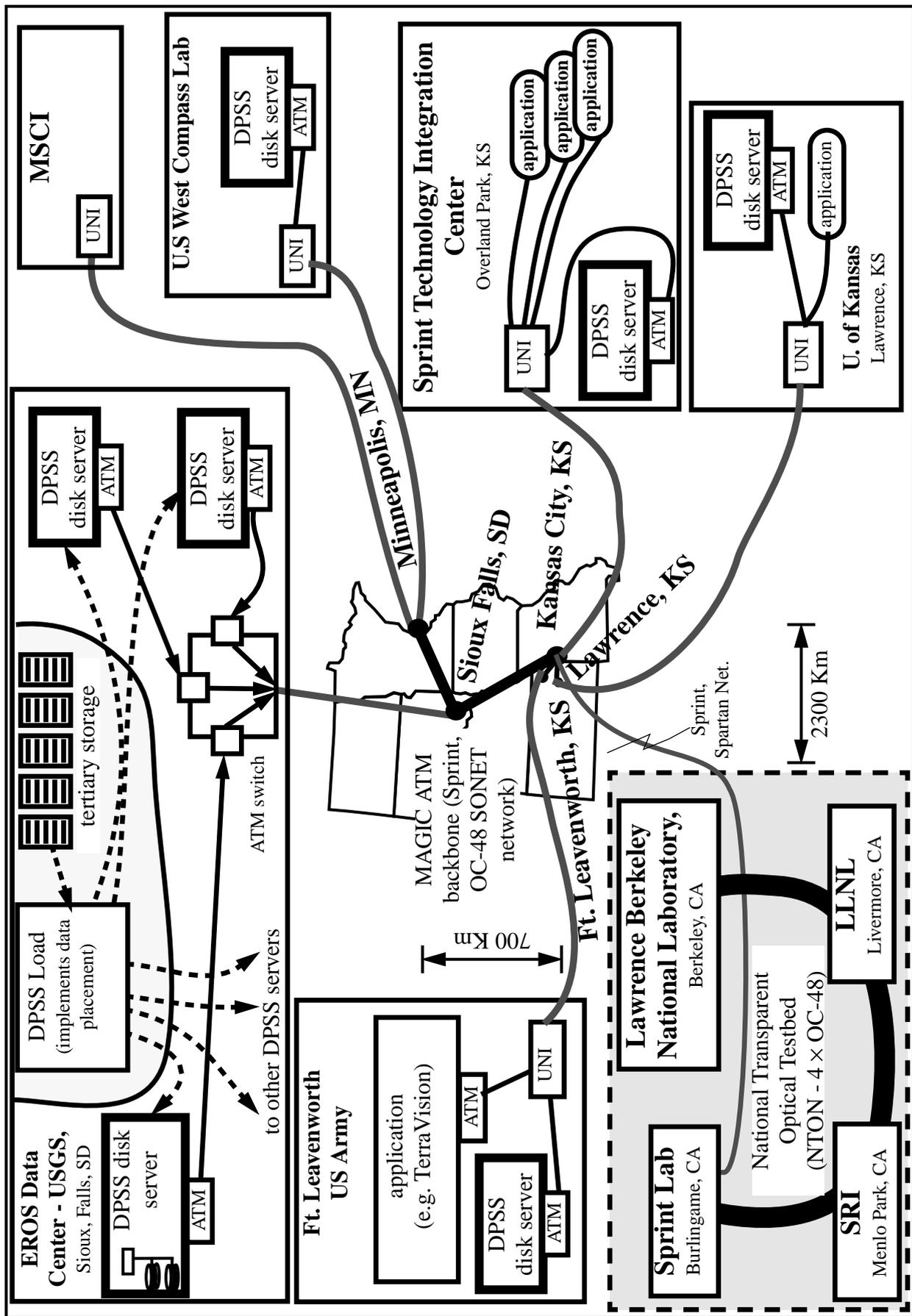


Figure 7 The MAGIC Testbed Distributed Application Environment

3.3 Data Management

In any scenario where data is generated in large volumes and with high throughput, and especially in a distributed environment where the people generating the data and the people cataloguing and using the data are geographically separated, a key issue is the automatic generation of at least minimal metadata, and the cataloguing of the data and the metadata as the data is received (or as close to real time as possible).

Such a capability is needed to help researchers organize, browse and search through data collections, especially collections where the original data are stored off-line in a mass storage system. The system should also facilitate co-operative research by allowing access to the data by specified users at local and remote sites, both by enabling immediate examination of the data, and thought incorporation of the data into other databases or documents.

Our approach to this capability (called “ImgLib” [ImgLib]) is to use the capabilities of the World Wide Web tools to provide a semi-automatic cataloguing of incoming data. This is done by extracting associated metadata and converting it into text records, together with generating auxiliary metadata and derived data, and combining these into Web documents.

In the case of image-like data, an ImgLib catalogue “entry” consists of (derived) thumbnail images and textual indexing material, together with pointers to the original images, video, or other data (whose permanent storage is typically off-line). The indexing material is kept on-line on a system that is running the ImgLib Web Server. The ImgLib server also directly manages a DPSS-based on-line cache, and manages migration of data from tertiary storage to cache.

A typical usage scenario involves creating a top level collection (a branch of a directory tree of html documents) on the server machine. Once the initial organization is established, the collection owner is able to manage the collection from his local workstation via Web form interfaces. ImgLib provides the following functions:

- Creates from the original data-objects, and stores, a set of “browsable” (e.g. “thumbnail”) images that can be easily viewed;
- Creates a text description file associated with each data-object and indexes over these files,
- Allows the collection to be laid out in hierarchical sub-collections which can be independent of how the original data-objects are organized or stored;
- Provides automatic and transparent management of the tertiary storage and the migration of data-objects from tertiary storage to on-line cache;
- Provides mechanisms to build, on-the-fly, subsets of the collection of data-objects based on keyword or metadata field searches or sub-collections defined by branches of the organizational hierarchy;
- Provides browsing tools to present a “view” of the subsets represented by thumbnail images and the associated textual metadata.

Figure 8 illustrates some of these points. It shows the results of a search on the textual metadata. The information about the data-objects that result from the search is shown as an associated collection of thumbnails, associated pointers to several other types of derived data (including a “movie” representation), and a pointer to the original data-object. The data-objects are kept on tertiary storage (a tape-robot based mass storage system in this case), and there is an option for forcing migration of data back to the on-line cache if the data of interest is not already there. In this example, the derived “movie” representation requires a special application to view.

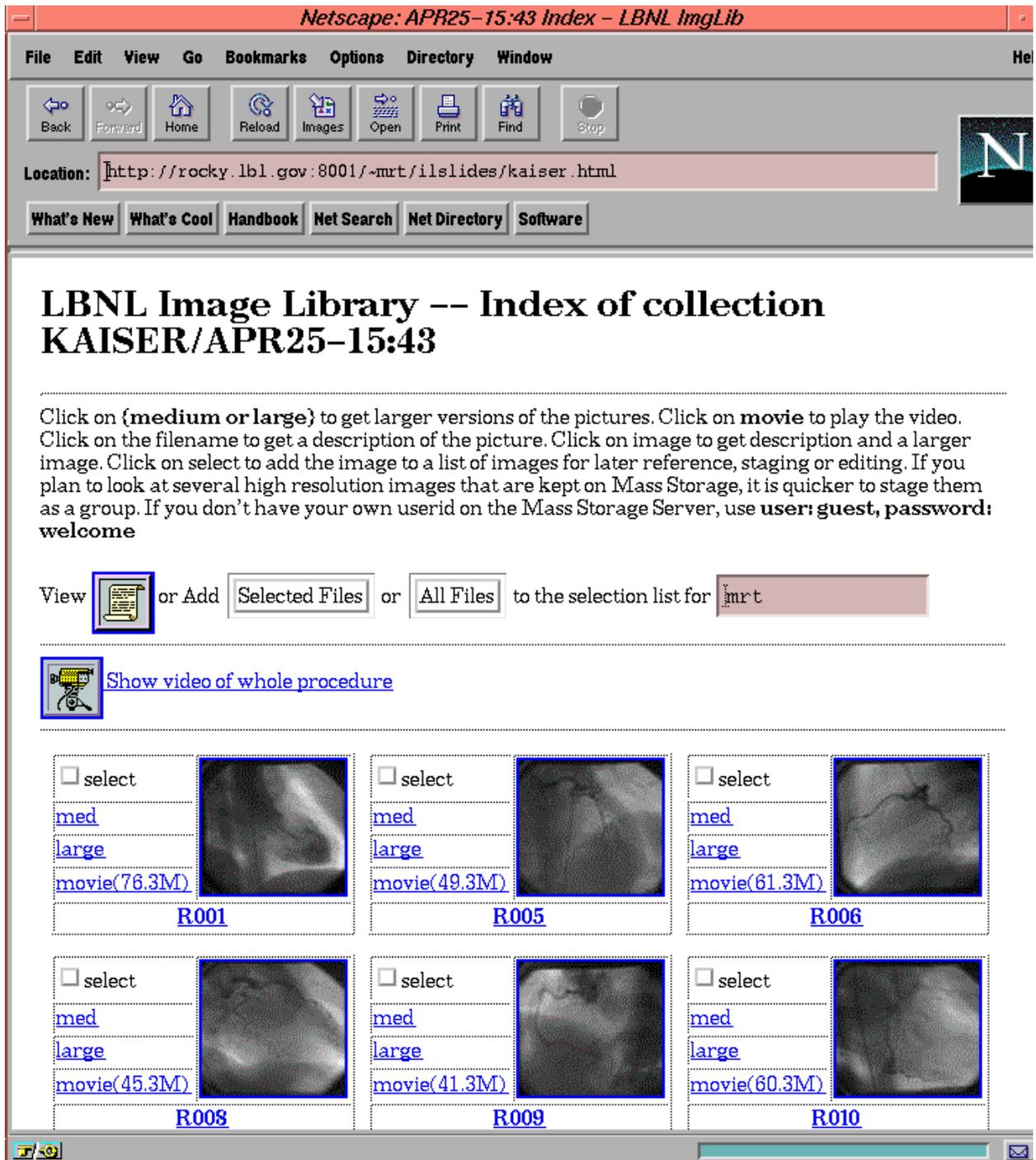


Figure 8 Large Data-Object Management for Video-like Objects

3.3.1 The Health Care Application

In the case of the cardio-angiography application, the video data is processed and then indexed using the method described above. Subsequently, physicians at Kaiser Oakland Hospital, Kaiser Division of Research, and Kaiser Vallejo Hospital can locate images and video sequences through the Web interface, and study those sequences from their workstations. The volume of original digital video data is sufficiently large that it is not practical sent it across the network for viewing at full frame rates on typical physician workstations. Therefore the video data is viewed with a

special video player that first uses a compressed version of the data to play the video-angiography sequences at video rates, and then display a subset of the original data. The compressed sequences are used as a “preview” to locate the angiography images of interest. As the video player is “slowed down” by the physician (for more detailed examination of the individual frames) and the frame rate reaches three or four frames per second, then the player switches to displaying the un-compressed digital video directly from the DPSS network storage system. (See Figure 1.) This approach uses the network to transport first the compressed video sequence, and then, in real-time, only those full resolution frames needed for close examination. This avoids having to move the very large amount of un-compressed digital video around the network many times and makes it possible for physicians to call up data “on the spur of the moment” and examine it on commonly available workstations that are not able to locally store the tens of Gigabytes of data needed for a full angiography study.

3.4 Performance

3.4.1 Performance Issues Related to ATM

There are virtually no behavioral aspects of an ATM “network” that can be taken for granted, even in an end-to-end ATM network. By “network” we mean the end-to-end data path from the transport API through the host network protocol (TCP/IP) software, the host network adaptors and their device drivers, the many different kinds of ATM switches and physical link bandwidths, and then up through the corresponding software stack on the receiver. Further, the behavior of different elements at similar places in the network architecture can be quite different because they are implemented in different ways. The combination of these aspects can lead to complex and unpredictable network behavior.

In the next two sections a provide some information and observations on this issue. We only address issues in the flow of data after the host-to-host connection is established. (That is, we do not address routing and signaling.)

3.4.2 ATM Network to Workstation Performance

The architecture and resulting environment that we are describing is a “data movement intensive” environment. That is, in the normal course of dealing with large data-objects in a fully distributed environment, large amounts of data are constantly on the move. This circumstance generates many the R&D issues from workstation network performance to security, and we discuss a few of them here.

In a scenario where you are counting on parallelism to provide high performance, then one of the important attributes is the independence of the parallel operating elements. One of the reasons that the DPSS has been so successful as a high performance storage system is that using wide area networks increases the independence of the servers (until you get to the final point of delivery).

One of the problems that we have encountered is that independence within the workstation platforms is much more difficult to achieve. Fortunately, platform hardware and software designers have put considerable effort into making key sub-systems operate independently. So, for example, we have found that the parallel operation of disks on controllers, controllers on busses, and memory sub-systems all do a reasonably good job of increasing performance through parallel operation.

On the other hand, this is not true for network interfaces. That is, through progress is being made, multiple network interfaces have, so far, generally done a poor job of providing increased throughput by parallel operation. For ATM interfaces, this is due in part to the newness of the

technology (and corresponding lack of experience in designing the supporting system hardware and software), and in part due to the fact that prior to the commercial interest in video file servers, striping across network interfaces was not a common practice.

The following synopsis describes the progress in this area (which is very important in high speed distributed applications). Some of the progress is due to the maturing of ATM network code stacks, and some is due to several years of interaction with the workstation vendors pointing out the problems.

The evolution of the IP over ATM infrastructure:

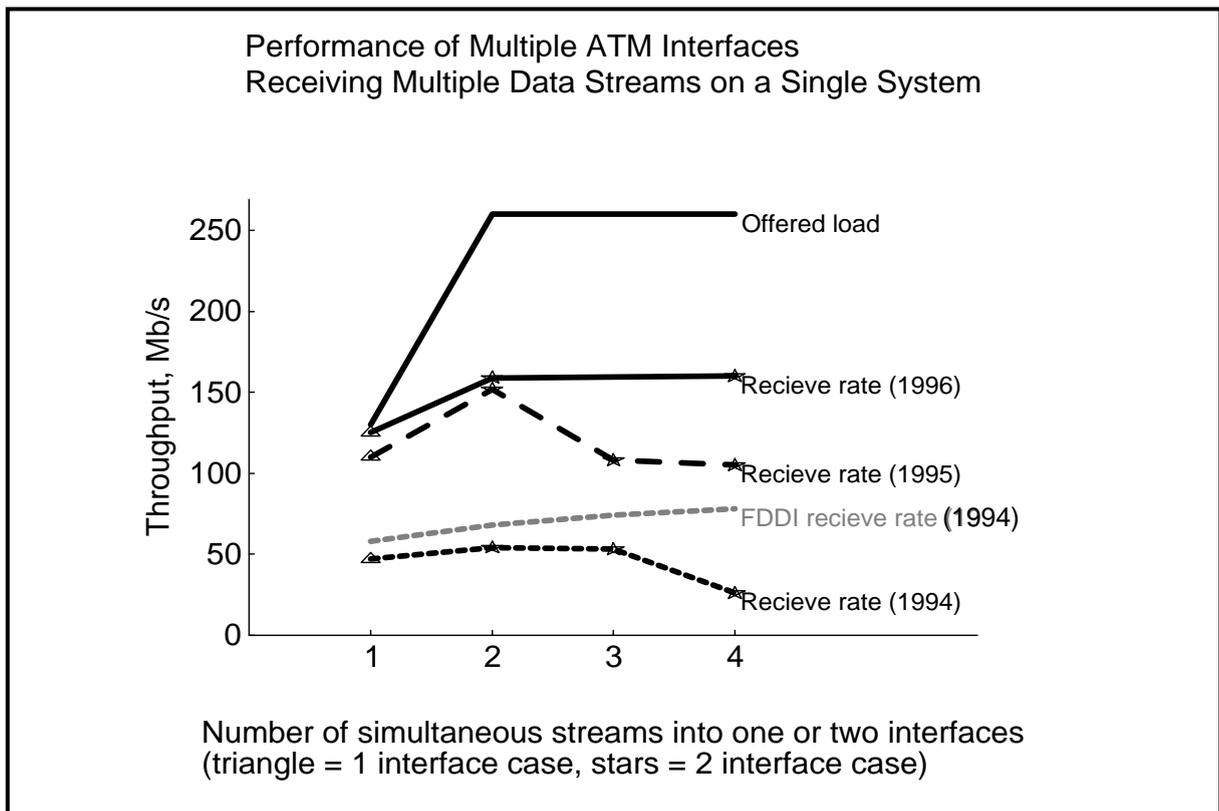


Figure 9 Evolution of the Performance of Workstation ATM Network Interfaces

- 1994
 - single interfaces were slow
 - multiple interfaces were no faster (poor independence of system data paths - nothing was multi-threaded)
 - switches dropped cells silently
- 1995
 - single interfaces are faster
 - multiple interface data path independence is improving
 - switches still drop cells silently
- 1996
 - single interfaces and systems are approaching wire speed
 - multiple interfaces better still, but not ideal

- switches now report cell loss

3.4.3 ATM Network Performance

3.4.3.1 Un-congested Behavior

Application Requirements:

To provide a specific example, we consider the Kaiser CalREN project and environment.

Medical imaging applications are a “systems” problem - no one requirement sets all of the necessary parameters). Distributed handling of video-angiography large data-objects is an example of an application constrained by the use of a public data network. Data is collected and stored locally during a short recording period, after which it is transmitted to a network-based storage system for subsequent processing, retrieval, and review at other locations. The quantitative characteristics are:

- video-angiography \equiv x-ray video camera in a bi-plane (two view) configuration (Phillips DCI Scanner)
- data is collected at $30 \text{ fr/sec} \times (512 \times 512 \times 8 \text{ bits} \times 2 \text{ views})$ (monochrome video, no compression) = 126 Mbits/sec (one of several relevant parameters)

Four GBytes of data is collected in a local buffer in about four minutes in the form of data-objects that range in size from 50-300 MBytes. Following the data collection, there are about twenty minutes in which to off-load local storage before the next collection period.

Four GBytes in 20 minutes is a rate of about 25Mbits/sec of user-level data on the network, or, accounting for the various protocol overheads, approximately $25 \times 1.4 \approx 34$ Mbits/sec at the ATM level.

34 Mbits/sec is about 23% of an OC-3 network, and it was our initial estimate of the maximum reasonable use of a shared network, in this case Pacific Bell's SF Bay Area ATM network - an ATM OC-3 network which has 30-40 attached sites.

Observations:

We ran a series of throughput tests to determine suitability of the Pacific Bell Metropolitan Area ATM for video-angiography transmission.

The experiments consisted of sending 25 Mbits/sec, with a duty cycle of 40 minutes on, 20 minutes off, 24 hours a day, and were conducted over several weeks.

Although this network was basically un-congested, several of the experiments were timed to correspond to some high resolution video and audio multicast sessions that generated 5-10 Mbits/sec of “continuous media” (steady state) traffic in the network, in addition to whatever else might have been going on in the network.

The experiments were done by putting the source host ATM adapter into “continuous bit rate” mode, but sending data via IP over AAL-5. This ensured that we transmitted a “true” 25 Mbits/sec, and did not inject high bandwidth “bursts” of data into the network (as would have been the case if we had rate-limited the TCP writes to an “average” of 25 Mbits/sec).

As can be seen in Figure 10, to the level of about 1%, there was no variation in our ability to send 25 Mbits/sec of data for any of the several week long periods during which the throughput

experiments were run. Figure 10 shows typical results that include some times with high rate multicast traffic was present (in particular, Friday afternoon).

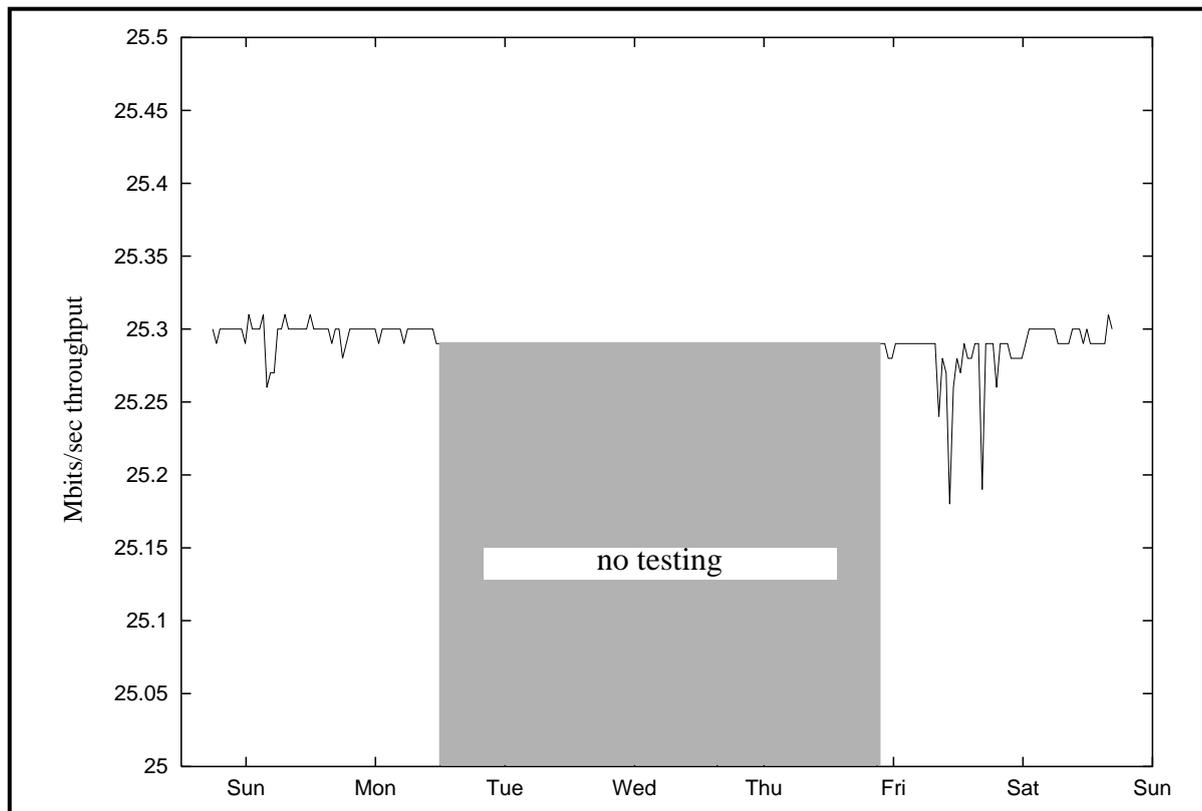


Figure 10 Throughput Testing to Determine the Suitability of the Pacific Bell, SF Metropolitan Area ATM for Video Angiography Transmission

3.4.3.2 Congested Behavior

The impact of congested ATM networks on TCP transport (which is used for all of the data transfers described here) can be very complex. Today's ATM networks and components are still relatively immature and with every new generation of switch software and hardware the congested behavior can change. The most severe impacts occur in circumstances when the characteristics of the various ATM links and network components conspire to defeat TCP's (by now otherwise well tuned) congestion control mechanisms. TCP has a suite of congestion detection and response algorithms [Stevens] that have evolved over the past 10-15 years in the Internet environment, and in that environment they work very well. (The Internet operates under conditions of very high load and high potential congestion, but congestion is generally avoided and the links usually operate at high utilization efficiency.) The problem in ATM networks is that the TCP control algorithms are based on certain assumptions about the relationships among parameters like the packet switch buffer sizes, link minimum transmission unit sizes, link bandwidth, etc. It is very easy for ATM networks to violate these assumptions, and when this happens both the throughput and the link utilization efficiency can be very low in the presence of congestion. (In current implementations of ATM networks, this congestion is frequently caused by small output port buffers in ATM switches. The good news is that most new switches are being designed with large output port buffers.) For a detailed analysis and case study of this situation,

the reader is referred to “Performance Analysis in High-Speed Wide Area ATM Networks: Top-to-Bottom End-to-End Monitoring” ([Tierney]).

3.5 Security Architectures for Large-Scale Remote Environments

The access to remote resources, whether for monitoring, control of experiments, or access to data, requires that security and access control mechanisms be in place in order to provide safeguards against unauthorized access, and privacy of proprietary or otherwise restricted data.

To provide security in large-scale, open access environments, we have evolved a decentralized security architecture that uses the concept of use-conditions imposed by the entity primarily responsible for a resource. This is, while the access restrictions are enforced at the resource, the conditions of access are specified by the principal responsible for specifying an access policy for the resource. The responsible principal is the only one who imposes and records use conditions, and those conditions are made available in a way that has nothing to do with the operation of the resource itself: The security mechanism at the resource is only responsible for checking that the use-conditions are satisfied, and enforcing access controls appropriately.

We believe that the combination of this model, and the generalized public-key certificate infrastructure that supports it, can provide a highly scalable and transparent infrastructure for both authentication and authorization of access to the full range of network based resources. The goal is a security architecture that can encode, distribute, and protect the information needed to enable routine secure availability of remote resources in a widely distributed environment.

3.5.1 The Problem and Objectives

The problem is to identify and demonstrate authentication and access control methodologies that are applicable to all user-level resources - data, storage and processing, control of external devices, etc. - in widely distributed systems, and that are also general, scalable, “strong”, non-intrusive (for both the administrator and the user), and easily managed.

Security is an essential part of remote access in open environments, and so awkward, intrusive, and expensive security will delay wide-spread use of these environments.

To address the problem, we have evolved the decentralized architecture noted above, and an implementation design that should provide “strong” security for applications that use the security architecture directly, and “good” security for a large class of unmodified applications.

As with many scalable distributed systems, the decentralized aspect of the architecture is the key to scalability: the administrative task of maintaining an appropriate expression of access policy is only done by those entities (principals, or their agents) that are directly responsible for the policy governing the use of a resource. The operational administration of that resource has only to deal with providing an access control mechanisms, not with access control information. Central administration or administration, by the resource operator of access control lists, or multiple roles that have to kept straight by resource users, are not a scalable approach in a global distributed environment, and are implementation of security that this model is intended to replace.

Our implementation approach is no more (and in the case of the capabilities described below, probably less) vulnerable to human error than other approaches, and it should be scalable and easily used. It uses “strong” cryptography to protect essential information, and so should not be any more or less vulnerable to “frontal” attack than conventional systems. In what we feel to be an essential trade-off favoring transparency of use, it is possible to use our approach in a way that permits security to be inherited by unmodified applications from a “secure shell”. This increases

potential exposure if the host operating system is compromised, however, by the same token, this is a great strength: It extends security to a wide range of unmodified applications in a way that is easily used by a novice user population.

3.5.2 Background

Public-key certificates (PKC) are an inherently reliable and secure mechanism for widely distributing the assured information needed for authentication and authorization. (See [RSA] for background information.) PKCs are also forming the foundation of the evolving commercial transaction security infrastructure. (See [MC] and [SET]).

Public-key certificates can encode, in a secure and verifiable way, one's identity, organizational / group affiliation, creditworthiness, level and scope of training, etc. Certificates that encode an attribute can enable automated processing of requests for operational control of remote equipment; remote, sensitive data access; use of resources; etc.

There are several other important aspects of a PK certificate infrastructure:

One is a mechanism to establish and represent trust relationships (policy and procedure agreements) among authorizing or verifying entities. The common ways of doing this are represented in a continuum from a centralized root, hierarchical structure of certification authorities (CA) at one end of the spectrum (see [Kent]), through "webs" of small (organization-scope) CAs, to the completely decentralized approach of PGP (where individuals attest to each other's attributes) at the other end. All of these approaches are in use, and we have focused on the middle ground of organization-level CAs that attest to specific facts and establish trust on a pair-wise basis among themselves.

A second crucial piece of the infrastructure is a mechanism for making certificates locatable, and widely and reliably available. There are several current approaches to this. X.500 directory servers provide, in principle, a standard way of searching for and distributing standard format certificates (X.509); There are already WWW sites and ftp sites that provide PGP certificates. Other distributed information mechanisms (e.g. whois++) are possible. Our current approach is based on using the "Lightweight Directory Access Protocol" [LDAP] to communicate with X.500 servers to obtain X.509 certificates.

A third important component is a mechanism for the representation and automatic manipulation of policies. While our first implementation will use ad-hoc representations, approaches such as PolicyMaker [Blaze] will be required for any general infrastructure.

3.5.3 The Security Model

The model that we are working toward is based on verifiable and securely represented resource controller / owner imposed use-conditions, and equally secure and verifiable satisfaction of those conditions: The individual identity that is so important in many security architectures is just one of many factors in our model.

The basic idea is that the entities that control resources will establish a set of conditions for the use of the resources. The use-conditions are whatever is appropriate to the resource: a level of training for operating an instrument; membership in a group (e.g. health care professionals working for a particular HMO) for access to patient data; creditworthiness for access to a resource that requires a commitment of payment; etc.

Paired with the resource owner's use conditions are the entities that can attest to the relevant attributes of a user or agent that is seeking access to a resource: A school registrar may certify a

level of training; a state examining board may certify a level of professional attainment; an organization can certify employment (and related group membership); a bank, a level of credit; and a certification authority may attest to individual identity.

In addition to conditions and certifications of those conditions, a collection of trust relationships must exist. So, for example, that a resource controller will accept a certificate that purports to represent satisfaction of a use-condition requires that the resource owner trusts the entity that is attesting to (certifying) an attribute of the user is qualified and/or certified to make such a claim. These trust relationships may be direct: the resource owner and the verifier of conditions have an explicit, pair-wise agreement; or they may be indirect: a resource owner accepts a statement of policy, perhaps policed by a third party, as an implicit guarantee of the satisfaction of a use-condition.

In the general case in our architecture, a broker / agent acting on behalf of the potential user, collects the certified use-conditions for a resource, and then collects the certificates / credentials that verify the user's satisfaction of the use conditions. The combination of resource identity, use-conditions, and matching credentials, all of which are cryptographically secured, are used to form a "capability": A capability to access, monitor, manipulate, use, etc., a remote resource.

However, even with all of these tools and infrastructure in place, they are still just that: tools and infrastructure. In each given situation where a resource is being protected, those responsible must evolve a security model that provides the level of assurances, protection, auditing, etc., that meets their requirements. Flexible tools and infrastructure make it possible to implement the specific security model in a usable and cost effective way.

Our general approach is motivated by the credential model used in most Western societies for consumer contracts. An example that illustrates many of the points is that of renting an automobile at an airport. The essential elements of this process include presenting and verifying a set of credentials, and then being given a capability (authorization certificate) in the form of a valid rental car contract, that is then used to gain access to the resource (the automobile).

In this model, a driving license is an attribute certificate that represents a guarantee by a trusted third party (the issuing state Dept. of Motor Vehicles) that an individual has had appropriate training and is currently authorized to drive a car. The driving license certificate also serves as proof of identity when the individual's signature must be verified in order to signify agreement to a contract. An individual's physical signature is a "private" or individual expression (like the private-key part of the public-private key pair). The photograph of the individual's face and a photo of a written signature together on a driving license certificate constitutes the certified "public" expression of identity (like a "classical" public-key certificate).

One's credit card provides the entree for on-line checking of creditworthiness. These two credentials (driving license and credit card) are verified (usually automatically), and a capability - to use the rental car, with restrictions - are represented in a third, combined authorization certificate (or "capability"): the rental car agreement.

The rental car company agents are the "access management agents" (AMA). They verify multiple use-conditions by checking a collection of attribute certificates and issue a (composite) authorization certificate. We call this new certificate a "capability" because it associates a set of otherwise unrelated certificates to represent the satisfaction of a set of use-conditions that, together, enable the use of a specific resource. The guard at the rental car lot gate is the resource controller, and when presented with a signed AMA certificate (the capability) enables the use of the rental car (the resource).

The success of this model is clear: most responsible people can rent an automobile - a dangerous and expensive resource - anywhere in the world by producing the commonly possessed set of attribute certificates: the driving license and the credit card.

This example embodies most, if not all, of the characteristics of a security architecture that can provide secured global capabilities in an open Internet and the abstraction of these operations are necessary, and probably sufficient, conditions to ensure general resource security in the Internet.

One of the expected advantages of the architecture described here is that it should provide a more natural (and distributed) responsibility for information validity, placing the basic responsibility and definition of policy with those entities that have the direct responsibility for the information (e.g. the home state Department of Motor Vehicles) or for setting use-conditions (the rental car company).

The combination of direct reference to those responsible for the generation and assurance of information (the resource owners or policy makers), and standardization of representations of the information that permit automated verification, should form the basis of an easily used and general purpose authentication and authorization system that provides such services as qualifying remote experimenters to use instrumentation systems in distributed collaboratories, as well as forming the basis for dynamically configured, large-scale, distributed computing resources.

Automated acquisition and verification of attribute credentials will increase efficiency and reliability, and will play a critical role in scaling remote environments by removing the need for human intervention in repetitive and trivial tasks such as credential checking, especially when the human participants are never in the same place at the same time, and may never physically visit the site where service is being requested.

3.5.4 Implementation Issues

Implementation of the model described above requires a number of operational components with capabilities beyond what we see today:

- ◆ Distributed brokers that understand the encoding of the information in corporate database objects represented in signed certificates will be needed to automatically build capabilities.
In the general model (see [Johnston96a]), when capabilities require a number of attribute certificates from different sources, a “broker” will locate the certificates required by the specified use-conditions and build the composite certificate (capability) on behalf of the user.
- ◆ Automated coupling between corporate databases and certificate servers will minimize the administrative overhead.
This will (eventually) provide the automated maintenance of attributes and automatic certificate generation based on information obtained from the primary corporate database (e.g. employee records, school records, etc.).
- ◆ Rules and mechanisms for encoding requirements (the “use-conditions”) for services will be required, as will mechanisms for determining the level of trust to place in other CAs (i.e. how are CA operating policies published and evaluated), etc. (For example, see [Mendez and Huitema] and [Blaze])

The basic technology is a generalized certificate understandable to a variety of distributed entities. However, while the certificate idea has been around for some time, there is no widespread use or understanding of many of the issues involved, which run the gamut from sociological to legal to technical. Even if we consider only the technical issues, we are faced with a considerable range of

topics to address: expression of attributes and certification policy in machine understandable form; distribution and revocation of certificates; usability and performance of certificates and their processing on the part of the servers, brokers, and clients, etc.

Our prototype implementation addresses:

- credential and key management via SSH and LDAP/X.500
- security context establishment via GSS/SPKM
- secure, authenticated, and integrity protected exchange via GSS-API

The SSH design [SSH] provides a convenient and transparent management of the credentials and public keys of servers and the private keys of the user. (Roles are possible through different “identities” of the user.) SSH also provides protection of unmodified applications by establishing a secure shell, and then redirecting communication with the outside world through secured ports.

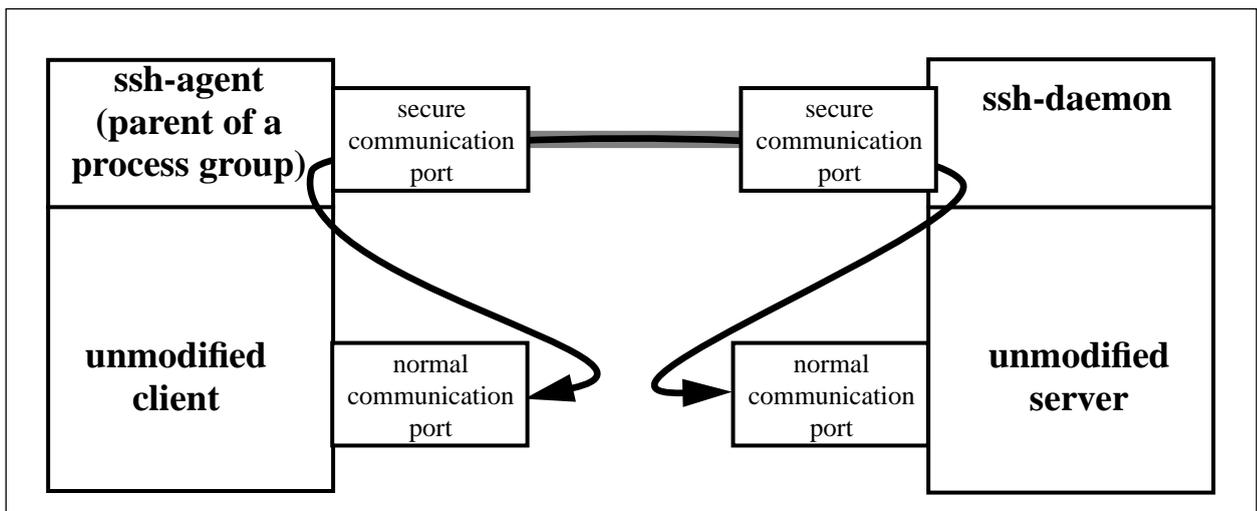


Figure 11 Protection of a Unmodified Client using a “Secure” Shell Approach

We have modified the SSH credential and key management facility to be able to set up the SPKM [GSS] security context for GSS. (See [Adams] for a good discussion of the rationale for SPKM and its relationship to GSS/Kerberos.)

SSH makes the user’s private key available to SPKM to do the session key exchange needed for GSS security context establishment. Once this is done, the GSS-API provides a simple application interface to send validated and/or confidential communication.

We currently use an ad-hoc use-condition encoding, but are examining the approach mentioned above.

3.5.5 Relationship to Other Work

Our general approach is an application of public key certificates, and as such it draws ideas from an active community of researchers.

The basic concepts of distributed systems security are laid out in various ISO (International Organization for Standardization) and ECMA (a Europe based international association for the standardization of information and communication systems) documents (see [SIRENE]). SESAME is an implementation that incorporates many of the [ECMA-219] and [ECMA-138] concepts. (See [SESAME].)

We have, however, taken a somewhat different direction than the ECMA and SESAME approaches with respect to attribute certificates, which are central to our approach. In the ECMA approach, most of the emphasis on the use of attribute certificates seems to be in the form of access control lists and roles (party A, who controls a resource, grants to party B the right to use that resource, perhaps under restricted conditions (roles)). Our approach treats use-conditions, rather than identity, as the fundamental concept. That is, we see resources primarily being protected by requirements on the user (e.g. level of training) rather than identity. (ACLs are, of course, easily implemented in this model.)

Our implementation is based on some refinements to the “secure shell” [SSH], combined with the “Generic Security Service Application Program Interface” [GSS] to provide an easily used application security mechanism.

3.5.6 Application of the Security Architecture

We are using the prototype implementation to provide security to a health care information application (see [Johnston95a]) that uses the Distributed-Parallel Storage System ([DPSS]). The application uses a WWW-based system for managing certain kinds of medical imaging records. The Vplayer application (“Vplayer”) is invoked via a WWW image database frontend that uses WWW security to protect the metadata. The application (a specialized, combined video and image browser) then accesses data stored on the DPSS. The Vplayer access to DPSS data is protected by the security architecture described here. (That is, credentialing is handled by the secure shell that runs Vplayer, and the DPSS access is controlled by checking authorization certificates.)

4.0 Conclusions

We have described several elements of an architecture the using high speed, wide area networks as the basis of a highly distributed environment for handling all aspects of large data-objects. Variations of the overall approach have been applied in several different projects, and we have described some of the results.

We believe this approach to be both practical and powerful given the increasing availability of high speed networks. Further work will be done on both the various functional elements and specific applications. Please see <http://www-itg.lbl.gov> for more information.

5.0 Acknowledgments

The overall distributed large data-object approach has evolved over several years with funding from both the Mathematical, Information, and Computational Sciences Division of the Office of Energy Research, U. S. Department of Energy, and ARPA Information Technology Office.

In the Kaiser CalREN project, E-J. Pol of Philips Research, Palo Alto, CA, designed and implemented the workstation interface for the cardio-angiography system. Dr.J. Terdiman, of the Kaiser Permanente Division of Research provided the leadership for the health care aspects of the project. Dr. Bob Lundstrum, Kaiser San Francisco Hospital, Cardiac Catheterization Laboratory, also provided guidance and information specific to the health care application.

Pacific Bell, through its CalREN program, provided funding for connecting the LBNL, Kaiser, and Philips sites to the SF Bay Area ATM network.

Sun Microsystems assisted in various ways, including loaning equipment to various Kaiser sites.

6.0 Notes and References

Adams Adams, C. "IDUP and SPKM: Developing Public-Key-Based APIs and Mechanisms for Communication Security Services". See <http://bilbo.isu.edu/sndss/sndss96.html>

Hudgins-Bonafield "Pipe Dreams And Lessons With TCP Over ATM". Published in Network Computing Online (<http://techweb.cmp.com:80/techweb/nc/docs/default.html>). See <http://techweb.cmp.com:80/techweb/nc/706/706hreportb.html>

Adams Adams, C. "IDUP and SPKM: Developing Public-Key-Based APIs and Mechanisms for Communication Security Services". See <http://bilbo.isu.edu/sndss/sndss96.html>

Blaze "Blaze, M., J. Feigenbaum, J. Lacy, "Decentralized Trust Management", to be presented at the 1996 IEEE Symposium on Research in Security and Privacy, Oakland, CA, May 1996. (<ftp://research.att.com/dist/mab/policymaker.ps>)

BAGnet-1 "Bay Area Gigabit Testbed (BAGNet)" See <http://www-itg.lbl.gov/BAGNet.html>

BAGNet-2 Pacific Bell Announces First CalREN Award"
<http://www-itg.lbl.gov/BAGNet.hm.pg.docs/PacBell.BAGigabit.Press.Release.html>

CRYPTOLOG ""The Internet Guide to Cryptography"
<http://www.enter.net/~chronos/cryptolog1.html>

DPSS The Distributed-Parallel Storage System (DPSS) Home Page
(<http://www-itg.lbl.gov/DPSS>)

ECMA-138 "Security in Open Systems - Data Elements and Service Definitions", 1st Edition (December, 1989). European Computer Manufacturers Association,
<http://www.ecma.ch/ecma-138.HTM>

ECMA-219 "Authentication and Privilege Attribute Security Application with related key distribution functions", 1st Edition (December 1994). European Computer Manufacturers Association, <http://www.ecma.ch/ecma-219.HTM>

GSS Linn, J., "Generic Security Service Application Program Interface"
(<http://ds.internic.net/rfc/rfc1508.txt>) Also see more recent and related drafts at the IETF Common Authentication Technology home page
(<http://www.ietf.cnri.reston.va.us/html.charters/cat-charter.html>) and at
<http://www.ietf.cnri.reston.va.us/ids.by.wg/cat.html>.

ImgLib See http://www-itg.lbl.gov/ImgLib/ImgLib_intro.html

Johnston96a Johnston, W. and C. Larsen, "Authenticated Global Capabilities for Distributed Systems" (draft)

Johnston95a Johnston, W., J. Terdiman, E. J. Pol, "The Use of Public ATM Networks for Health Care Imaging Information" (see <http://www-itg.lbl.gov/Kaiser/LKP>)

Johnston95b W. Johnston, and D. Agarwal, "The Virtual Laboratory: Using Networks to Enable Widely Distributed Collaboratory Science" A NSF Workshop Virtual Laboratory whitepaper. (See <http://www-itg.lbl.gov/~johnston/Virtual.Labs.html>)

Johnston95c Johnston, W., "BAGNet: A High Speed, Metropolitan Area, IP over ATM Network Testbed", IEEE Comcon'95, March, 1995, San Francisco, CA. Also see <http://www-itg.lbl.gov/BAGNet.hm.pg.docs/CompCon.95.3.html> .

Johnston93 W. Johnston and A. Allen, "Regional Health Care Information Systems: Motivation, Architecture, and Implementation," LBL Technical Report 34770 (draft), Dec. 1993.

Kent Steve Kent, "Privacy Enhancement for Internet Electronic Mail: Part II: Certificate-Based Key Management", RFC-1422, <http://ds.internic.net/rfc/rfc1422.txt>.

Kaiser-CalREN Kaiser - LBNL - Philips CalREN project. See <http://www-itg.lbl.gov/Kaiser/LKP/homepage.html> .

LDAP "Lightweight Directory Access Protocol", <http://www.umich.edu/~rsug/ldap/ldap.html>.

MAGIC testbed, [http:](http://)

Mendez and Huitema "A New Approach to the X.509 Framework: Allowing a Global Authentication Infrastructure without a Global Trust Model", S. Mendez and C. Huitema, in "Proceedings: Symposium on Network and Distributed System Security", IEEE, San Diego, CA, February 16-17, 1995.

MC "Visa & MasterCard Combine Secure Specifications For Card Transactions On The Internet Into One Standard" <http://www.mastercard.com/Press/release-960201.htm>

NTON The National Transparent Optical Network (<http://www.cais.net/nton/>) connects a number of organizations in the San Francisco Bay Area)

RSA "RSA Labs' Frequently Asked Questions About Today's Cryptography" http://www.rsa.com/rsalabs/faq/faq_home.html

SESAME

SET "Secure Electronic Transactions" <http://www.mastercard.com/set/set.htm>

SIRENE "Security in Computer Networks"
(<http://www.zurich.ibm.com/pub/sti/www/g-kk/sirene/> also see <http://www.zurich.ibm.com/pub/sti/www/g-kk/sirene/pointers.html>)

Spartan Sprint test network. T-3 / OC-3 ATM from Kansas City, Kansas to the Sprint Test Laboratory in Burlingame, CA.

SSH Ylonen, T., "The SSH (Secure Shell) Remote Login Protocol"
(<ftp://ietf.cnri.reston.va.us/internet-drafts/draft-ylonen-ssh-protocol-00.txt>). Also see "Ssh (Secure Shell) Remote Login Program" (<http://www.cs.hut.fi/ssh/>).

SSL Leay, E. "SSLLeay". (SSLLeay is a free implementation of Netscape's Secure Socket Layer - the software encryption protocol behind the Netsite Secure Server and the Netscape Browser.)

See <http://www.psy.uq.edu.au:8080/~ftp/Crypto> and <http://www.home-port.org/~adam/crypto/ssl.html>

"What started as an effort to implement the SSL protocol has turned into a fairly complete cryptographic library (which Eric is still working on). There is also quite a bit of ASN.1 support, with routines to convert and manipulate the base ASN.1 types, X509v3

certificates, certificate requests, certificate revocation lists (CRL), RSA private keys and DH parameters. There are routines to load and write these objects in base64 encoding and routines to convert ASN.1 object identifiers to/from ASCII representations and an internal form. There are functions for verification of X509 certificates and for specifying where to look for certificates to 'climb' the x509 'tree'. This last part of the library is still evolving.

The big number library is quite complete and has no restrictions on the size of the numbers manipulated. RSA and Diffie-Hellman routines have been layered on top of this library.”

Stevens Stevens, W. “TCP Slow Start, Congestion Avoidance, Fast Retransmit, and Fast Recovery Algorithms”. IETF Draft. See <ftp://ietf.cnri.reston.va.us/internet-drafts/draft-stevens-tcpca-spec-01.txt> .

TerraVision A terrain visualization application that uses the DPSS to supply high speed data streams for realtime, “virtual reality” type navigation of real landscape. (See <http://www.ai.sri.com/~magic/terravision.html>)

Tierney Tierney, B., W. Johnston, G. Hoo, J. Lee, “Performance Analysis in High-Speed Wide Area ATM Networks: Top-to-Bottom End-to-End Monitoring”, To be published in IEEE Networking Journal, LBL Report 38246, 1996. (Also see <http://www-itg.lbl.gov/DPSS/papers.html> .)

Wiltzius Wiltzius, D., L. Berc, and S. Devadhar, “BAGNet: Experiences with an ATM metropolitan-area network”, ConneXions, Volume 10, No. 3, March 1996. Also see <http://www.llnl.gov/bagnet/article.html> .