

# CS162 Operating Systems and Systems Programming Lecture 22

## TCP/IP (Continued)

November 18<sup>th</sup>, 2015  
Prof. John Kubiatowicz  
<http://cs162.eecs.Berkeley.edu>

## Recall: Two-Phase Commit

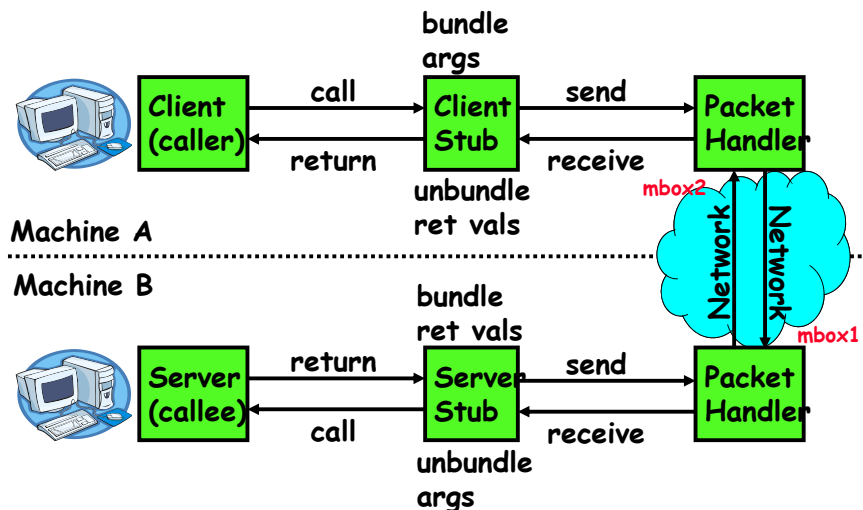
- Since we can't solve the General's Paradox (i.e. simultaneous action), let's solve a related problem
  - Distributed transaction: Two machines agree to do something, or not do it, **atomically**
- Two-Phase Commit protocol:
  - **Persistent stable log on each machine**: keep track of whether commit has happened
    - » If a machine crashes, when it wakes up it first checks its log to recover state of world at time of crash
  - **Prepare Phase**:
    - » The global coordinator requests that all participants will promise to commit or rollback the transaction
    - » Participants record promise in log, then acknowledge
    - » If anyone votes to abort, coordinator writes "Abort" in its log and tells everyone to abort; each records "Abort" in log
  - **Commit Phase**:
    - » After all participants respond that they are prepared, then the coordinator writes "Commit" to its log
    - » Then asks all nodes to commit; they respond with ack
    - » After receive acks, coordinator writes "Got Commit" to log
  - Log can be used to complete this process such that all machines either commit or don't commit

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.2

## Recall: RPC Information Flow



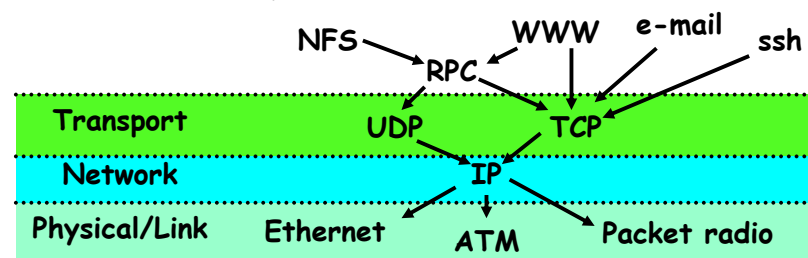
11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.3

## Network Protocols

- Networking protocols: many levels
  - Physical level: mechanical and electrical network (e.g. how are 0 and 1 represented)
  - Link level: packet formats/error control (for instance, the CSMA/CD protocol)
  - Network level: network routing, addressing
  - Transport Level: reliable message delivery
- Protocols on today's Internet:



11/18/15

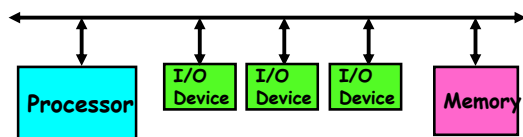
Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.4

## Broadcast Networks



### • Broadcast Network: Shared Communication Medium



- Shared Medium can be a set of wires
  - » Inside a computer, this is called a bus
  - » All devices simultaneously connected to devices



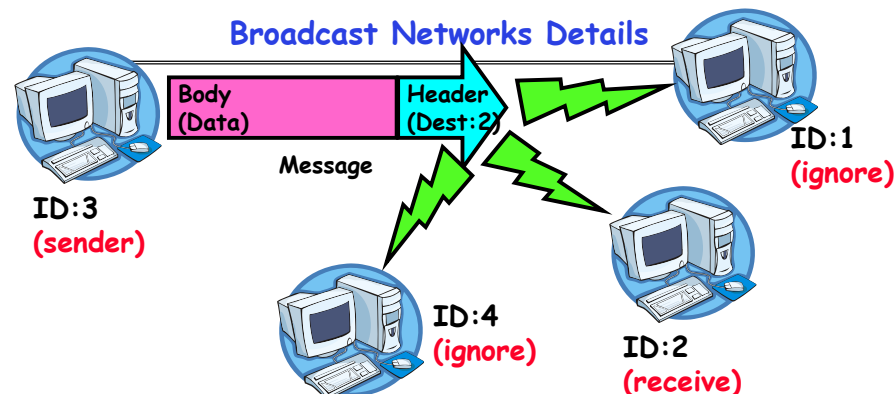
- Originally, Ethernet was a broadcast network
  - » All computers on local subnet connected to one another
- More examples (wireless: medium is air): cellular phones, GSM GPRS, EDGE, CDMA 1xRTT, and 1EvDO

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.5

## Broadcast Networks Details



- **Delivery:** When you broadcast a packet, how does a receiver know who it is for? (packet goes to everyone!)
  - Put header on front of packet: [ Destination | Packet ]
  - Everyone gets packet, discards if not the target
  - In Ethernet, this check is done in hardware
    - » No OS interrupt if not for particular destination
  - This is layering: we're going to build complex network protocols by layering on top of the packet

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.6

## Carrier Sense, Multiple Access/Collision Detection

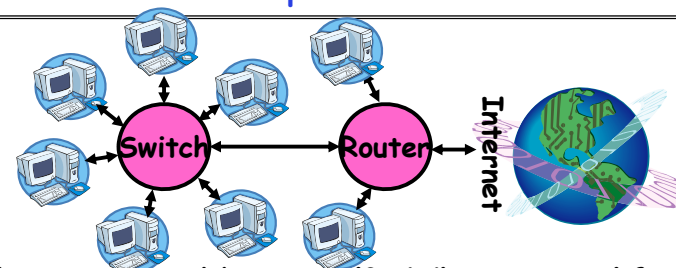
- Ethernet (early 80's): first practical local area network
  - It is the most common LAN for UNIX, PC, and Mac
  - Use wire instead of radio, but still broadcast medium
- Key advance was in arbitration called CSMA/CD: Carrier sense, multiple access/collision detection
  - **Carrier Sense:** don't send unless idle
    - » Don't mess up communications already in process
  - **Collision Detect:** sender checks if packet trampled.
    - » If so, abort, wait, and retry.
  - **Backoff Scheme:** Choose wait time before trying again
- How long to wait after trying to send and failing?
  - What if everyone waits the same length of time? Then, they all collide again at some time!
  - Must find way to break up shared behavior with nothing more than shared communication channel
- Adaptive randomized waiting strategy:
  - **Adaptive and Random:** First time, pick random wait time with some initial mean. If collide again, pick random value from bigger mean wait time. Etc.
  - Randomness is important to decouple colliding senders
  - Scheme figures out how many people are trying to send!

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.7

## Point-to-point networks



- Why have a shared bus at all? Why not simplify and only have point-to-point links + routers/switches?
  - Originally wasn't cost-effective
  - Now, easy to make high-speed switches and routers that can forward packets from a sender to a receiver.
- **Point-to-point network:** a network in which every physical wire is connected to only two computers
- **Switch:** a bridge that transforms a shared-bus (broadcast) configuration into a point-to-point network.
- **Router:** a device that acts as a junction between two networks to transfer data packets among them.

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.8

## The Internet Protocol: "IP"

- The Internet is a large network of computers spread across the globe
  - According to the Internet Systems Consortium, there were over 681 million computers as of July 2009
  - In principle, every host can speak with every other one under the right circumstances
- **IP Packet:** a network packet on the internet
- **IP Address:** a 32-bit integer used as the destination of an IP packet
  - Often written as four dot-separated integers, with each integer from 0–255 (thus representing  $8 \times 4 = 32$  bits)
  - Example CS file server is: 169.229.60.83  $\equiv 0xA9E53C53$
- **Internet Host:** a computer connected to the Internet
  - Host has one or more IP addresses used for routing
    - » Some of these may be private and unavailable for routing
  - Not every computer has a unique IP address
    - » Groups of machines may share a single IP address
    - » In this case, machines have private addresses behind a "Network Address Translation" (NAT) gateway

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.9

## Address Subnets

- **Subnet:** A network connecting a set of hosts with related destination addresses
- With IP, all the addresses in subnet are related by a prefix of bits
  - **Mask:** The number of matching prefix bits
    - » Expressed as a single value (e.g., 24) or a set of ones in a 32-bit value (e.g., 255.255.255.0)
- A subnet is identified by 32-bit value, with the bits which differ set to zero, followed by a slash and a mask
  - Example: 128.32.131.0/24 designates a subnet in which all the addresses look like 128.32.131.XX
  - Same subnet: 128.32.131.0/255.255.255.0
- Difference between subnet and complete network range
  - Subnet is always a subset of address range
  - Once, subnet meant single physical broadcast wire; now, less clear exactly what it means (virtualized by switches)

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.10

## Address Ranges in IP

- IP address space divided into prefix-delimited ranges:
  - Class A: NN.0.0.0/8
    - » NN is 1-126 (126 of these networks)
    - » 16,777,214 IP addresses per network
    - » 10.xx.yy.zz is private
    - » 127.xx.yy.zz is loopback
  - Class B: NN.MM.0.0/16
    - » NN is 128-191, MM is 0-255 (16,384 of these networks)
    - » 65,534 IP addresses per network
    - » 172.[16-31].xx.yy are private
  - Class C: NN.MM.LL.0/24
    - » NN is 192-223, MM and LL 0-255 (2,097,151 of these networks)
    - » 254 IP addresses per networks
    - » 192.168.xx.yy are private
- Address ranges are often owned by organizations
  - Can be further divided into subnets

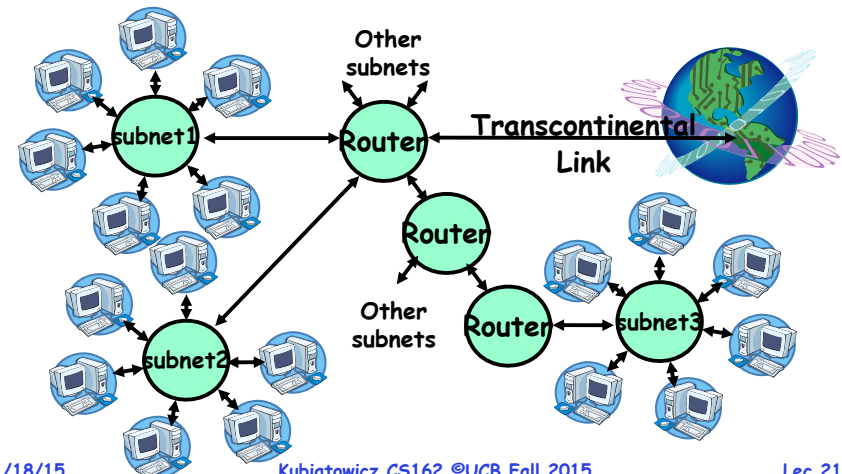
11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.11

## Hierarchical Networking: The Internet

- How can we build a network with millions of hosts?
  - Hierarchy! Not every host connected to every other one
  - Use a network of Routers to connect subnets together
    - » Routing is often by prefix: e.g. first router matches first 8 bits of address, next router matches more, etc.



11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.12

## Simple Network Terminology

- **Local-Area Network (LAN)** – designed to cover small geographical area
  - Multi-access bus, ring, or star network
  - Speed  $\approx 10 - 1000$  Megabits/second (even 40-100GB/s)
  - Broadcast is fast and cheap
  - In small organization, a LAN could consist of a single subnet. In large organizations (like UC Berkeley), a LAN contains many subnets
- **Wide-Area Network (WAN)** – links geographically separated sites
  - Point-to-point connections over long-haul lines (often leased from a phone company)
  - Speed  $\approx 1.544 - 155$  Megabits/second
  - Broadcast usually requires multiple messages

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.13

## Administrivia

- **Midterm II: Next Monday (11/23)**
  - Time: 7:00PM-10:00PM
  - Location: VLSB: 2040/2050/2060
  - Division of people between rooms will be posted
  - All topics from Midterm I, up to this Wednesday, including:
    - » Address Translation/TLBs/Paging
    - » I/O subsystems, Storage Layers, Disks/SSD
    - » Performance and Queueing Theory
    - » File systems
    - » Distributed systems, TCP/IP, RPC
- **Closed book, one page of notes** – both sides
- **Makeup: Make sure to contact William**
  - We need to know everyone who needs a makeup exam
  - Also – we need want to know why: Data Sciences exam not valid reason
- **Review session:**
  - 306 Soda Hall
  - Sunday TBA

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.14

## Routing

- **Routing: the process of forwarding packets hop-by-hop through routers to reach their destination**
  - Need more than just a destination address!
    - » Need a path
  - **Post Office Analogy:**
    - » Destination address on each letter is not sufficient to get it to the destination
    - » To get a letter from here to Florida, must route to local post office, sorted and sent on plane to somewhere in Florida, be routed to post office, sorted and sent with carrier who knows where street and house is...
- **Internet routing mechanism: routing tables**
  - Each router does table lookup to decide which link to use to get packet closer to destination
  - Don't need 4 billion entries in table: routing is by subnet
  - Could packets be sent in a loop? Yes, if tables incorrect
- **Routing table contains:**
  - Destination address range  $\rightarrow$  output link closer to destination
  - Default entry (for subnets without explicit entries)



11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.15

## Setting up Routing Tables

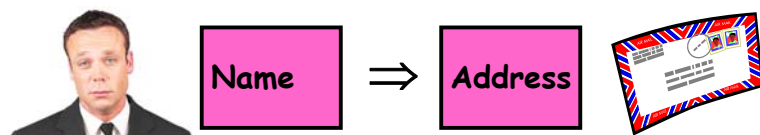
- **How do you set up routing tables?**
  - Internet has no centralized state!
    - » No single machine knows entire topology
    - » Topology constantly changing (faults, reconfiguration, etc)
  - Need dynamic algorithm that acquires routing tables
    - » Ideally, have one entry per subnet or portion of address
    - » Could have "default" routes that send packets for unknown subnets to a different router that has more information
- **Possible algorithm for acquiring routing table**
  - Routing table has "cost" for each entry
    - » Includes number of hops to destination, congestion, etc.
    - » Entries for unknown subnets have infinite cost
  - Neighbors periodically exchange routing tables
    - » If neighbor knows cheaper route to a subnet, replace your entry with neighbors entry (+1 for hop to neighbor)
- **In reality:**
  - Internet has networks of many different scales
  - Different algorithms run at different scales

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.16

## Naming in the Internet



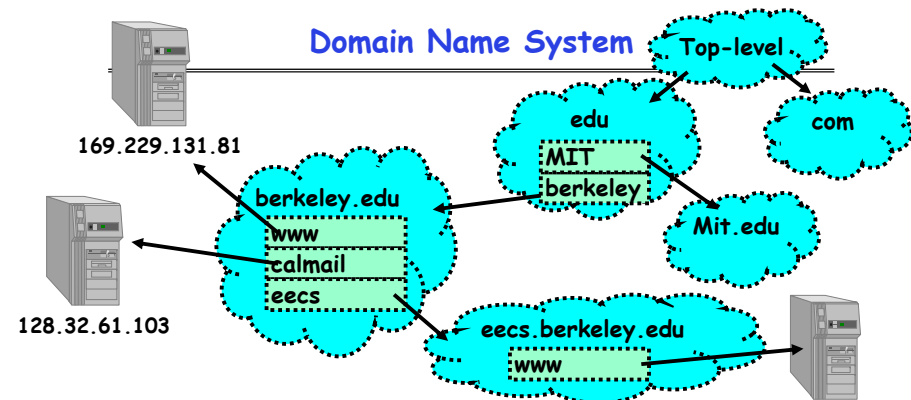
- How to map human-readable names to IP addresses?
  - E.g. `www.berkeley.edu`  $\Rightarrow$  `128.32.139.48`
  - E.g. `www.google.com`  $\Rightarrow$  different addresses depending on location, and load
- Why is this necessary?
  - IP addresses are hard to remember
  - IP addresses change:
    - » Say, Server 1 crashes gets replaced by Server 2
    - » Or - `google.com` handled by different servers
- Mechanism: Domain Naming System (DNS)

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.17

## Domain Name System



- DNS is a hierarchical mechanism for naming
  - Name divided in domains, right to left: `www.eecs.berkeley.edu`
- Each domain owned by a particular organization
  - Top level handled by ICANN (Internet Corporation for Assigned Numbers and Names)
  - Subsequent levels owned by organizations
- Resolution: series of queries to successive servers
- Caching: queries take time, so results cached for period of time

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.18

## How Important is Correct Resolution?

- If attacker manages to give incorrect mapping:
  - Can get someone to route to server, thinking that they are routing to a different server
    - » Get them to log into "bank" - give up username and password
- Is DNS Secure?
  - Definitely a weak link
    - » What if "response" returned from different server than original query?
    - » Get person to use incorrect IP address!
  - Attempt to avoid substitution attacks:
    - » Query includes random number which must be returned
- In July 2008, hole in DNS security located!
  - Dan Kaminsky (security researcher) discovered an attack that broke DNS globally
    - » One person in an ISP convinced to load particular web page, then *all* users of that ISP end up pointing at wrong address
  - High profile, highly advertised need for patching DNS
    - » Big press release, lots of mystery
    - » Security researchers told no speculation until patches applied

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.19

## Network Layering

- **Layering:** building complex services from simpler ones
  - Each layer provides services needed by higher layers by utilizing services provided by lower layers
- The physical/link layer is pretty limited
  - Packets are of limited size (called the "Maximum Transfer Unit or MTU: often 200-1500 bytes in size)
  - Routing is limited to within a physical link (wire) or perhaps through a switch
- Our goal in the following is to show how to construct a secure, ordered, message service routed to anywhere:

Physical Reality: Packets	Abstraction: Messages
Limited Size	Arbitrary Size
Unordered (sometimes)	Ordered
Unreliable	Reliable
Machine-to-machine	Process-to-process
Only on local area net	Routed anywhere
Asynchronous	Synchronous
Insecure	Secure

11/18/15

Lec 21.20

## Building a messaging service

- Handling Arbitrary Sized Messages:
  - Must deal with limited physical packet size
  - Split big message into smaller ones (called fragments)
    - » Must be reassembled at destination
  - Checksum computed on each fragment or whole message
- Internet Protocol (IP): Must find way to send packets to arbitrary destination in network
  - Deliver messages unreliably ("best effort") from one machine in Internet to another
  - Since intermediate links may have limited size, must be able to fragment/reassemble packets on demand
  - Includes 256 different "sub-protocols" build on top of IP
    - » Examples: ICMP(1), TCP(6), UDP (17), IPSEC(50,51)

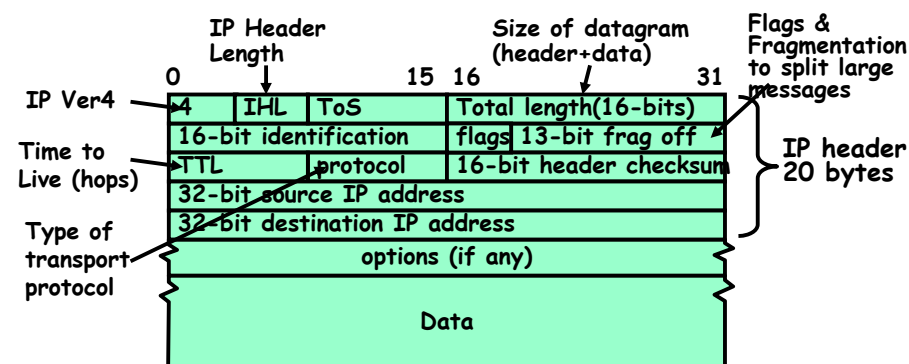
11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.21

## IP Packet Format

- IP Packet Format:



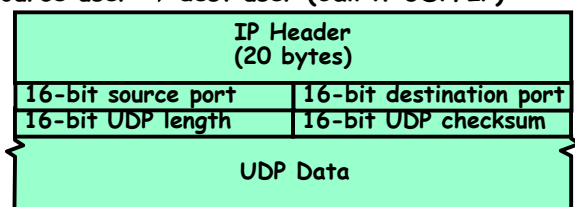
11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.22

## Building a messaging service

- Process to process communication
  - Basic routing gets packets from machine→machine
  - What we really want is routing from process→process
    - » Add "ports", which are 16-bit identifiers
    - » A communication channel (connection) defined by 5 items: [source addr, source port, dest addr, dest port, protocol]
- UDP: The Unreliable Datagram Protocol
  - Layered on top of basic IP (IP Protocol 17)
    - » Datagram: an unreliable, unordered, packet sent from source user → dest user (Call it UDP/IP)



- Important aspect: low overhead!
  - » Often used for high-bandwidth video streams
  - » Many uses of UDP considered "anti-social" - none of the "well-behaved" aspects of (say) TCP/IP

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.23

## Ordered Messages

- Ordered Messages
  - Several network services are best constructed by ordered messaging
    - » Ask remote machine to first do x, then do y, etc.
  - Unfortunately, underlying network is packet based:
    - » Packets are routed one at a time through the network
    - » Can take different paths or be delayed individually
  - IP can reorder packets!  $P_0, P_1$  might arrive as  $P_1, P_0$
- Solution requires queuing at destination
  - Need to hold onto packets to undo misordering
  - Total degree of reordering impacts queue size
- Ordered messages on top of unordered ones:
  - Assign sequence numbers to packets
    - » 0, 1, 2, 3, 4, ...
    - » If packets arrive out of order, reorder before delivering to user application
    - » For instance, hold onto #3 until #2 arrives, etc.
  - Sequence numbers are specific to particular connection
    - » Reordering among connections normally doesn't matter
  - If restart connection, need to make sure use different range of sequence numbers than previously...

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.24

## Reliable Message Delivery: the Problem

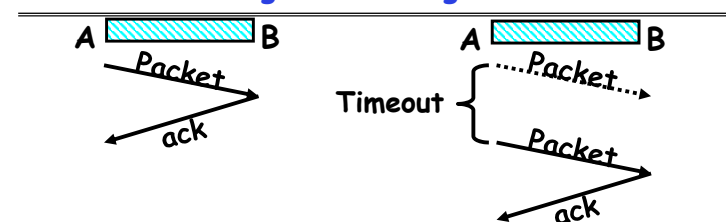
- All physical networks can garble and/or drop packets
  - Physical media: packet not transmitted/received
    - » If transmit close to maximum rate, get more throughput - even if some packets get lost
    - » If transmit at lowest voltage such that error correction just starts correcting errors, get best power/bit
  - Congestion: no place to put incoming packet
    - » Point-to-point network: insufficient queue at switch/router
    - » Broadcast link: two host try to use same link
    - » In any network: insufficient buffer space at destination
    - » Rate mismatch: what if sender send faster than receiver can process?
- Reliable Message Delivery on top of Unreliable Packets
  - Need some way to make sure that packets actually make it to receiver
    - » Every packet received at least once
    - » Every packet received at most once
  - Can combine with ordering: every packet received by process at destination exactly once and in order

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.25

## Using Acknowledgements



- How to ensure transmission of packets?
  - Detect garbling at receiver via checksum, discard if bad
  - Receiver acknowledges (by sending "ack") when packet received properly at destination
  - Timeout at sender: if no ack, retransmit
- Some questions:
  - If the sender doesn't get an ack, does that mean the receiver didn't get the original message?
    - » No
  - What if ack gets dropped? Or if message gets delayed?
    - » Sender doesn't get ack, retransmits. Receiver gets message twice, acks each.

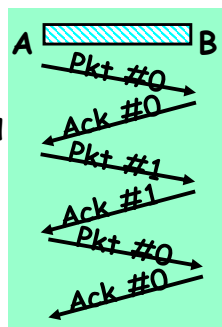
11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.26

## How to deal with message duplication?

- Solution: put sequence number in message to identify re-transmitted packets
  - Receiver checks for duplicate #'s; Discard if detected
- Requirements:
  - Sender keeps copy of unack'ed messages
    - » Easy: only need to buffer messages
  - Receiver tracks possible duplicate messages
    - » Hard: when ok to forget about received message?
- Alternating-bit protocol:
  - Send one message at a time; don't send next message until ack received
  - Sender keeps last message; receiver tracks sequence # of last message received
- Pros: simple, small overhead
- Con: Poor performance
  - Wire can hold multiple messages; want to fill up at (wire latency  $\times$  throughput)
- Con: doesn't work if network can delay or duplicate messages arbitrarily



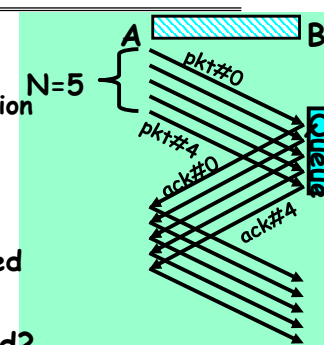
11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.27

## Better messaging: Window-based acknowledgements

- Windowing protocol (not quite TCP):
  - Send up to N packets without ack
    - » Allows pipelining of packets
    - » Window size (N) < queue at destination
  - Each packet has sequence number
    - » Receiver acknowledges each packet
    - » Ack says "received all packets up to sequence number X"/send more
- Acks serve dual purpose:
  - Reliability: Confirming packet received
  - Ordering: Packets can be reordered at destination
- What if packet gets garbled/dropped?
  - Sender will timeout waiting for ack packet
    - » Resend missing packets  $\Rightarrow$  Receiver gets packets out of order!
  - Should receiver discard packets that arrive out of order?
    - » Simple, but poor performance
  - Alternative: Keep copy until sender fills in missing pieces?
    - » Reduces # of retransmits, but more complex
- What if ack gets garbled/dropped?
  - Timeout and resend just the un-acknowledged packets

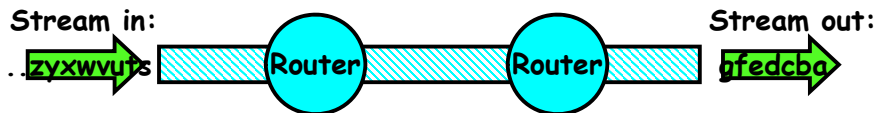


11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.28

## Transmission Control Protocol (TCP)



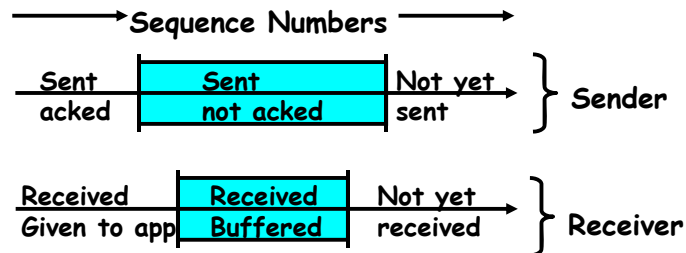
- Transmission Control Protocol (TCP)
  - TCP (**IP Protocol 6**) layered on top of IP
  - Reliable byte stream between two processes on different machines over Internet (read, write, flush)
- TCP Details
  - Fragments byte stream into packets, hands packets to IP
    - » IP may also fragment by itself
  - Uses window-based acknowledgement protocol (to minimize state at sender and receiver)
    - » "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"

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.29

## TCP Windows and Sequence Numbers



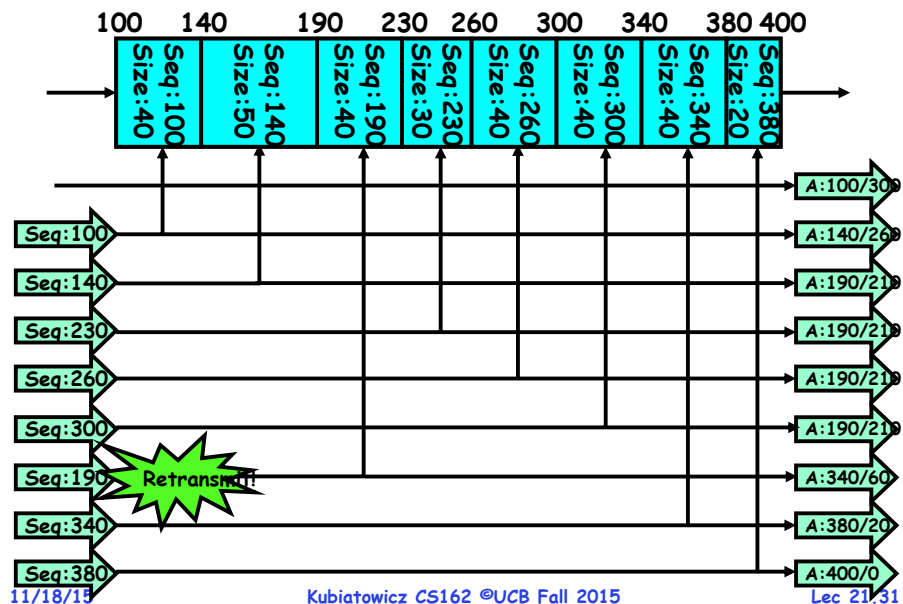
- Sender has three regions:
  - Sequence regions
    - » sent and ack'd
    - » Sent and not ack'd
    - » not yet sent
  - Window (colored region) adjusted by sender
- Receiver has three regions:
  - Sequence regions
    - » received and ack'd (given to application)
    - » received and buffered
    - » not yet received (or discarded because out of order)

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.30

## Window-Based Acknowledgements (TCP)

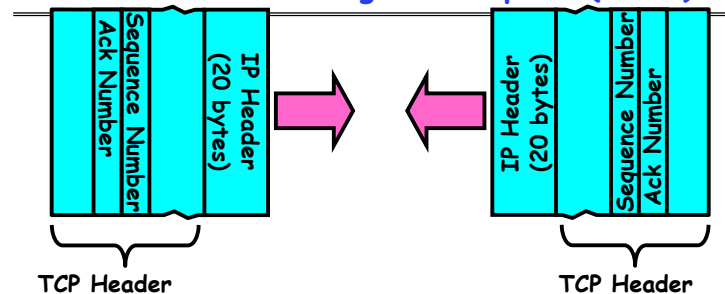


11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.31

## Selective Acknowledgement Option (SACK)



- Vanilla TCP Acknowledgement
  - Every message encodes Sequence number and Ack
  - Can include data for forward stream and/or ack for reverse stream
- Selective Acknowledgement
  - Acknowledgement information includes not just one number, but rather ranges of received packets
  - Must be specially negotiated at beginning of TCP setup
    - » Not widely in use (although in Windows since Windows 98)

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.32



## 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 ⇒ ack is delayed ⇒ unnecessary timeout ⇒ more traffic ⇒ more congestion
    - » 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 ⇒ congestion, so cut window size in half
  - "Additive Increase, Multiplicative Decrease"

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.33

## Sequence-Number Initialization

- How do you choose an initial sequence number?
  - When machine boots, ok to start with sequence #0?
    - » No: could send two messages with same sequence #!
    - » Receiver might end up discarding valid packets, or duplicate ack from original transmission might hide lost packet
  - Also, if it is possible to predict sequence numbers, might be possible for attacker to hijack TCP connection
- Some ways of choosing an initial sequence number:
  - Time to live: each packet has a deadline.
    - » If not delivered in X seconds, then is dropped
    - » Thus, can re-use sequence numbers if wait for all packets in flight to be delivered or to expire
  - Epoch #: uniquely identifies *which* set of sequence numbers are currently being used
    - » Epoch # stored on disk, Put in every message
    - » Epoch # incremented on crash and/or when run out of sequence #
  - Pseudo-random increment to previous sequence number
    - » Used by several protocol implementations

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.34

## Open Connection: 3-Way Handshaking

- Goal: agree on a set of parameters, i.e., the start sequence number for each side
  - Starting sequence number: sequence of first byte in stream
  - Starting sequence numbers are random

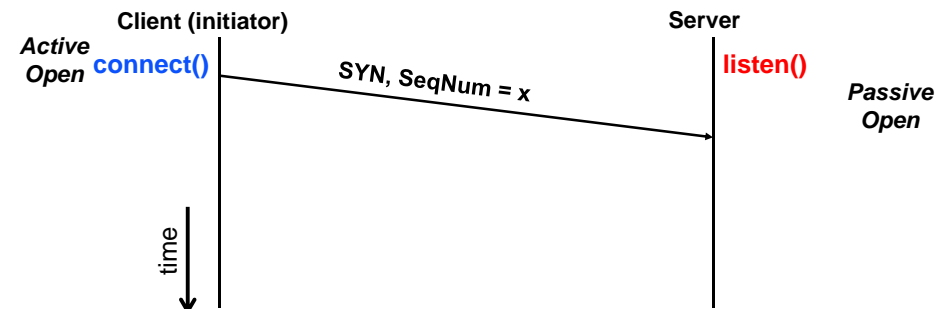
11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.35

## Open Connection: 3-Way Handshaking

- Server waits for new connection calling **listen()**
- Sender call **connect()** passing socket which contains server's IP address and port number
  - OS sends a special packet (SYN) containing a proposal for first sequence number,  $x$



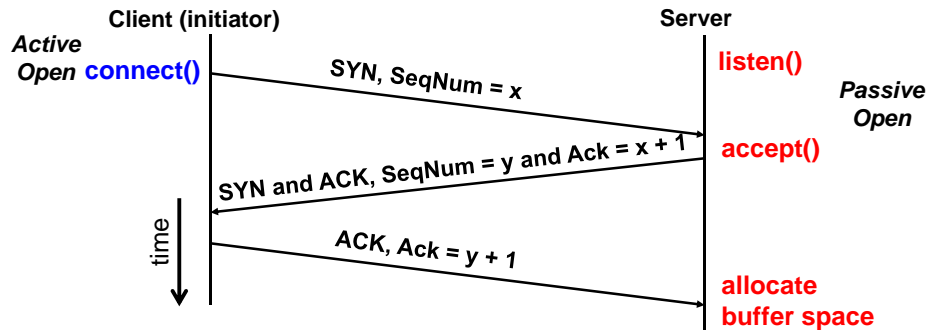
11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.36

## Open Connection: 3-Way Handshaking

- If it has enough resources, server calls `accept()` to accept connection, and sends back a SYN ACK packet containing
  - Client's sequence number incremented by one,  $(x + 1)$ 
    - » Why is this needed?
  - A sequence number proposal,  $y$ , for first byte server will send



11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.37

## 3-Way Handshaking (cont'd)

- Three-way handshake adds 1 RTT delay
- Why do it this way?
  - Congestion control: SYN (40 byte) acts as cheap probe
  - Protects against delayed packets from other connection (would confuse receiver)

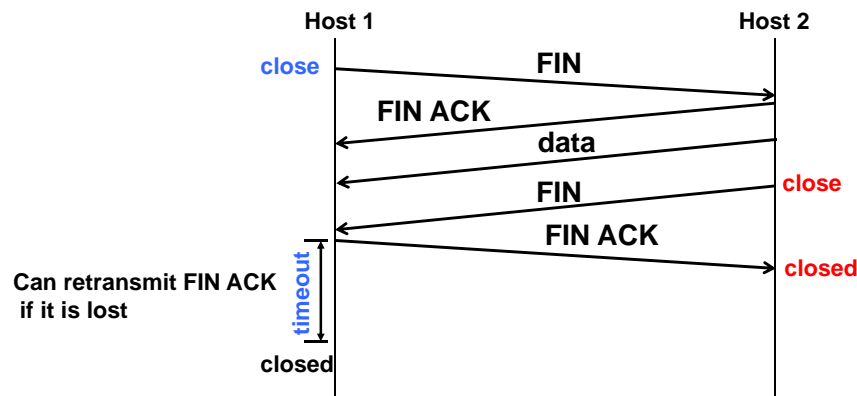
11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.38

## Close Connection

- Goal: both sides agree to close the connection
- 4-way connection tear down



11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.39

## Use of TCP: Sockets

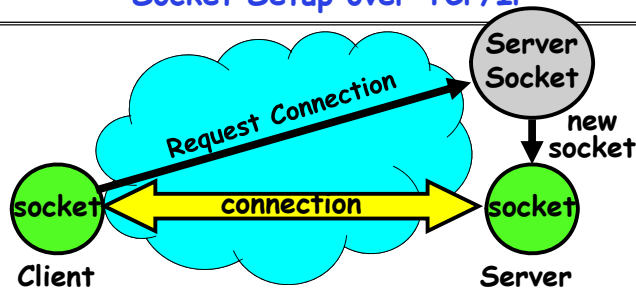
- **Socket:** an abstraction of a network I/O queue
  - Embodies one side of a communication channel
    - » Same interface regardless of location of other end
    - » Could be local machine (called "UNIX socket") or remote machine (called "network socket")
  - First introduced in 4.2 BSD UNIX: big innovation at time
    - » Now most operating systems provide some notion of socket
- Using Sockets for Client-Server (C/C++ interface):
  - On server: set up "server-socket"
    - » Create socket, Bind to protocol (TCP), local address, port
    - » Call `listen()`: tells server socket to accept incoming requests
    - » Perform multiple `accept()` calls on socket to accept incoming connection request
    - » Each successful `accept()` returns a new socket for a new connection; can pass this off to handler thread
  - On client:
    - » Create socket, Bind to protocol (TCP), remote address, port
    - » Perform `connect()` on socket to make connection
    - » If `connect()` successful, have socket connected to server

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.40

## Socket Setup over TCP/IP



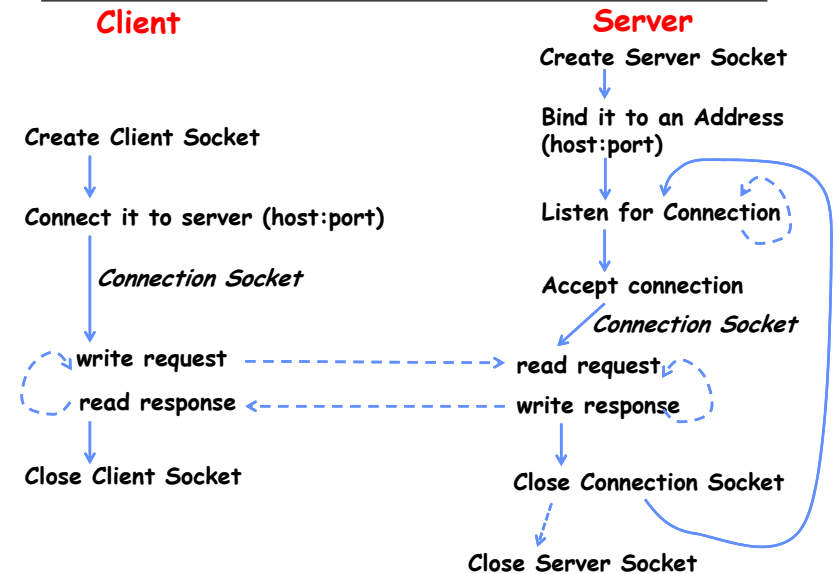
- **Server Socket:** Listens for new connections
  - Produces new sockets for each unique connection
- **Things to remember:**
  - Connection involves 5 values: [ Client Addr, Client Port, Server Addr, Server Port, Protocol ]
  - Often, Client Port "randomly" assigned
    - » Done by OS during client socket setup
  - Server Port often "well known"
    - » 80 (web), 443 (secure web), 25 (sendmail), etc
    - » Well-known ports from 0–1023

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.41

## Recall: Sockets in concept



11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.42

## Recall: Client Protocol

```
char *hostname;
int sockfd, portno;
struct sockaddr_in serv_addr;
struct hostent *server;

server = buildServerAddr(&serv_addr, hostname, portno);

/* Create a TCP socket */
sockfd = socket(AF_INET, SOCK_STREAM, 0)

/* Connect to server on port */
connect(sockfd, (struct sockaddr *) &serv_addr, sizeof(serv_addr))
printf("Connected to %s:%d\n", server->h_name, portno);

/* Carry out Client-Server protocol */
client(sockfd);

/* Clean up on termination */
close(sockfd);
```

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.43

## Recall: Server Protocol (v1)

```
/* Create Socket to receive requests*/
lstnsckfd = socket(AF_INET, SOCK_STREAM, 0);

/* Bind socket to port */
bind(lstnsckfd, (struct sockaddr *)&serv_addr, sizeof(serv_addr));
while (1) {
/* Listen for incoming connections */
listen(lstnsckfd, MAXQUEUE);

/* Accept incoming connection, obtaining a new socket for it */
consockfd = accept(lstnsckfd, (struct sockaddr *) &cli_addr,
&clilen);

server(consockfd);

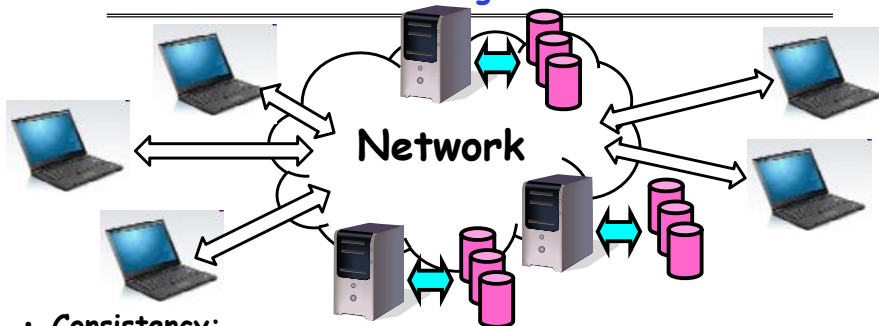
close(consockfd);
}
close(lstnsckfd);
```

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.44

## Network-Attached Storage and the CAP Theorem



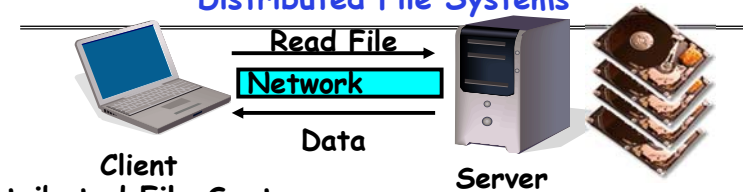
- **Consistency:**
  - Changes appear to everyone in the same serial order
- **Availability:**
  - Can get a result at any time
- **Partition-Tolerance**
  - System continues to work even when network becomes partitioned
- **Consistency, Availability, Partition-Tolerance (CAP) Theorem:**
  - **Cannot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

11/18/15

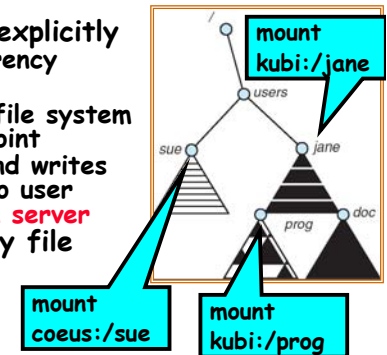
Kubiawicz CS162 ©UCB Fall 2015

Lec 21.45

## 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

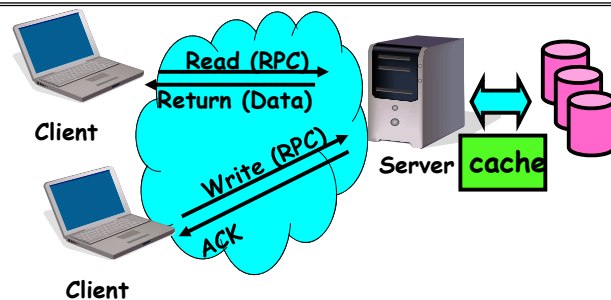


11/18/15

Kubiawicz CS162 ©UCB Fall 2015

Lec 21.46

## Simple Distributed File System



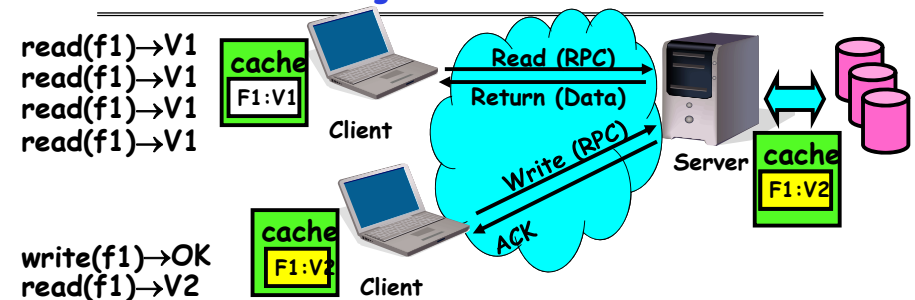
- **Remote Disk:** Reads and writes forwarded to server
  - Use Remote Procedure Calls (RPC) to translate file system calls into remote requests
  - No local caching/can be caching at server-side
- **Advantage:** Server provides completely consistent view of file system to multiple clients
- **Problems?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

11/18/15

Kubiawicz CS162 ©UCB Fall 2015

Lec 21.47

## 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

11/18/15

Kubiawicz CS162 ©UCB Fall 2015

Lec 21.48

## 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

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.49

## 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)

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.50

## NFS Continued

- NFS servers are **stateless**; each request provides all arguments required 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)

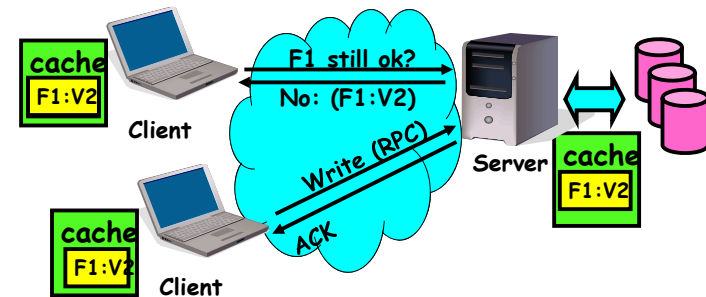
11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.51

## 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.



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

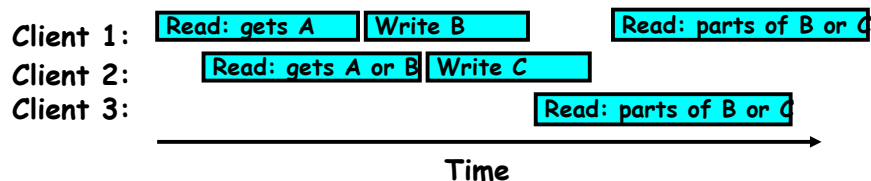
11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.52

## 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

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.53

## 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

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.54

## Summary

- Internet Protocol (IP)
  - Used to route messages through routes across globe
  - 32-bit addresses, 16-bit ports
- DNS: System for mapping from names⇒IP addresses
  - Hierarchical mapping from authoritative domains
  - Recent flaws discovered
- Datagram: a self-contained message whose arrival, arrival time, and content are not guaranteed
- Reliable/Ordered messages:
  - Use sequence numbers and reorder at destination
  - Use Acknowledgements
- 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
- Distributed File System:
  - Transparent access to files stored on a remote disk
  - Caching for performance

11/18/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.55