

CS162 Operating Systems and Systems Programming Lecture 25

RPC, NFS and AFS

April 26th, 2022

Prof. Anthony Joseph and John Kubiatowicz

<http://cs162.eecs.Berkeley.edu>



Recall: Transmission Control Protocol (TCP)

- Transmission Control Protocol (TCP)
 - TCP ([IP Protocol 6](#)) layered on top of IP
 - Reliable byte stream between two processes on different machines over Internet (read, write, flush)
- TCP Details
 - Fragments byte stream into packets, hands packets to IP
 - » IP may also fragment by itself
 - » "Window" reflects storage at receiver – sender shouldn't overrun receiver's buffer space
 - » Also, window should reflect speed/capacity of network – sender shouldn't overload network
 - Automatically retransmits lost packets
 - » **Adjusts rate of transmission to avoid congestion**
 - » A "good citizen"

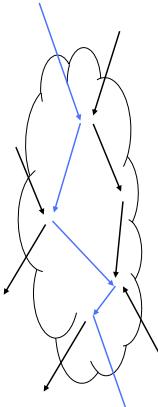
4/26/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 25.2

Congestion

- Too much data trying to flow through some part of the network



Congestion Avoidance

- Congestion
 - How long should timeout be for re-sending messages?
 - » Too long → wastes time if message lost
 - » Too short → retransmit even though ACK will arrive shortly
 - Stability problem: more congestion \Rightarrow ACK is delayed \Rightarrow unnecessary timeout \Rightarrow more traffic
 - Closely related to window size at sender: too big means putting too much data into network
- How does the sender's window size get chosen?
 - Must be less than receiver's advertised buffer size
 - Try to match the rate of sending packets with the rate that the slowest link can accommodate
 - Sender uses an adaptive algorithm to decide size of N
 - » Goal: fill network between sender and receiver
 - » Basic technique: slowly increase size of window until acknowledgements start being delayed/lost
- TCP solution: "slow start" (start sending slowly)
 - If no timeout, slowly increase window size (throughput) by 1 for each ACK received
 - Timeout \Rightarrow congestion, so cut window size in half
- "Additive Increase, Multiplicative Decrease"

4/26/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

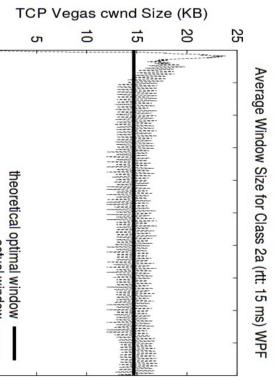
Lec 25.3

4/26/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 25.4

Congestion Management

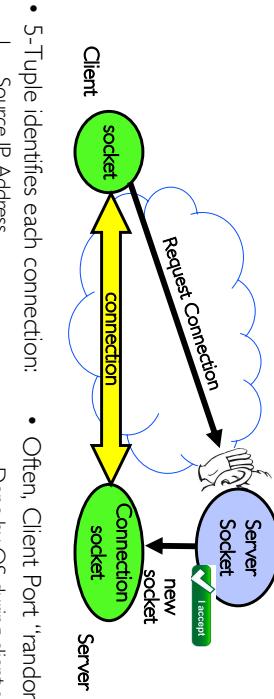


- TCP artificially restricts the window size if it sees packet loss
 - Careful control loop to make sure:
 1. We don't send too fast and overwhelm the network
 2. We utilize most of the bandwidth the network has available
 - In general, these are conflicting goals!
- From Low, Peterson, and Wang, "Understanding vegas: Duality Model". J. ACM, March 2002.

4/26/22

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

Lec 25.5



- 5-Tuple identifies each connection:
 1. Source IP Address
 2. Destination IP Address
 3. Source Port Number
 4. Destination Port Number
 5. Protocol (always TCP here)
 - Well-known ports from 0—1023
 - Done by OS during client socket setup
 - Server Port often “well known”
 - 80 (web), 443 (secure web), 25 (sendmail), etc

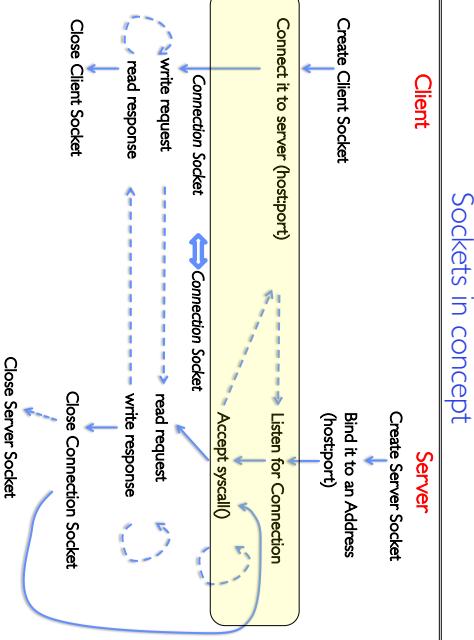
4/26/22

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

Lec 25.6

Establishing TCP Service

1. Open connection: 3-way handshaking
2. Reliable byte stream transfer from ($(\text{IP}_a, \text{TCP_Port}_1)$) to ($(\text{IP}_b, \text{TCP_Port}_2)$)
 - Indication if connection fails: Reset
3. Close (tear-down) connection



Lec 25.7

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

Lec 25.8

Recall: Connection Setup over TCP/IP

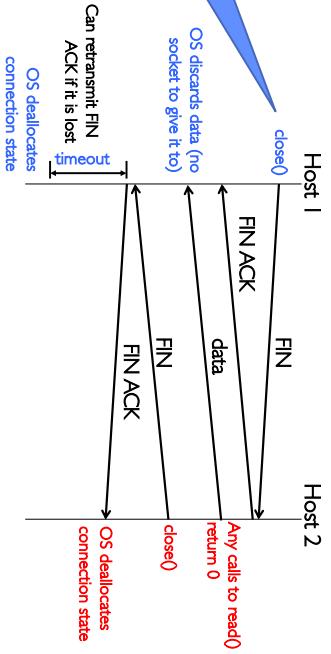
Open Connection: 3-Way Handshake

- Server calls `listen()` to wait for a new connection
 - Client calls `connect()` providing server's IP address and port number
 - Each side sends SYN packet proposing an initial sequence number (one for each sender) and ACKs the other
-
- $\text{SYN, SeqNum} = x$
- $\text{SYN and ACK, SeqNum} = y \text{ and } \text{Ack} = x + 1$
- $\text{ACK, Ack} = y + 1$
- time

4/26/22

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

Lec 25.9

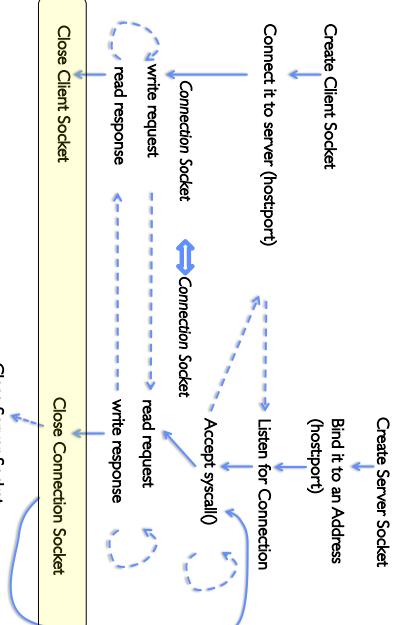


Lec 25.11

4/26/22

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

Sockets in concept



4/26/22

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

Lec 25.10

Recall: Distributed Applications Build With Messages



- How do you actually program a distributed application?
 - Need to synchronize multiple threads, running on different machines
 - » No shared memory, so cannot use test&set
- Interface:
 - Mailbox (mbox): temporary holding area for 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

Lec 25.12

4/26/22

Question: Data Representation

- An object in memory has a machine-specific binary representation
 - Threads within a single process have the same view of what's in memory
 - Easy to compute offsets into fields, follow pointers, etc.
- In the absence of shared memory, externalizing an object requires us to turn it into a sequential sequence of bytes
 - **Serialization/Marshalling**: Express an object as a sequence of bytes
 - **Deserialization/Unmarshalling**: Reconstructing the original object from its marshalled form at destination

- Neither one is "wrong" but sender and receiver should be consistent!

4/26/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 25 | 13

4/26/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 25 | 14

Machine Representation

- Consider using the machine representation:
 - `fwrite(&x, sizeof(uint32_t), 1, f);`
- How do we know if the recipient represents `x` in the same way?
 - For pipes, is this a problem?
 - What about for sockets?

```
int main(int argc, char **argv[])
{
    int val = 0x12345678;
    int i;
    printf("%x\n", val);
    for (i = 0; i < sizeof(val); i++) {
        printf("%x[%d] = %x\n", i, ((uint8_t *) &val)[i]);
    }
}
```

Processor	Endianness
Motorola 68000	Big Endian
PowerPC (PPC)	Big Endian
Sun Sparc	Big Endian
IBM S/390	Big Endian
Intel x86 (32 bit)	Little Endian
Intel x86_64 (64 bit)	Little Endian
DEC VAX	Little Endian
Alpha	Big Endian
ARM	Big Endian
Intel x86 (64 bit)	Little Endian
MIPS	Big Endian

Simple Data Types

- Suppose I want to write a `x` to a file
 - First, open the file: `FILE* f = fopen("foo.txt", "w");`
 - Then, I have two choices:
 1. `fprintf(f, "%lu", x);`
 2. `fwrite(&x, sizeof(uint32_t), 1, f);`
 - » Or equivalently, `write(fd, &x, sizeof(uint32_t));` (perhaps with a loop to be safe)

4/26/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 25 | 15

4/26/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 25 | 16

What Endian is the Internet?

- Big Endian
 - Network byte order
 - Vs. “host byte order”

	NAME
	applyPatch - definitions for internal operations
	SYNOPSIS
	#include <arpa/inet.h>
	DESCRIPTION

The .NET ADDRESSER [10] and NERF ADDRESSER Gracings shall be defined as described in [Section 6.1](#).
 The following shall either be declared as functions, defined as macros, or both. If functions are declared, function prototypes shall be defined as follows:

```

  uint32_t hex2int(uint32_t);
  uint16_t hex2int(uint16_t);
  uint32_t hex2int(uint32_t, t);
  uint16_t hex2int(uint16_t, t);

  The .NET2 and .NET3 types shall be defined as described in Section 6.2.
```

The following shall be declared as functions and may also be defined as macros. Function prototypes shall be provided in .NET2:

```

  int32_t lntc_hex2int(const char *);
```

char * lntc_hex2str(int32_t);

const char * lntc_hex2str(int32_t);

```

  int32_t lntc_hex2int(const void *restrict, char *restrict,
  const void *restrict, t);
```

```

  lntc_hex2int(int32_t, const char *restrict, void *restrict);
```

Inclusion of the `restrict` keyword may also break symbols from [Section 6.1](#) and [Section 6.2](#).

47622

Joseph & Kubiatowicz CS162 @ UCB Spring 2022

卷之三

What About Richer Objects?

- Consider `word_count_t` of Homework 0 and | ...
 - Each element contains:
 - An `int`
 - A pointer to a string (of some length)
 - A pointer to the next element
 - `fprintf_f_words` writes these as a sequence of lines (character strings with `\n`) to a file like stream
 - What if you wanted to write the whole list as a binary object (and read it back as one)?
 - How do you represent the string?
 - Does it make any sense to write the pointer?

Lec 25.17

4/26/22

Joseph & Kubiatowicz CS162 @ UCB Spring 2022

Data Serialization Formats

- JSON and XML are commonly used in web applications
 - Lots of ad-hoc formats

Consequently, the first step in the process of developing a new product is to identify the needs of the target market. This involves conducting market research to understand the wants and requirements of potential customers. Once these needs are identified, the next step is to develop a product concept that addresses them. This involves creating a detailed product specification, including features, benefits, and price points. The final step is to prototype the product and test it in the market to gather feedback and make any necessary improvements.

- Decide on an “on-wire” endianness
- Convert from native endianness to “on-wire” endianness before sending out data
(serialization/marshalling)
 - `uint32_t htonl(uint32_t)` and `uint16_t htons(uint16_t)` convert from native endianness to network endianness (big endian)
- Convert from “on-wire” endianness to native endianness when receiving data
(deserialization/unmarshalling)
 - `uint32_t ntohl(uint32_t)` and `uint16_t ntohs(uint16_t)` convert from network endianness to native endianness (big endian)

Dealing with Endianness

- ## Data Serialization Formats

 - Decide on an “on-wire” endianness
 - Convert from native endianness to “on-wire” endianness before sending out data (**serialization/marshalling**)
 - `uint32_t htonl(uint32_t)` and `uint16_t htons(uint16_t)` convert from native endianness to network endianness (big endian)
 - Convert from “on-wire” endianness to native endianness when receiving data (**deserialization/unmarshalling**)
 - `uint32_t ntohl(uint32_t)` and `uint16_t ntohs(uint16_t)` convert from network endianness to native endianness (big endian)
 - JSON and XML are commonly used in web applications
 - Lots of ad-hoc formats

Lec 25.18

4/2622

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 25.19

4/2622

Joseph & Kubiatowicz CS162 @ UCB Spring 2022

Lec 25.2C

Data Serialization Formats

4/2622

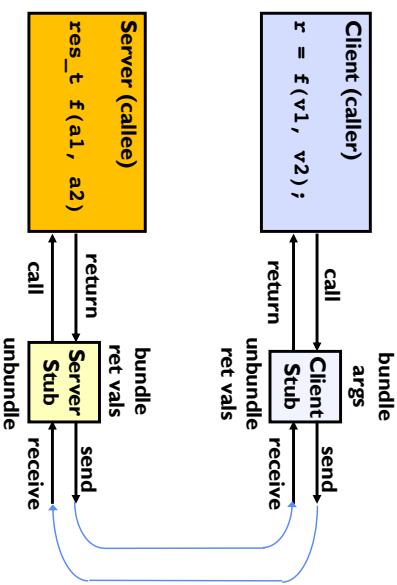
Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 25.2 |

4/2622

- 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
 - **And must deal with machine representation by hand**
 - Another option: Remote Procedure Call (RPC)
 - Calls a procedure on a remote machine
 - Idea: Make communication look like an ordinary function call
 - Automate all of the complexity of translating between representations
 - Client calls:
`remoteFileSystem->Read("rutabaga");`
 - Translated automatically into call on server:
`fileSys->Read("rutabaga");`

RPC Concept



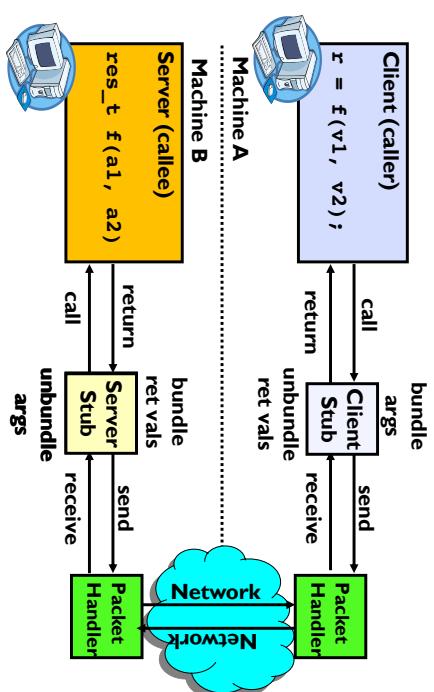
Lec 25.23

4/2622

& Kubiatowicz CS162 © UCB Spring 2022

Lec 25.24

RPC Information Flow



Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 25.22

Remote Procedure Call (RPC)

- 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
 - **And must deal with machine representation by hand**

4/2622

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 25.2 |

4/2622

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 25.22

RPC 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.

4/26/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

RPC Details (2/3)

- Cross-Platform issues:
 - What if client/server machines are different architectures/ languages?
 - » Convert everything to/from some canonical form
 - » Tag every item with an indication of how it is encoded (avoids unnecessary conversions)
 - How does client know which mbox (destination queue) 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
 - 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

4/26/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

RPC Details (1/3)

- 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

Lec 25.25

4/26/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

RPC Details (3/3)

- 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

Lec 25.27

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 25.28

4/26/22

Problems with RPC: Non-Atomic Failures

- Different failure modes in dist. 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

4/26/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 25.29

Problems with RPC: Performance

- RPC is not performance transparent:
 - Cost of Procedure call « same-machine RPC « network RPC
 - Overheads: Marshalling, Stubs, Kernel-Crossing, Communication
- Programmers must be aware that RPC is not free
 - Caching can help, but may make failure handling complex

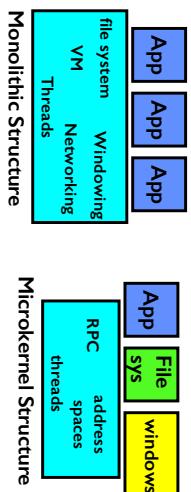
4/26/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 25.30

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 or the same machine
 - Services can be run wherever it’s most appropriate
 - Access to local and remote services looks the same
- Examples of RPC systems:
 - CORBA (Common Object Request Broker Architecture)
 - DCOM (Distributed COM)
 - RMI (Java Remote Method Invocation)



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.

Lec 25.31

Joseph & Kubiatowicz CS162 © UCB Spring 2022

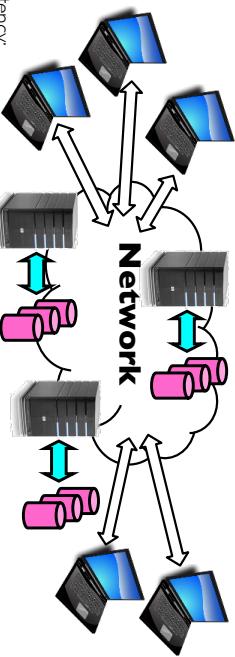
4/26/22

Lec 25.32

4/26/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Network-Attached Storage and the CAP Theorem



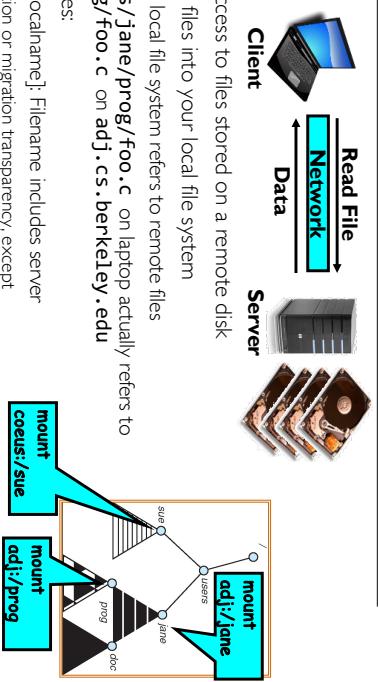
- Consistency:
 - Changes appear to everyone in the same serial order
- Availability:
 - Can get a result at any time
- Partition-Tolerance
 - System continues to work even when network becomes partitioned
 - Consistency, Availability, Partition-Tolerance (CAP) Theorem **Cannot have all three at same time**
 - Otherwise known as "Brewer's Theorem"

4/26/22

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

Lec 25.33

Distributed File Systems



- Transparent access to files stored on a remote disk
- Mount remote files into your local file system
 - Directory in local file system refers to remote files
 - e.g., /users/Jane/prog/foo.c on laptop actually refers to /prog/foo.c on adj.cs.berkeley.edu
- Naming Choices:
 - [Hostname.localname]: Filename includes server
 - » No location or migration transparency, except through DNS remapping
 - A global name space: Filename unique in "world"
 - » Can be served by any server

4/26/22

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

Lec 25.35

Administrivia

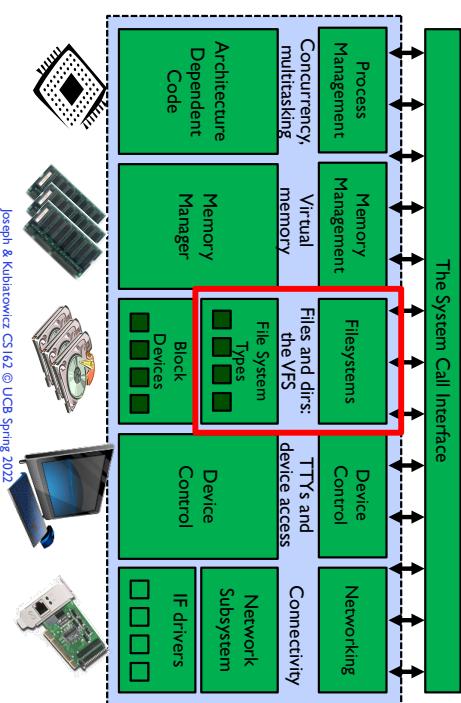
- Midterm 3: Thursday 4/28: 7-9PM
 - All course material
 - Review session Monday 4/25 1-3PM

4/26/22

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

Lec 25.34

Enabling Design: VFS



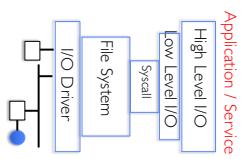
4/26/22

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

Lec 25.36

Recall: Layers of I/O...

User App:
User library:
Application / Service



```

length = read(input_fd, buffer, BUFFER_SIZE);

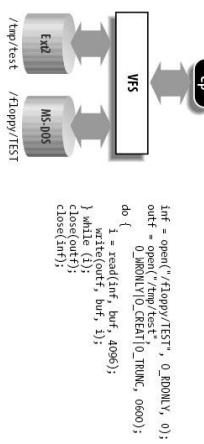
User App:
ssize_t read(int, void *, size_t) {
    marshal args into registers
    issue syscall
    register result of syscall to rtn value
}

Exception U>K, interrupt processing
void syscall_handler (struct intr_frame *f) {
    size_t count, loff_t *pos) {
        unmarshal args from regs
        dispatch : handlers[call#](args)
        marshal results to syscall ret
    }
}
  
```

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 25.37

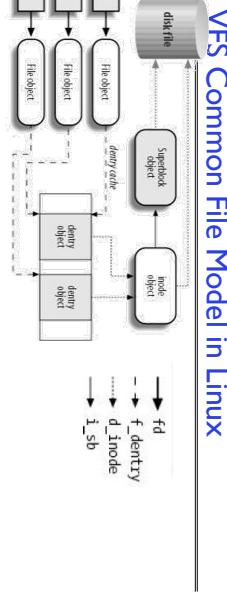
- VFS: Virtual abstraction similar to local file system
 - Provides virtual superblocks, inodes, files, etc
 - Compatible with a variety of local and remote file systems
 - » provides object-oriented way of implementing file systems
 - VFS allows the same system call interface (the API) to be used for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



Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 25.39

Virtual Filesystem Switch



- Four primary object types for VFS:

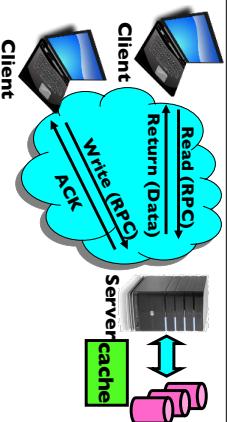
- superblock object: represents a specific mounted filesystem
- inode object: represents a specific file
- dentry object: represents a directory entry
- file object: represents open file associated with process
- There is no specific directory object (VFS treats directories as files)
- May need to fit the model by faking it
 - Example: make it look like directories are files
 - Example: make it look like have inodes, superblocks, etc.

4/26/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 25.39

Simple Distributed File System



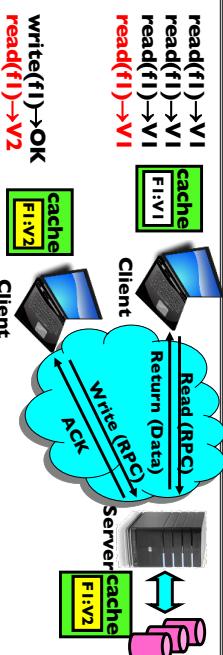
- Remote Disk: Reads and writes forwarded to server
 - Use **Remote Procedure Calls (RPC)** to translate file system calls into remote requests
- No local caching, but can be cache at server-side
 - Advantage: Server provides consistent view of file system to multiple clients
 - Problems? Performance!
 - Going over network is slower than going to local memory
 - Lots of network traffic not well pipelined
 - Server can be a bottleneck

4/26/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 25.40

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
 - Client caches not consistent with server/each other

4/26/22

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

Lec 25.41

Dealing with Failures

- What if server crashes? Can client wait until it comes back and just continue making requests?
 - Changes in server's cache but not in disk are lost
- What if there is shared state across RPC's?
 - Client opens file, then does a seek
 - Server crashes
 - Similar problem: What if client removes a file but server crashes before acknowledgement?

4/26/22

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

Lec 25.42

Stateless Protocol

- Stateless Protocol:** A protocol in which all information required to service a request is included with the request

- Even better: Idempotent Operations – repeating an operation multiple times is same as executing it just once (e.g. storing to a mem addr.)
- Client: timeout expires without reply, just run the operation again (safe regardless of first attempt)
- Recall HTTP: Also a stateless protocol
 - Include cookies with request to simulate a session

Case Study: 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
 - NFS service layer: bottom layer of the architecture
 - Implements the NFS protocol
- NFS Protocol: RPC for file operations on server
 - XDR: Serialization standard for data format independence
 - 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)

4/26/22

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

Lec 25.43

4/26/22

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

Lec 25.44

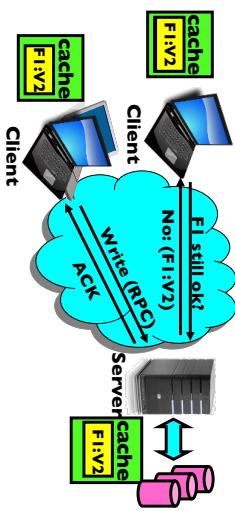
NFS Continued

- NFS servers are **stateless**; each request provides all arguments required for execution
 - E.g., reads include information for entire operation, such as **Readt (inumber, position)**, not **Read (open file)**
 - 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 them 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 other 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)

4/26/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 25.45



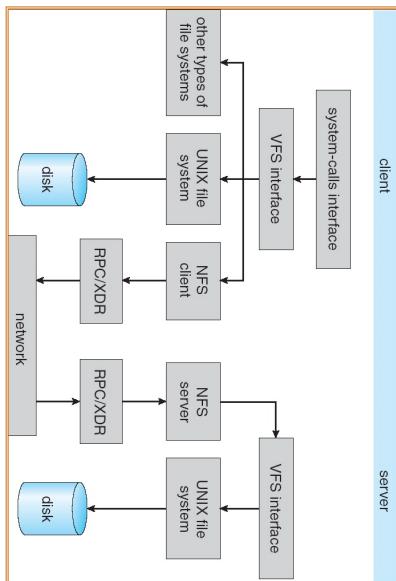
- What if multiple clients write to same file?
- » In NFS, can get either version (or parts of both)
- » Completely arbitrary!

4/26/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 25.47

NFS Architecture



4/26/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

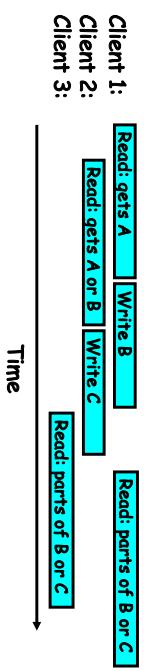
Lec 25.46

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 is tunable parameter).
 - » Thus, when file is changed on one client, server is notified, but other clients use old version of file until timeout.

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

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 25.48

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
 - » Server becomes bottleneck due to polling traffic

4/26/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

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 \Rightarrow more files can be cached locally
 - Callbacks \Rightarrow server not involved if file is read-only
- For both AFS and NFS, central server is bottleneck!
 - Performance: all writes \rightarrow server, cache misses \rightarrow server
 - Availability: Server is single point of failure
 - Cost: server machine's high cost relative to workstation

Lec 25.49

4/26/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Andrew File System

- Andrew File System (AFS, late 80's) \rightarrow 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

4/26/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Summary (1/2)

- **TCP:** Reliable byte stream between two processes on different machines over Internet (read, write, flush)
 - Uses window-based acknowledgement protocol
 - Congestion-avoidance dynamically adapts sender window to account for congestion in network
- **Remote Procedure Call (RPC):** Call procedure on remote machine or in remote domain
 - Provides same interface as procedure
 - Automatic packing and unpacking of arguments without user programming (in stub)
 - Adapts automatically to different hardware and software architectures at remote end

Lec 25.50

4/26/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

4/26/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

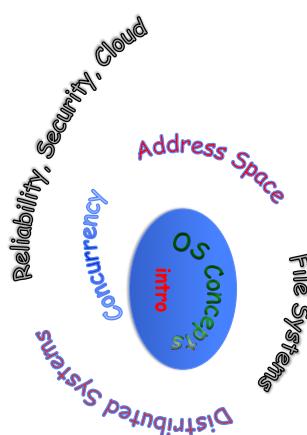
Lec 25.51

Lec 25.52

Summary (2/2)

- **Distributed File System:**
 - Transparent access to files stored on a remote disk
 - Caching for performance
- **VFS:** Virtual File System layer (Or Virtual Filesystem Switch)
 - Provides mechanism which gives same system call interface for different types of file systems
- **Cache Consistency:** Keeping 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 to be notified by server of changes

Thank you!



- Thanks for all your great questions!
- Good Bye! You have all been great!