

Slides for Chapter 8: Distributed File Systems



From **Coulouris, Dollimore and Kindberg**
**Distributed Systems:
Concepts and Design**
Edition 4, © Pearson Education 2005

DISTRIBUTED SYSTEMS
CONCEPTS AND DESIGN
George Coulouris
Jean Dollimore
Tim Kindberg

Learning Objectives

- ⌘ Understand the requirements that affect the design of distributed storage services
- ⌘ Case study on NFS: understand how a relatively simple, widely used service is designed

Instructor's Guide for Coulouris, Dollimore and Kindberg Distributed Systems: Concepts and Design, Edn. 4 © Pearson Education 2005

Distributed storage systems

- ⌘ Earlier storage systems are file systems (e.g. NFS); units are files.
- ⌘ More recently, distributed object systems (e.g. CORBA, Java); units are objects.

Instructor's Guide for Coulouris, Dollimore and Kindberg Distributed Systems: Concepts and Design, Edn. 4 © Pearson Education 2005

Storage systems and their properties

	Sharing	Persistence	Distributed cache/replicas	Consistency maintenance	Example
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.

Instructor's Guide for Coulouris, Dollimore and Kindberg Distributed Systems: Concepts and Design, Edn. 4 © Pearson Education 2005

Characteristics of (non-distributed) file systems

- ⌘ data and attributes (Fig 8.3)
- ⌘ directory: mapping from text names to internal file identifiers
- ⌘ layers of modules in file systems (Fig 8.2)
- ⌘ file operation system calls in UNIX (Fig. 8.4)

Instructor's Guide for Coulouris, Dollimore and Kindberg Distributed Systems: Concepts and Design, Edn. 4 © Pearson Education 2005

File attributes

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

Instructor's Guide for Coulouris, Dollimore and Kindberg Distributed Systems: Concepts and Design, Edn. 4 © Pearson Education 2005

File system modules

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

Instructor's Guide for Conradin, Dellinger and Kinberg Distributed Systems: Concepts and Design Eds. 4
© Pearson Education 2005

UNIX file system operations

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

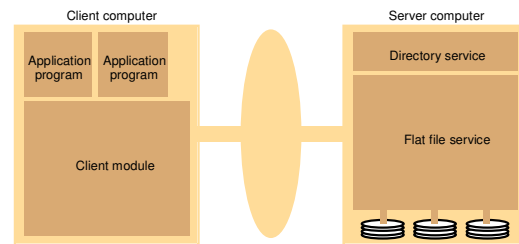
Instructor's Guide for Conradin, Dellinger and Kinberg Distributed Systems: Concepts and Design Eds. 4
© Pearson Education 2005

Distributed file system requirements

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

Instructor's Guide for Conradin, Dellinger and Kinberg Distributed Systems: Concepts and Design Eds. 4
© Pearson Education 2005

File service architecture (Author's model)



Instructor's Guide for Conradin, Dellinger and Kinberg Distributed Systems: Concepts and Design Eds. 4
© Pearson Education 2005

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

Instructor's Guide for Conradin, Dellinger and Kinberg Distributed Systems: Concepts and Design Eds. 4
© Pearson Education 2005

Flat file service interface

- ⌘ RPC used by client modules
 - ☑ not by user-level programs (which use client modules)
- ⌘ More fault tolerant compared to UNIX
 - ☑ Repeatable (idempotent) operations
 - ☑ [except for Create: re-execution gets a new file]
 - ☑ 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

Instructor's Guide for Conradin, Dellinger and Kinberg Distributed Systems: Concepts and Design Eds. 4
© Pearson Education 2005

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).

Instructor's Guide for Conrad, DeLamater and Kinberg Distributed Systems: Concepts and Design, Edn. 4
© Pearson Education 2005

Access control

- ⌘ UNIX checks access rights when a file is opened
 - ☑ subsequent checks during read/write are not necessary
- ⌘ distributed environment
 - ☑ server has to check
 - ☑ stateless approaches
 1. access check once when 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.
 - ☑ second is more common

Instructor's Guide for Conrad, DeLamater and Kinberg Distributed Systems: Concepts and Design, Edn. 4
© Pearson Education 2005

Directory service operations

<i>Lookup(Dir, Name) -> FileId</i> — throws <i>NotFound</i>	Locates the text name in the directory and returns the relevant UFID. If <i>Name</i> is not in the directory, throws an exception.
<i>AddName(Dir, Name, FileId)</i> — throws <i>NameDuplicate</i>	If <i>Name</i> is not in the directory, adds (<i>Name, File</i>) to the directory and updates the file's attribute record. If <i>Name</i> is already in the directory: throws an exception.
<i>UnName(Dir, Name)</i> — throws <i>NotFound</i>	If <i>Name</i> is in the directory: the entry containing <i>Name</i> is removed from the directory. If <i>Name</i> is not in the directory: throws an exception.
<i>GetNames(Dir, Pattern) -> NameSeq</i>	Returns all the text names in the directory that match the regular expression <i>Pattern</i> .

Instructor's Guide for Conrad, DeLamater and Kinberg Distributed Systems: Concepts and Design, Edn. 4
© Pearson Education 2005

Hierarchical file system

- ⌘ Directories containing other directories and files
- ⌘ Each file can have more than one name (pathnames)
 - ☑ how in UNIX, Windows?

Instructor's Guide for Conrad, DeLamater and Kinberg Distributed Systems: Concepts and Design, Edn. 4
© Pearson Education 2005

File groups

- ⌘ a logical collection of files on one server
 - ☑ a server can have more than one group
 - ☑ a group can change server
 - ☑ filesystems in unix
 - ☑ different devices for non-distributed
 - ☑ different hosts for distributed

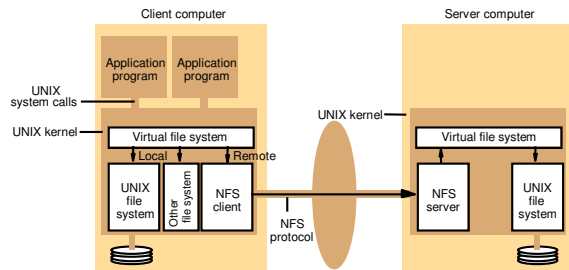
Instructor's Guide for Conrad, DeLamater and Kinberg Distributed Systems: Concepts and Design, Edn. 4
© Pearson Education 2005

Case Study: Sun NFS

- ⌘ Industry standard for local networks since the 1980's
- ⌘ OS independent
- ⌘ unix implementation
- ⌘ rpc
- ⌘ udp or tcp

Instructor's Guide for Conrad, DeLamater and Kinberg Distributed Systems: Concepts and Design, Edn. 4
© Pearson Education 2005

NFS architecture: virtual file system



Instructor's Guide for Condorin, Dellinger and Kinberg Distributed Systems: Concepts and Design Edn. 4
© Pearson Education 2005

Virtual file system

- ⌘ 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
- ⌘ VFS
 - ☐ struct for each file system
 - ☐ v-node for each open file
 - ☐ file handle for remote file
 - ☐ i-node number for local file

Instructor's Guide for Condorin, Dellinger and Kinberg Distributed Systems: Concepts and Design Edn. 4
© Pearson Education 2005

Client integration

- ⌘ 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

Instructor's Guide for Condorin, Dellinger and Kinberg Distributed Systems: Concepts and Design Edn. 4
© Pearson Education 2005

Access control

- ⌘ nfs server is stateless, doesn't keep open files for clients
- ⌘ server checks identity each time (uid and gid)

Instructor's Guide for Condorin, Dellinger and Kinberg Distributed Systems: Concepts and Design Edn. 4
© Pearson Education 2005

NFS server operations (simplified)

What do you notice about read() and write()?

<code>lookup(dirfh, name) -> fh, attr</code>	Returns file handle and attributes for the file <i>name</i> in the directory <i>dirfh</i> .
<code>create(dirfh, name, attr) -> newfh, attr</code>	Creates a new file <i>name</i> in directory <i>dirfh</i> with attributes <i>attr</i> and returns the new file handle and attributes.
<code>remove(dirfh, name) status</code>	Removes file <i>name</i> from directory <i>dirfh</i> .
<code>getattr(fh) -> attr</code>	Returns file attributes of file <i>fh</i> . (Similar to the UNIX <i>stat</i> system call.)
<code>setattr(fh, attr) -> attr</code>	Sets the attributes (mode, user id, group id, size, access time and modify time) of a file.
<code>read(fh, offset, count) -> attr, data</code>	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.
<code>write(fh, offset, count, data) -> attr</code>	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.
<code>rename(dirfh, name, todirfh, toname) -> status</code>	Changes the name of file <i>name</i> in directory <i>dirfh</i> to <i>toname</i> in directory <i>todirfh</i> .
<code>link(newdirfh, newname, dirfh, name) -> status</code>	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> . Continues on next slide ...

Instructor's Guide for Condorin, Dellinger and Kinberg Distributed Systems: Concepts and Design Edn. 4
© Pearson Education 2005

NFS server operations (simplified) – 2

<code>symlink(newdirfh, newname, string) -> status</code>	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.
<code>readlink(fh) -> string</code>	Returns the string that is associated with the symbolic link file identified by <i>fh</i> .
<code>mkdir(dirfh, name, attr) -> newfh, attr</code>	Creates a new directory <i>name</i> with attributes <i>attr</i> and returns the new file handle and attributes.
<code>rmdir(dirfh, name) -> status</code>	Removes the empty directory <i>name</i> from the parent directory <i>dirfh</i> . Fails if the directory is not empty.
<code>readdir(dirfh, cookie, count) -> entries</code>	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.
<code>statvfs(fh) -> fstats</code>	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> .

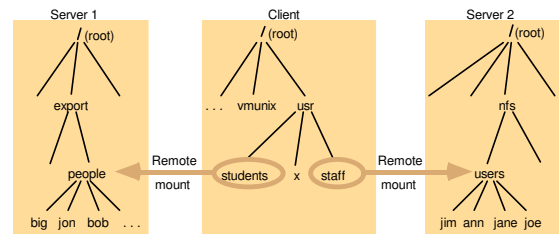
Instructor's Guide for Condorin, Dellinger and Kinberg Distributed Systems: Concepts and Design Edn. 4
© Pearson Education 2005

Mount service

- ⌘ 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

Instructor's Guide for Condorin, DeLineiro and Kinberg Distributed Systems: Concepts and Design, Edn. 4
© Pearson Education 2005

Local and remote file systems



Note: The file system mounted at `/usr/students` in the client is actually the sub-tree located at `/export/people` in Server 1; the file system mounted at `/usr/staff` in the client is actually the sub-tree located at `/nfs/users` in Server 2.

Instructor's Guide for Condorin, DeLineiro and Kinberg Distributed Systems: Concepts and Design, Edn. 4
© Pearson Education 2005

Pathname translation

- ⌘ 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

Instructor's Guide for Condorin, DeLineiro and Kinberg Distributed Systems: Concepts and Design, Edn. 4
© Pearson Education 2005

Automounter

- ⌘ what if a user process reference 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 check find the first server that is up
- ⌘ mount it at some location and set a symbolic link (original impl)
- ⌘ mount it at the mount point (later impl)
- ⌘ could help fault tolerance, the same mount point with multiple replicated servers.

Instructor's Guide for Condorin, DeLineiro and Kinberg Distributed Systems: Concepts and Design, Edn. 4
© Pearson Education 2005

Server caching

- ⌘ caching file pages, directory/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 option
 1. write-through: write to disk before replying to the client
 2. cache and commit:
 - ☑ stored in memory cache
 - ☑ write to disk before replying to a "commit" request from the client

Instructor's Guide for Condorin, DeLineiro and Kinberg Distributed Systems: Concepts and Design, Edn. 4
© Pearson Education 2005

Client caching

- ⌘ caches results of read, write, getattr, lookup, readdir
- ⌘ clients responsibility to poll the server for consistency

Instructor's Guide for Condorin, DeLineiro and Kinberg Distributed Systems: Concepts and Design, Edn. 4
© Pearson Education 2005

Client caching: reading (1)

- ⌘ timestamp-based methods for 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:
 1. $T - T_c < t$, where t is the freshness interval
 - ☒ t is adaptively adjusted:
 - files: 3 to 30 seconds depending on freq of updates
 - directories: 30 to 60 seconds
 2. $T_{m_{client}} = T_{m_{server}}$

Instructor's Guide for Condorin, DeLineiro and Kinberg Distributed Systems: Concepts and Design, Edn. 4
© Pearson Education 2005

Client caching: reading (2)

- ⌘ need validation for all cache accesses
- ⌘ condition #1 can be determined by the client alone--performed first
- ⌘ Reducing getattr() to the server [for getting $T_{m_{server}}$]
 1. new value of $T_{m_{server}}$ is received, apply to all cache entries from the same file
 2. piggyback getattr() on file operations
 3. adaptive alg for update t (condition #1)
- ⌘ validation doesn't guarantee the same level of consistency as one-copy

Instructor's Guide for Condorin, DeLineiro and Kinberg Distributed Systems: Concepts and Design, Edn. 4
© Pearson Education 2005

Client caching: writing

- ⌘ 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/write

Instructor's Guide for Condorin, DeLineiro and Kinberg Distributed Systems: Concepts and Design, Edn. 4
© Pearson Education 2005

Other optimization

- ⌘ UDP, RPC
- ⌘ extended to 9 kilobytes--entire block in a single packet

Instructor's Guide for Condorin, DeLineiro and Kinberg Distributed Systems: Concepts and Design, Edn. 4
© Pearson Education 2005

Security

- ⌘ stateless nfs server
- ⌘ user's identity in each request
- ⌘ Kerberos authentication during the mount process, which includes uid and host address
- ⌘ server maintain authentication info for the mount
- ⌘ on each file request, nfs checks the uid and address
- ⌘ one user per client

Instructor's Guide for Condorin, DeLineiro and Kinberg Distributed Systems: Concepts and Design, Edn. 4
© Pearson Education 2005

Performance

- ⌘ overhead/penalty is low
- ⌘ main problems
 - ☒ frequent getattr() for cache validation (piggybacking)
 - ☒ relatively poor performance if write-through is used on the server (delay-write/commit in current versions)
- ⌘ write < 5%
- ⌘ lookup is almost 50% (step by step pathname translation)

Instructor's Guide for Condorin, DeLineiro and Kinberg Distributed Systems: Concepts and Design, Edn. 4
© Pearson Education 2005

Summary for NFS

- ⌘ 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

Instructor's Guide for Coulouris, Dillman and Kinberg Distributed Systems: Concepts and Design, Edn. 4
© Pearson Education 2005

WebNFS (p. 360)

- ⌘ HTTP and NFS both can read files, what NFS can do more than HTTP in terms of reading? [what can read() do in a file system that HTTP can't?]
- ⌘ NFS is designed to be on "fast" LANs
- ⌘ WebNFS is designed to be on "slower" WANs
- ⌘ WebNFS clients talk to NFS servers
 - ☑ Small set up cost (thinner client)
 - ☑ public files, mostly read access, no authentication
 - ☑ no mounting
 - ☑ access portions of a file
- ⌘ `nfs://xyz.com/someFile`

Instructor's Guide for Coulouris, Dillman and Kinberg Distributed Systems: Concepts and Design, Edn. 4
© Pearson Education 2005