

The Architecture of the High Performance Storage System (HPSS)

Danny Teaff

IBM Federal
3700 Bay Area Blvd.
Houston, TX 77058
teaff@vnet.ibm.com
+1-713-282-8137
Fax +1-713- 282-8074

Dick Watson

Lawrence Livermore National Laboratory
PO Box 808, L-560
Livermore, CA 94550
dwatson@llnl.gov
+1-510-422-9216
Fax +1-510-423-7997

Bob Coyne

IBM Federal
3700 Bay Area Blvd., 5th Floor
Houston, TX 77058
coyne@vnet.ibm.com
+1-713-282-8039
Fax +1-713-282-8074

Abstract

The rapid growth in the size of datasets has caused a serious imbalance in I/O and storage system performance and functionality relative to application requirements and the capabilities of other system components. The High Performance Storage System (HPSS) is a scalable, next-generation storage system that will meet the functionality and performance requirements of large-scale scientific and commercial computing environments.

Our goal is to improve the performance and capacity of storage systems by two orders of magnitude or more over what is available in the general or mass marketplace today. We are also providing corresponding improvements in architecture and functionality. This paper describes the architecture and functionality of HPSS.

Introduction

The rapid improvement in computational science, processing capability, main memory sizes, data collection devices, multimedia capabilities, and integration of enterprise data are producing very large datasets. These datasets range from tens to hundreds of gigabytes up to terabytes. In the near future, storage systems must manage total capacities, both distributed and at single sites, scalable into the petabyte range. We expect these large datasets and capacities to be common in high-performance and large-scale national information infrastructure scientific and commercial environments. The result of this rapid growth of data is a serious imbalance in I/O and storage system performance and functionality relative to application requirements and the capabilities of other system components.

To deal with these issues, the performance and capacity of large-scale storage systems must be improved by two orders of magnitude or more over what is available in the general or mass marketplace today, with corresponding improvements in architecture and functionality. The goal of the HPSS collaboration is to provide such improvements. HPSS is the major development project within the National Storage Laboratory (NSL). The NSL was established to investigate, demonstrate, and commercialize new mass storage system architecture to meet the needs above [5,7,21]. The NSL and closely related projects involve more than 20 participating organization from industry, Department of Energy (DOE) and other federal laboratories, universities, and National Science Foundation (NSF) supercomputer centers. The current HPSS development team consists of IBM U.S. Federal, four DOE laboratories (Lawrence Livermore, Los Alamos, Oak Ridge, and Sandia), Cornell University, and NASA Langley and Lewis Research Centers. Ampex, IBM, Maximum Strategy Inc., Network Systems Corp., PsiTech, Sony Precision Graphics, Storage Technology, and Zitel have supplied hardware in support of HPSS development and demonstration. Cray Research, Intel, IBM, and Meiko are cooperating in the development of high-performance access for supercomputers and MPP clients.

The HPSS commercialization plan includes availability and support by IBM as a high-end Service offering through IBM U.S. Federal. HPSS source code can also be licensed and marketed by any US. company.

Architectural Overview

The HPSS architecture is based on the IEEE Mass Storage Reference Model: version 5 [6,9] and is network-centered, including a high speed network for data transfer and a separate network for control (*Figure 1*) [4,7,13,16]. The control network uses the Open Software Foundation's (OSF) Distributed Computing Environment DCE Remote Procedure Call technology [17]. In actual implementation, the control and data transfer networks may be physically separate or shared. An important feature of HPSS is its

support for both parallel and sequential input/output (I/O) and standard interfaces for communication between processors (parallel or otherwise) and storage devices. In typical use, clients direct a request for data to an HPSS server. The HPSS server directs the network-attached storage devices or servers to transfer data directly, sequentially or in parallel to the client node(s) through the high speed data transfer network. TCP/IP sockets and IPI-3 over High Performance Parallel Interface (HIPPI) are being utilized today; Fibre Channel Standard (FCS) with IPI-3 or SCSI, or Asynchronous Transfer Mode (ATM) will also be supported in the future [3,20,22]. Through its parallel storage support by data striping HPSS will continue to scale upward as additional storage devices and controllers are added to a site installation.

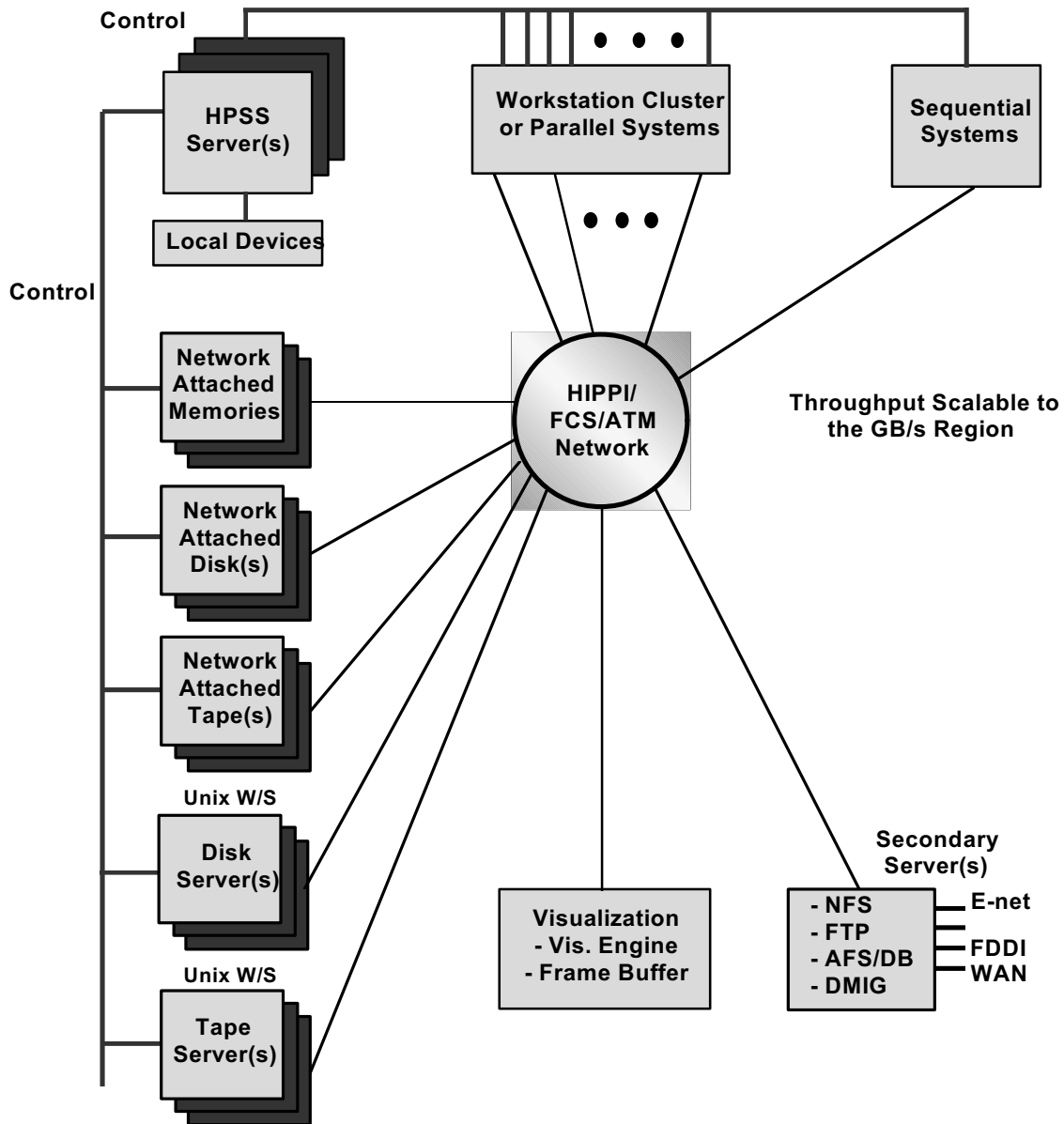


Figure 1 - Example of the type of configuration HPSS is designed to support

The key objectives of HPSS are now described.

Scalability

A major driver for HPSS is to develop a scalable, distributed, high performance storage management system. HPSS is designed to scale in several dimensions.

The HPSS I/O architecture is designed to provide I/O performance scaling by supporting parallel I/O through software striping [1]. The system will support application data transfers from megabytes to gigabytes per second with total system throughput of many gigabytes per second. Data object number and size must scale to support billions of data objects, each potentially terabytes or larger in size, for total storage capacities in petabytes. This is accomplished through 64-bit metadata fields and scalable organization of system metadata. The system also is required to scale geographically to support distributed systems with hierarchies of hierarchical storage systems. Multiple storage systems located in different areas must integrate into a single logical system accessible by personal computers, workstations, and supercomputers. These requirements are accomplished using a client/server architecture, the use of OSF's DCE as its distributed infrastructure, support for distributed file system interfaces and multiple servers. HPSS also supports a scalable storage object name service capable of managing millions of directories and the ability to support hundreds to thousands of simultaneous clients. The latter is achieved through the ability to multitask, multiprocess and replicate the HPSS servers.

Modularity and APIs

The HPSS architecture is highly modular. Each replicable software component is responsible for a set of storage objects, and acts as a service provider for those objects. The IEEE Reference Model, on which the HPSS design is based, provides the modular layered functionality (see *Figure 2*) [6,9]. The HPSS software components are loosely coupled, with open application program interfaces (APIs) defined at each component level. Most users will access HPSS at its high level interfaces—currently client API, FTP (both parallel and sequential), NFS, Parallel File System (PFS), with AFS/DFS, Unix Virtual File System (VFS), and Data Management Interface Group (DMIG) interfaces in the future) [11,15,18,19]. However, APIs are available to the underlying software components for applications, such as large scale data management, digital library or video-on-demand requiring high performance or special services. This layered architecture affords the following advantages:

- **Replacement of selected software components**—As new and better commercial software and hardware components became available, an installation can add or replace existing components. For example, an installation might add or replace

Physical Volume Repositories, Movers or the HPSS Physical Volume Library with other commercially available products.

- **Support of applications direct access to lower level services**—The layered architecture is designed to accommodate efficient integration of different applications such as digital library, object store, multimedia, and data management systems. Its modularity will enable HPSS to be embedded transparently into the large distributed information management systems that will form the information services in the emerging national information infrastructure. Support for different name spaces or data organizations is enabled through introduction of new Name Servers and data management applications.

Portability and Standards

Another important design goal is portability to many vendor's platforms to enable OEM and multivendor support of HPSS. HPSS has been designed to run under Unix requiring no kernel modifications, and to use standards based protocols, interfaces, and services where applicable. HPSS is written in ANSI C, and uses POSIX functions to enhance software portability. Use of existing commercial products for many of the infrastructure services supported on multiple-vendor platforms enables portability, while also providing market proven dependability. Open Software Foundation (OSF) Distributed Computing Environment (DCE), Transarc's Encina transaction manager [8], Kinesix SAMMI and X-windows are being used by HPSS because of their support across multiple vendor platforms, in addition to the rich set of functionality provided. The HPSS component APIs have been turned over to the IEEE Storage System Standards Working Group as a basis for its standards activities.

Reliability and Recovery

Reliable and recoverable storage of data is mandatory for any storage system. HPSS supports several mechanisms to facilitate this goal. The client-server interactions between HPSS software components have been designed to be based on atomic transactions in order to maintain system state consistency [14]. Within the scope of a given request, a transaction may be established so that an abort (or commit) in one component will cause the other participating components to abort (or commit). The HPSS Metadata Manager is fully integrated with its Transaction Manager. Following an abort, the non-volatile file and name space metadata changes within the scope of the transactions will automatically be rolled back. For recovery purposes, mirroring of the storage object and name space metadata is supported. The HPSS architecture will also support data mirroring if desired in a future release.

Support is also provided to recover from failed devices and bad media. An administrator interface is provided to place a device off line. Once the device has been repaired, it may

then be placed back on line. For bad media, an application interface is provided to move storage segments from a virtual volume to a new virtual volume.

The HPSS software components execute in a distributed manner. Should a processor fail, any of the HPSS software components may be moved to another platform. Component services are registered with the DCE Cell Directory Service (CDS) so that components may locate the services. Each component has also been designed to perform reconnect logic when a connection to a peer component fails. Connection context is maintained by selected components. When a connection context is established, a keep-alive activity is started to detect broken connections. A server may use the context information associated with a broken connection to perform any necessary clean up.

Security and Privacy

HPSS uses DCE and POSIX security and privacy mechanisms for authentication, access control lists, permissions and security labels. Security policy is handled by a separate policy module. Audit trails are also supported. Further, HPSS design and implementation use a rigorous software engineering methodology which support its reliability and maintainability.

Storage System Management

HPSS has a rich set of storage system management services for operators and system administrators based on managed object definitions. The application programming interface supports monitoring, reporting and controlling operations (see Appendix A).

Software Components

The HPSS software components are shown *Figure 2*. The shaded boxes are defined in the IEEE Mass Storage Reference Model: version 5 [9].

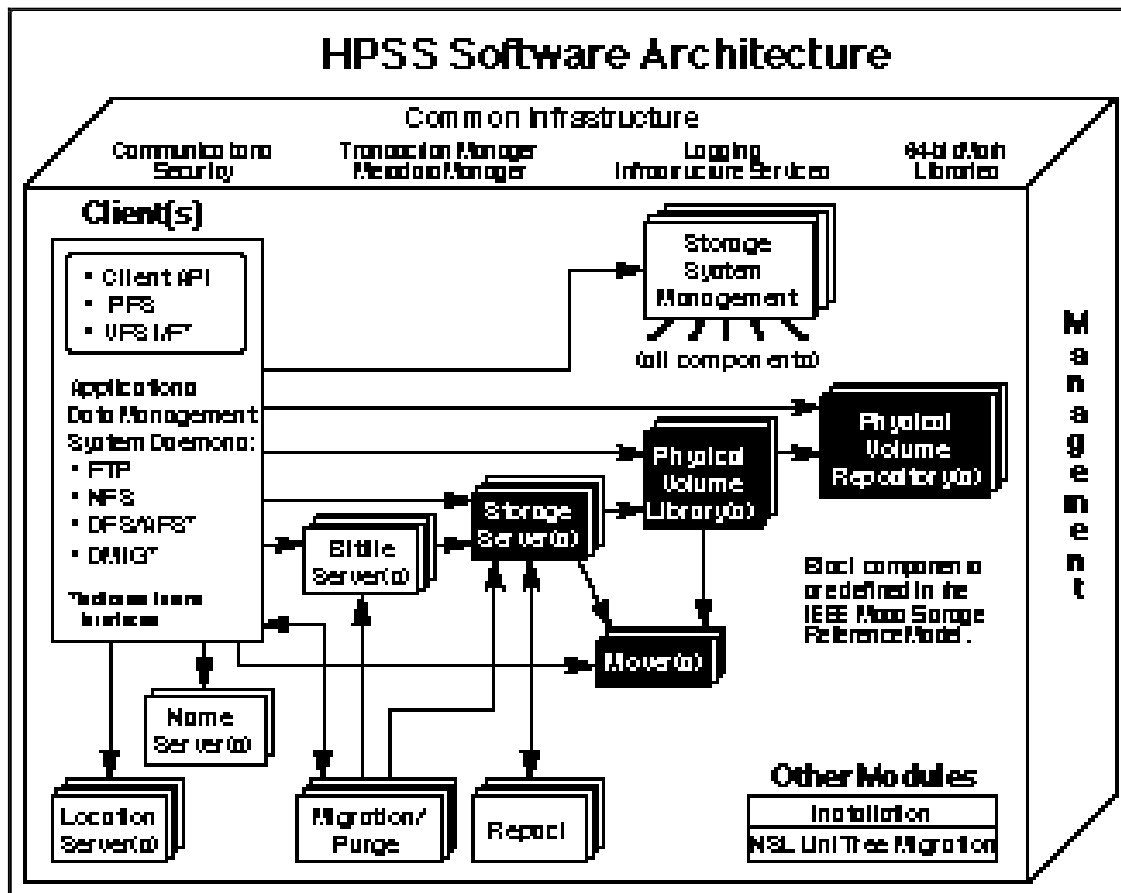


Figure 2 - Software Model Diagram

This section outlines the function of each component.

Infrastructure

HPSS design is based upon a well-formed industry standard infrastructure. The key infrastructure components are now outlined.

Distributed Computing Environment

HPSS uses OSF's DCE as the base infrastructure for its distributed architecture [17]. This standards-based framework will enable the creation of distributed storage systems for a national information infrastructure capable of handling gigabyte-terabyte-class files at gigabyte per second data transfer rates.

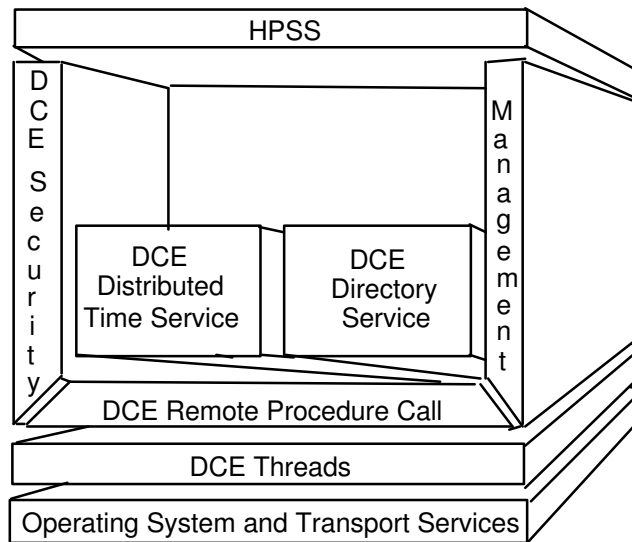


Figure 3 - HPSS DCE Architecture Infrastructure

DCE was selected because of its wide adoption among vendors and its near industry-standard status. HPSS uses the DCE Remote Procedure Call (RPC) mechanism for control messages and DCE Threads for multitasking. The DCE threads package is vital for HPSS to serve large numbers of concurrent users and to enable multiprocessing of its servers. HPSS also uses the DCE Security, Cell Directory, and Time services. A library of DCE convenience functions was developed for use in HPSS.

Transaction Management

Requests to HPSS to perform actions such as creating bitfiles or accessing file data results in client/server interactions between software components. Transaction integrity is required to guarantee consistency of server state and metadata in case a particular component should fail. As a result, a transaction manager was required by HPSS. Encina, from Transarc, was selected by the HPSS project as its transaction manager [8]. This selection was based on functionality, its use of DCE, and multi-platform vendor support. Encina provides begin-commit-abort semantics, distributed two-phase commit, and nested transactions. In addition, Transaction RPCs (TRPCs), which extend DCE RPCs with transaction semantics, are provided. For recovery purposes, Encina uses a write-ahead log for storing transaction outcomes and updates to recoverable metadata. Mirroring of data is also provided.

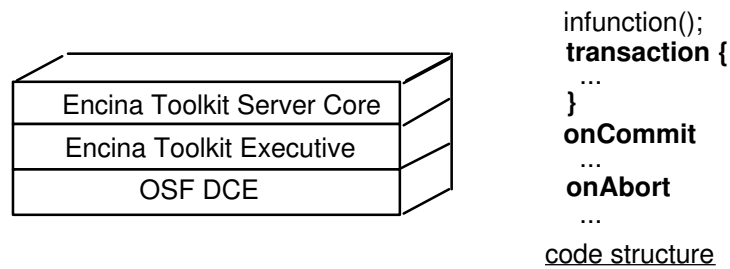


Figure 4 - Encina Components

Metadata Management

Each HPSS software component has system metadata associated with the objects it manages. Each server with non-volatile metadata requires the ability to reliably store its metadata. It is also required that metadata management performance be scalable as the number of object instances grow. In addition, access to metadata by primary and secondary keys is required. The Structured File Server (SFS), an Encina optional product, was selected by the HPSS project as its metadata manager. SFS provides B-tree clustered file records, record and field level access, primary and secondary keys, and automatic byte ordering between machines. SFS is also fully integrated with the Encina transaction manager. As a result, SFS provides transaction consistency and data recovery from transaction aborts. For reliability purposes, HPSS metadata stored in SFS is mirrored. A library of metadata manager convenience functions for retrieving, adding, updating, and deleting metadata for each of the HPSS components was developed.

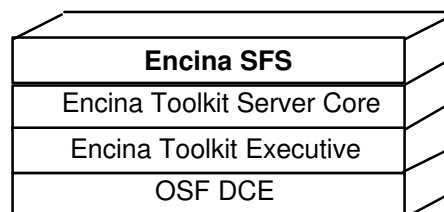


Figure 5 - Structured File Server (SFS)

Security

The security components of HPSS provide authentication, authorization, enforcement, and audit capabilities for the HPSS components. Authentication is responsible for guaranteeing that a principal is the entity that is claimed, and that information received from an entity is from that entity. Authorization is responsible for enabling an authenticated entity access to an allowed set of resources and objects. Authorization enables end user access to HPSS directories and bitfiles. Enforcement is responsible for guaranteeing that operations are restricted to the authorized set of operations.

Enforcement applies to end user access to bitfiles. Audit is responsible for generating a log of security relevant activity. HPSS security libraries utilize DCE and DCE security. The authentication service, which is part of DCE, is based on Kerberos v5. The following figure depicts how HPSS security fits with DCE and Kerberos.

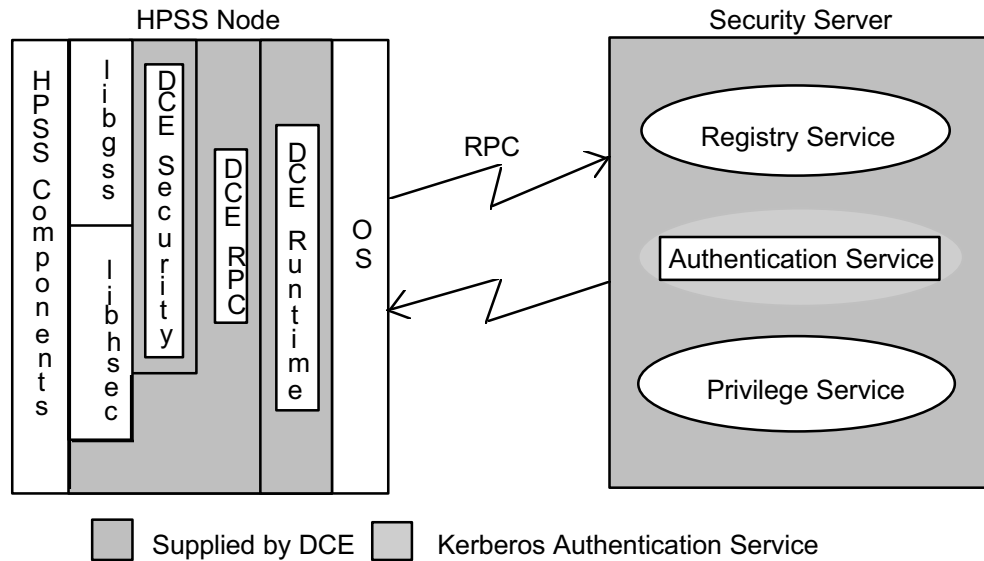


Figure 6 - HPSS Security

Communication

The control path communications between HPSS components is through DCE RPCs or Encina transaction RPCs. For data path communication, the HPSS Mover(s) currently utilize either Sockets or IPI-3 (over HIPPI) libraries. Future support is planned for IPI-3 and SCSI over Fibre Channel Standard and TCP/IP over ATM. A special parallel data transfer library has been developed. This library allows data to be transferred across many parallel data connections. The library transfers data headers that identify the data that follows. This allows data to be sent and arrive in any order on the parallel paths.

Logging

The HPSS logger is used to record alarms, events, requests, security audit records, accounting records, and trace information from the HPSS components. A central log is maintained which contains records from all HPSS components. A local log of activity from components on each HPSS node is also supported. When the central log fills, it will switch to a secondary log file. A configuration option allows the filled log to be automatically archived to HPSS. A delog function is provided to extract and format log records. Delog options support filtering by time interval, record type, server, and user.

64 Bit Arithmetic Libraries

HPSS supports file sizes up to 2**64 bytes. Many vendor platforms support only 32 bit integer arithmetic. In order to support large file sizes and large numbers of objects on 32 bit platforms, a library of 64 bit arithmetic functions has been developed. The functions support both big endian and little endian I/O architectures.

Interfaces

HPSS supports several high-level interfaces: currently Client API, FTP (both standard and parallel), and NFS, with DFS/AFS, DMIG, and VFS planned for future releases.

Client API

The HPSS Client file server API mirrors the POSIX file system interface specification where possible. The Client API also supports extensions to allow the programmer to take advantage of the specific features provided by HPSS (e.g., class-of-service, storage/access hints passed at file creation and support for parallel data transfers).

FTP (standard and parallel)

HPSS provides a standard FTP server interface to transfer files from HPSS to a local file system. Parallel FTP, an extension and superset of standard FTP, has been implemented to provide high performance data transfers to client systems. The standard FTP protocol supports third-party data transfer through separation of the data transfer and control paths, but it does not offer parallel data paths [11]. HPSS modified and augmented the standard client FTP file retrieval and storage functions to offer parallel data paths for HPSS data transfers. This approach provides high performance FTP transfers to the client while still supporting the FTP command set. Additional commands have been added to support parallel transfer. This work will be submitted to the Internet Engineering Task Force for standardization.

NFS

The NFS V2 Server interface for HPSS provides transparent access to HPSS name space objects and bitfile data for client systems from both the native HPSS and the Network File System V2 service. The NFS V2 Server translates standard NFS calls into HPSS control calls and provides data transfers for NFS read and write requests. The NFS V2 Server handles optimization of data movement requests by the caching of data and control information. If the server machine crashes, the NFS V2 Server is in charge of recovery of all cached data at the time of the crash. The NFS V2 Server will also recover

when HPSS crashes. Before NFS clients can request NFS services, they must mount an exported HPSS directory by calling the Mount daemon mount API. Support for NFS V3 is planned for a future release.

Parallel File System

HPSS provides the capability to act as an external hierarchical file system to vendor Parallel File Systems (PFS). The first implementation supports the IBM SPx PIOFS. Early deployment is also planned for Intel Paragon and Meiko PFS integration with HPSS.

Name Server (NS)

The Name Server maps a file name to an HPSS object. The Name Server provides a POSIX view of the name space which is a hierarchical structure consisting of directories, files, and links. File names are human readable ASCII strings. Namable objects are any object identified by HPSS Storage Object IDs. The commonly named objects are bitfiles, directories, or links. In addition to mapping names to unique object identifiers, the Name Server provides access verification to objects. POSIX Access Control Lists (ACLs) are supported for the name space objects. A key requirement of the Name Server is to be able to scale to millions of directories and greater than a billion name space entries.

Bitfile Server (BFS)

The Bitfile Server provides the POSIX file abstraction to its clients. A logical bitfile is an uninterpreted bit string. HPSS supports bitfile sizes up to 2^{64} bytes. A bitfile is identified by a Bitfile Server generated name called a bitfile-id. Mapping of a human readable name to the bitfile id is provided by a Name Server external to the Bitfile Server. Clients may reference portions of a bitfile by specifying the bitfile-id and a starting address and length. The writes and reads to a bitfile are random and the writes may leave holes where no data has been written. The Bitfile Server supports both sequential and parallel read and write of data to bitfiles. In conjunction with Storage Servers, the Bitfile Server maps logical portions of bitfiles onto physical storage devices.

Storage Server (SS)

The Storage Server provides a hierarchy of storage objects: logical storage segments, virtual volumes and physical volumes. All three layers of the Storage Server can be accessed by appropriately privileged clients. The server translates references to storage segments into references to virtual volume and finally into physical volume references. It also schedules the mounting and dismounting of removable media through the Physical Volume Library. The Storage Server in conjunction with the Mover have the main responsibility for orchestration of HPSS's parallel I/O operations.

The storage segment service is the conventional method for obtaining and accessing HPSS storage resources. The Storage Server maps an abstract storage space, the storage segment, onto a virtual volume, resolving segment addresses as required. The client is presented with a storage segment address space, with addresses from 0 to N-1, where N is the byte length of the segment. Segments can be opened, created, read, written, closed and deleted. Characteristics and information about segments can be retrieved and changed.

The virtual volume service is the method provided by the Storage Server to group physical storage volumes. The virtual volume service supports striped volumes today and mirrored volume in a future release. Thus, a virtual volume can span multiple physical volumes. The Storage Server maps the virtual volume address space onto the component physical volumes in a fashion appropriate to the grouping. The client is presented with a virtual volume that can be addressed from 0 to N-1, where N is the byte length of the virtual volume. Virtual volumes can be mounted, created, read, written, unmounted and deleted. Characteristics of the volume can be retrieved and in some cases, changed.

The physical volume service is the method provided by the storage server to access the physical storage volumes in HPSS. Physical volumes can be mounted, created, read, written, unmounted and deleted. Characteristics of the volume can be retrieved and in some cases, changed.

Repack runs as a separate process. It provides defragmentation of physical volumes. Repack utilizes a Storage Server provided function which moves storage segments to a different virtual volume.

Mover (Mvr)

The Mover is responsible for transferring data from a source device(s) to a sink device(s). A device can be a standard I/O device with geometry (e.g., a tape or disk), or a device without geometry (e.g., network, memory). The Mover also performs a set of device control operations. Movers perform the control and transfer of both sequential and parallel data transfers.

The Mover consists of several major parts: Mover parent task, Mover listen task/request processing task, Data Movement, Device control, and System Management.

The Mover parent task performs Mover initialization functions, and spawns processes to handle the Mover's DCE communication, data transfer connections, as well as the Mover's functional interface. The Mover listen task listens on a well-known TCP port for incoming connections to the Mover, spawns request processing tasks, and monitors completion of those tasks. The request processing task performs initialization and return functions common to all Mover requests. Data movement supports client requests to transfer data to or from HPSS. Device control supports querying the current device

read/write position, changing the current device read/write position, loading a physical volume into a drive, unloading a physical volume from a drive, flushing data to the media, writing a tape mark, loading a message to a device's display area, reading a media label, writing a media label, and zeroing a portion of disk. System management supports querying and altering device characteristics and overall Mover state.

Physical Volume Library (PVL)

The PVL manages all HPSS physical volumes. Clients can ask the PVL to mount and dismount sets of physical volumes. Clients can also query the status and characteristics of physical volumes. The PVL maintains a mapping of physical volume to cartridge and a mapping of cartridge to PVR. The PVL also controls all allocation of drives. When the PVL accepts client requests for volume mounts, the PVL allocates resources to satisfy the request. When all resources are available, the PVL issues commands to the PVR(s) to mount cartridges in drives. The client is notified when the mount has completed.

The Physical Volume Library consists of two major parts: Volume mount service and Storage system management service.

The volume mount service is provided to clients such as a Storage Server. Multiple physical volumes belonging to a virtual volume may be specified as part of a single request. All of the volumes will be mounted before the request is satisfied. All volume mount requests from all clients are handled by the PVL. This allows the PVL to prevent multiple clients from deadlocking when trying to mount intersecting sets of volumes. The standard mount interface is asynchronous. A notification is provided to the client when the entire set of volumes has been mounted. A synchronous mount interface is also provided. The synchronous interface can only be used to mount a single volume, not sets of volumes. The synchronous interface might be used by a non-HPSS process to mount cartridges which are in a tape library, but not part of the HPSS system.

The storage system management service is provided to allow a management client control over HPSS tape repositories. Interfaces are provided to import, export, and move volumes. When volumes are imported into HPSS, the PVL is responsible for writing a label to the volume. This label can be used to confirm the identity of the volume every time it is mounted. Management interfaces are also provided to query and set the status of all hardware managed by the PVL (volumes, drives, and repositories).

Physical Volume Repository (PVR)

The PVR manages all HPSS supported robotics devices and their media such as cartridges. Clients can ask the PVR to mount and dismount cartridges. Every cartridge in HPSS must be managed by exactly one PVR. Clients can also query the status and characteristics of cartridges.

The Physical Volume Repository consists of these major parts: Generic PVR service, and support for devices such as Ampex, STK, and 3494/3495 robot services, as well as an operator mounted device service.

The generic PVR service provides a common set of APIs to the client regardless of the type of robotic device being managed. Functions to mount, dismount, inject and eject cartridges are provided. Additional functions to query and set cartridge metadata are provided. The mount function is asynchronous. The PVR calls a well-known API in the client when the mount has completed. For certain devices, like operator mounted repositories, the PVR will not know when the mount has completed. In this case it is up to the client to determine when the mount has completed. The client may poll the devices or use some other method. When the client determines a mount has completed, the client should notify the PVR using one of the PVR's APIs. All other PVR functions are synchronous. The generic PVR maintains metadata for each cartridge managed by the PVR. The generic PVR interface calls robotics vendor supplied code to manage specific robotic devices.

The operator mounted device service manages a set of cartridges that are not under the control of a robotics device. These cartridges are mounted to a set of drives by operators. The Storage System Manager is used to inform the operators when mount operations are required.

Storage System Management (SSM)

The HPSS SSM architecture is based on the ISO managed object architecture [10,12]. The Storage System Manager (SSM) monitors and controls the available resources of the HPSS storage system in ways that conform to the particular management policies of a given site. Monitoring capabilities include the ability to query the values of important management attributes of storage system resources as well as an ability to receive notifications of alarms and other significant system events. Controlling capabilities include the ability to set the values of management attributes of storage system resources and storage system policy parameters. Additionally, SSM can request that specific operations be performed on resources within the storage system, such as adding and deleting logical or physical resources. The operations performed by SSM are usually accomplished through standard HPSS server APIs.

SSM management roles cover a wide spectrum, including configuration aspects of installation, creating new volumes, initialization, operations, and termination tasks. SSM can provide management capabilities to a range of clients, including site administrators, systems administrators, operations personnel, complex graphical user interface (GUI) management environments, and independent management applications responsible for tasks such as purges, migration, and reclamation. Some of the functional areas of SSM include fault management, configuration management, security management, accounting management, and performance management.

SSM consists of these major parts: SSM Graphical User Interface (SAMMI GUI Displays), SAMMI Data Server, and System Manager.

The SSM Graphical User Interface allows operators, administrators, and users to interactively monitor and control the HPSS storage system. Kinesix's SAMMI product is used to provide the HPSS GUI services. SAMMI is built on X-windows and OSF's Motif. It provides mechanisms to simplify screen design and data management services for screen fields. Standard Motif widgets such as menus, scrollbar lists, and buttons are used. In addition SAMMI specific widgets such as dials, gauges, and bar charts are used for informational and statistical data.

The SAMMI Data Server is a client to the System Manager and a server to the SAMMI Runtime Display Manager. The SAMMI Data Server is the means by which data is acquired and fed to the SAMMI Displays.

The Storage System Manager is a client to the HPSS servers and a server to the SAMMI Data Server and other external clients wishing to perform management specific operations. It interfaces to the managed objects defined by the HPSS servers.

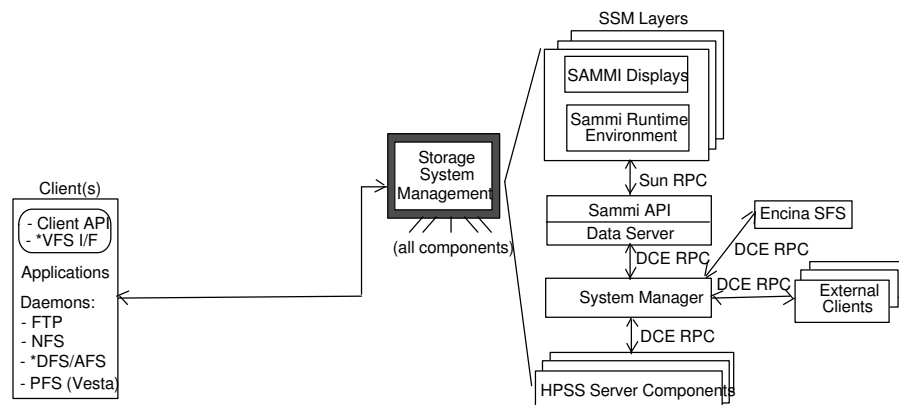


Figure 7 - Storage System Management

Migration - Purge

The Migration-Purge server provides hierarchical storage management for HPSS through migration and caching of data between devices. There are two types of migration and caching: disk migration and caching and tape migration and caching. Multiple storage hierarchies are supported by HPSS [2]. Data is cached to the highest level (fastest) device in a given hierarchy when accessed and migrated when inactive and space is required.

The main purpose of disk migration is to free up the disk storage. This type of migration contains two functions; migration and purge. Migration selects the qualified bitfiles and copies these bitfiles to the next storage level defined in the hierarchy. Purge later frees the original bitfiles from the disk storage.

The main purpose of tape migration is to free up tape volumes, and not just migrate bitfiles. The active bitfiles in the target virtual volumes are moved laterally to the free tape volumes in the same storage level. The inactive bitfiles in the target virtual volumes are migrated to the free tape volumes in the next storage level.

The HPSS component client APIs provide the vehicle for the Storage System Manager to request the server to start migration and purge whenever it is necessary. The migration-purge server is set up to run migration periodically with the time interval specified in the migration policy. In addition, the server will start the migration and purge to run automatically if the free space of a storage class is below the percentage specified in the migration-purge policy.

Other

Installation

Installation software is provided for system administrators to install/update HPSS, and perform the initial configuration of HPSS following installation. The full HPSS system is first installed to an installation node. Selected HPSS software components may then be installed (using the remote installation feature) from the installation node to the other nodes where HPSS components will be executed.

NSL-UniTree Migration

HPSS, through its support of parallel storage, provides significant improvements in I/O rates and storage capacity over existing storage systems software. In transitioning from existing systems, a migration path is required. The migration path should be transparent to end users of the storage system. The capability to migrate from NSL UniTree to HPSS is provided. The migration software handles both file metadata and actual data. Utilities convert the file metadata (e.g., storage maps, virtual volume data, physical volume data), and name space metadata from UniTree format to HPSS format. Actual data is not moved. The HPSS Mover software contains additional read logic to recognize NSL UniTree data formats when an NSL UniTree file is accessed. Utilities to support migration from other legacy storage systems will also be provided as required.

Accounting

HPSS provides interfaces to collect accounting information (initially storage space utilization). These interfaces may be used by site specific programs to charge for data storage. SSM provides user interfaces to run the accounting collection utility, change account numbers and change the account code assigned to storage objects.

Summary and Status

We have described the key objectives, features and components of the HPSS architecture. At the time this paper is being written, December 1994, HPSS Release 1 (R1) is in integration testing and planning for its early deployment at several sites has begun. R1 contains all the basic HPSS components and services and supports parallel tape. It is targeted at MPP environments with existing parallel disk services. Much of the coding for Release 2 (R2) has been completed also. R2 adds support for parallel disks, migration and caching between levels of the hierarchy and other functionality. R2 will be a complete stand-alone system and is targeted for third quarter 1995.

We demonstrated, HPSS at Supercomputing 1994 with R1 and early R2 capabilities of parallel disks, and tape access (Ampex D2, IBM NTP and 3490), to an IBM SP2, IBM RS 6000, PsiTech framebuffer, and Sony high-resolution monitor over a NSC HIPPI switch. HPSS R1 is on order 95K lines of executable source code and R2 is expected to add on another 50K lines of executable source code.

Our experience indicates that the architectural choices of basing the system on the IEEE Reference Model, use of an industry defacto standard infrastructure based on OSF DCE and Transarc Encina, and use of other industry standards such as POSIX, C, Unix, ISO managed object model for Storage System Management and standard communication protocols is sound. This foundation plus the software engineering methodology employed, we believe, positions HPSS for a long and useful life for both scientific and commercial high performance environments.

Acknowledgments

We wish to acknowledge the many discussions and shared design, implementation, and operation experiences with our colleagues in the National Storage Laboratory collaboration, the IEEE Mass Storage Systems and Technology Technical Committee, the IEEE Storage System Standards Working Group, and in the storage community. Specifically we wish to acknowledge the people on the HPSS Technical Committee and Development Teams. At the risk of leaving out a key colleague in this ever-growing collaboration, the authors wish to acknowledge Dwight Barrus, Ling-Ling Chen, Ron Christman, Danny Cook, Lynn Kluegel, Tyce McLarty, Christina Mercier, and Bart Parlman from LANL; Larry Berdahl, Jim Daveler, Dave Fisher, Mark Gary, Steve Louis, Donna Mecozzi, Jim Minton, and Norm Samuelson from LLNL; Marty Barnaby, Rena Haynes, Hilary Jones, Sue Kelly, and Bill Rahe from SNL; Randy Burris, Dan Million, Daryl Steinert, Vicky White, and John Wingenbach from ORNL; Donald Creig Humes, Juliet Pao, Travis Priest and Tim Starrin from NASA LaRC; Andy Hanushevsky, Lenny Silver, and Andrew Wyatt from Cornell; and Paul Chang, Jeff Deutsch, Kurt Everson, Rich Ruef, Tracy Tran, Terry Tyler, and Benny Wilbanks from IBM U.S. Federal and its contractors.

This work was, in part, performed by the Lawrence Livermore National Laboratory, Los Alamos National Laboratory, Oak Ridge National Laboratory, and Sandia National Laboratories, under auspices of the U.S. Department of Energy Cooperative Research and Development Agreements, by Cornell, Lewis Research Center and Langley Research Center under auspices of the National Aeronautics and Space Agency and by IBM U.S. Federal under Independent Research and Development and other internal funding.

References

1. Berdahl, L., ed., "Parallel Transport Protocol," draft proposal, available from Lawrence Livermore National Laboratory, Dec. 1994.
2. Buck, A. L., and R. A. Coyne, Jr., "Dynamic Hierarchies and Optimization in Distributed Storage System," Digest of Papers, Eleventh IEEE Symposium on Mass Storage Systems, Oct. 7-10, 1991, IEEE Computer Society Press, pp. 85-91.
3. Christensen, G. S., W. R. Franta, and W. A. Petersen, "Future Directions of High-speed Networks for Distributed Storage Environments," Digest of Papers, Eleventh IEEE Symposium on Mass Storage Systems, Oct. 7-10, 1991, IEEE Computer Society Press, pp. 145-148.
4. Collins, B., et al., "Los Alamos HPDS: High-Speed Data Transfer," Proc. Twelfth IEEE Symposium on Mass Storage Systems, Monterey, April 1993.
5. Coyne, R. A., H. Hulen, and R. W. Watson, "The High Performance Storage System," Proc. Supercomputing 93, Portland, IEEE Computer Society Press, Nov. 1993.
6. Coyne, R. A. and H. Hulen, "An Introduction to the Mass Storage System Reference Model, Version 5," Proc. Twelfth IEEE Symposium on Mass Storage Systems, Monterey, April 1993.
7. Coyne, R. A., H. Hulen, and R. W. Watson, "Storage Systems for National Information Assets," Proc. Supercomputing 92, Minneapolis, Nov. 1992, pp. 626-633.
8. Dietzen, Scott, Transarc Corporation, "Distributed Transaction Processing with Encina and the OSF/DCE", Sept. 1992, 22 pages.
9. IEEE Storage System Standards Working Group (SSSWG) (Project 1244), "Reference Model for Open Storage Systems Interconnection, Mass Storage Reference Model Version 5," Sept. 1994. Available from the IEEE SSSWG Technical Editor Richard Garrison, Martin Marietta (215) 532-6746

10. "Information Technology - Open Systems Interconnection - Structure of Management Information - Part 4: Guidelines for the Definition of Management Objects," ISO/IEC 10165-4, 1991.
11. Internet Standards. The official Internet standards are defined by RFC's (TCP protocol suite). RFC 783; TCP standard defined. RFC 959; FTP protocol standard. RFC 1068; FTP use in third-party transfers. RFC 1094; NFS standard defined. RFC 1057; RPC standard defined.
12. ISO/IEC DIS 10040 Information Processing Systems - Open Systems Interconnection - Systems Management Overview, 1991.
13. Katz, R. H., "High Performance Network and Channel-Based Storage," *Proceedings of the IEEE*, Vol. 80, No. 8, pp. 1238-1262, August 1992.
14. Lampson, B. W., M. Paul, and H. J. Siegart (eds.), "Distributed Systems - Architecture and Implementation," Berlin and New York: Springer-Verlag, 1981.
15. Morris, J. H., et al., "Andrew: A Distributed Personal Computing Environment," *Comm. of the ACM*, Vol. 29, No. 3, March 1986.
16. Nelson, M., et al., "The National Center for Atmospheric Research Mass Storage System," *Digest of Papers, Eighth IEEE Symposium on Mass Storage Systems*, May 1987, pp. 12-20.
17. Open Software Foundation, *Distributed Computing Environment Version 1.0 Documentation Set*. Open Software Foundation, Cambridge, Mass. 1992.
18. OSF, *File Systems in a Distributed Computing Environment, White Paper, Open Software Foundation*, Cambridge, MA, July 1991.
19. Sandberg, R., et al., "Design and Implementation of the SUN Network Filesystem," *Proc. USENIX Summer Conf.*, June 1989, pp. 119-130.
20. Tolmie, D. E., "Local Area Gigabit Networking," *Digest of Papers, Eleventh IEEE Symposium on Mass Storage Systems*, Oct. 7-10, 1991, IEEE Computer Society Press, pp. 11-16.
21. Watson, R. W., R. A. Coyne, "The National Storage Laboratory: Overview and Status," *Proc. Thirteenth IEEE Symposium on Mass Storage Systems*, Annecy France, June 12-15, 1994, pp. 39-43.
22. Witte, L. D., "Computer Networks and Distributed Systems," *IEEE Computer*, Vol. 24, No. 9, Sept. 1991, pp. 67-77.