

Distributed File Systems

Vaidé Narváez

Computer Information Systems

July 21st, 2010

	<i>Sharing</i>	<i>Persistence</i>	<i>Distributed caches/replicas</i>	<i>Consistency maintenance</i>	<i>Example</i>
Main memory	✗	✗	✗	1	RAM
File system	✗	✓	✗	1	UNIX file system
Distributed file system	✓	✓	✓	✓	Sun NFS
Web	✓	✓	✓	✗	Web server
Distributed shared memory	✓	✗	✓	✓	Ivy (DSM, Ch. 18)
Remote objects (RMI/ORB)	✓	✗	✗	1	CORBA
Persistent object store	✓	✓	✗	1	CORBA Persistent Object Service
Peer-to-peer storage system	✓	✓	✓	2	OceanStore (Ch. 10)

Types of consistency:

1: strict one-copy. : ✓ slightly weaker guarantees. 2: considerably weaker guarantees.

Characteristics of file systems

- Organization
- Storage
- Retrieval
- Naming
- Sharing
- Protection of files

Characteristics of file systems

- Data and attributes

■ Data and attributes

File length
Creation timestamp
Read timestamp
Write timestamp
Attribute timestamp
Reference count
Owner
File type
Access control list

- Data and attributes

File length
Creation timestamp
Read timestamp
Write timestamp
Attribute timestamp
Reference count
Owner
File type
Access control list

© Addison-Wesley 2005

- A **directory** is a file, often of a special type, that provides mapping from text names to internal file identifiers

Layered module structure of FS

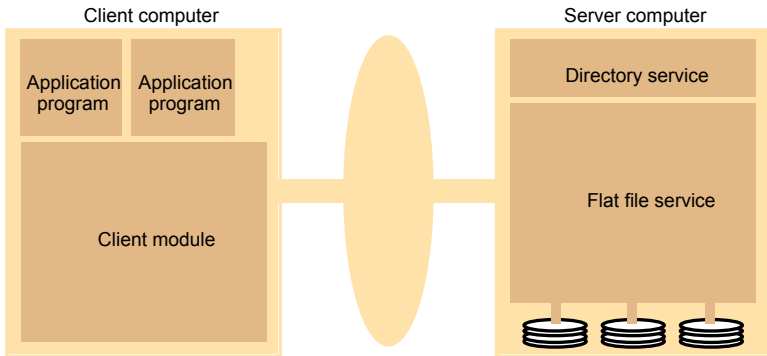
Directory module:	relates file names to file IDs
File module:	relates file IDs to particular files
Access control module:	checks permission for operation requested
File access module:	reads or writes file data or attributes
Block module:	accesses and allocates disk blocks
Device module:	disk I/O and buffering

© Addison-Wesley 2005

<i>filedes</i> = <i>open</i> (<i>name</i> , <i>mode</i>)	Opens an existing file with the given <i>name</i> .
<i>filedes</i> = <i>creat</i> (<i>name</i> , <i>mode</i>)	Creates a new file with the given <i>name</i> .
	Both operations deliver a file descriptor referencing the open file. The <i>mode</i> is <i>read</i> , <i>write</i> or both.
<i>status</i> = <i>close</i> (<i>filedes</i>)	Closes the open file <i>filedes</i> .
<i>count</i> = <i>read</i> (<i>filedes</i> , <i>buffer</i> , <i>n</i>)	Transfers <i>n</i> bytes from the file referenced by <i>filedes</i> to <i>buffer</i> .
<i>count</i> = <i>write</i> (<i>filedes</i> , <i>buffer</i> , <i>n</i>)	Transfers <i>n</i> bytes to the file referenced by <i>filedes</i> from <i>buffer</i> .
	Both operations deliver the number of bytes actually transferred and advance the read-write pointer.
<i>pos</i> = <i>lseek</i> (<i>filedes</i> , <i>offset</i> , <i>whence</i>)	Moves the read-write pointer to <i>offset</i> (relative or absolute, depending on <i>whence</i>).
<i>status</i> = <i>unlink</i> (<i>name</i>)	Removes the file <i>name</i> from the directory structure. If the file has no other names, it is deleted.
<i>status</i> = <i>link</i> (<i>name1</i> , <i>name2</i>)	Adds a new name (<i>name2</i>) for a file (<i>name1</i>).
<i>status</i> = <i>stat</i> (<i>name</i> , <i>buffer</i>)	Gets the file attributes for file <i>name</i> into <i>buffer</i> .

- transparency:
 - ▶ access
 - ▶ location
 - ▶ mobility
 - ▶ performance
 - ▶ scaling
- concurrent file updates
- file replication
- consistency
- fault tolerance
- hardware and os heterogeneity
- security
- efficiency

File service architecture



- Flat file service
 - ▶ unique file identifiers (UFID)
- Directory service
 - ▶ map names to UFIDs
- Client module
 - ▶ integrate/extend flat file and directory services
 - ▶ provide a common application programming interface (can emulate different file interfaces)
 - ▶ stores location of flat file and directory services

- RPC used by client modules
 - ▶ not by user-level programs (which use client modules)
- More fault tolerant compared to UNIX
 - ▶ Repeatable (idempotent) operations
 - ▶ at-least-once semantics
 - ▶ no open (hence close) so no state to remember
 - ▶ specify starting location and UFID (from directory service) in Read/Write
 - ▶ Stateless server
 - ▶ Can be restarted without the server or client restoring any state information

Flat file service operations

<i>Read(FileId, i, n) -> Data</i> – throws <i>BadPosition</i>	If $1 \leq i \leq \text{Length}(\text{File})$: Reads a sequence of up to n items from a file starting at item i and returns it in <i>Data</i> .
<i>Write(FileId, i, Data)</i> – throws <i>BadPosition</i>	If $1 \leq i \leq \text{Length}(\text{File})+1$: Writes a sequence of <i>Data</i> to a file, starting at item i , extending the file if necessary.
<i>Create() -> FileId</i>	Creates a new file of length 0 and delivers a UFID for it.
<i>Delete(FileId)</i>	Removes the file from the file store.
<i>GetAttributes(FileId) -> Attr</i>	Returns the file attributes for the file.
<i>SetAttributes(FileId, Attr)</i>	Sets the file attributes (only those attributes that are not shaded in Figure 8.3).

- UNIX checks access rights when a file is opened
 - ▶ subsequent checks during read/write are not necessary

- distributed environment
 - ▶ access rights are checked at the server

 - ▶ stateless approaches
 1. access check whenever UFID is issued
 - client gets an encoded "capability" (who can access and how)
 - capability is submitted with each subsequent request
 2. access check for each request.

Lookup(Dir, Name) -> FileId
 – throws *NotFound*

Locates the text name in the directory and returns the relevant UFID. If *Name* is not in the directory, throws an exception.

AddName(Dir, Name, FileId)
 – throws *NameDuplicate*

If *Name* is not in the directory, adds (*Name, File*) to the directory and updates the file's attribute record.
 If *Name* is already in the directory: throws an exception.

UnName(Dir, Name)
 – throws *NotFound*

If *Name* is in the directory: the entry containing *Name* is removed from the directory.
 If *Name* is not in the directory: throws an exception.

GetNames(Dir, Pattern) -> NameSeq

Returns all the text names in the directory that match the regular expression *Pattern*.

- Directories are arranged in a tree structure
- Any file or directory is reference by a *pathname*
- Files can have more than one name (*link*)

A **file group** is a logical collection of files located on a given server

- a server can have more than one group

- a group can change server

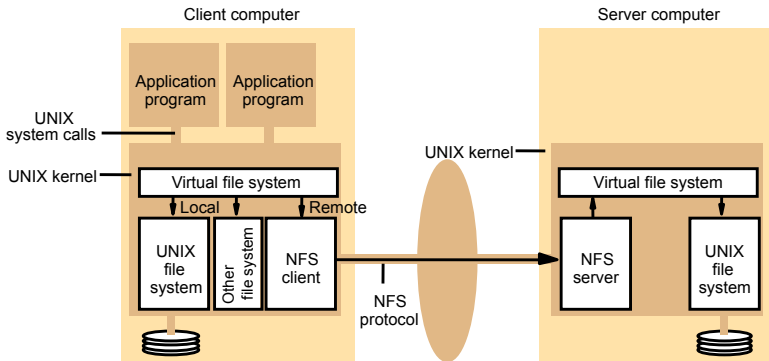
- *filesystems* in unix

- Industry standard for local networks since the 1985
- OS independent
- unix implementation
- rpc
- udp or tcp

NFS Illustrated

Brent Callaghan





- access transparency

- part of unix kernel

- NFS *file handle*, 3 components:
 - ▶ filesystem identifier
 - different groups of files
 - ▶ i-node (index node)
 - structure for finding the file
 - ▶ i-node generation number
 - i-nodes are reused
 - incremented when reused

<i>lookup(dirfh, name) -> fh, attr</i>	Returns file handle and attributes for the file <i>name</i> in the directory <i>dirfh</i> .
<i>create(dirfh, name, attr) -> newfh, attr</i>	Creates a new file name in directory <i>dirfh</i> with attributes <i>attr</i> and returns the new file handle and attributes.
<i>remove(dirfh, name) status</i>	Removes file name from directory <i>dirfh</i> .
<i>getattr(fh) -> attr</i>	Returns file attributes of file <i>fh</i> . (Similar to the UNIX <i>stat</i> system call.)
<i>setattr(fh, attr) -> attr</i>	Sets the attributes (mode, user id, group id, size, access time and modify time of a file). Setting the size to 0 truncates the file.
<i>read(fh, offset, count) -> attr, data</i>	Returns up to <i>count</i> bytes of data from a file starting at <i>offset</i> . Also returns the latest attributes of the file.
<i>write(fh, offset, count, data) -> attr</i>	Writes <i>count</i> bytes of data to a file starting at <i>offset</i> . Returns the attributes of the file after the write has taken place.
<i>rename(dirfh, name, todirfh, toname) -> status</i>	Changes the name of file <i>name</i> in directory <i>dirfh</i> to <i>toname</i> in directory to <i>todirfh</i>
<i>link(newdirfh, newname, dirfh, name) -> status</i>	Creates an entry <i>newname</i> in the directory <i>newdirfh</i> which refers to file <i>name</i> in the directory <i>dirfh</i> .

<i>symlink(newdirfh, newname, string)</i> -> <i>status</i>	Creates an entry <i>newname</i> in the directory <i>newdirfh</i> of type symbolic link with the value <i>string</i> . The server does not interpret the <i>string</i> but makes a symbolic link file to hold it.
<i>readlink(fh)</i> -> <i>string</i>	Returns the string that is associated with the symbolic link file identified by <i>fh</i> .
<i>mkdir(dirfh, name, attr)</i> -> <i>newfh, attr</i>	Creates a new directory <i>name</i> with attributes <i>attr</i> and returns the new file handle and attributes.
<i>rmdir(dirfh, name)</i> -> <i>status</i>	Removes the empty directory <i>name</i> from the parent directory <i>dirfh</i> . Fails if the directory is not empty.
<i>readdir(dirfh, cookie, count)</i> -> <i>entries</i>	Returns up to <i>count</i> bytes of directory entries from the directory <i>dirfh</i> . Each entry contains a file name, a file handle, and an opaque pointer to the next directory entry, called a <i>cookie</i> . The <i>cookie</i> is used in subsequent <i>readdir</i> calls to start reading from the following entry. If the value of <i>cookie</i> is 0, reads from the first entry in the directory.
<i>statfs(fh)</i> -> <i>fsstats</i>	Returns file system information (such as block size, number of free blocks and so on) for the file system containing a file <i>fh</i> .

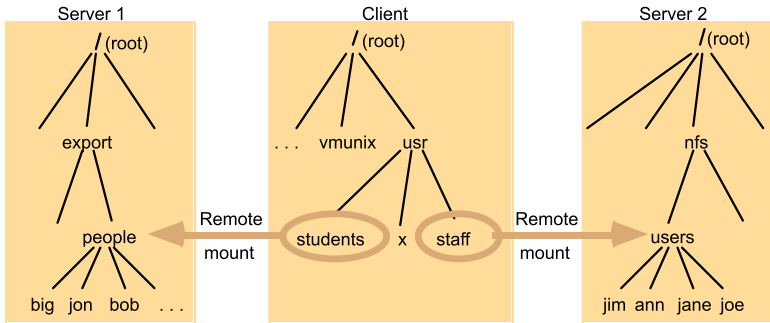
- NFS client emulates Unix file semantics

- in the kernel, not in a library, because:
 - ▶ access files via system calls
 - ▶ single client module for multiple user processes
 - ▶ encryption can be done in the kernel

- NFS server is stateless, doesn't keep open files for clients

- server checks identity each time (uid and gid)

- the process of including a new filesystem is called mounting
- `/etc/exports` has filesystems that can be mounted by others
- clients use a modified mount command for remote filesystems
- communicates with the mount process on the server in a mount protocol
- hard-mounted
 - ▶ user process is suspended until request is successful
 - ▶ when server is not responding
 - ▶ request is retried until it's satisfied
- soft-mounted
 - ▶ if server fails, client returns failure after a small # of retries
 - ▶ user process handles the failure



- pathname: /users/students/dc/abc
- server doesn't receive the entire pathname for translation **why?**
- client breaks down the pathnames into parts
- iteratively translate each part
- translation is cached

- what if a user process references a file on a remote filesystem that is not mounted?

- what if a user process references a file on a remote filesystem that is not mounted?
- table of mount points (pathname) and servers
- NFS client sends the reference to the automounter
- automounter finds the first server that is up
- mounts it at some location and sets a symbolic link (original impl)
- mounts it at the mount point (later impl)
- could help for fault tolerance, the same mount point with multiple replicated servers.

- caching file pages, directory and file attributes
- read-ahead: prefetch pages following the most-recently read file pages
- delayed-write: write to disk when the page in memory is needed for other purposes
- "sync" flushes "dirty" pages to disk every 30 seconds
- two write options:
 - ▶ write-through: write to disk before replying to the client
 - ▶ cache and commit: - stored in memory cache
- write to disk before replying to a "commit" request from the client

- caches results of `read`, `write`, `getattr`, `lookup`, `readdir`
- clients responsibility to poll the server for consistency
- Reading:
 - ▶ timestamp-based consistency validation
 - T_c : time when the cache entry was last validated
 - T_m : time when the block was last modified at the server
 - cache entry is valid if: $(T - T_c < t) \vee (T_{mclient} = T_{mserver})$, where t is the freshness interval
 - ▶ validation doesn't guarantee the same level of consistency as one-copy

- dirty: modified page in cache
- flush to disk: file is closed or sync from client
- bio-daemon (block input-output)
 - ▶ read-ahead: after each read request, request the next file block from the server as well
 - ▶ delayed write: after a block is filled, it's sent to the server
 - ▶ reduce the time to wait for read or write

- overhead is low

- main problems:
 - ▶ frequent `getattr()` for cache validation (piggybacking)
 - ▶ relatively poor performance if write-through is used on the server (delay-write or commit in current versions)

- `write < 5%`

- lookup is almost 50% (step by step pathname translation)

- access transparency: same system calls for local or remote files
- location transparency: could have a single name space for all files (depending on all the clients to agree the same name space)
- mobility transparency: mount table need to be updated on each client (not transparent)
- scalability: can usually support large loads, add processors, disks, servers...
- file replication: read-only replication, no support for replication of files with updates
- hardware and OS: many ports
- fault tolerance: stateless and idempotent
- consistency: not quite one-copy for efficiency
- security: added encryption–Kerberos
- efficiency: pretty efficient, wide-spread use

Google FS, Hadoop, Amazon S3...