

CS162
Operating Systems and
Systems Programming
Lecture 21

Distributed Systems,
Networking, TCP/IP, RPC

April 15th, 2015
Prof. John Kubiawicz
<http://cs162.eecs.Berkeley.edu>

Recall: The ACID properties of Transactions

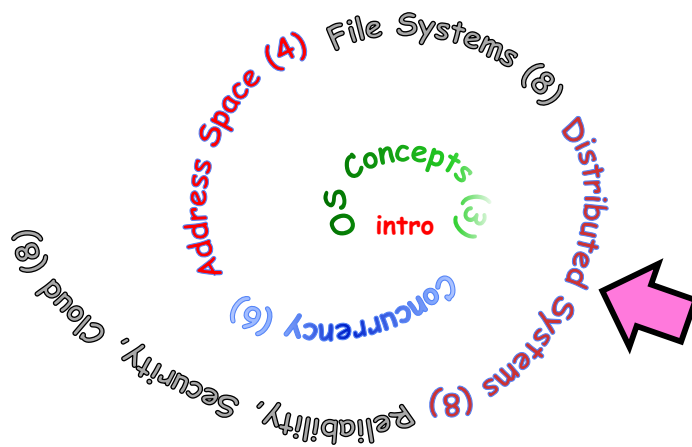
- **Atomicity:** all actions in the transaction happen, or none happen
- **Consistency:** transactions maintain data integrity, e.g.,
 - Balance cannot be negative
 - Cannot reschedule meeting on February 30
- **Isolation:** execution of one transaction is isolated from that of all others; no problems from concurrency
- **Durability:** if a transaction commits, its effects persist despite crashes

4/15/15

Kubiawicz CS162 ©UCB Spring 2015

Lec 21.2

Course Structure: Spiral



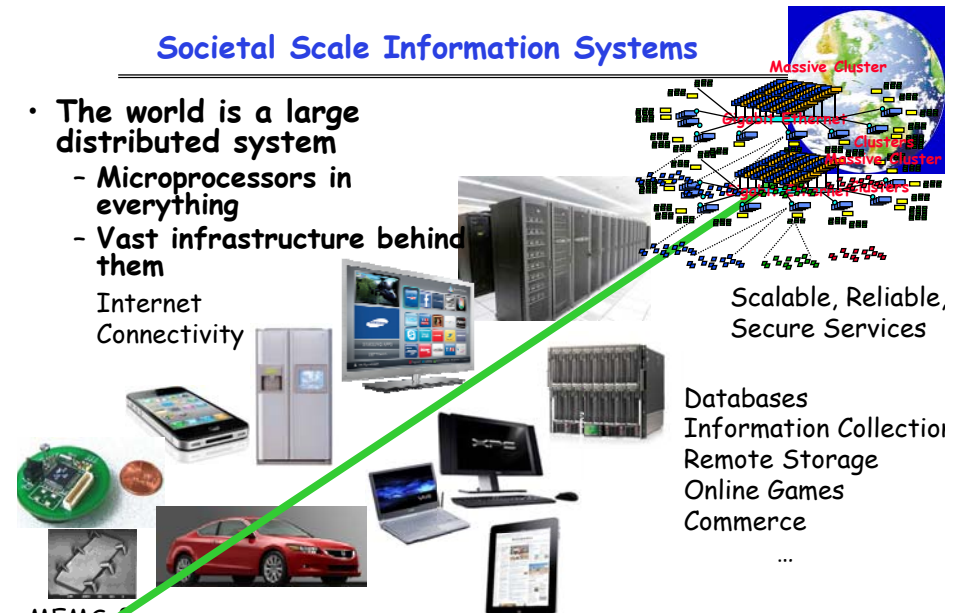
4/15/15

Kubiawicz CS162 ©UCB Spring 2015

Lec 21.3

Societal Scale Information Systems

- The world is a large distributed system
 - Microprocessors in everything
 - Vast infrastructure behind them
- Internet Connectivity



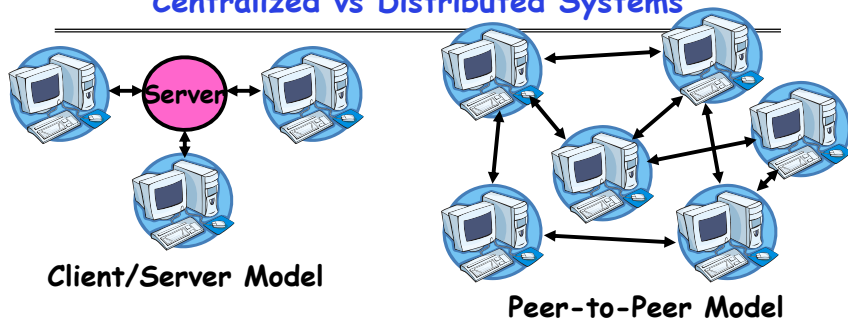
MEMS for
Sensor Nets

4/15/15

Kubiawicz CS162 ©UCB Spring 2015

Lec 21.4

Centralized vs Distributed Systems



- **Centralized System:** System in which major functions are performed by a single physical computer
 - Originally, everything on single computer
 - Later: client/server model
- **Distributed System:** physically separate computers working together on some task
 - Early model: multiple servers working together
 - » Probably in the same room or building
 - » Often called a "cluster"
 - Later models: peer-to-peer/wide-spread collaboration

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.5

Distributed Systems: Motivation/Issues

- Why do we want distributed systems?
 - Cheaper and easier to build lots of simple computers
 - Easier to add power incrementally
 - Users can have complete control over some components
 - Collaboration: Much easier for users to collaborate through network resources (such as network file systems)
- The *promise* of distributed systems:
 - Higher availability: one machine goes down, use another
 - Better durability: store data in multiple locations
 - More security: each piece easier to make secure
- Reality has been disappointing
 - Worse availability: depend on every machine being up
 - » Lamport: "a distributed system is one where I can't do work because some machine I've never heard of isn't working!"
 - Worse reliability: can lose data if any machine crashes
 - Worse security: anyone in world can break into system
- Coordination is more difficult
 - Must coordinate multiple copies of shared state information (using only a network)
 - What would be easy in a centralized system becomes a lot more difficult

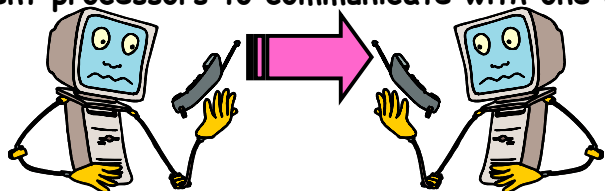
4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.6

Distributed Systems: Goals/Requirements

- **Transparency:** the ability of the system to mask its complexity behind a simple interface
- Possible transparencies:
 - **Location:** Can't tell where resources are located
 - **Migration:** Resources may move without the user knowing
 - **Replication:** Can't tell how many copies of resource exist
 - **Concurrency:** Can't tell how many users there are
 - **Parallelism:** System may speed up large jobs by splitting them into smaller pieces
 - **Fault Tolerance:** System may hide various things that go wrong in the system
- Transparency and collaboration require some way for different processors to communicate with one another



4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.7

What Is A Protocol?

- A protocol is an **agreement on how to communicate**
- Includes
 - **Syntax:** how a communication is specified & structured
 - » Format, order messages are sent and received
 - **Semantics:** what a communication means
 - » Actions taken when transmitting, receiving, or when a timer expires
- Described formally by a state machine
 - Often represented as a message transaction diagram

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.8

Examples of Protocols in Human Interactions

• Telephone

1. (Pick up / open up the phone)
2. Listen for a dial tone / see that you have service
3. Dial
4. Should hear ringing ...
5. Callee: "Hello?"
6. Caller: "Hi, it's John...."
Or: "Hi, it's me" (← what's *that* about?)
7. Caller: "Hey, do you think ... blah blah blah ..." pause
1. Callee: "Yeah, blah blah blah ..." pause
2. Caller: Bye
3. Callee: Bye
4. Hang up

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.9

Protocols in Human Interactions

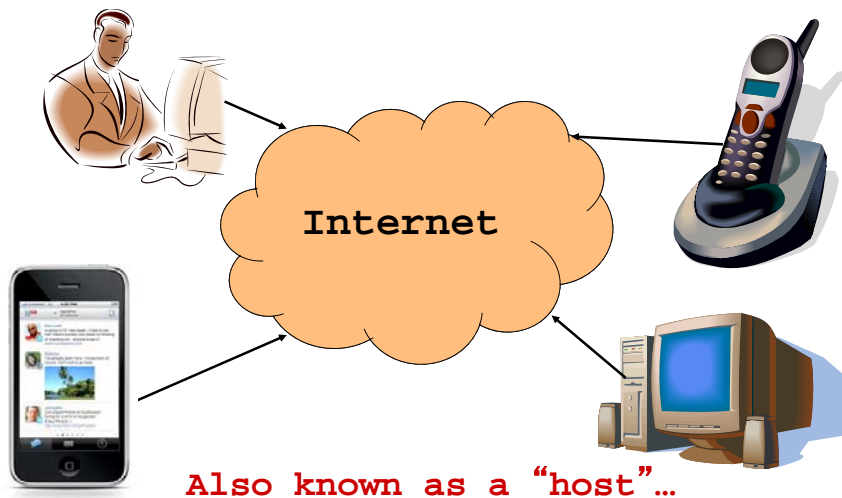
- Asking a question
 - Raise your hand
 - Wait to be called on
 - Or: wait for speaker to pause and vocalize

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.10

End System: Computer on the 'Net



4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.11

Recall: Namespaces for communication over IP

- Hostname
 - `www.eecs.berkeley.edu`
- IP address
 - `128.32.244.172` (ipv6?)
- Port Number
 - 0-1023 are "**well known**" or "system" ports
 - » Superuser privileges to bind to one
 - 1024 - 49151 are "registered" ports (**registry**)
 - » Assigned by IANA for specific services
 - 49152-65535 ($2^{15}+2^{14}$ to $2^{16}-1$) are "dynamic" or "private"
 - » Automatically allocated as "ephemeral Ports"

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.12

Recall: Client: getting the server address

```
struct hostent *buildServerAddr(struct sockaddr_in *serv_addr,
                               char *hostname, int portno) {
    struct hostent *server;

    /* Get host entry associated with a hostname or IP address */
    server = gethostbyname(hostname);
    if (server == NULL) {
        fprintf(stderr, "ERROR, no such host\n");
        exit(1);
    }

    /* Construct an address for remote server */
    memset((char *) serv_addr, 0, sizeof(struct sockaddr_in));
    serv_addr->sin_family = AF_INET;
    bcopy((char *)server->h_addr,
          (char *)&(serv_addr->sin_addr.s_addr), server->h_length);
    serv_addr->sin_port = htons(portno);

    return server;
}
```

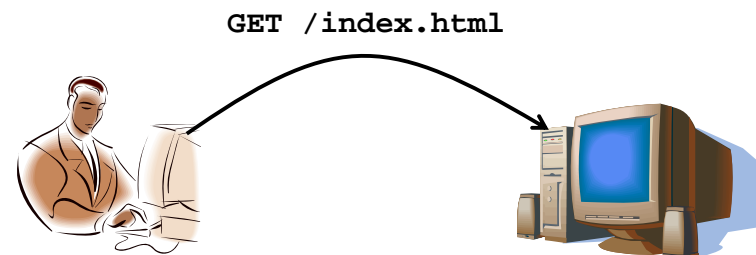
4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.13

Clients and Servers

- Client program
 - Running on end host
 - Requests service
 - E.g., Web browser



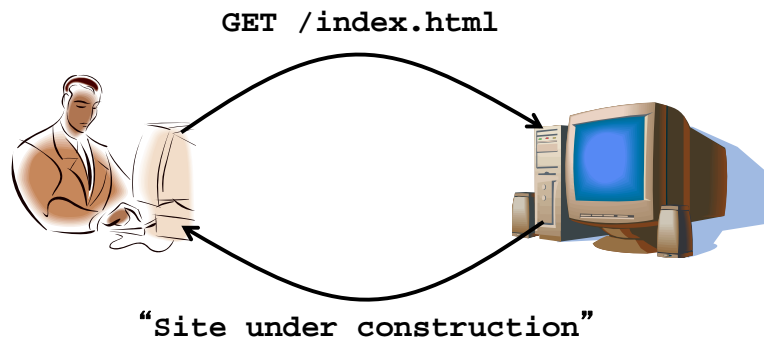
4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.14

Clients and Servers

- Client program
 - Running on end host
 - Requests service
 - E.g., Web browser
- Server program
 - Running on end host
 - Provides service
 - E.g., Web server



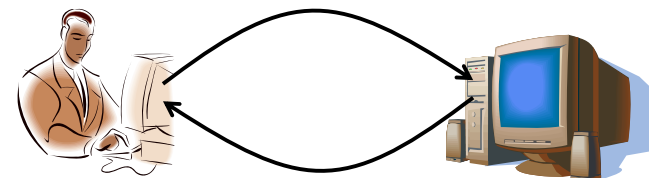
4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.15

Client-Server Communication

- Client "sometimes on"
 - Initiates a request to the server when interested
 - E.g., Web browser on your laptop or cell phone
 - Doesn't communicate directly with other clients
 - Needs to know the server's address
- Server is "always on"
 - Services requests from many client hosts
 - E.g., Web server for the *www.cnn.com* Web site
 - Doesn't initiate contact with the clients
 - Needs a fixed, well-known address



4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.16

Peer-to-Peer Communication

- No always-on server at the center of it all
 - Hosts can come and go, and change addresses
 - Hosts may have a different address each time
- Example: peer-to-peer file sharing (e.g., BitTorrent)
 - Any host can request files, send files, query to find where a file is located, respond to queries, and forward queries
 - Scalability by harnessing millions of peers
 - Each peer acting as **both a client and server**

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.17

The Problem

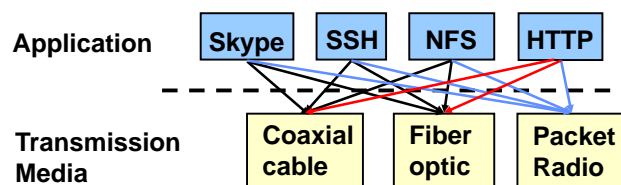
- Many different applications
 - email, web, P2P, etc.
- Many different network styles and technologies
 - Wireless vs. wired vs. optical, etc.
- How do we organize this mess?

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.18

The Problem (cont'd)



- Re-implement every application for every technology?
- No! But how does the Internet design avoid this?

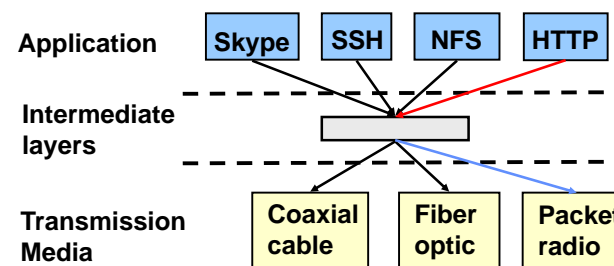
4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.19

Solution: Intermediate Layers

- Introduce intermediate layers that provide **set of abstractions** for various network functionality & technologies
 - A new app/media implemented only once
 - Variation on “add another level of indirection”

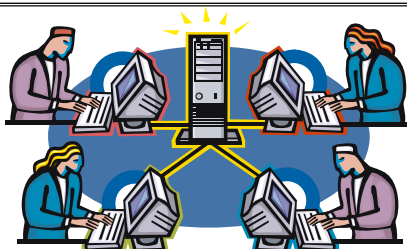


4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.20

Networking Definitions



- **Network:** physical connection that allows two computers to communicate
- **Packet:** unit of transfer, sequence of bits carried over the network
 - Network carries packets from one CPU to another
 - Destination gets interrupt when packet arrives
- **Protocol:** agreement between two parties as to how information is to be transmitted

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.21

Administrivia

- Midterm II: Next Wednesday (4/22)
 - Time: 6:30PM - 9:30PM
 - Location: Dwinelle: 145/155
 - Division of people between rooms will be posted
 - All topics from Midterm I, up to next Monday, including:
 - » Address Translation/TLBs/Paging
 - » I/O subsystems, Storage Layers, Disks/SSD
 - » Performance and Queuing Theory
 - » File systems
 - » Distributed systems, TCP/IP, RPC
 - » NFS/AFS, Key-Value Store
- Closed book, one page of notes - both sides
- Review session:
 - 306 Soda Hall
 - Sunday 4/19. 4:00pm—6:00pm

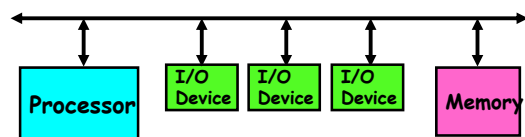
4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.22

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



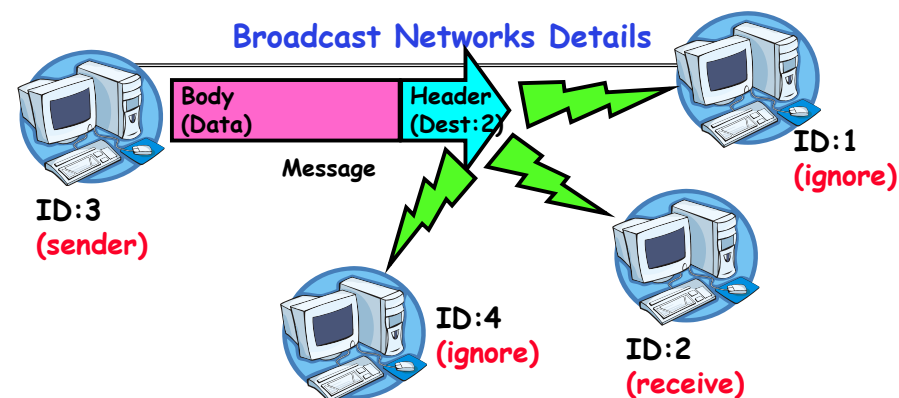
- Originally, Ethernet was a broadcast network
 - » All computers on local subnet connected to one another
- More examples (wireless: medium is air): cellular phones, GSM GPRS, EDGE, CDMA 1xRTT, and 1EvDO

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.23

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

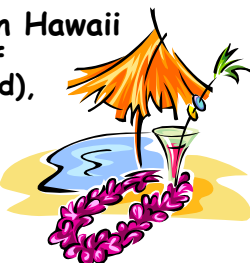
4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.24

Broadcast Network Arbitration

- **Arbitration:** Act of negotiating use of shared medium
 - What if two senders try to broadcast at same time?
 - Concurrent activity but can't use shared memory to coordinate!
- Aloha network (70's): packet radio within Hawaii
 - Blind broadcast, with checksum at end of packet. If received correctly (not garbled), send back an acknowledgement. If not received correctly, discard.
 - » Need checksum anyway - in case airplane flies overhead
 - Sender waits for a while, and if doesn't get an acknowledgement, re-transmits.
 - If two senders try to send at same time, both get garbled, both simply re-send later.
 - Problem: Stability: what if load increases?
 - » More collisions \Rightarrow less gets through \Rightarrow more resent \Rightarrow more load... \Rightarrow More collisions...
 - » Unfortunately: some sender may have started in clear, get scrambled without finishing



4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.25

Carrier Sense, Multiple Access/Collision Detection

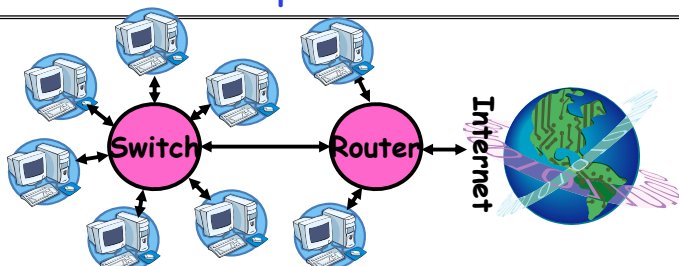
- Ethernet (early 80's): first practical local area network
 - It is the most common LAN for UNIX, PC, and Mac
 - Use wire instead of radio, but still broadcast medium
- Key advance was in arbitration called CSMA/CD: Carrier sense, multiple access/collision detection
 - **Carrier Sense:** don't send 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!

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.26

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.

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.27

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

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.28

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)

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.29

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

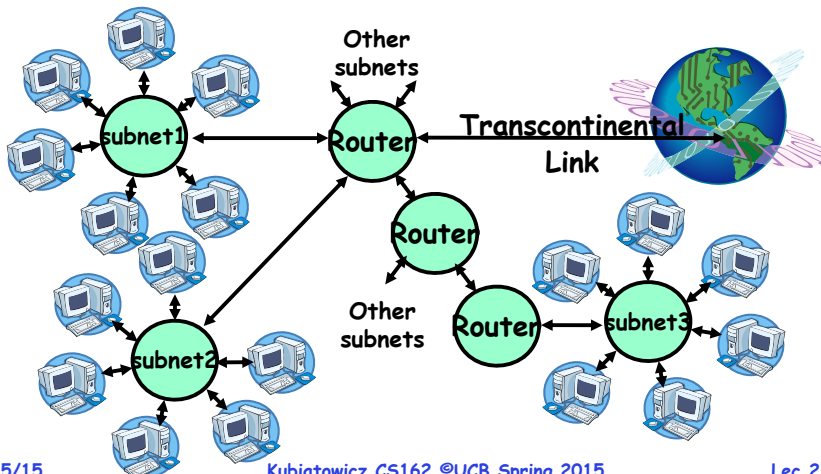
4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.30

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.



4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.31

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 - 150 Megabits/second
 - Broadcast usually requires multiple messages

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.32

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 → output link closer to destination
 - Default entry (for subnets without explicit entries)



4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.33

Setting up Routing Tables

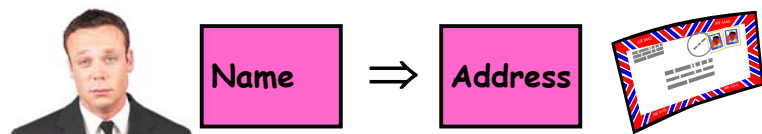
- **How do you set up routing tables?**
 - Internet has no centralized state!
 - » No single machine knows entire topology
 - » Topology constantly changing (faults, reconfiguration, etc)
 - Need dynamic algorithm that acquires routing tables
 - » Ideally, have one entry per subnet or portion of address
 - » 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

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.34

Naming in the Internet



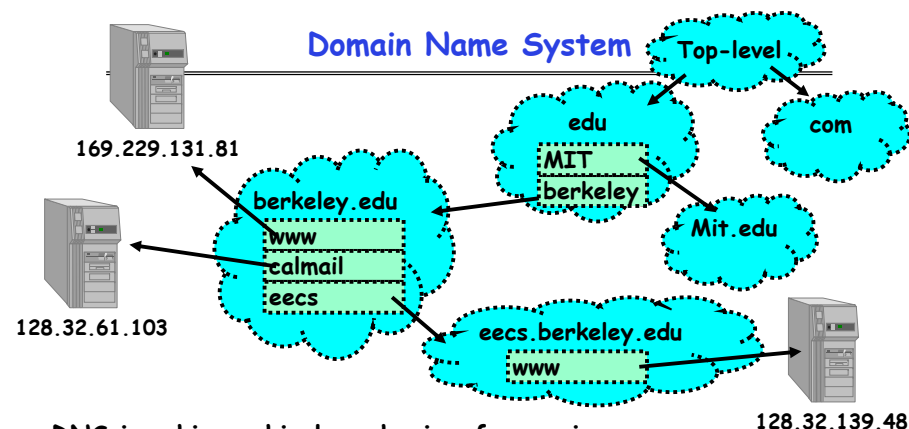
- **How to map human-readable names to IP addresses?**
 - E.g. `www.berkeley.edu` ⇒ `128.32.139.48`
 - E.g. `www.google.com` ⇒ different addresses depending on location, and load
- **Why is this necessary?**
 - IP addresses are hard to remember
 - IP addresses change:
 - » Say, Server 1 crashes gets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

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.35

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

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.36

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

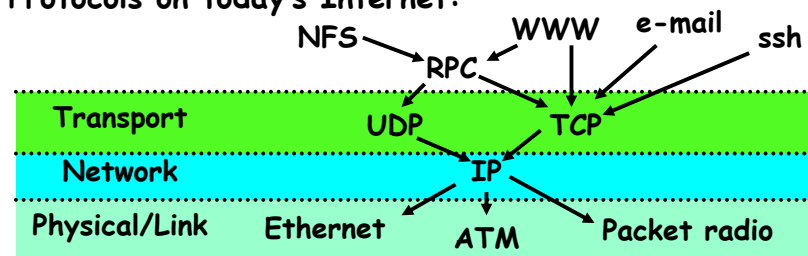
4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.37

Network Protocols

- **Protocol:** Agreement between two parties as to how information is to be transmitted
 - Example: system calls are the protocol between the operating system and application
 - Networking examples: 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:



4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.38

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 |

4/15/15

Lec 21.39

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)

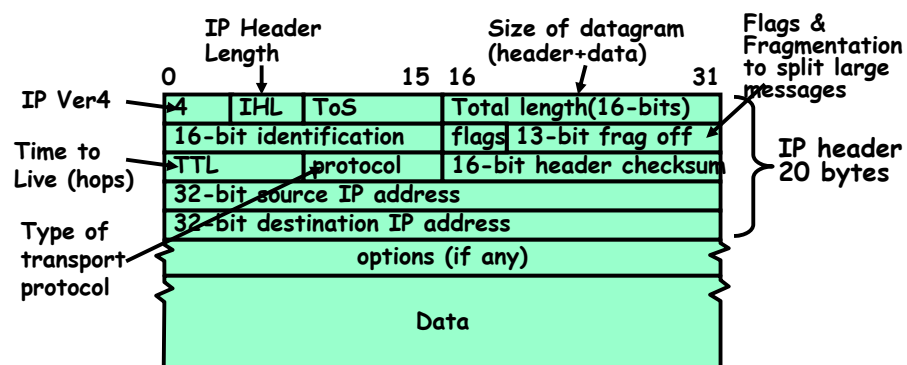
4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.40

IP Packet Format

IP Packet Format:



4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.41

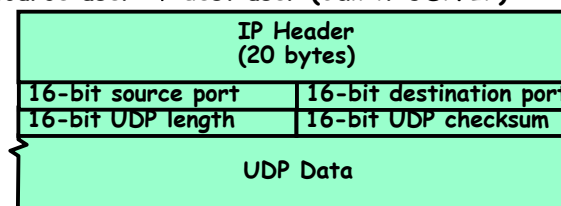
Building a messaging service

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

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.42

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...

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.43

Reliable Message Delivery: the Problem

All physical networks can garble and/or drop packets

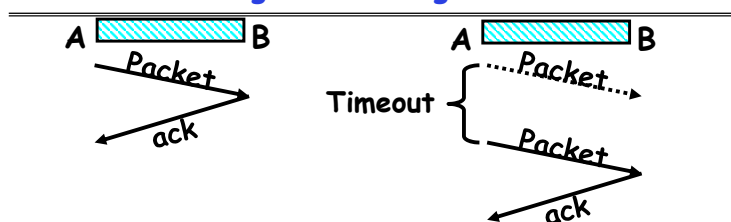
- Physical media: packet not transmitted/received
 - » If transmit close to maximum rate, get more throughput - even if some packets get lost
 - » If transmit at lowest voltage such that error correction just starts correcting errors, get best power/bit
- Congestion: no place to put incoming packet
 - » 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

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.44

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.

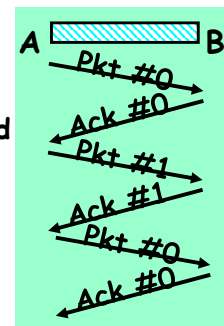
4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.45

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 × throughput)
- Con: doesn't work if network can delay or duplicate messages arbitrarily



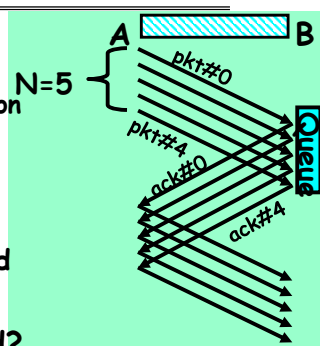
4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.46

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 ⇒ 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

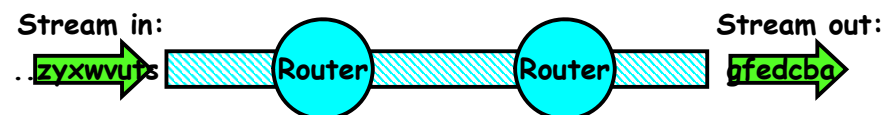


4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.47

Transmission Control Protocol (TCP)



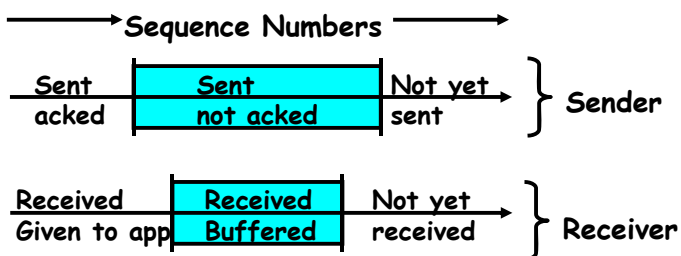
- Transmission Control Protocol (TCP)
 - TCP (IP Protocol 6) layered on top of IP
 - Reliable byte stream between two processes on different machines over Internet (read, write, flush)
- TCP Details
 - Fragments byte stream into packets, hands packets to IP
 - » IP may also fragment by itself
 - Uses 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"

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.48

TCP Windows and Sequence Numbers



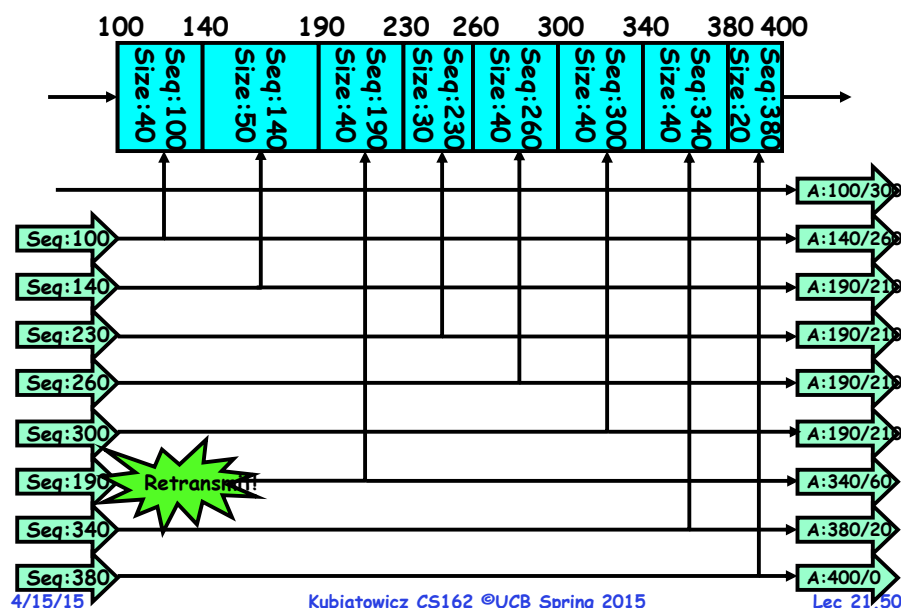
- Sender has three regions:
 - Sequence regions
 - » sent and ack'd
 - » Sent and not ack'd
 - » not yet sent
 - Window (colored region) adjusted by sender
- Receiver has three regions:
 - Sequence regions
 - » received and ack'd (given to application)
 - » received and buffered
 - » not yet received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

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.49

Window-Based Acknowledgements (TCP)

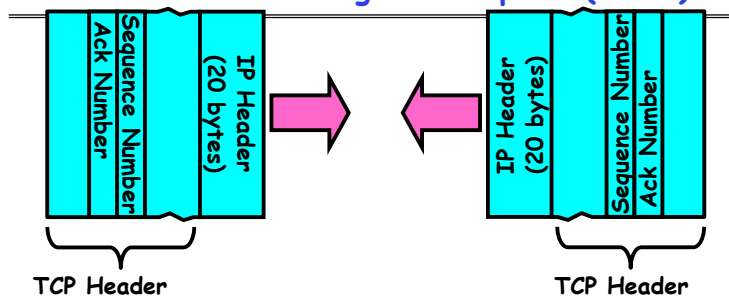


4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.50

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)

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.51

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"

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.52

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

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.53

Use of TCP: Sockets

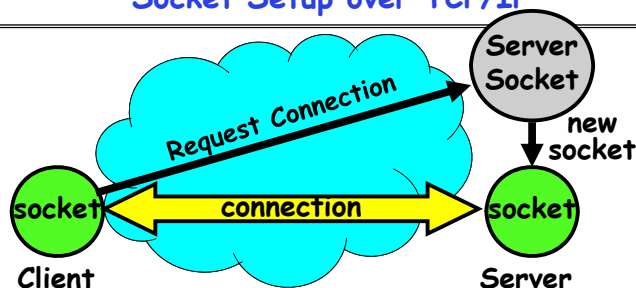
- **Socket:** an abstraction of a network I/O queue
 - Embodies one side of a communication channel
 - » Same interface regardless of location of other end
 - » Could be local machine (called "UNIX socket") or remote machine (called "network socket")
 - First introduced in 4.2 BSD UNIX: big innovation 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

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.54

Socket Setup over TCP/IP



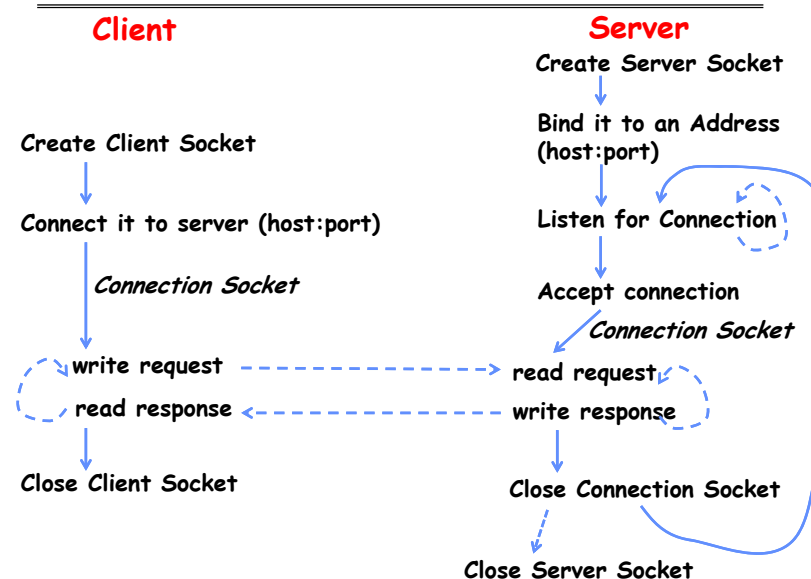
- Server Socket: Listens for new connections
 - Produces new sockets for each unique connection
- Things to remember:
 - Connection involves 5 values: [Client Addr, Client Port, Server Addr, Server Port, Protocol]
 - Often, Client Port "randomly" assigned
 - » Done by OS during client socket setup
 -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

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.55

Recall: Sockets in concept



4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.56

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);
```

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.57

Recall: Server Protocol (v1)

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

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

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

server(consockfd);

close(consockfd);
}
close(ltnsockfd);
```

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.58

Distributed Applications

- 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



- One 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

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.59

Using Messages: Send/Receive behavior

- When should send(message, mbox) return?
 - When receiver gets message? (i.e. ack received)
 - When message is safely buffered on destination?
 - Right away, if message is buffered on source node?
- Actually two questions here:
 - When can the sender be sure that receiver actually received the message?
 - When can sender reuse the memory containing message?
- Mailbox provides 1-way communication from T1→T2
 - T1→buffer→T2
 - Very similar to producer/consumer
 - » Send = V, Receive = P
 - » However, can't tell if sender/receiver is local or not!

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.60

Messaging for Producer-Consumer Style

- Using send/receive for producer-consumer style:

```

Producer:
int msg1[1000];
while(1) {
    prepare message;
    send(msg1,mbox);
}
    
```

Send Message

```

Consumer:
int buffer[1000];
while(1) {
    receive(buffer,mbox);
    process message;
}
    
```

Receive Message

- No need for producer/consumer to keep track of space in mailbox: handled by send/receive
 - One of the roles of the window in TCP: window is size of buffer on far end
 - Restricts sender to forward only what will fit in buffer

Messaging for Request/Response communication

- What about two-way communication?
 - Request/Response
 - Read a file stored on a remote machine
 - Request a web page from a remote web server
 - Also called: **client-server**
 - Client ≡ requester, Server ≡ responder
 - Server provides "service" (file storage) to the client
- Example: File service

```

Client: (requesting the file)
char response[1000];
    
```

Request File

```

send("read rutabaga", server_mbox);
receive(response, client_mbox);
    
```

Get Response

```

Server: (responding with the file)
char command[1000], answer[1000];
    
```

```

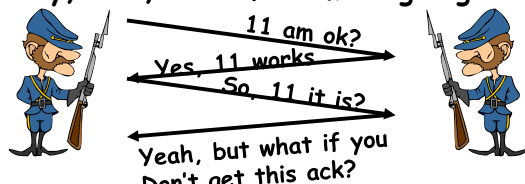
receive(command, server_mbox);
decode command;
read file into answer;
send(answer, client_mbox);
    
```

Receive Request

Send Response

General's Paradox

- General's paradox:
 - Constraints of problem:
 - Two generals, on separate mountains
 - Can only communicate via messengers
 - Messengers can be captured
 - Problem: need to coordinate attack
 - If they attack at different times, they all die
 - If they attack at same time, they win
 - Named after Custer, who died at Little Big Horn because he arrived a couple of days too early
- Can messages over an unreliable network be used to guarantee two entities do something simultaneously?
 - Remarkably, "no", even if all messages get through



- No way to be sure last message gets through!

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 does this
 - Use a persistent, stable log on each machine to 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

Two phase commit example

- Simple Example: A=WellsFargo Bank, B=Bank of America
 - Phase 1: **Prepare** Phase
 - » A writes "Begin transaction" to log
 - A→B: OK to transfer funds to me?
 - » Not enough funds:
 - B→A: transaction aborted; A writes "Abort" to log
 - » Enough funds:
 - B: Write new account balance & promise to commit to log
 - B→A: OK, I can commit
 - Phase 2: A can decide for both whether they will **commit**
 - » A: write new account balance to log
 - » Write "Commit" to log
 - » Send message to B that commit occurred; wait for ack
 - » Write "Got Commit" to log
- What if B crashes at beginning?
- Wakes up, does nothing; A will timeout, abort and retry
- What if A crashes at beginning of phase 2?
- Wakes up, sees that there is a transaction in progress; sends "Abort" to B
- What if B crashes at beginning of phase 2?
- B comes back up, looks at log; when A sends it "Commit" message, it will say, "oh, ok, commit"

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.65

Distributed Decision Making Discussion

- Why is distributed decision making desirable?
 - Fault Tolerance!
 - A group of machines can come to a decision even if one or more of them fail during the process
 - » Simple failure mode called "failstop" (different modes later)
 - After decision made, result recorded in multiple places
- Undesirable feature of Two-Phase Commit: **Blocking**
 - One machine can be stalled until another site recovers:
 - » Site B writes "prepared to commit" record to its log, sends a "yes" vote to the coordinator (site A) and crashes
 - » Site A crashes
 - » Site B wakes up, check its log, and realizes that it has voted "yes" on the update. It sends a message to site A asking what happened. At this point, B cannot decide to abort, because update may have committed
 - » B is blocked until A comes back
 - A blocked site holds resources (locks on updated items, pages pinned in memory, etc) until learns fate of update
- **PAXOS**: An alternative used by GOOGLE and others that does not have this blocking problem
- What happens if one or more of the nodes is malicious?
 - **Malicious**: attempting to compromise the decision making

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.66

Byzantine General's Problem



- Byzantine General's Problem (n players):
 - One General
 - n-1 Lieutenants
 - Some number of these (f) 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

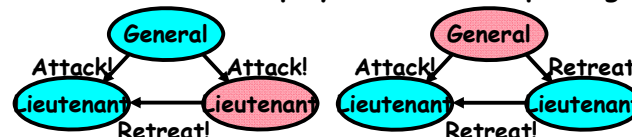
4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

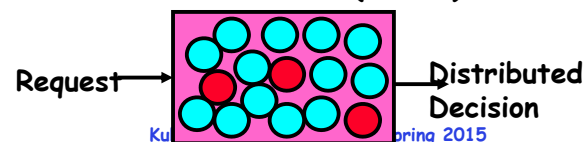
Lec 21.67

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



4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.68

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:

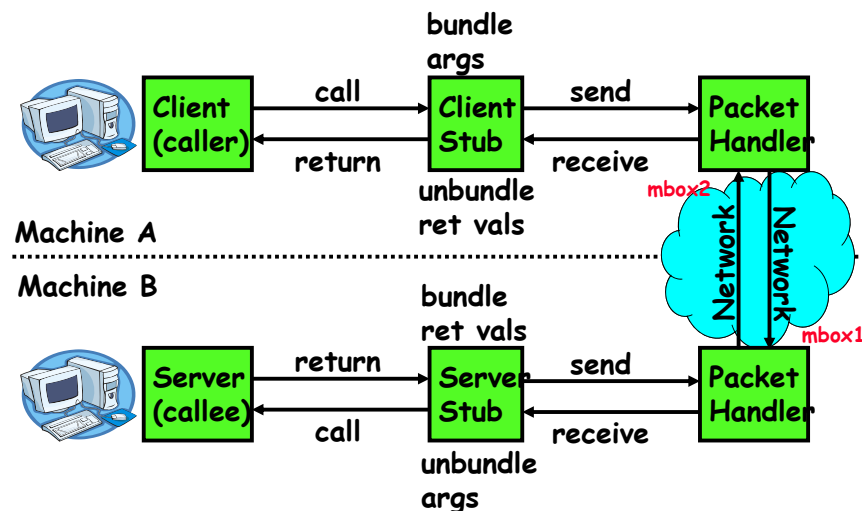

```
fileSys→Read("rutabaga");
```
- Implementation:
 - Request-response message passing (under covers!)
 - "Stub" provides glue on client/server
 - » Client stub is responsible for "marshalling" arguments and "unmarshalling" the return values
 - » Server-side stub is responsible for "unmarshalling" arguments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

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.69

RPC Information Flow



4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.70

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).

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.71

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

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.72

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

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.73

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 modern RPC systems:
 - CORBA (Common Object Request Broker Architecture)
 - DCOM (Distributed COM)
 - RMI (Java Remote Method Invocation)

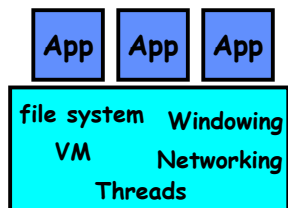
4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

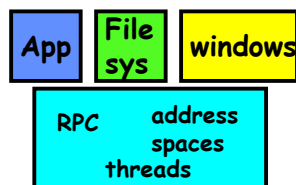
Lec 21.74

Microkernel operating systems

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



Monolithic Structure



Microkernel Structure

- 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.

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.75

Summary (1/2)

- Network: physical connection that allows two computers to communicate
 - Packet: sequence of bits carried over the network
- **Broadcast Network**: Shared Communication Medium
 - Transmitted packets sent to all receivers
 - Arbitration: act of negotiating use of shared medium
 - » Ethernet: Carrier Sense, Multiple Access, Collision Detect
- **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.
- **Protocol**: Agreement between two parties as to how information is to be transmitted
- Internet Protocol (IP)
 - Used to route messages through routes across globe
 - 32-bit addresses, 16-bit ports
- **DNS**: System for mapping from names \Rightarrow IP addresses
 - Hierarchical mapping from authoritative domains
 - Recent flaws discovered

4/15/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 21.76

Summary (2/2)

- **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
- **Two-phase commit**: distributed decision making
 - First, make sure everyone guarantees that they will commit if asked (prepare)
 - Next, ask everyone to commit
- **Byzantine General's Problem**: distributed decision making with malicious failures
 - One general, $n-1$ lieutenants: some number of them may be malicious (often f of them)
 - All non-malicious lieutenants must come to same decision
 - If general not malicious, lieutenants must follow general
 - Only solvable if $n \geq 3f+1$
- **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)