

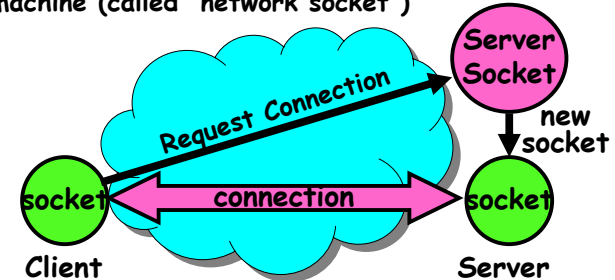
CS162 Operating Systems and Systems Programming Lecture 24

Distributed File Systems

November 23, 2005
Prof. John Kubiatowicz
<http://inst.eecs.berkeley.edu/~cs162>

Review: Network Communication

- **TCP**: Reliable byte stream between two processes on different machines over Internet (read, write, flush)
- **Socket**: an abstraction of a network I/O queue
 - Embodies one side of a communication channel
 - » Same interface regardless of location of other end
 - » Could be local machine (called "UNIX socket") or remote machine (called "network socket")



- **Two-phase commit**: distributed decision making
 - First, make sure everyone guarantees that they will commit if asked (prepare)
 - Next, ask everyone to commit

11/23/05

Kubiatowicz CS162 ©UCB Fall 2005

Lec 24.2

Review: Distributed Applications



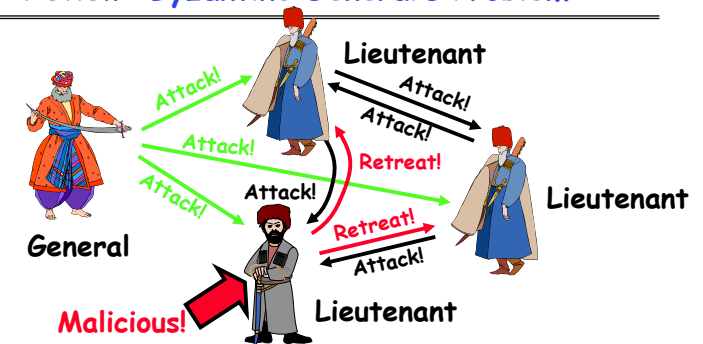
- **Message Abstraction**: send/receive messages
 - Already atomic: no receiver gets portion of a message and two receivers cannot get same message
- **Interface**:
 - Mailbox (mbox): temporary holding area for messages
 - » Includes both destination location and queue
 - Send (message, mbox)
 - » Send message to remote mailbox identified by mbox
 - Receive (buffer, mbox)
 - » Wait until mbox has message, copy into buffer, and return
 - » If threads sleeping on this mbox, wake up one of them

11/23/05

Kubiatowicz CS162 ©UCB Fall 2005

Lec 24.3

Review: Byzantine General's Problem



- **Byzantine General's Problem** (n players):
 - One General
 - n-1 Lieutenants
 - Some number of these ($f < n/3$) can be insane or malicious
- The commanding general must send an order to his n-1 lieutenants such that:
 - IC1: All loyal lieutenants obey the same order
 - IC2: If the commanding general is loyal, then all loyal lieutenants obey the order he sends

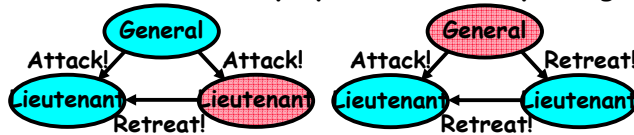
11/23/05

Kubiatowicz CS162 ©UCB Fall 2005

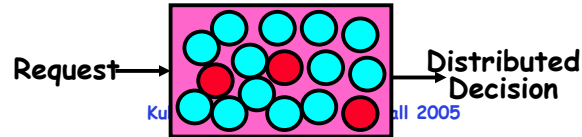
Lec 24.4

Review: Byzantine General's Problem (con't)

- **Impossibility Results:**
 - Cannot solve Byzantine General's Problem with $n=3$ because one malicious player can mess up things



- With f faults, need $n > 3f$ to solve problem
- Various algorithms exist to solve problem
 - Original algorithm has #messages exponential in n
 - Newer algorithms have message complexity $O(n^2)$
 - » One from MIT, for instance (Castro and Liskov, 1999)
- Use of BFT (Byzantine Fault Tolerance) algorithm
 - Allow multiple machines to make a coordinated decision even if some subset of them ($< n/3$) are malicious



11/23/05

Ku

all 2005

Lec 24.5

Review: Remote Procedure Call

- Raw messaging is a bit too low-level for programming
 - Must wrap up information into message at source
 - Must decide what to do with message at destination
 - May need to sit and wait for multiple messages to arrive
- Better option: Remote Procedure Call (RPC)
 - Calls a procedure on a remote machine
 - Client calls:


```
remoteFileSystem->Read("rutabaga");
```
 - Translated automatically into call on server:


```
fileSys->Read("rutabaga");
```
- Implementation:
 - Request-response message passing (under covers!)
 - "Stub" provides glue on client/server
 - » Client stub is responsible for "marshalling" arguments and "unmarshalling" the return values
 - » Server-side stub is responsible for "unmarshalling" arguments and "marshalling" the return values.
- **Marshalling** involves (depending on system)
 - Converting values to a canonical form, serializing objects, copying arguments passed by reference, etc.

11/23/05

Kubiatowicz CS162 ©UCB Fall 2005

Lec 24.6

Goals for Today

- Finish RPC
- Examples of Distributed File Systems
- Cache Coherence Protocols

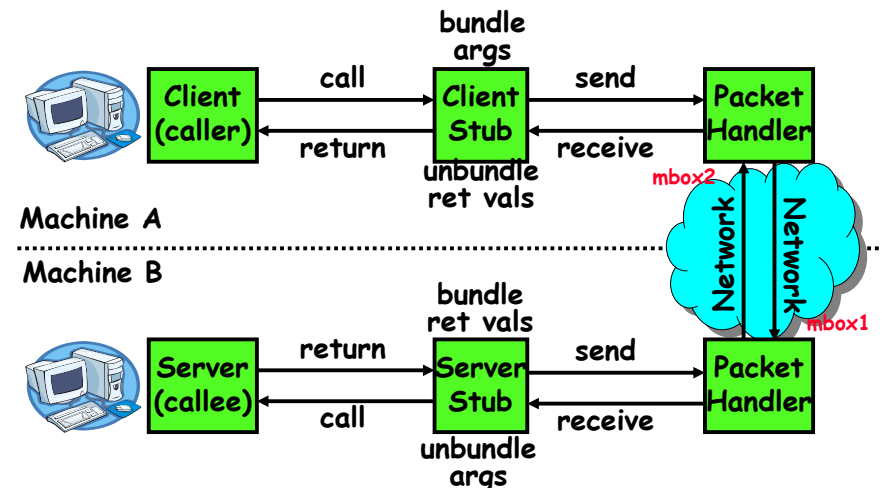
Note: Some slides and/or pictures in the following are adapted from slides ©2005 Silberschatz, Galvin, and Gagne

11/23/05

Kubiatowicz CS162 ©UCB Fall 2005

Lec 24.7

RPC Information Flow



11/23/05

Kubiatowicz CS162 ©UCB Fall 2005

Lec 24.8

RPC Details

- **Equivalence with regular procedure call**
 - Parameters \leftrightarrow Request Message
 - Result \leftrightarrow Reply message
 - Name of Procedure: Passed in request message
 - Return Address: mbox2 (client return mail box)
- **Stub generator: Compiler that generates stubs**
 - Input: interface definitions in an "interface definition language (IDL)"
 - » Contains, among other things, types of arguments/return
 - Output: stub code in the appropriate source language
 - » Code for client to pack message, send it off, wait for result, unpack result and return to caller
 - » Code for server to unpack message, call procedure, pack results, send them off
- **Cross-platform issues:**
 - What if client/server machines are different architectures or in different languages?
 - » Convert everything to/from some canonical form
 - » Tag every item with an indication of how it is encoded (avoids unnecessary conversions).

11/23/05

Kubiatowicz CS162 ©UCB Fall 2005

Lec 24.9

RPC Details (continued)

- **How does client know which mbox to send to?**
 - Need to translate name of remote service into network endpoint (Remote machine, port, possibly other info)
 - **Binding:** the process of converting a user-visible name into a network endpoint
 - » This is another word for "naming" at network level
 - » Static: fixed at compile time
 - » Dynamic: performed at runtime
- **Dynamic Binding**
 - Most RPC systems use dynamic binding via name service
 - » Name service provides dynamic translation of service \rightarrow mbox
 - **Why dynamic binding?**
 - » Access control: check who is permitted to access service
 - » Fail-over: If server fails, use a different one
- **What if there are multiple servers?**
 - Could give flexibility at binding time
 - » Choose unloaded server for each new client
 - Could provide same mbox (router level redirect)
 - » Choose unloaded server for each new request
 - » Only works if no state carried from one call to next
- **What if multiple clients?**
 - Pass pointer to client-specific return mbox in request

11/23/05

Kubiatowicz CS162 ©UCB Fall 2005

Lec 24.10

Problems with RPC

- **Non-Atomic failures**
 - Different failure modes in distributed system than on a single machine
 - Consider many different types of failures
 - » User-level bug causes address space to crash
 - » Machine failure, kernel bug causes all processes on same machine to fail
 - » Some machine is compromised by malicious party
 - Before RPC: whole system would crash/die
 - After RPC: One machine crashes/compromised while others keep working
 - Can easily result in inconsistent view of the world
 - » Did my cached data get written back or not?
 - » Did server do what I requested or not?
 - Answer? Distributed transactions/Byzantine Commit
- **Performance**
 - Cost of Procedure call \ll same-machine RPC \ll network RPC
 - Means programmers must be aware that RPC is not free
 - » Caching can help, but may make failure handling complex

11/23/05

Kubiatowicz CS162 ©UCB Fall 2005

Lec 24.11

Administrivia

- **My office hours**
 - No office hours Thursday (Thanksgiving)
- **Project 4 design document**
 - Due Tuesday November 29th
- **MIDTERM II: Monday December 5th!**
 - 5:30-8:30pm, 10 Evans
 - All material from last midterm and up to previous class
 - Includes virtual memory
- **Final Exam**
 - December 17th, 12:30 - 3:30, 220 Hearst Gym
- **Final Topics: Any suggestions?**

11/23/05

Kubiatowicz CS162 ©UCB Fall 2005

Lec 24.12

Cross-Domain Communication/Location Transparency

- How do address spaces communicate with one another?
 - Shared Memory with Semaphores, monitors, etc...
 - File System
 - Pipes (1-way communication)
 - "Remote" procedure call (2-way communication)
- RPC's can be used to communicate between address spaces on different machines on the same machine
 - Services can be run wherever it's most appropriate
 - Access to local and remote services looks the same
- Examples of modern RPC systems:
 - CORBA (Common Object Request Broker Architecture)
 - DCOM (Distributed COM)
 - RMI (Java Remote Method Invocation)

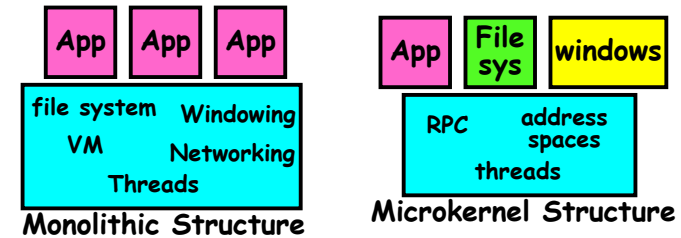
11/23/05

Kubiatowicz CS162 ©UCB Fall 2005

Lec 24.13

Microkernel operating systems

- Example: split kernel into application-level servers.
 - File system looks remote, even though on same machine



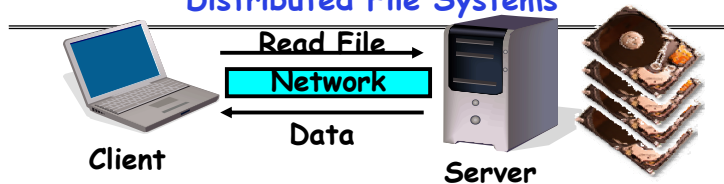
- Why split the OS into separate domains?
 - Fault isolation: bugs are more isolated (build a firewall)
 - Enforces modularity: allows incremental upgrades of pieces of software (client or server)
 - Location transparent: service can be local or remote
 - » For example in the X windowing system: Each X client can be on a separate machine from X server; Neither has to run on the machine with the frame buffer.

11/23/05

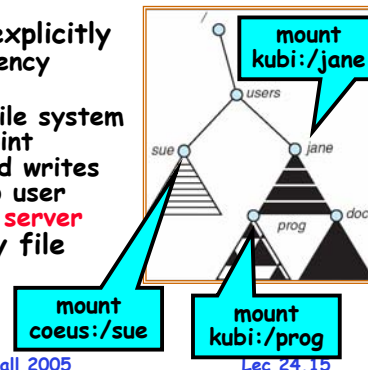
Kubiatowicz CS162 ©UCB Fall 2005

Lec 24.14

Distributed File Systems



- Distributed File System:
 - Transparent access to files stored on a remote disk
- Naming choices (always an issue):
 - *Hostname:localname*: Name files explicitly
 - » No location or migration transparency
 - *Mounting of remote file systems*
 - » System manager mounts remote file system by giving name and local mount point
 - » Transparent to user: all reads and writes look like local reads and writes to user e.g. `/users/sue/foo` → `/sue/foo` on server
 - *A single, global name space*: every file in the world has unique name
 - » Location Transparency: servers can change and files can move without involving user

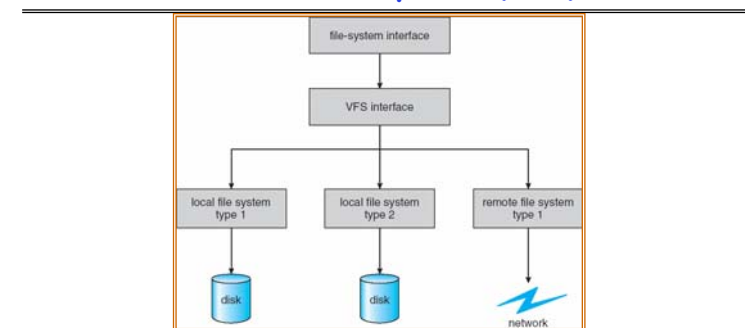


11/23/05

Kubiatowicz CS162 ©UCB Fall 2005

Lec 24.15

Virtual File System (VFS)



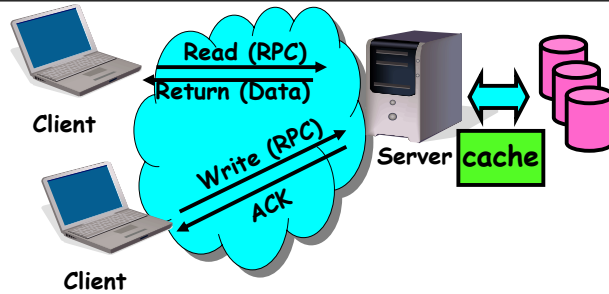
- **VFS**: Virtual abstraction similar to local file system
 - Instead of "inodes" has "vnodes"
 - Compatible with a variety of local and remote file systems
 - » provides object-oriented way of implementing file systems
- VFS allows the same system call interface (the API) to be used for different types of file systems
 - The API is to the VFS interface, rather than any specific type of file system

11/23/05

Kubiatowicz CS162 ©UCB Fall 2005

Lec 24.16

Simple Distributed File System



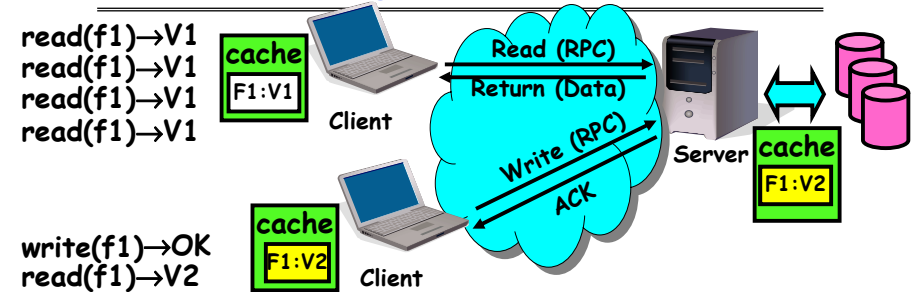
- Remote Disk: Reads and writes forwarded to server
 - Use RPC to translate file system calls
 - No local caching/can be caching at server-side
- Advantage: Server provides completely consistent view of file system to multiple clients
- Problems? Performance!
 - Going over network is slower than going to local memory
 - Lots of network traffic/not well pipelined
 - Server can be a bottleneck

11/23/05

Kubiatowicz CS162 ©UCB Fall 2005

Lec 24.17

Use of caching to reduce network load



- Idea: Use caching to reduce network load
 - In practice: use buffer cache at source and destination
- Advantage: if open/read/write/close can be done locally, don't need to do any network traffic...fast!
- Problems:
 - Failure:
 - » Client caches have data not committed at server
 - Cache consistency!
 - » Client caches not consistent with server/each other

11/23/05

Kubiatowicz CS162 ©UCB Fall 2005

Lec 24.18

Failures



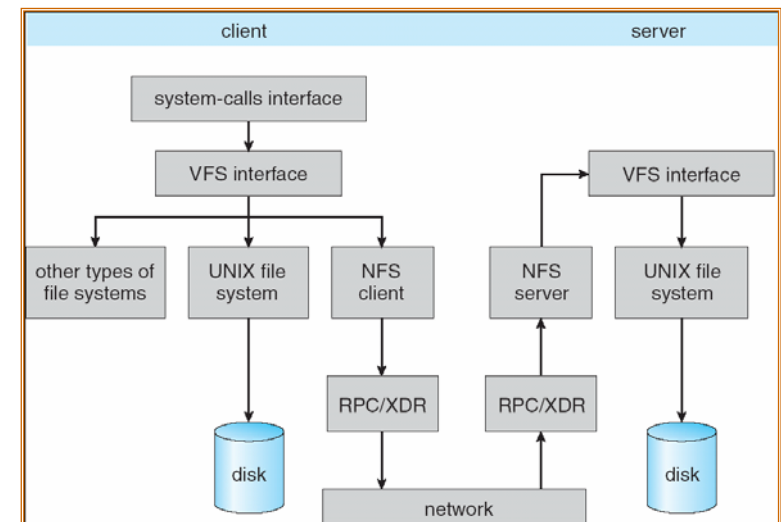
- What if server crashes? Can client wait until server comes back up and continue as before?
 - Any data in server memory but not on disk can be lost
 - Shared state across RPC: What if server crashes after seek? Then, when client does "read", it will fail
 - Message retries: suppose server crashes after it does UNIX "rm foo", but before acknowledgment?
 - » Message system will retry: send it again
 - » How does it know not to delete it again? (could solve with two-phase commit protocol, but NFS takes a more ad hoc approach)
- **Stateless protocol:** A protocol in which all information required to process a request is passed with request
 - Server keeps no state about client, except as hints to help improve performance (e.g. a cache)
 - Thus, if server crashes and restarted, requests can continue where left off (in many cases)
- What if client crashes?
 - Might lose modified data in client cache

11/23/05

Kubiatowicz CS162 ©UCB Fall 2005

Lec 24.19

Schematic View of NFS Architecture



11/23/05

Kubiatowicz CS162 ©UCB Fall 2005

Lec 24.20

Network File System (NFS)

- Three Layers for NFS system
 - **UNIX file-system interface**: open, read, write, close calls + file descriptors
 - **VFS layer**: distinguishes local from remote files
 - » Calls the NFS protocol procedures for remote requests
 - **NFS service layer**: bottom layer of the architecture
 - » Implements the NFS protocol
- NFS Protocol: RPC for file operations on server
 - Reading/searching a directory
 - manipulating links and directories
 - accessing file attributes/reading and writing files
- **Write-through caching**: Modified data committed to server's disk before results are returned to the client
 - lose some of the advantages of caching
 - time to perform write() can be long
 - Need some mechanism for readers to eventually notice changes! (more on this later)

11/23/05

Kubiatowicz CS162 ©UCB Fall 2005

Lec 24.21

NFS Continued

- NFS servers are **stateless**; each request provides all arguments require for execution
 - E.g. reads include information for entire operation, such as ReadAt(inumber, position), not Read(openfile)
 - No need to perform network open() or close() on file - each operation stands on its own
- **Idempotent**: Performing requests multiple times has same effect as performing it exactly once
 - Example: Server crashes between disk I/O and message send, client resend read, server does operation again
 - Example: Read and write file blocks: just re-read or re-write file block - no side effects
 - Example: What about "remove"? NFS does operation twice and second time returns an advisory error
- **Failure Model**: Transparent to client system
 - Is this a good idea? What if you are in the middle of reading a file and server crashes?
 - Options (NFS Provides both):
 - » Hang until server comes back up (next week?)
 - » Return an error. (Of course, most applications don't know they are talking over network)

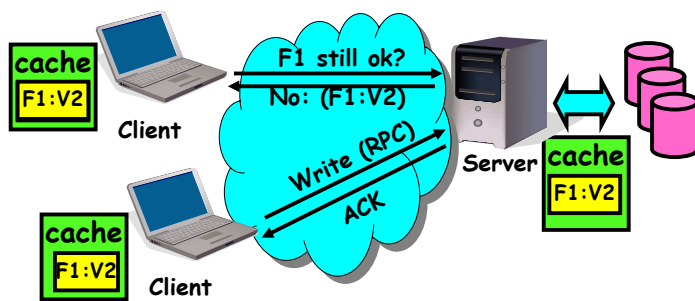
11/23/05

Kubiatowicz CS162 ©UCB Fall 2005

Lec 24.22

NFS Cache consistency

- NFS protocol: weak consistency
 - Client polls server periodically to check for changes
 - » Polls server if data hasn't been checked in last 3-30 seconds (exact timeout it tunable parameter).
 - » Thus, when file is changed on one client, server is notified, but other clients use old version of file until timeout.



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

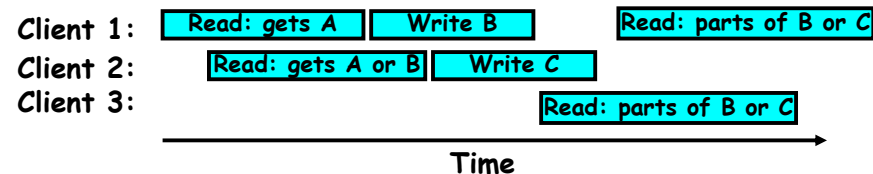
11/23/05

Kubiatowicz CS162 ©UCB Fall 2005

Lec 24.23

Sequential Ordering Constraints

- What sort of cache coherence might we expect?
 - i.e. what if one CPU changes file, and before it's done, another CPU reads file?
- Example: Start with file contents = "A"



- What would we actually want?
 - Assume we want distributed system to behave exactly the same as if all processes are running on single system
 - » If read finishes before write starts, get old copy
 - » If read starts after write finishes, get new copy
 - » Otherwise, get either new or old copy
 - For NFS:
 - » if read starts more than 30 seconds after write, get new copy; otherwise, could get partial update

11/23/05

Kubiatowicz CS162 ©UCB Fall 2005

Lec 24.24

NFS Pros and Cons

- NFS Pros:
 - Simple, Highly portable
- NFS Cons:
 - Sometimes inconsistent!
 - Doesn't scale to large # clients
 - » Must keep checking to see if caches out of date

11/23/05

Kubiatowicz CS162 ©UCB Fall 2005

Lec 24.25

Andrew File System

- Andrew File System (AFS, late 80's) → DCE DFS (commercial product)
- **Callbacks:** Server records who has copy of file
 - On changes, server immediately tells all with old copy
 - No polling bandwidth (continuous checking) needed
- Write through on close
 - Changes not propagated to server until close()
 - Session semantics: updates visible to other clients only after the file is closed
 - » As a result, do not get partial writes: all or nothing!
 - » Although, for processes on local machine, updates visible immediately to other programs who have file open
- In AFS, everyone who has file open sees old version
 - Don't get newer versions until reopen file

11/23/05

Kubiatowicz CS162 ©UCB Fall 2005

Lec 24.26

Andrew File System (con't)

- Data cached on local disk of client as well as memory
 - On open with a cache miss (file not on local disk):
 - » Get file from server, set up callback with server
 - On write followed by close:
 - » Send copy to server; tells all clients with copies to fetch new version from server on next open (using callbacks)
- What if server crashes? Lose all callback state!
 - Reconstruct callback information from client: go ask everyone "who has which files cached?"
- AFS Pro: Relative to NFS, less server load:
 - Disk as cache ⇒ more files can be cached locally
 - Callbacks ⇒ server not involved if file is read-only
- For both AFS and NFS: central server is bottleneck!
 - Performance: all writes→server, cache misses→server
 - Availability: Server is single point of failure
 - Cost: server machine's high cost relative to workstation

11/23/05

Kubiatowicz CS162 ©UCB Fall 2005

Lec 24.27

Conclusion

- **Remote Procedure Call (RPC):** Call procedure on remote machine
 - Provides same interface as procedure
 - Automatic packing and unpacking of arguments without user programming (in stub)
- **VFS:** Virtual File System layer
 - Provides mechanism which gives same system call interface for different types of file systems
- **Distributed File System:**
 - Transparent access to files stored on a remote disk
 - » NFS: Network File System
 - » AFS: Andrew File System
 - Caching for performance
- **Cache Consistency:** Keeping contents of client caches consistent with one another
 - If multiple clients, some reading and some writing, how do stale cached copies get updated?
 - NFS: check periodically for changes
 - AFS: clients register callbacks so can be notified by server of changes

11/23/05

Kubiatowicz CS162 ©UCB Fall 2005

Lec 24.28