

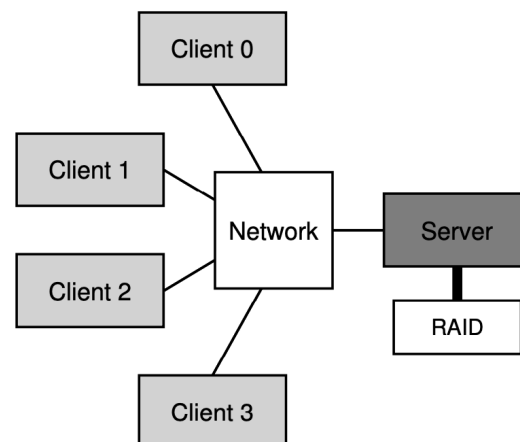
# Distributed Systems

- 48 - *Communication Basics*
- 49 - NFS
- 50 - AFS
- GFS

121

## NFS *Sun Microsystems*

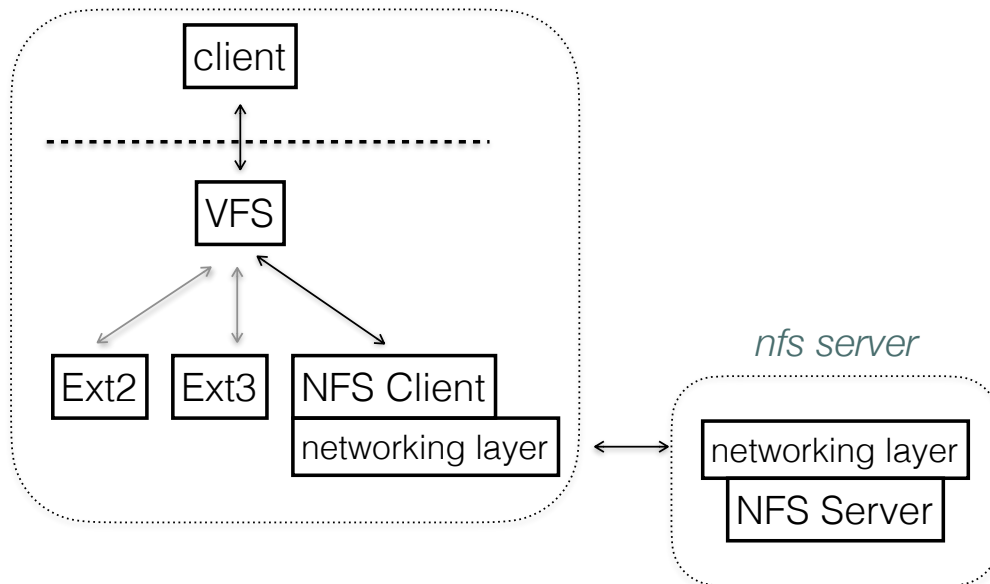
- first widely used distributed file system
  - clients diskless
    - easy sharing
    - centralized admin
    - security



122

# NFS

- distributed file system should be *transparent*
  - except possibly in performance
  - client issues same file-system calls as standalone system



123

# NFS *actually NFSv2*

“a distributed system is one where a machine I've never heard of goes down and I can't read my email”

- Leslie Lamport: Turing Award Winner for his work on distributed systems

- NFS goals:
  - simple and fast file recovery
  - *stateless* protocol : *server keeps no client state*
    - server scales well
    - client crashes transparent
    - server crashes transparent
    - client must maintain all state the the server needs for any communication

124

# NFS *actually NFSv2*

- file handle : uniquely describe file or directory
  - volume ID
  - inode number
  - generation number (in numbers get re-used)

NFSPROC_GETATTR	file handle returns: attributes
NFSPROC_SETATTR	file handle, attributes returns: -
NFSPROC_LOOKUP	directory file handle, name of file/dir to look up returns: file handle
NFSPROC_READ	file handle, offset, count data, attributes
NFSPROC_WRITE	file handle, offset, count, data attributes
NFSPROC_CREATE	directory file handle, name of file, attributes -
NFSPROC_REMOVE	directory file handle, name of file to be removed -
NFSPROC_MKDIR	directory file handle, name of directory, attributes file handle
NFSPROC_RMDIR	directory file handle, name of directory to be removed -
NFSPROC_READDIR	directory handle, count of bytes to read, cookie returns: directory entries, cookie (to get more entries)

125

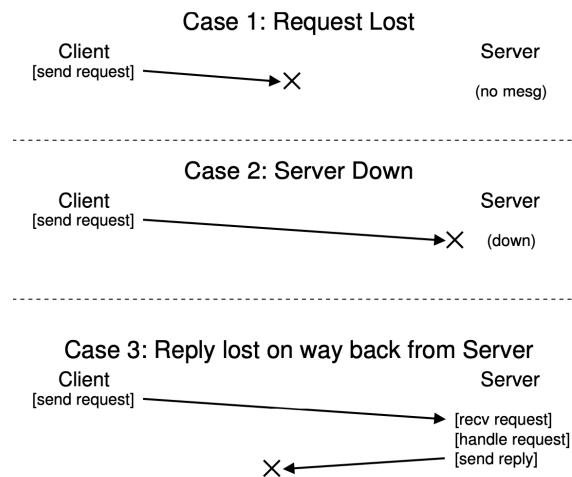
# NFS *reading a file :*

Client	Server
<b>fd = open("/foo", ...);</b> Send LOOKUP (rootdir FH, "foo")	Receive LOOKUP request look for "foo" in root dir return foo's FH + attributes
Receive LOOKUP reply allocate file desc in open file table store foo's FH in table store current file position (0) return file descriptor to application	
<b>read(fd, buffer, MAX);</b> Index into open file table with fd get NFS file handle (FH) use current file position as offset Send READ (FH, offset=0, count=MAX)	Receive READ request use FH to get volume/inode num read inode from disk (or cache) compute block location (using offset) read data from disk (or cache) return data to client
Receive READ reply update file position (+bytes read) set current file position = MAX return data/error code to app	
<b>read(fd, buffer, MAX);</b> Same except offset=MAX and set current file position = 2*MAX	
<b>read(fd, buffer, MAX);</b> Same except offset=2*MAX and set current file position = 3*MAX	
<b>close(fd);</b> Just need to clean up local structures Free descriptor "fd" in open file table (No need to talk to server)	

126

## NFS *server failures*

- server crashes / restarts, knowing nothing about clients
  - because most client requests are *idempotent*
    - lookups, reads don't change server state
    - writes contain data and exact offset to write to
- client handles all timeouts in the same way



127

## NFS *performance*

- client-side caching
  - read file data (and metadata) cached by client
  - all good unless the file changes on the server
- client-side write buffers
  - coalescing
  - aggregating disparate messages
  - writes sent back to server asynchronously (but before close())
- However : cache consistency!

128

# NFS *cache consistency*

Problems:

- **update visibility**
  - $C_1$  writes `foo.c`, but does not immediately push to server
  - $C_2$  reads, **sees old version**
  - $C_1$  flushes to server
- **stale cache**
  - $C_2$  closes and reads again, sees old version (`foo.c` locally cached)

Fixes:

- **close-to-open consistency**
  - every open guaranteed to see every prior write to the server
    - must validate cache (GETATTR)
    - but maybe not all the time

*NFS consistency is weak... (so are most other FSs)*

129

# NFS *server caching*

- **tons of memory**
  - wants to use it for disk cache (satisfy reads)
  - wants to use it for write buffer (quickly ack writes)
    - what could go wrong?
- **server could ack a write before writing to disk!**

- say file initially has three 4k blocks of data:

```
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
YYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYY
ZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZ
```

- client overwrites with:

```
write(aaa..., 0)., write(bbb..., 4k), write(ccc..., 8k):
```

- server crashes after acking second block, before writing:

```
aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa
YYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYY <--- oops
ccccccccccccccccccccccccccccccccccccccccccccccccccccccccccc
```

- client never evens knows that the server crashed

130

## NFS *cache consistency*

Problem: *poor performance for client<sub>i</sub>; the same file again*

- *fix: allow client<sub>i</sub> to cache data and attributes on client*
  - *but when client<sub>i</sub> re-opens not guaranteed most recent version*
- *fix: have clients re-validate on open*
  - *but slow*
- *fix: time out the cached attributes*
  - *means data can all be cached, attributes sometimes validated w/ server before accesses*
  - *but when client<sub>i</sub> re-opens not guaranteed most recent version (still)*

non-fix: *NFS consistency is weak... (same true for other FS's)*

131

## NFS *innovations*

- *stateless protocol*
  - *minimizes state server needs to track*
  - *server can crash and recover w/o clients being aware*
- *idempotent requests*
  - *necessary for statelessness*
  - *client treats network message drops, server failure the same*
  - *client does not need to know which is which*
- *client and server buffering*
  - *essential for performance*
  - *cache consistency issues*
    - *server flushes writes before acking*
    - *client attribute cache times out*
- *VFS interface*
  - *makes application API independent of underlying FS*

132

# NFS *later versions*

- version v3
  - 64-bit sizes and offsets (large files)
  - synchronous server writes
  - readdirplus (reads dir, also includes the file handles)
  - tcp
- version v4
  - strong security (kerberos, public key protocols)
  - performance improvements
  - stateful protocol (mostly for file consistency)
  - open standard (IETF)

133

# Distributed Systems

- 48 - *Communication Basics*
- 49 - *NFS*
- 50 - *AFS*
- *Review*

134

# Andrew File System *AFS v1*

- primary motivation was *scale*
  - how many clients could a single server accommodate?
  - user-visible behavior well-defined
  - whole-file (not block) caching

TestAuth	Test whether a file has changed (used to validate cached entries)
GetFileStat	Get the stat info for a file
Fetch	Fetch the contents of file
Store	Store this file on the server
SetFileStat	Set the stat info for a file
ListDir	List the contents of a directory

Figure 50.1: AFSv1 Protocol Highlights

135

# Andrew File System *AFS v2*

- Problems w/ v1:
  - full path traversal costs (on the server!)
  - client issues too many TestAuth msgs
  - also:
    - load not balanced across servers (fix using volumes)
    - server has a process per client (fix using threads)
- Improving the protocol:
  - client callbacks:
    - promise from the server to notify client if cached file changed
  - file identifier (FID)
    - volume id
    - file id
    - “uniquifier” (usually called epochs elsewhere)

136



# AFS

*example*

Client (C <sub>1</sub> )	Server
<code>fd = open("/home/remzi/notes.txt", ...);</code> Send Fetch (home FID, "remzi")	Receive Fetch request look for remzi in home dir establish callback(C <sub>1</sub> ) on remzi return remzi's content and FID
Receive Fetch reply write remzi to local disk cache record callback status of remzi Send Fetch (remzi FID, "notes.txt")	
Receive Fetch reply write notes.txt to local disk cache record callback status of notes.txt local <code>open()</code> of cached notes.txt return file descriptor to application	Receive Fetch request look for notes.txt in remzi dir establish callback(C <sub>1</sub> ) on notes.txt return notes.txt's content and FID
<code>read(fd, buffer, MAX);</code> perform local <code>read()</code> on cached copy	
<code>close(fd);</code> do local <code>close()</code> on cached copy if file has changed, flush to server	
<code>fd = open("/home/remzi/notes.txt", ...);</code> Foreach dir (home, remzi) if (callback(dir) == VALID) use local copy for lookup(dir) else Fetch (as above) if (callback(remzi) == VALID) open local cached copy return file descriptor to it else Fetch (as above) then open and return fd	

## Andrew File System *cache consistency*

Mentioned two issues w/ NFS:

- *update visibility*
  - when will server be updated w/ client write?
- *cache staleness:*
  - when will clients be informed their versions are out of date?
- *AFS procedure:*
  - client writes, possibly many times
  - *closes*
    - writes complete file back to server, becomes visible
    - server *breaks callback*
      - contact each server w/ a callback and invalidate its copy

*all apps on single machine see same copy*

# Andrew File System *cache consistency*

P <sub>1</sub>	Client <sub>1</sub>		Client <sub>2</sub>		Server Disk	Comments
	P <sub>2</sub>	Cache	P <sub>3</sub>	Cache		
open(F)		-		-	-	File created
write(A)		A		-	-	
close()		A		-	A	
	open(F)	A		-	A	
	read() → A	A		-	A	
	close()	A		-	A	
	open(F)	A		-	A	
	write(B)	B		-	A	
	open(F)	B		-	A	Local processes see writes immediately
	read() → B	B		-	A	
	close()	B		-	A	
		B	open(F)	A	A	Remote processes do not see writes...
		B	read() → A	A	A	
		B	close()	A	A	
close()		B		<del>A</del>	B	... until close() has taken place
		B	open(F)	B	B	
		B	read() → B	B	B	
		B	close()	B	B	
		B	open(F)	B	B	
open(F)		B		B	B	
write(D)		D		B	B	
		D	write(C)	C	B	
		D	close()	C	C	
close()		D		<del>C</del>	D	
		D	open(F)	D	D	Unfortunately for P <sub>3</sub> the last writer wins
		D	read() → D	D	D	
		D	close()	D	D	

139

# Andrew File System *cache consistency*

- AFS provides also *close-to-open consistency*
  - whole-file caching and updating
    - never see concurrent writes diff clients in same version of a file
  - “last writer wins” (really last *closer* wins)
- Crash recovery complicated
  - crashing client might miss callback (*client treats cache as suspect after crash*)
  - crashing server loses callbacks table
    - server might inform all clients after recovery
    - or clients constantly check for server liveness w/ heartbeats
  - *there is a cost to building a more sensible and scalable caching model*

140

# NFS vs AFS

- primarily differ in caching
  - What to cache?
    - NFS caches blocks
    - AFS entire files (on disk)
  - When to push writes to server?
    - Loosely defined for NFS:
      - any time from *right away*, to *when file is closed*
      - (only modified blocks)
    - If any part modified, AFS pushes entire file at `close()`
  - Final contents after concurrent merges by different clients:
    - NFS: writes by the different clients might be intermingled
    - AFS: final version reflects *the last write*; other write is lost

141

## Exam 2 *review*

- exam topics:
  - disk performance: perf from latencies
  - disk scheduling: SSTF, CSCAN, SCAN
  - RAID 0, 1, 5
  - FFS: advantages, write order
  - journaling, meta-data journaling
  - SSDS: simple mapping table, hybrid mapping table,
  - end-to-end argument
  - LFS: structure, write cost calculations
  - NFS and AFS: structure, cache consistency
- what to review
  - quizzes 7-10
  - lectures after spring break

142

# RH 7 review

Q1

5 Points

Given disk:

- 6000 RPM
- 200 sectors/track
- sector is 8KB
- avg. seek time 2 msec
- read/write no difference
- no track caching

**Write only a number in each of the boxes for this question: no explanation, no units, no nothing.**

Q1.1

1 Point

In msec, what is the average rotational latency?

5

Explanation

6000 RPM  $\Rightarrow$  100/sec  $\Rightarrow$  10 msec per rotation.  
Average latency would be half of this, i.e. 5 msec.

Q1.2

1 Point

In msec, what is the average sector transfer time?

0.05

Explanation

10 msec/rot, 1/200th of a rotation

Q1.3

1 Point

In msec, what is the average cost of a random 4k read?

7.05

Explanation

$= \text{seek} + \text{latency}_{\text{rotational}} + \text{transfer}_{\text{sector}}$   
 $= 2 + 5 + 0.05 = 7.05 \text{ msec}$   
you have to read an entire sector

# RH 7 review

Q1.5

1 Point

In msec, what would be the minimum expected cost of reading 10 sequentially ordered sectors?

7.5

Explanation

The minimum expected would be if they were all laid out in the same track, so we only pay seek time and rotational latency once. After that we just pay the transfer time for the rest of the sectors.

$$7.05 + 9 * 0.05 = 7.5 \text{ msec}$$

# RH 8 *review*

## Q2

5 Points

List the writes that should occur when creating a 100-byte file `bar.c` in the directory `/foo` for a generic, non-journaing file system, in a correct order (there may be more than one):

### Explanation

write data block bitmap (async)

write `bar.c` data (async)

write inode bitmap (async)

write `bar.c` inode

write `/foo` data

write `/foo` inode

## Q3

5 Points

Instead, how many disk writes would a log-structured file system ideally issue?

(enter just an integer)

1