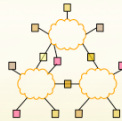


15-446 Distributed Systems Spring 2009



L-17 Distributed File Systems

1

Outline

- Why Distributed File Systems?
- Basic mechanisms for building DFSs
 - Using NFS and AFS as examples
- Design choices and their implications
 - Naming
 - Authentication and Access Control
 - Caching
 - Concurrency Control
 - Locking

2

What Distributed File Systems Provide

- Access to data stored at servers using file system interfaces
- What are the file system interfaces?
 - Open a file, check status of 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, create/delete a directory
 - Delete a file, rename a file, add a symlink to a file
 - etc

3

Why DFSs are Useful

- Data sharing among multiple users
- User mobility
- Location transparency
- Backups and centralized management

4

Outline

- Why Distributed File Systems?
- Basic mechanisms for building DFSs
 - Using NFS and AFS as examples
- Design choices and their implications
 - Naming
 - Authentication and Access Control
 - Caching
 - Concurrency Control
 - Locking

5

Components in a DFS Implementation

- Client side:
 - What has to happen to enable applications to access a remote file the same way a local file is accessed?
 - Accessing remote files in the same way as accessing local files → kernel support
- Communication layer:
 - Just TCP/IP or a protocol at a higher level of abstraction?
- Server side:
 - How are requests from clients serviced?

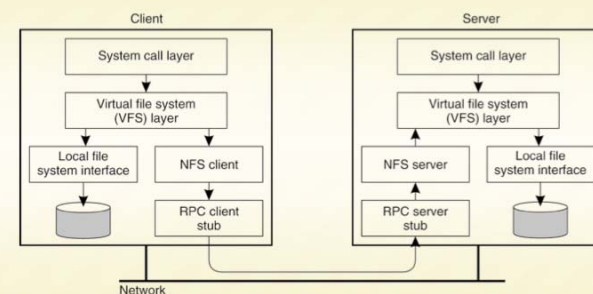
6

VFS interception

- VFS provides “pluggable” file systems
- Standard flow of remote access
 - User process calls read()
 - Kernel dispatches to VOP_READ() in some VFS
 - nfs_read()
 - check local cache
 - send RPC to remote NFS server
 - put process to sleep
 - server interaction handled by kernel process
 - retransmit if necessary
 - convert RPC response to file system buffer
 - store in local cache
 - wake up user process
 - nfs_read()
 - copy bytes to user memory

7

VFS Interception



8

Communication Layer Example: Remote Procedure Calls (RPC)

RPC call

xid
"call"
service
version
procedure
auth-info
arguments
....

RPC reply

xid
"reply"
reply_stat
auth-info
results
...

- Failure handling: timeout and re-issue

9

Extended Data Representation (XDR)

- Argument data and response data in RPC are packaged in XDR format
 - Integers are encoded in big-endian format
 - 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
- Marshalling and un-marshalling data
- Extra overhead in data conversion to/from XDR

10

Some NFS V2 RPC Calls

- NFS RPCs using XDR over, e.g., 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

- fhandle: 32-byte opaque data (64-byte in v3)

11

Server Side Example: mountd and nfsd

- 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 should be exported to the client
- nfsd: answers the RPC calls, gets reply from local file system, and sends reply via RPC
 - Usually listening at port 2049
- Both mountd and nfsd use underlying RPC implementation

12

NFS V2 Design

- “Dumb”, “Stateless” servers
- Smart clients
- Portable across different OSs
- Immediate commitment and idempotency of operations
- Low implementation cost
- Small number of clients
- Single administrative domain

13

Stateless File Server?

- Statelessness
 - Files are state, but...
 - Server **exports** files without creating extra state
 - No list of “who has this file open” (permission check on each operation on open file!)
 - No “pending transactions” across crash
- Results
 - Crash recovery is “fast”
 - Reboot, let clients figure out what happened
 - Protocol is “simple”
- State stashed elsewhere
 - Separate MOUNT protocol
 - Separate NLM locking protocol

14

NFS V2 Operations

- V2:
 - NULL, GETATTR, SETATTR
 - LOOKUP, READLINK, READ
 - CREATE, WRITE, REMOVE, RENAME
 - LINK, SYMLINK
 - READIR, MKDIR, RMDIR
 - STATFS (get file system attributes)

15

NFS V3 and V4 Operations

- V3 added:
 - REaddirPLUS, COMMIT (server cache!)
 - FSSTAT, FSINFO, PATHCONF
- V4 added:
 - COMPOUND (bundle operations)
 - LOCK (server becomes more stateful!)
 - PUTROOTFH, PUTPUBFH (no separate MOUNT)
 - Better security and authentication
 - Very different than V2/V3 → stateful

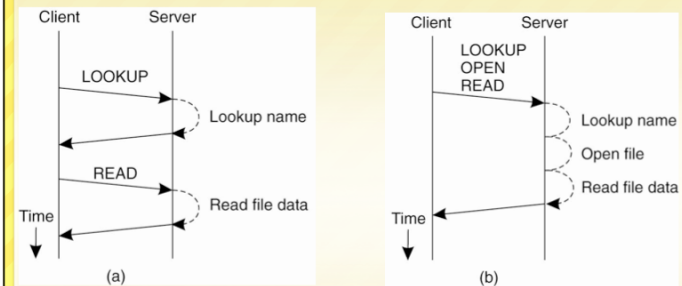
16

Operator Batching

- Should each client/server interaction accomplish one file system operation or multiple operations?
 - Advantage of batched operations?
 - How to define batched operations
- Examples of Batched Operators
 - NFS v3:
 - REaddirPLUS
 - NFS v4:
 - COMPOUND RPC calls

17

Remote Procedure Calls in NFS



- (a) Reading data from a file in NFS version 3
- (b) Reading data using a compound procedure in version 4.

18

AFS Goals

- Global distributed file system
 - "One AFS", like "one Internet"
 - Why would you want more than one?
- LARGE** numbers of clients, servers
 - 1000 machines could cache a single file,
 - some local, some (very) remote
- Goal: $O(0)$ work per client operation
 - $O(1)$ may just be too expensive!

19

AFS Assumptions

- Client machines are un-trusted
 - Must **prove** they act for a specific user
 - Secure RPC layer
 - Anonymous "system: anyuser"
- Client machines have disks(!!)
 - Can cache whole files over long periods
- Write/write and write/read sharing are rare
 - Most files updated by one user, on one machine

20

AFS Cell/Volume Architecture

- Cells correspond to administrative groups
 - /afs/andrew.cmu.edu is a **cell**
- Client machine has cell-server database
 - **protection server** handles authentication
 - **volume location server** maps volumes to servers
- Cells are broken into **volumes** (miniature file systems)
 - One user's files, project source tree, ...
 - Typically stored on one server
 - Unit of disk quota administration, backup

21

Outline

- Why Distributed File Systems?
- Basic mechanisms for building DFSs
 - Using NFS and AFS as examples
- Design choices and their implications
 - Naming
 - Authentication and Access Control
 - Caching
 - Concurrency Control
 - Locking

22

Topic 1: Name-Space Construction and Organization

- NFS: per-client linkage
 - Server: export /root/fs1/
 - Client: mount server:/root/fs1 /fs1 → fhandle
- 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 fid = <vol_id, vnode #, uniquifier>
 - All AFS servers keep a copy of "volume location database", which is a table of vol_id → server_ip mappings

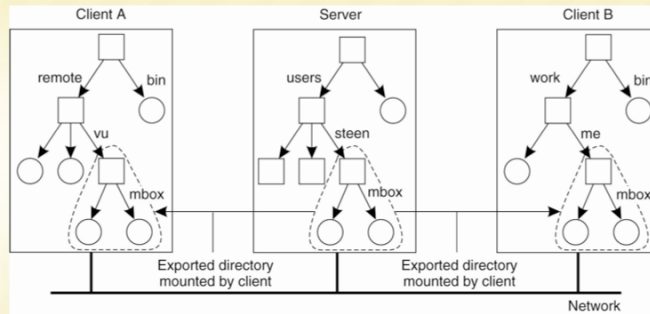
23

Implications on Location Transparency

- NFS: no transparency
 - If a directory is moved from one server to another, client must remount
- AFS: transparency
 - If a volume is moved from one server to another, only the volume location database on the servers needs to be updated

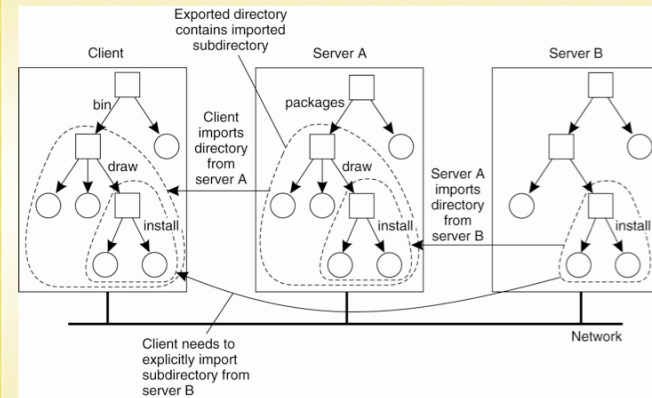
24

Naming in NFS (1)



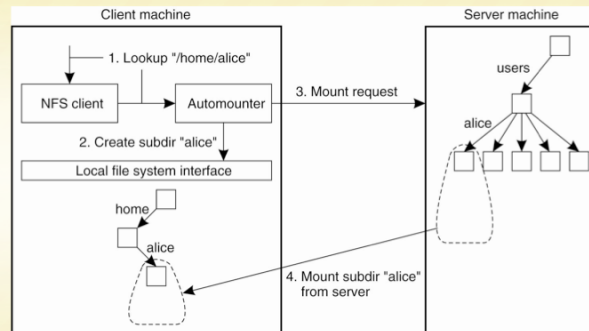
29

Naming in NFS (2)



26

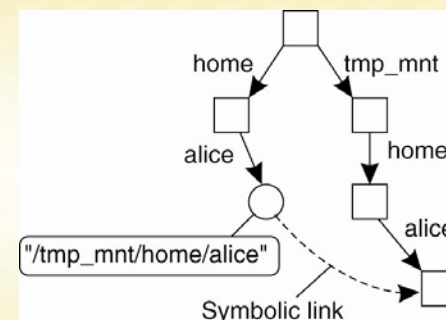
Automounting (1)



- A simple automounter for NFS.

27

Automounting (2)



- Using symbolic links with automounting.

Topic 2: User Authentication and Access Control

- 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?
- 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
 - 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
 - Problem: if any client workstation is broken into, the protection of data on the server is lost;
 - $\langle \text{uid}, \text{gid} \rangle$ sent in clear-text over wire \rightarrow request packets can be faked easily

29

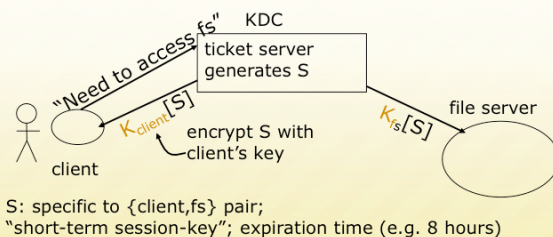
User Authentication (cont'd)

- How do we fix the problems in NFS v2
 - Hack 1: root remapping \rightarrow strange behavior
 - Hack 2: UID remapping \rightarrow no user mobility
 - Real Solution: use a centralized Authentication/Authorization/Access-control (AAA) system

30

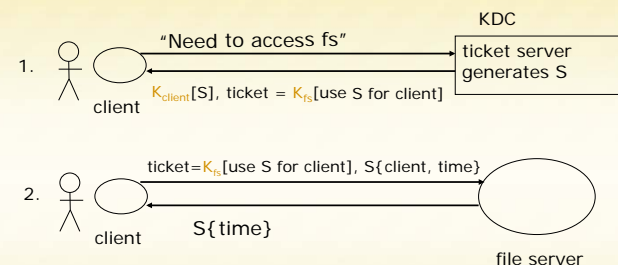
A Better AAA System: Kerberos

- Basic idea: shared secrets
 - User proves to KDC who he is; KDC generates shared secret between client and file server



31

Kerberos Interactions



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

32

AFS Security (Kerberos)

- Kerberos has multiple administrative domains (realms)
 - principal@realm
 - srini@cs.cmu.edu sseshan@andrew.cmu.edu
- Client machine presents Kerberos ticket
 - Arbitrary binding of (user,machine) to Kerberos (principal,realm)
 - dongsuh on grad.pc.cs.cmu.edu machine can be srini@cs.cmu.edu
- Server checks against access control list (ACL)

33

AFS ACLs

- Apply to directory, not to file
- Format:
 - sseshan rlidwka
 - srini@cs.cmu.edu rl
 - sseshan:friends rl
- Default realm is typically the cell name (here andrew.cmu.edu)
- Negative rights
 - Disallow "joe rl" even though joe is in sseshan:friends

34

Topic 3: Client-Side Caching

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

35

Client Caching in NFS v2

- Cache both clean and dirty file data and file attributes
- File attributes in the client cache expire after 60 seconds (file data doesn't expire)
- File data is checked against the modified-time in file attributes (which could be a cached copy)
 - Changes made on one machine can take up to 60 seconds to be reflected on another machine
- Dirty data are buffered on the client machine until 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

36

Implication of NFS v2 Client Caching

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

37

Client Caching in AFS v2

- 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
 - Why? What’s the disadvantage of doing so?
- However, when a client caches file data, it obtains a “callback” on the file
- In case another client writes to the file, the server “breaks” the callback
 - Similar to invalidations in distributed shared memory implementations
- Implication: file server must keep state!

38

AFS v2 RPC Procedures

- 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
 - What should the client do if a callback is revoked?
 - Store: store the status and optionally data of a file
- Rest are similar to NFS calls

39

Failure Recovery in AFS v2

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

40

Key to Simple Failure Recovery

- Try not to keep any state on the server
- If you must keep some state on the server
 - Understand why and what state 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

41

Topic 4: File Access Consistency

- 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
 - The kernel locks the file vnode
 - Each file write is immediately visible to all file readers
- Neither NFS nor AFS provides such concurrency control
 - NFS: “sometime within 30 seconds”
 - AFS: session semantics for consistency

42

Semantics of File Sharing

Method	Comment
UNIX semantics	Every operation on a file is instantly visible to all processes
Session semantics	No changes are visible to other processes until the file is closed
Immutable files	No updates are possible; simplifies sharing and replication
Transactions	All changes occur atomically

- Four ways of dealing with the shared files in a distributed system.

43

Session Semantics in AFS v2

- What it means:
 - A file write is visible to processes on the same box immediately, but not visible to processes on other machines 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
- 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

44

AFS Write Policy

- Data transfer is by chunks
 - Minimally 64 KB
 - May be whole-file
- Writeback cache
 - Opposite of NFS "every write is sacred"
 - Store chunk back to server
 - When cache overflows
 - On last user close()
 - ...or don't (if client machine crashes)
- Is writeback crazy?
 - Write conflicts "assumed rare"
 - Who wants to see a half-written file?

45

Access Consistency in the "Sprite" File System

- Sprite: a research file system developed in UC Berkeley in late 80's
- 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?

46

Implementing Sequential Consistency

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

47

Implication of "Sprite" Caching

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

48

“Tokens” in DCE DFS

- How does one implement sequential consistency in a file system that spans multiple sites over WAN
- 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
 - “read” token: cached data are valid
 - “write” token: can write to data and keep dirty data at client
 - Status tokens: provide guarantee of file attributes
 - “read” status token: cached attribute is valid
 - “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

49

Compatibility Rules for Tokens

- 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?
- Data tokens: R/W and W/W are incompatible if the byte range overlaps
- Status tokens: R/W and W/W are incompatible
- Data token and status token: compatible or incompatible?

50

Token Manager

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

51

Topic 5: File Locking for Concurrency Control

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

52

Issues in Locking Implementations

- 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?

53

AFS Locking

- Locking
 - Server refuses to keep a waiting-client list
 - Client cache manager refuses to poll server
 - User program must invent polling strategy

54

Wrap up: Design Issues

- Name space
- Authentication
- Caching
- Consistency
- Locking

55

AFS Retrospective

- Small AFS installations are hard
 - Step 1: Install Kerberos
 - 2-3 servers
 - Inside locked boxes!
 - Step 2: Install ~4 AFS servers (2 data, 2 pt/vldb)
 - Step 3: Explain Kerberos to your users
 - Ticket expiration!
 - Step 4: Explain ACLs to your users

56

AFS Retrospective

- Worldwide file system
- Good security, scaling
- Global namespace
- "Professional" server infrastructure per cell
 - Don't try this at home
 - Only ~190 AFS cells (2002-03)
 - 8 are cmu.edu, 14 are in Pittsburgh
- "No write conflict" model only partial success

57