



Distributed File Systems

Chien-Min Wang
Institute of Information Science
Academia Sinica

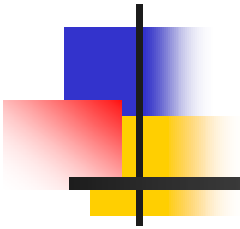


Contents

- n File System Overview
- n Distributed File Systems: Issues
- n Distributed File Systems: Case Studies
- n Distributed File Systems for Clouds

Lecture 2

Distributed File Systems: Issues





Outline

- n Introduction
- n Basic Implementation Mechanisms
- n Design Choices



Why distributed file systems?

- n The data may be much larger than the storage space of a computer.
- n The data may survive much longer than the life of a computer.
- n A user may access his/her data from different machines at different geographic locations.
- n A user may want to share his/her data with users around the world.



Accessing Files on Remote Sites

- n FTP

- ┆ Explicit access
- ┆ User-directed connection to access remote resources

- n We want more transparency

- ┆ Allow user to access remote files just as local ones



File Service Types₁

n Upload/Download Model

- | Read file: copy file from server to client
- | Write file: copy file from client to server

n Advantage

- | Simple

n Problems

- | Wasteful: what if client needs small piece?
- | Problematic: what if client doesn't have enough space?
- | Consistency: what if others need to modify the same file?



File Service Types₂

- n Remote Access Model
- n File service provides functional interface:
 - | create, delete, read bytes, write bytes, etc...
- n Advantages:
 - | Client gets only what's needed
 - | Server can manage coherent view of the file system
- n Problem:
 - | Possible server and network congestion
 - └ Servers are accessed for duration of file access
 - └ Same data may be requested repeatedly



File Servers

- n File Directory Service

- ┆ Maps textual names for file to internal locations that can be used by file service

- n File service

- ┆ Provides file access interface to clients

- n Client module (driver)

- ┆ Client side interface for file and directory service
- ┆ if done right, helps provide access transparency
 - ┆ e.g. under vnode layer of Linux virtual file system



Distributed File Systems

- n Provide accesses to data stored at servers using *file system interfaces*.
- n What are the file system interfaces?
 - | Open a file, check status on a file, close a file;
 - | Read data from a file;
 - | Write data to a file;
 - | Lock a file or part of a file;
 - | List files in a directory, delete a directory;
 - | Delete a file, rename a file, add a symbolic link to a file;
 - | etc;



Why is DFS useful?

- n Data sharing of multiple users
- n User mobility
- n Location transparency
- n Location independence
- n Backups and centralized management

- n Not all DFS are the same:
 - ▀ High-speed network DFS vs. low-speed network DFS



Interface: File vs. Block

- n Data are organized in files, which in turn are organized in directories
- n Compare these with disk-level access or “block” access interface: [Read/Write, LUN, block#]
- n Key differences:
 - ┆ Implementation of the directory/file structure and semantics
 - ┆ Synchronization



Digression: Buzz Word Discussion

	NAS	SAN
Access Methods	File access	Disk block access
Access Medium	Ethernet	Fiber Channel and Ethernet
Transport Protocol	Layer over TCP/IP	SCSI/FC and SCSI/IP
Efficiency	Less	More
Sharing and Access Control	Good	Poor
Integrity demands	Strong	Very strong
Clients	Workstations	Database servers



Sequential Semantics of File Sharing

- n Read returns result of last write
- n Easily achieved if
 - | Only one server
 - | Clients do not cache data
- n BUT
 - | Performance problems if no cache
 - u Obsolete data
 - | We can write-through
 - u Must notify clients holding copies
 - u Requires extra state, generates extra traffic



Session Semantics of File Sharing

- n Relax the rules
- n Changes to an open file are initially visible only to the process (or machine) that modified it.
- n Last process to modify the file wins.



Other solutions

- n Make files immutable
 - | Aids in replication
 - | Does not help with detecting modification
- n Use atomic transactions
 - | Each file access is an atomic transaction
 - | If multiple transactions start concurrently
 - u Resulting modification is serial



File Usage Patterns

- n We can't have the best of all worlds
- n Where to compromise?
 - l Semantics vs. efficiency
 - l Efficiency = client performance, network traffic, server load
- n Understand how files are used



File Usage

- n Most files are <10 Kbytes
 - | 2005: average size of 385,341 files =197 KB
 - | 2007: average size of 440,519 files =451 KB
 - | Feasible to transfer entire files (simpler)
 - | Still have to support long files
- n Most files have short lifetimes
 - | Perhaps keep them local
- n Few files are shared
 - | Overstated problem
 - | Session semantics will cause no problem most of the time



Outline

- n Introduction
- n Basic Implementation Mechanisms
- n Design Choices

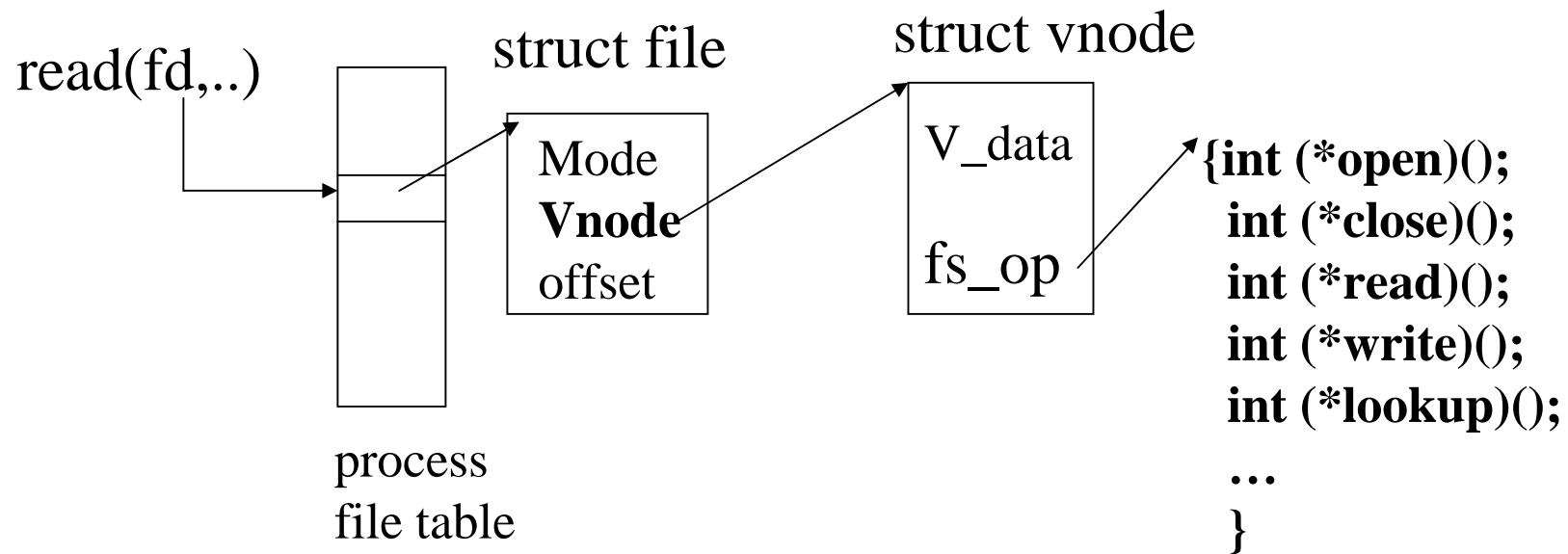


Components in a DFS

- n Client side:
 - | What has to happen to enable applications access a remote file in the same way as accessing a local file
- n Communication layer:
 - | Just TCP/IP or some protocol at higher abstraction
- n Server side:
 - | How does it service requests from the client

Client Side Example: UNIX

- n Accessing remote files in the same way as accessing local files à kernel support
 - i Vnode interface





VFS Interception

- n VFS provides “pluggable” file systems
- n Standard flow of remote access
 - | User process calls read()
 - | Kernel dispatches to VOP_READ() in some VFS
 - | nfs_read()
 - u check local cache
 - u send RPC to remote NFS server
 - u put process to sleep



VFS Interception

- n Standard flow of remote access (continued)
 - | server interaction handled by kernel process
 - u retransmit if necessary
 - u convert RPC response to file system buffer
 - u store in local cache
 - u wake up user process
 - | `nfs_read()`
 - u copy bytes to user memory



Communication Layer Example: RPC

RPC call

xid
“call”
service
version
procedure
auth-info
arguments
...

RPC reply

xid
“reply”
reply_stat
auth-info
results
...

- n Failure handling: timeout and re-issuance
- n RPC over UDP vs. RPC over TCP



Extended Data Representation (XDR)

- n Argument data and response data in RPC are packaged in XDR format
 - | Integers are encoded in big-endian
 - | Strings: len followed by ascii bytes with NULL padded to four-byte boundaries
 - | Arrays: 4-byte size followed by array entries
 - | Opaque: 4-byte len followed by binary data
- n Marshalling and un-marshalling
- n Extra overhead in data conversion to/from XDR



NFS RPC Calls

n NFS / RPC using XDR / TCP/IP

Proc.	Input args	Results
lookup	dirfh, name	status, fhandle, fattr
read	fhandle, offset, count	status, fattr, data
create	dirfh, name, fattr	status, fhandle, fattr
write	fhandle, offset, count, data	status, fattr

- n fhandle: 32-byte opaque data (64-byte in v3)
 - ┆ What's in the file handle



NFS Operations

n V2:

- | NULL, GETATTR, SETATTR
- | LOOKUP, READLINK, READ
- | CREATE, WRITE, REMOVE, RENAME
- | LINK, SYMLINK
- | READDIR, MKDIR, RMDIR
- | STATFS

n V3: add

- | *READDIRPLUS, COMMIT*
- | FSSTAT, FSINFO, PATHCONF



Server Side Example: mountd and nfsd

- n Mountd: provides the initial file handle for the exported directory
 - | Client issues nfs_mount request to mountd
 - | Mountd checks if the pathname is a directory and if the directory is exported to the client
- n nfsd: answers the rpc calls, gets reply from local file system, and sends reply via rpc
 - | Usually listening at port 2049
- n Both mountd and nfsd use underlying RPC implementation



NFS Client Server Interactions

- n Client machine:

- | Application à nfs_vnops à nfs client code à
rpc client interface

- n Server machine:

- | rpc server interface à nfs server code à
ufs_vnops à ufs code à disks



NFS File Server Failure Issues

- n Semantics of file write in V2
 - | Bypass UFS file buffer cache
- n Semantics of file write in V3
 - | Provide “COMMIT” procedure
- n Server-side retransmission cache
 - | Idempotent vs. non-idempotent requests



Outline

- n Introduction
- n Basic Implementation Mechanisms
- n Design Choices



Topic 1: Naming

- n NFS: per-client linkage
 - | Server: `export /root/fs1/`
 - | Client: `mount server:/root/fs1 /fs1 à fhandle`
- n AFS: global name space
 - | Name space is organized into Volumes
 - Global directory `/afs`;
 - `/afs/cs.wisc.edu/vol1/...`; `/afs/cs.stanford.edu/vol1/...`
 - | Each file is identified as `<vol_id, vnode#, vnode_gen>`
 - | All AFS servers keep a copy of “volume location database”, which is a table of `vol_id à server_ip` mappings



Location Transparency

- n NFS: no transparency
 - | If a directory is moved from one server to another, client must remount
- n AFS: transparency
 - | If a volume is moved from one server to another, only the volume location database on the servers needs to be updated
 - | Implementation of volume migration
 - | File lookup efficiency
- n Are there other ways to provide location transparency?



Topic 2: User Authentication and Access Control

- n User X logs onto workstation A, wants to access files on server B
 - | How does A tell B who X is
 - | Should B believe A
- n Choices made in NFS v2
 - | All servers and all client workstations share the same $\langle \text{uid}, \text{gid} \rangle$ name space \Rightarrow B send X's $\langle \text{uid}, \text{gid} \rangle$ to A
 - u Problem: root access on any client workstation can lead to creation of users of arbitrary $\langle \text{uid}, \text{gid} \rangle$
 - | Server believes client workstation unconditionally
 - u Problem: if any client workstation is broken into, the protection of data on the server is lost;
 - u $\langle \text{uid}, \text{gid} \rangle$ sent in clear-text over wire \Rightarrow request packets can be faked easily

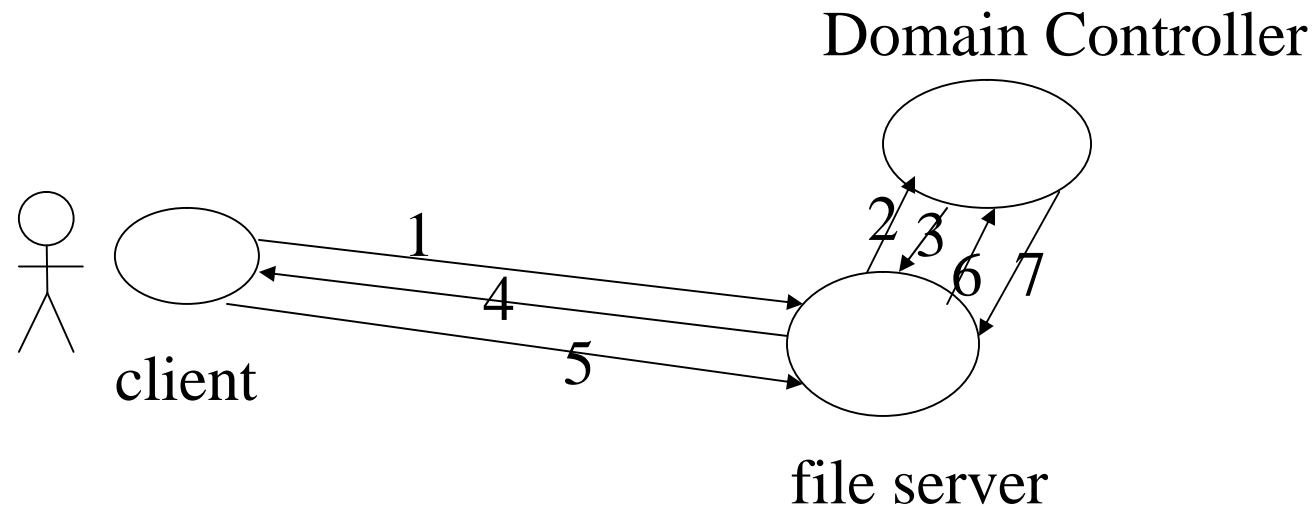


User Authentication

- n How do we fix the problems in NFS v2
 - | Hack1: root remapping → strange behavior
 - | Hack 2: UID remapping → no user mobility
 - | Real Solution: use a centralized Authentication/Authorization/Access-control (AAA) system

Example AAA System: NTLM

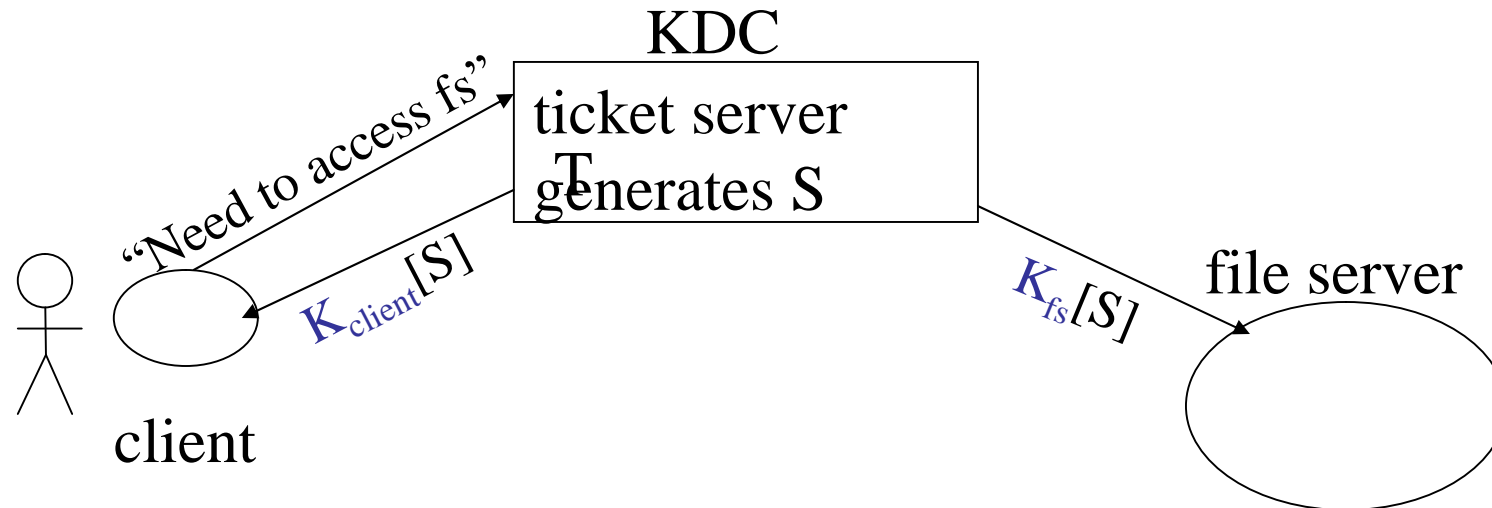
- n Microsoft Windows Domain Controller
 - ┆ Centralized AAA server
 - ┆ NTLM v2: per-connection authentication



A Better AAA System: Kerberos

n Basic idea: shared secrets

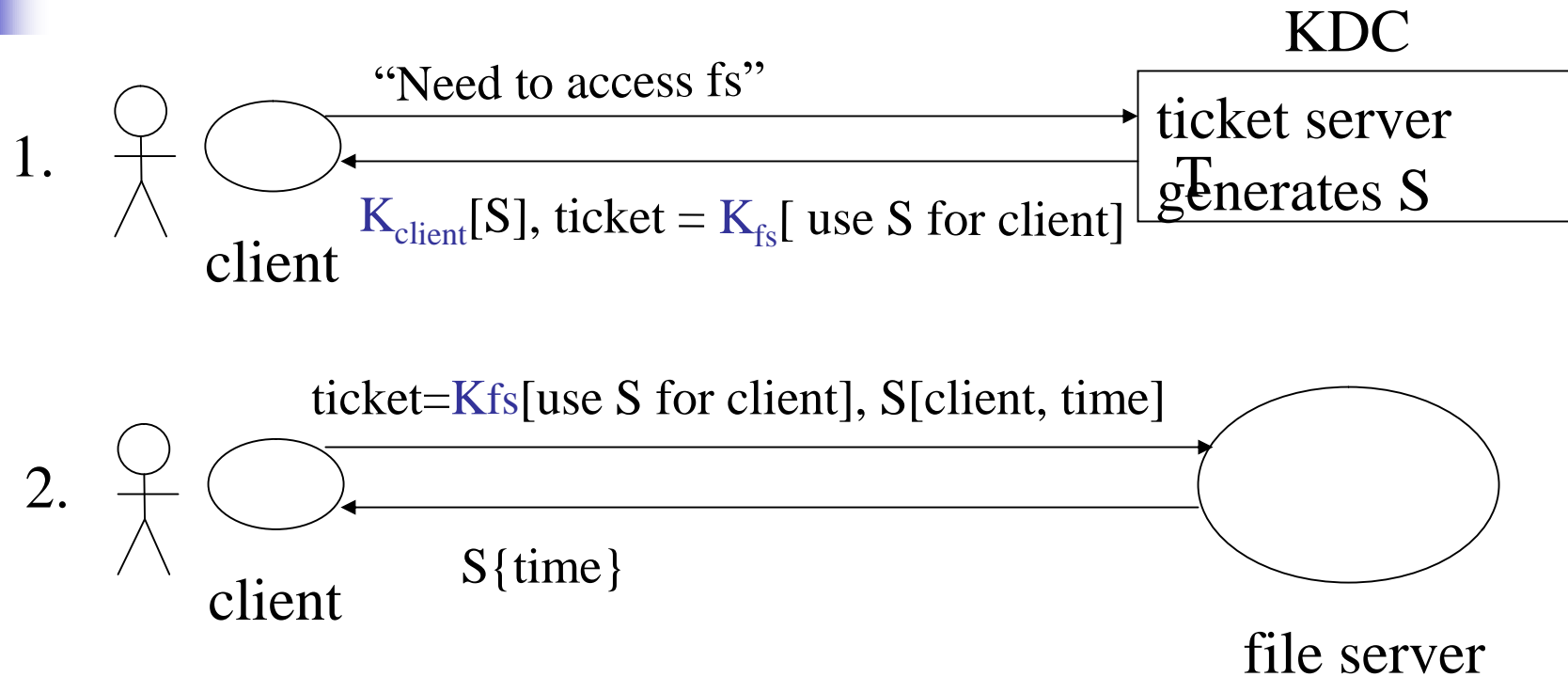
- User prove to KDC who he is; KDC generates shared secret between client and file server



S: specific to {client,fs} pair;

“short-term session-key”; has expiration time (e.g. 8 hours);

Kerberos Interactions



- why "time": guard against replay attack
- mutual authentication
- File server doesn't store S , which is specific to $\{\text{client}, \text{fs}\}$
- Client doesn't contact "ticket server" every time it contacts fs



Kerberos: User Log-on Process

- n How does user prove to KDC who the user is
 - | Long-term key: $1\text{-way-hash-func}(\text{passwd})$
 - | Long-term key comparison happens once only, at which point the KDC generates a shared secret for the user and the KDC itself → ticket-granting ticket, or “logon session key”
 - | The “ticket-granting ticket” is encrypted in KDC’s long-term key



Topic 3: Operator Batching

- n Should each client/server interaction accomplish one file system operation or multiple operations?
- n Advantage of batched operations
- n How to define batched operations



Examples of Batched Operators

- n NFS v3:
 - | Readdirplus
- n NFS v4:
 - | Compound RPC calls
- n CIFS:
 - | “AND-X” requests



Topic 4: Client-Side Caching

- n Why is client-side caching necessary
- n What are cached
 - ┆ Read-only file data and directory data → easy
 - ┆ Data written by the client machine → when are data written to the server? What happens if the client machine goes down?
 - ┆ Data that are written by other machines → how to know that the data have been changed? How to ensure data consistency?
 - ┆ Is there any pre-fetching?



Client Caching in NFS v2

- n Cache both clean and dirty file data and file attributes
- n File attributes in the client cache are expired after 60 seconds
- n File data are checked against the modified-time in file attributes (which could be a cached copy)
 - | Changes made on one machine can take up to 60 secs to be reflected on another machine
- n Dirty data are buffered on the client machine till file close or up to 30 seconds
 - | If the machine crashes before then, the changes are lost
 - | Similar to UNIX FFS local file system behavior



Implication of NFS v2 Client Caching

- n Data consistency guarantee is very poor
 - | Simply unacceptable for some distributed applications
 - | Productivity apps tend to tolerate such loose consistency
- n Different client implementations implement the “prefetching” part differently
- n Generally clients do not cache data on local disks



Client Caching in AFS

- n Client caches both clean and dirty file data and attributes
 - | The client machine uses local disks to cache data
 - | When a file is opened for read, the whole file is fetched and cached on disk
 - u Why? What's the disadvantage of doing so?
- n However, when a client caches file data, it obtains a “callback” on the file
- n In case another client writes to the file, the server “breaks” the callback
 - | Similar to invalidations in distributed shared memory implementations
- n Implications: file server must keep states!



AFS RPC Procedures

- n Procedures that are not in NFS
 - | Fetch: return status and optionally data of a file or directory, and place a callback on it
 - | RemoveCallBack: specify a file that the client has flushed from the local machine
 - | BreakCallBack: from server to client, revoke the callback on a file or directory
 - u What should the client do if a callback is revoked?
 - | Store: store the status and optionally data of a file
- n Rest are similar to NFS calls



Failure Recovery in AFS

- n What if the file server fails
 - | Two candidate approaches to failure recovery
- n What if the client fails
- n What if both the server and the client fail
- n Network partition
 - | How to detect it? How to recover from it?
 - | Is there anyway to ensure absolute consistency in the presence of network partition?
 - u Reads
 - u Writes
- n What if all three fail: network partition, server, client



Key to Simple Failure Recovery

- n Try not to keep any state on the server
- n If you must keep some states on the server
 - ┆ Understand why and what states the server is keeping
 - ┆ Understand the worst case scenario of no state on the server and see if there are still ways to meet the correctness goals
 - ┆ Revert to this worst case in each combination of failure cases



Topic 5: File Access Consistency

- n In UNIX local file system, concurrent file reads and writes have “sequential” consistency semantics
 - | Each file read/write from user-level app is an atomic operation
 - u The kernel locks the file vnode
 - | Each file write is immediately visible to all file readers
- n Neither NFS nor AFS provides such concurrency control
 - | NFS: “sometime within 30 seconds”
 - | AFS: session semantics for consistency



Session Semantics in AFS

n What it means:

- | A file write is visible to processes on the same box immediately, but not visible to processes on other machine until the file is closed
- | When a file is closed, changes are visible to new opens, but are not visible to “old” opens
- | All other file operations are visible everywhere immediately

n Implementation

- | Dirty data are buffered at the client machine until file close, then flushed back to server, which leads the server to send “break callback” to other clients
- | Problems with this implementation



Access Consistency in “Sprite”

- n Sprite: a research file system developed in UC Berkeley in late 80’s
- n Implements “sequential” consistency
 - ┆ Caches only file data, not file metadata
 - ┆ When server detects a file is open on multiple machines but is written by some client, client caching of the file is disabled; all reads and writes go through the server
 - ┆ “Write-back” policy otherwise
 - ┆ Why?



Implementing Sequential Consistency

- n How to identify out-of-date data blocks
 - ┆ Use file version number
 - ┆ No invalidation
 - ┆ No issue with network partition
- n How to get the latest data when read-write sharing occurs
 - ┆ Server keeps track of last writer



Implication of “Sprite” Caching

- n Server must keep states!
 - | Recovery from power failure
 - | Server failure doesn't impact consistency
 - | Network failure doesn't impact consistency
- n Price of sequential consistency: no client caching of file metadata; all file opens go through server
 - | Performance impact
 - | Suited for wide-area network?



Access Consistency in AFS v3

n Motivation

- | How does one implement sequential consistency in a file system that spans multiple sites over WAN
 - u Why Sprite's approach won't work
 - u Why AFS v2 approach won't work
 - u Why NFS approach won't work

n What should be the design guidelines?

- | What are the common share patterns?



“Tokens” in AFS v3

- n Callbacks are evolved into 4 kinds of “Tokens”
 - | Open tokens: allow holder to open a file; submodes: read, write, execute, exclusive-write
 - | Data tokens: apply to a range of bytes
 - u “read” token: cached data are valid
 - u “write” token: can write to data and keep dirty data at client
 - | Status tokens: provide guarantee of file attributes
 - u “read” status token: cached attribute is valid
 - u “write” status token: can change the attribute and keep the change at the client
 - | Lock tokens: allow holder to lock byte ranges in the file



Compatibility Rules for Tokens

- n Open tokens:
 - | Open for exclusive writes are incompatible with any other open, and “open for execute” are incompatible with “open for write”
 - | But “open for write” can be compatible with “open for write” --- why?
- n Data tokens: R/W and W/W are incompatible if the byte range overlaps
- n Status tokens: R/W and W/W are incompatible
- n Data token and status token: compatible or incompatible?



Token Manager

- n Resolve conflicts: block the new requester and send notification to other clients' tokens
- n Handle operations that request multiple tokens
 - l Example: rename
 - l How to avoid deadlocks



Failure Recovery in Token Manager

- n What if the server fails
- n What if a client fails
- n What if network partition happens



Topic 6: File Locking

n Issues

- | Whole file locking or byte-range locking
- | Mandatory or advisory
 - u UNIX: advisory
 - u Windows: if a lock is granted, it's mandatory on all other accesses

n NFS: network lock manager (NLM)

- | NLM is not part of NFS v2, because NLM is stateful
- | Provides both whole file and byte-range locking
- | Advisory
- | Relies on “network status monitor” for server monitoring



Issues in Locking Implementations

- n Synchronous and Asynchronous calls
 - | NLM provides both
- n Failure recovery
 - | What if server fails
 - Lock holders are expected to re-establish the locks during the “grace period”, during which no other locks are granted
 - | What if a client holding the lock fails
 - | What if network partition occurs



Wrap up: Comparing the File Systems

- n Caching:
 - | NFS
 - | AFS
 - | Sprite
- n Consistency
 - | NFS
 - | AFS
 - | Sprite
 - | AFS v3
- n Locking



Wrap up: Comparison with the Web

n Differences:

- | Web offers HTML, etc. DFS offers binary data only
- | Web has a few but universal clients; DFS is implemented in the kernel

n Similarities:

- | Caching with TTL is similar to NFS consistency
- | Caching with IMS-every-time is similar to Sprite consistency
 - u As predicted in AFS studies, there is a scalability problem here

n Security mechanisms

- | AAA similar
- | Encryption?



Topic 7: Stateful or stateless design?

- n Stateful
 - | Server maintains client-specific states
- n Shorter requests
- n Better performance in processing requests
- n Cache coherence is possible
 - | Server can know who's accessing what
- n File locking is possible



Topic 7: Stateful or stateless design?

- n Stateless
 - | Server maintains no information on client accesses
- n Each request must identify file and offsets
- n Server can crash and recover
 - | No state to lose
- n Client can crash and recover
- n No open/close needed
 - | They only establish state
- n No server space used for state
 - | Don't worry about supporting many clients
- n Problems if file is deleted on server
- n File locking not possible