Review

- Congestion control
 - AIMD additive increase, multiplicative decrease
 - Internet stability relies on politeness (!)
- Send/Recv, RPC
 - RPC seems like magic bullet, yet...

Outline:

- NW Performance/LogP
- Distributed FS intro

Network (I/O) performance: LogP model.

--> If you want good network performance, need to **pipeline** requests [picture]

Pipeline more complex than for CPU because (a) not fixed number of stages, (b) not balanced stages... essentially, CPU designed around pipeline --> simpler to think about it's pipline. In distributed systems, need a more complete model.

N3:



FIGURE 1. LogP parameters in a generic platform

LogP model – similar to Patterson/Hennessy "Iron triangle of performance": CPI, IPS (frequency), nInstructions

Latency – Elapsed wallclock time from X to Y.

note: "latency" alone is ambiguous; need to say latency from what to what; "latency from x to y". E.g., 1-way latency v. round-trip latency v. ...

e.g., how long from when one byte packet sent to when it is received

b/c pipeline is complex, need to specify which pipeline stages you are talking about (sometimes care about different subsets of pipelines)

NOTE: can hide latency with **pipelining**; latency does not imply resource is busy. Latency tells you how deep pipeline needs to be.

Overhead – Time for first pipeline stage.

Bottleneck time to initiate operations. (Can't be overlapped.) e.g., Cpu time to put packet on wire.

Throughput, bandwidth – (1/gap) -- time for bottleneck stage. Maximum steady state rate.

Time consumed by slowest pipeline stage.

e.g., maximum bytes per second

How does overhead differ from latency? Overhead: resource usage Latency: real-time end-to-end delay How would you measure latency of a network request? How does overhead differ from latency?

How does overhead differ from bandwidth? How would you measure overhead of sending a packet? How would you measure bandwidth of a network?

Example

E.g., suppose you open a TCP connection and start sending 1KB messages to another node on a 10Mbit/s Ethernet



1) What is "bottleneck rate"? (for overhead, BW)

The only tricky thing about this is that you have complex pipelining models (e.g., a disk request "occupies" CPU, bus, scsi controller, scsi bus, disk arm)

Which one is the bottleneck depends on configuration (how many disks? How many SCSI busses? How fast CPU?)

Which one is the bottleneck depends on how question is asked: E.g., "For a Seagate Barracuda 5100 disk, what is the average overhead per 1-sector disk request?" v. "For a Dell Dimension 5100, what is the overhead per 1-sector disk request?" The first is asking how long a disk seek and rotation take; the second is asking how long the CPU is busy to set up a request.

Need to consider: What bottleneck is the question asking about? For throughput, steady state bottleneck is the same in both cases. For overhead, first stage overhead differs.

Examples:

Latency – significant fraction of the speed of light (1 foot/ns) \rightarrow <1us anywhere in building

Overhead (network send/receive)-- 10's-100's us to send/receive TCP/IP packet; 1-10us for streamlined protocols [note: old #'s]

Bandwidth -- 1Mbit/s 3G phone, 1-10 Mbit/s home internet connection, 10-50Mbit/s WiFi, 100-1000Mbit/s desktop, 1Gbit/s-10Gbit/s data center network

	Throughput	Overhead	100 byte	4KB	Remote
					4kB read
TCP/IP	10 Mbit/s	0.1 ms	.1ms +	.1ms +	3.3ms
Wireless			.08ms	3ms	
TCP/IP	100 Mbit/s	0.1 ms	.1ms +	.1ms +	0.5 ms
Ethernet			.008ms	.3ms	
TCP/IP	1000	0.1 ms	.1ms +	.1ms +	0.2 ms
Gigabit	Mbit/s		.0008ms	.03ms	
Ethernet					
AM/	1200Mbit/s	.007ms	.007ms +	.007ms +	.04ms
Myrinet			.001ms	.03ms	

Example

Create 1MB file using NFSv3. Close --> client needs to write back 1MB to server. How long?

Assume client sends one 4KB block at a time, waits for server to get block safely to disk.

[[See NFSExample.tex for better picture]]

```
100Mbit/s network; o_send = o_recv = 100us
each block:
100us (o_send) + 4KB/100Mbit/s + 1us (latency from last byte off NIC to last byte arrives) + 100us (o_recv) + 10ms (disk) + 100us (o_send) + .5KB/100Mbit/s + 1us + 100us (o_recv)
= 100us + 300us + 1us + 100us + 10000us + 100us + 35us + 1us + 100us = 10702us
1MB/4KB = 256 blocks 256 * 10702us = 2.75s
1 Gbit/s network
each block:
100us + 4KB/1Gbit/s + 1us + 100us + 10ms + 100us + .5KB/1Gbit/s + 1us + 100us = 100us + 30us + 1us + 100us + 100us + 100us + 100us + .5KB/1Gbit/s + 1us + 100us + 100us + 100us + 100us + .5KB/1Gbit/s + 1us + 100us = 100us + 30us + 1us + 100us + 100us + 100us + 3.5us + 1us + 100us = 100us + 30us + 1us + 100us + 100us + 3.5us + 1us + 100us = 100us + 30us + 1us + 100us + 100us + 3.5us + 1us + 100us = 100us + 30us + 1us + 100us + 100us + 3.5us + 1us + 100us = 10436us
```

256 * 10436 = 2.74s

MORAL: fast network doesn't buy you much if you haven't paid attention to latency and overhead.

Example (cont)

---> Instead of sending one block at time, send all blocks. Then wait for server to send "Ack" saying all on disk

1 Gbit/s network:

-- bottleneck is o_send and o_recv

-- picture

- -- Time: 256*100us (now last packet starting to go on wire)
 - + 4KB/1000Mbit/s (now last packet entirely on wire)
 - + 1us (now last byte of last packet arriving at receiver)
 - + 100us (now last packet received at receiver)
- + 4KB/100MB/s (now last packet on disk; assuming end of streaming, sequential write)
 - aming, sequential write)
 - + 100us (now ack is on wire)
 - + 256B/1000Mbit/s (now ack on wire)
 - + lus (now last byte of ack is at receiver)
 - + 100us (now ack received at receiver)
 - = 25600us
 - + 30us
 - + 1us
 - + 100us
 - + 40us
 - + 100us
 - + 1.5us
 - + 1us
 - + 100us
 - = 26236us
 - = 26ms

MORAL: Need to pipeline to get good performance from IO systems

Notes on example

-- NFS v3 kind of works like first example (except 5-10 outstanding requests at a time instead of 1); NFS v4 adds support for something like second approach

-- Send/receiver overheads in example are clearly too high. Any modern machine can keep a 100Mbit/s or even 1Gbit/s network "full" of 4KB messages.

(What does o_send need to be?)

What is latency if go cross-country? 3000 miles * 5000 ft/mile \rightarrow 15ms now 4KB read dominated by latency for all networks

Key to good performance:

- \blacksquare (1) in LAN minimize overhead
- (2) in WAN keep pipeline full

Example: LogP network benchmark

LogP Performance Assessment of Fast Network Interfaces www.cs.rutgers.edu/~rmartin/papers/papers/micropaper.ps

> Sender: (uniprocessor) Thread1: Start timer repeat M times send msg spin for delta Stop timer

Thread2: while (1) receive msg

Receiver: while(1) receive msg send reply



FIGURE 3. Request issue time-line



FIGURE 4. Expected microbenchmark signature



FIGURE 6. Myrinet microbenchmark signature

Example: Batching

General rule of thumb: OS provides abstraction of byte transfers, but batch into block I/O for efficiency (pro-rates overhead and latency over larger unit)

Example

- Suppose CPU takes 100us of processing to issue one 512 byte write request
- Each request is to a random sector on disk
- Disk has parameters as above (4ms avg seek, 3ms ½ rot, transfer .02ms)
- 32KB write buffer on disk (producer/consumer bounded buffer)

• Writes are issued asynchronously (CPU can issue k+1 as soon as k is in write buffer)



(1) Suppose CPU issues k back-to-back requests, when does CPU complete?



(2) When does first write to disk complete at the disk?

(e.g., latency from when first write starts at CPU until done at disk?)

7.1ms

(3) Suppose there are 500 writes in a burst, when does the last write complete at the disk?

100us + 500 * 7ms

Distributed file system

Outline Distributed File Systems 2 Case studies: NFS, AFS Crosscutting issues Performance Failures Cache coherence/consistency Distributed commit

> A **distributed file system** provides transparent access to files stored on a remote disk

Themes:

failures: what happens when a server crashes, but a client doesn't? Or vice versa?

Performance \rightarrow caching; use caching at both clients and server to improve performance

cache coherence – how do we make sure each client sees most up-todate copy?

Atomic update – how to update state at two or more machines

These issues and strategies we will discuss are much more general than file system – arise in many distributed systems.

Simple: no caching

use RPC to forward every file system request to remote server (e.g. Novell Netware)

Example operations: open, seek, read, write, close

Server implements each operation as it would for a local request and sends back result to client *straightforward utilization of RPC*



Advantage: server provides consistent view of file system to both A and B

issues: Failures, performance Failures – see NFS (below)

Performance can be lousy: going over network is slower than going to local memory! lots of network traffic server can be a bottleneck – what if lots of clients?

NFS (Sun Network File System)

(I'll talk about "NFS v.3" to illustrate issues in a simple system; NFS v. 4 makes significant changes, including some of the state-of the art techniques I'll talk about later this week...)

Idea: use caching to reduce network load

Cache file blocks, file headers, etc at both clients and servers



Advantage: if open/read/write/close can be done locally, no network traffic

Issues:

(1) no longer have automatic stub generation \rightarrow lose one advantage of "RPC" over message passing

(2) helps performance; challenges failures and cache consistency

Issues: part 1: cache consistency

What if multiple clients are sharing same files? Easy if they are both reading – each gets a copy of the file

What if one writing? How do updates happen?

At writer – NFS has hybrid delayed write/write through policy

• write through within 30 seconds or immediately when file closed

How does other client find out about change (it has cached copy, so doesn't see any reason to talk to the server)

NFS protocol, part 1: weak consistency

In NFS, client polls server periodically, to check if file has changed. Poll server if data hasn't been checked in last 3-30 seconds (exact timeout is tunable parameter)

Thus, when file is changed on one client, server is notified, but other clients use old version of file until timeout. They then check server, and get new version.



What if multiple clients write the same file? In NFS, can get either version (or parts of both). Completely arbitrary!

HTTP uses essentially same protocol

If rule #1 in CS is "any problem can be solved with an additional level of indirection", Dahlin's rule #2 is "I can make it go as fast as you want, as long as you don't need the right answer"

We'll talk about better ways to enforce consistency next week.

Issues, part 2: Failures

What if server crashes? Can client wait until server comes back up, and continue as before?

- 1) any data in server memory but not yet on disk can be lost
- 2) shared state across RPCs. Ex: open, seek, read. What if server crashes after seek? Then when client does "read", it will fail.
- 3) Message retries: suppose server crashes after it does UNIX "rm foo", but before acknowledgement?

Message system will retry – send it again. How does it know not to delete it again? (Could solve this with two-phase commit protocol, but NFS takes a more ad hoc approach – sound familiar?)

What if client crashes?1) Might lose modified data in client cache

NFS: Solve problems in protocol (ad hoc?)

NFS Protocol (part 2): solutions

Key idea: Server is **stateless.** Client not allowed to rely on any server state

1) write through caching – when a file is closed, all modified blocks are sent immediately to server disk. To the client "close" doesn't return until all bytes are stored on server disk.

Client caches dirty data until close. Client failure --> data loss. Network write (to server) -- block until data safely on disk.

- 2) Stateless protocol server keeps no state about client (except as hints to help improve performance; e.g. a cache)
 - each read request gives enough information to do entire operation ReadAt(inumber, position) not Read(openFile)
 - when server crashes and restarts, can start processing requests immediately, as if nothing happened

3) Timeout and repeat requests to mask lost messages

Standard RPC technique.

Simple solution:

Request/acknowledge protocol

Common case:

1) Sender sends message (msg, msgId) and sets timer

2) Receiver receives message and sends (ack, msgId)

3) Sender receives (ack, msgId) and clears timer

If timer goes off, goto (1)

How does this work? Local procedure call guarantes *exactly once* semantics. What does retransmission guarantee?

- What if msg 1 lost?
- What if ack lost?

Guarantees *at least once* semantics *assuming no machines crash or otherwise discontinue protocol*

■ Receiver guaranteed to recv message at least once

3)

 Operations are "idempotent": all requests are OK to repeat (all requests are done *at least once*). So, if server crashes between disk I/O and message send, client can resend message, server just does operation all over again

- read and write file block are easy just re-read or re-write file block; no side effects
- What about "remove"? NFS just ignores this problem does the remove twice; second time returns an error if file not found

5) Failures are transparent to client system

Is this a good idea? What should happen if server crashes? Suppose you are an application, in middle of reading a file, and server crashes?

Options;

- a) hang until server comes back up (next week)?
- b) return an error? Problem is: most applications don't know they are talking over the network we're transparent, right?

Many UNIX apps simply ignore errors! Crash if there is a problem. (Network \rightarrow many more errors than before)

NFS does both options – can select which one. Usually, hang and only return error if really must – if see "NFS stale file handle" that's why

NFS Summary

NFS pros & cons

+ simple

- + highly portable
- sometimes inconsistent

■ doesn't scale up to large # of clients

Might think NFS is really stupid, but Netscape/WWW does something similar: cache recently seen pages, and refetch them if they are too old. Nothing in WWW to help with cache coherence

Notice: what happened to "RPC \rightarrow transparent distributed system"?

- performance \rightarrow add caching
- failures \rightarrow change all public methods to (mostly) idempotent
- performance v. failures \rightarrow write through cache
- performance v. failures \rightarrow weak consistency

Basically ended up rearchitecting and rewriting everything!

Next 2 weeks -- address fundamental problems in distributed systems. You see them in NFS

- -- performance
- -- consistency
- -- distributed commit
- -- security

After we talk about (some of) these, we'll revisit in context of a scalable cluster file system: The Google File System

Performance

Cost of a procedure call << same machine RPC << network RPC

means programmer must be aware that RPC is cheap, but not free

Caching can help, but

- generally gets rid of "transparent stub generation" advantage of RPC
- makes failure handling more complex, raises consistency issues

Not work for all worlkloads, all cases. (E.g., web caching -- data changes, zipf distribution --> client caches have 20-50% hit rate --> network performance dominates (amdhal's law)

--> Caching alone can't fully mask slow network.

NFS Example:

File close needs to write back all dirty sectors from client cache to server disk.

Network performance

"How fast is your network?"

Bandwidth isn't whole story. Bandwidth is the MIPS of I/O In architecture, MIPS is one of three factors (cycles per instruction, instruction count, instructions per second) -- only looking at one is misleading Similar issues for IO

Example

Suppose I have a 100Mbps and 1000Mbps network. Is second network 10x faster? Not if I use it to do a "remote read" (50 byte request, 50 byte

response) Graph: (lab) 510us (100Mbps), 501us (1000 Mbps) (Graph: fixed portion + variable portion...) Cross-country: 50.5ms (10Mbps), 50.5ms (100Mbps) What's going on?

Example

e.g., Suppose I replace load/store from local memory with load/store from remote machine via network.

Bandwidth not *that* different -- maybe 10-100 GB/s v. 10 Gbit/s (2011) --> 10-100x

But slowdown would probably be many times that (1000x-100,000x)

Other factors

-- Latency. Speed of light to get across building (~us)/campus(100us)/country(10's of ms) (v. 100ns to memory)

-- Overhead. Thousands of instructions to send/receive a packet (100us to send/recv a packet)

Result: Even if network bandwidth is 10Gbit/s, if I only access one remote word at a time, I'll probably see an effective bandwidth of 1 word per 100us or 1ms (100-1000x slower)

So, if bandwidth alone can get you off by a factor of 1000x, how do you reason about performance?

Cache consistency

Today: cache consistency – callbacks, leases Wednesday: reliability

Recall -- NFS caching sometimes gives wrong answer

-- client caching data checks with server to see if still valid if it has been more than 30s since last check

--> Window of vulnerability

--> My compiles occasionally fail and I tear my hair out

"I can make any system run fast as long as you don't insist on the right answer."

Sequential ordering constraints

Cache coherence – what should happen? What if one Cpu changes file and before it's done, another CPU reads file?

"right answer" turns out to be more subtle than one might hope...

- Essentially same problem as reasoning about synchronization of multi-threaded programs we have several programs running on (potentially) multiple processors (at arbitrary speeds), what can they see as they read and write memory (or files)?
- But now even load/store may not be atomic operations
 - Caching, write buffering, multipath routing through network
 - $\circ \rightarrow$ write by one thread may not immediately be seen by another!

- Consistency/coherence/staleness semantics define how "non-atomic" memory can be (and give you a basis for reasoning about distributed programs)
 - Essentially ask "can a distributed program tell that memory is 'playing tricks on it' compared to case where all threads run on uniprocessor with single memory?"
 - Memory system semantics restrict/define which "new" behaviors a memory system (or file system) can expose to a program

consistency v. coherence v. staleness

Coherence restricts *order* of reads and writes to *one location* - Can you tell memory system is playing tricks on you by looking at one location? – Example P1: P2: for(ii = 0; ii < 100; ii++){</pre> while(1){ write(A, ii); printf(''%d '', read(A)); } } - Where is incoherence? 1233349109111213... - Why might a system exhibit incoherence? e.g., 2 nodes, writer sends updates via Internet; updates get reordered en route...

e.g., cooperative caching -- read cached value from two different peers, could get out-of-order answer

```
e.g., client switching between two servers (e.g., on Internet, get redirected to different Akamai node)
```

Staleness bounds the maximum (real-time) *delay* between writes and reads to one location. – Can you tell memory system is playing tricks on you by looking at clock?

```
At 1:00:01 price is 10.55
At 1:00:02 price is 10.65
At 1:00:02 price is 10.65
At 1:00:02 price is 10.65
At 1:00:05 price is 13.18
...
```

-- Why might a system exhibit staleness?

e.g., NFS polling interval

e.g., network delay prevents update/invalidation from reaching cache for a while...

Consistency restricts order of reads and writes *across locations* - Can you tell memory system is playing tricks on you by looking at multiple locations? - Example 1 P1: P2: for(ii = 0; ii < 100; ii++){</pre> while(1){ printf(''(%d, %d), '', write(A, ii); write(B, ii); read(A), read(B)); } } - Where is inconsistency?: (0,0), (0,1), (1,2), (4,3), (4,8), (8,9), (9,9), (9,10), (9,10), (10,10), (11,10), (11,11), (12,12), ... - Is there also incoherence? - Example 2 (a classic) P1: P2: write(A, 0); write(B, 0); write(A, 1); write(B, 1); if(read(B) == 0){ if(read(A) == 0){ printf(''P1.''); printf(''P2.''); } } - Which outputs are legal under strict coherence? Under sequential consistency? "P1." "P2." w "P1.P2." "P2.P1." - Why might a system exhibit inconsistency? - Notice In first example, order between writes must be maintained...fairly obvious notion of causality

In second example, order between writes and reads must be maintained. (Less obvious?) Consistency involves ordering both writes and reads.



Semantics for non-atomic memory

Above defines 3 axes/dimensions:



consistency

Now we can start talking about particular design points in this space.

One option: Insist that distributed memory look "just like" local memory

Definitions (from Tannenbaum *Distributed Systems (with slight modifications and additions)*)

- sequential consistency – The result of any execution is the same as if the (read and write) operations by all processes on the data store were executed in some sequential order and the operations of each individual process appear in the sequence in the order specified by its program

Sounds pretty strong (and it is). But it is not "perfect" - e.g.,

A B

output "A=0 B=0" is legal under sequential consistency!

 \rightarrow Expect certain staleness guarantees

- delta coherence - the maximum real-time delay between when a write completes and when a subsequent read begins such that the read must return a value at least as new as that write

- **strict coherence** - any read on a data item x returns a value corresponding to the most recent write on x

strict coherence = delta coherence, delta = 0

linearizability = sequential consistency + delta coherence + delta = 0 (formal definition is essentialy -- sequential consistency + the global sequence is consistent with real time)

--> linearizability is essentially the origin in the design space -- the strongest consistency we typically ask for

Even this is a bit less tight than you might hope... Simple to see what strict coherence means if reads and writes are instantaneous. But they are not!

Note that every operation takes time: actual read could occur anytime between when system call is started, and when system call returns



Assume what we want is distributed system to behave exactly the same as if all processes are running on a single UNIX system if read finishes before write starts, then get old copy if read starts after write finishes, then get new copy

Otherwise indeterminant can get either new or old copy

Similarly, if write starts before another write finishes, may get either old or new version. (Hence, in above diagram, non-deterministic as to which you end up with!)

In NFS, if read starts more than 30 seconds after write finishes, get new copy. Othewise, who knows? Could get partial update.

Regular v. atomic semantics

Regular semantics -- return either the value of the last completed write or that of one of the writes which are concurrent with the read.

Atomic semantics --- guarantee that the read and write operations to the variable behave exactly as if they happened instantaneously in some point in time which is within the actual time where the operation took place. (Usually this is what is assumed for strict coherence)

Strict coherence = delta coherence with delta = 0

Limitations of strong consistency

So, we can define "perfect" consistency/coherence/staleness.

Are we done? "Distribibuted systems should implement linearizability"???

Unfortunately, no. Implementing these semantics has costs. Some of these costs are fundamental (and sometimes they are unacceptable.)

- Sequential consistency has fundamental performance cost: fast reads or fast writes but not both

r + **w** >= **t** (where r is read time, \$w is the write time, and t is the minimal packet transfer time between nodes.) [Lipton and Sandberg]

- Sequential consistency has a fundamental **CAP dilemma** (Brewer): A system can not have sequential **C**onsistency and maintain 100% **A**vailability in the presence of **P**artitions.

 \rightarrow develop weaker models

- causal consistency – writes that are potentially causally related must be seen by all
processes in the same order. Concurrent writes may be seen in a different order on different
machines.

Basic idea -- if I see a write that you issued, then I can also see (at least) all writes you could have seen when you issued the write

- if \$P_1\$ reads a write \$A\$ before issuing a write \$B\$, then any process that sees \$B\$ cannot subsequently see the old value of \$A\$
- Hard to see why this is useful if you assume a centralized consistency server. Think about a world where machines can send writes to one another. If \$A\$ reads a bunch of writes from \$B\$ and then creates some writes of its own. Then \$C\$ synchronizes with \$A\$, \$A\$ must send \$B\$'s writes to \$C\$ before sending its own.



e.g., while(1) write(A, I++)	while(1) write(B, j++)	while(1) print(A, B)	while(1) print(A, B)
		(1,1) (1,2) (1,3) (1,4) (2,5)	(1,1) (2,1) (3,1) (4,1) (5,3)

This is causally consistent, but not sequentially consistent.

This would be useful if A and B were on different nodes on internet – I might see the closer node's updates before the more distant nodes, and you might see a different order...

e.g., while(1) write(A, I++)	while(1) write(B, readA)	while(1) print(A, B)
		(1,1) (1,2) (1,3) (1,4) (2,5)

No longer causally consistent - if I see new value of B, then I need to see new value of A

Theorem: Causally consistent is the strongest consistency you can provide without giving up availability. (Dahlin, Alvisi, Mahajan -- April 2011)

There are also weaker options

- **FIFO** consistency (aka **PRAM** consistency) – writes done by a single process are seen by all other processes in the order in which they were issued, but writes from different processes may be seen in different orders by different processes

- Is FIFO stronger or weaker than causal?

(A weaker semantic allows more legal orderings than a stronger semantic. Consistency semantic A is stronger than consistency semantic B if any sequence of read and write results that are legal in A are also legal in B but there is at least one sequence that is legal in B but that is not legal in A.)

Why would you ever want this? Requires no coordination at all. E.g., 2 web servers....

How to provide improved consistency across clients?

(1) Poll each read – send every read to central server → get centralized semantics (e.g., can get sequential consistency this way – see the global order?)

We can optimize this with getattr(), but still slow...

Callbacks (e.g., Sprite, Andrew File System)

AFS (CMU late 80's) → DCE DFS (commercial product)

Notify client if data they are caching is no longer valid



Callbacks:

- (1) When a client reads data from server, server remembers that client has data
- (2) When client writes data, server notifies all other clients (that are caching the data) that they must contact server on next read

Write begins Tell server "I want to write foo" Server tells all clients "discard current copy of foo" All clients acknowledge Server tells writer "ok to write foo" Write completes

What semantics does this provide?

- Inearizability if client issues one operation at a time and blocks until completion
- (weaker if I can, say, read from cache while my write is pending)

Fault tolerance 1: Recovery of callback state

AFS approach: protocol level design (e.g., ad-hoc)

Challenge: improved caching + consistency increases failure handling complexity:

What if server crashes? Lose all callback state

QUESTION: Why is this a problem?

QUESTION: What can you do? Reconstruct callback information from clients – go ask everyone "who has which files cached?"

QUESTION: What if client crashes?

Fault tolerance 2: CAP -- consistency v. availability during partitions

Key idea: Leases

CAP says sequential consistency must give up availability during partitions

How does this manifest in AFS?

Write completes when all caching clients have acknowledged QUESTION: why do I have to wait? [answer – you can return early if you are willing to weaken semantics... but if you want linearizability, you have to wait] Naïve solution: client blocks indefinitely if any client crashes

■ How does this scale as we increase # clients?

Solution: lease -- combine polling and callbacks

Lease: cache has the right to access cached object X for Y seconds; after Y seconds, must renew lease before accessing cached object

Server does callbacks for X seconds after lease

New solution:

- (1) Write waits until all caching nodes acknowledge or leases expire (sequential coherence)
- (2) Write returns immediately (delta coherence)

Enhancement: Volume lease...

Other AFS features

files cached on local disk
 NFS caches only in memory
 → reduce server load

- 2) more precise consistency model
 - 1) callbacks
 - \circ server records who has copy of file
 - o send "callback" on each update
 - 2) write-through on close

If file changes, server is updated (on close) Server then immediately tells those with old copy

 session semantics – updates visible only on close In UNIX (single machine) updates visible immediately to other programs who have file open In AFS, everyone who has file open sees old version; anyone who opens file again will see new version

In AFS slight variation: session semantics

- a) on open and cache miss get file from server; set up callback
- b) on write close: send copy to server; tells all clients with copies to fetch new version on next open

Essentially – think of all reads happening when file opened and all writes happening when file closed...

AFS pros & cons

Relative to NFS, less server load:

- + disk as cache \rightarrow more files can be cached locally
- + callbacks \rightarrow server not involved if file is read-only
- more complex recovery

Fault tolerance 3: Disconnected operation

Leases do a pretty good engineering job on CAP dilemma. If I can talk to server, I can access data. Clients disconnected from server are stuck. (Notice -- they are stuck even if they have the data they want to read in their cache.)

AFS stores data on local disk

Suppose server crashes – can client keep going?

almost – except renewing callbacks on open/close; writing though on close Support *disconnected operation* – allow client to access cached data even when it cannot contact server.

- Improve availability
- Support mobility

Coda (and NTFS)

- (1) Reads -- prefetch "hoard" data into local cache Want to make sure you have everything in cache you need. What should you do? (Hoard list)
- (2) Writes -- write updates to local log; send log to server when reconnect

Need to make sure that updates you did when disconnected make it back to server. What should you do? (Log writes + reconciliation)

CAP dilemma: Cannot provide sequential consistency and 100% availability in a system that can be partitioned.

■ What consistency does this provide? (causal?)

Problem: Conflicting writes...

What happens if two disconnected nodes both write same file? Is this OK?

Coda solution:

- (1) Detect
- (2) Regular files: manual selection of "right" version to keep
- (3) Directories: automatically correct most cases (manual for the rest)

Avoiding central server

Coda lets me write when disconnected, but all updates go through server

What if you don't want to have to synchronize through a server

Basic idea

- Each node's writes are a log [picture]
- Version vector index of highest known write from each node
- Log exchange you send me your VV, I send you all updates you have not yet seen
- \rightarrow Eventual consistency
- Lamport clock my accept stamp = max(VV) + 1
- Send elements from my log sorted by accept stamp
- \rightarrow Causal consistency
- Still need to deal with conflicting concurrent writes (how can you detect?)

Google file system [see gradOS notes]

Consistency in memcached

memcached: reading from a database is slow --> have another set of machines act as a cache (distributed hash table) [picture]

```
basic idea:
read(x)
data = memcached->read(x)
if(data) return data
else
data = db->read(x)
memcached->set(x, data)
return x
```

write(x, data) db->write(x, data) memcached->set(x, data)

```
What incoherence might you observe?
```

How long can incoherence last (how much staleness)

```
What if writer crashes after setting DB but before setting memcached?
```

```
Obvious fix (?)
```

write(x, data)
 memcached->clear(x);
 db->write(x, data);
 memcached->set(x, data);

Does this solve the problem?

Next time: improve consistency, 2 phase commit \rightarrow atomic distributed updates