

CS162

Operating Systems and Systems Programming Lecture 24

Networking and TCP/IP (Cont'd), RPC, Distributed File Systems

April 21st, 2022

Prof. Anthony Joseph and John Kubiatowicz

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

- Consensus problem
 - All nodes propose a value
 - Some nodes might crash and stop responding
 - Eventually, all remaining nodes decide on the same value from set of proposed values
- Distributed Decision Making
 - Choose between "true" and "false"
 - Or Choose between "commit" and "abort"
- Equally important (but often forgotten!): make it durable!
 - How do we make sure that decisions cannot be forgotten?
 - » This is the "D" of "ACID" in a regular database
 - In a global-scale system?
 - » What about erasure coding or massive replication?
 - » Like **BlockChain** applications!

4/21/22

Joseph & Kubiatowicz CS162 @ UCB Spring 2022

Lec 24.2

Recall: Distributed Consensus Making

- 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"
 - Then asks all nodes to commit; they respond with ACK
 - After receive ACKs, coordinator writes "Got Commit" to log
 - Log used to guarantee that all machines either commit or don't

4/21/22

Joseph & Kubiatowicz CS162 @ UCB Spring 2022

Lec 24.2

Distributed Decision Making Discussion (1/2)

- 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 "fail-stop" (different modes later)
 - After decision made, result recorded in multiple places
- Why is 2PC not subject to the General's paradox?
 - Because 2PC is about *all nodes eventually coming to the same decision – not necessarily at the same time!*
 - Allowing us to reboot and continue allows time for collecting and collating decisions

4/21/22

Joseph & Kubiatowicz CS162 @ UCB Spring 2022

Lec 24.3

4/21/22

Joseph & Kubiatowicz CS162 @ UCB Spring 2022

Lec 24.4

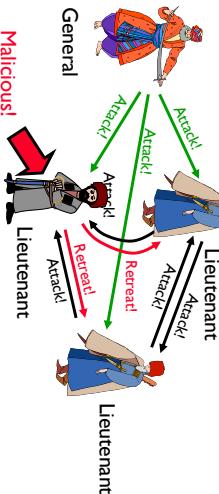
Distributed Decision Making Discussion (2/2)

- 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, checks 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

4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.5



Byzantine General's Problem

- Byzantine General's Problem (n players):
 - One General and 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 the following Integrity Constraints apply:
 - [C1]: All loyal lieutenants obey the same order
 - [C2]: If the commanding general is loyal, then all loyal lieutenants obey the order he sends

4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.7

Alternatives to 2PC

- Three-Phase Commit:** One more phase, allows nodes to fail or block and still make progress.
- PAXOS:** An alternative used by Google and others that does not have 2PC blocking problem
 - Developed by Leslie Lamport (Turing Award Winner)
 - No fixed leader; can choose new leader on the fly, deal with failure
 - Some think this is extremely complex!
- RAFT:** PAXOS alternative from John Osterhout (Stanford)
 - Simpler to describe complete protocol

- What happens if one or more of the nodes is malicious?
 - Malicious: attempting to compromise the decision making
 - Use a more hardened decision making process;

Byzantine Agreement and Block Chains

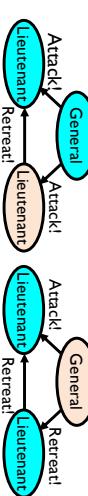
4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.6

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
 - Never algorithms have message complexity $\mathcal{O}(n^2)$
 - 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



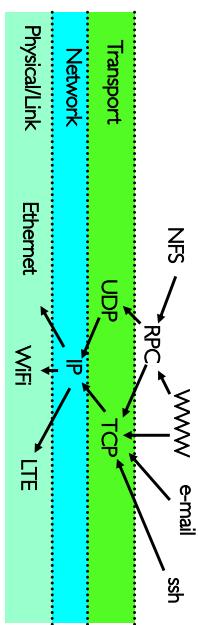
4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.8

Network Protocols

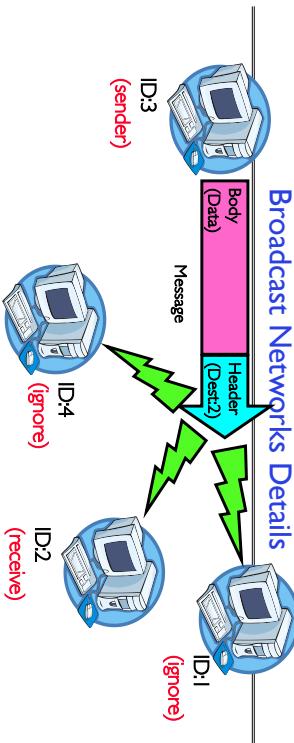
- Networking protocols: many levels
 - Physical level: mechanical and electrical network (e.g., how are 0 and 1 represented)
 - Link level: packet formats/error control (for instance, the CSMA/CD protocol)
 - Network level: network routing, addressing
 - Transport Level: reliable message delivery
- Protocols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



4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.11

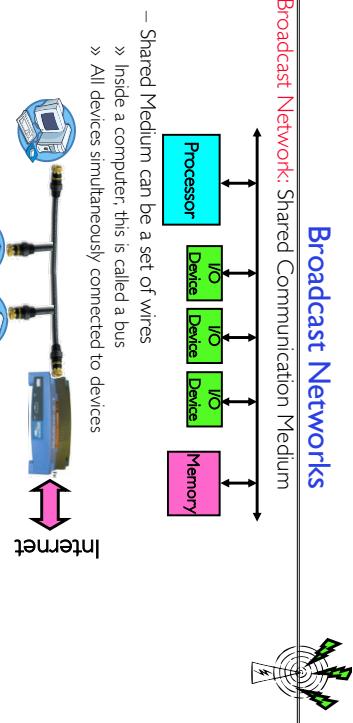


4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.13

Point-to-point networks



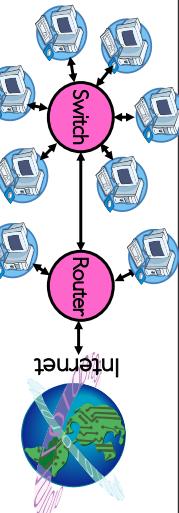
4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.12

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



4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.16

- Media Access Control (MAC) Address:**
 - 48-bit physical address for hardware interface
 - Every device (in the world!) has a unique address
 - Delivery:** When you broadcast a packet, how does a receiver know who it is for?
- Point-to-point network:** a network in which every physical wire is connected to only two computers
 - 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 hardware is cheap!
- Switch:** a bridge that transforms a shared-bus (broadcast) configuration into a point-to-point network
 - Adaptively figures out which ports have which MAC addresses
- Router:** a device that acts as a junction between physical networks to transfer data packets among them
 - Routes between switching domains using (for instance) IP addresses
- No OS interrupt if not for particular destination

4/21/22

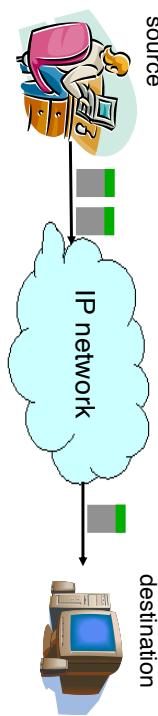
Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.12

The Internet Protocol (IP)

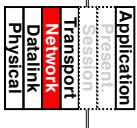
- Internet Protocol: Internet's network layer
 - Service it provides: "Best-Effort" Packet Delivery
 - Tries its "best" to deliver packet to its destination
 - Packets may be lost
 - Packets may be corrupted
 - Packets may be delivered out of order

- [IP] Is a Datagram service!
 - Routes across many physical switching domains (subnets)



Lec 24.12

Joseph & Kubiatowicz CS 162 © UCB Spring 2022



- **IP Address:** a 32-bit integer used as destination of IP packet
 - Often written as four dot-separated integers, with each integer representing 8 bits (thus representing 32 bits)
 - Example CS file server is: 169.229.60.83 ≡ 0xa9xE53C53

- **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
 - » Groups of machines may share a single IP address
 - » In this case, machines have private addresses behind a "Network Address Translation" (NAT) gateway
 - Not every computer has a unique IP address

Lec 24.17

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

IPv4 Address Space

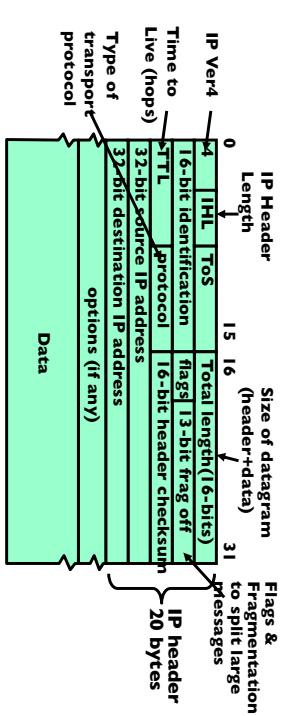
- **Subnet:** network connecting hosts with related IP addresses
 - A subnet is identified by 32-bit value, with the bits which differ set to zero, followed by a slash and a mask
 - » Example: 128.32.1.0/24 designates a subnet in which all the addresses look like 128.32.1.31.XX
 - » Same subnet: 128.32.1.31/25.255.255.2550
 - **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)
 - **Often routing within subnet is by MAC address (smart switches)**

Lec 24.18

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

IPv4 Packet Format

- IP Packet Format:



- IP Datagram: an unreliable, unordered, packet sent from source to destination
 - Function of network – deliver datagrams!

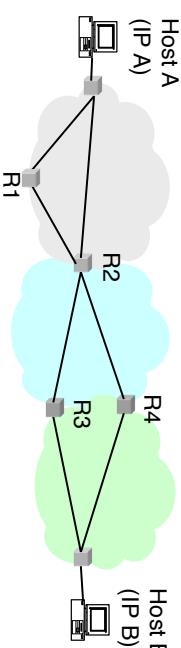
Lec 24.20

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

Wide Area Network

- **Wide Area Network (WAN):** network that covers a broad area (e.g., city, state, country, entire world)
 - E.g., Internet is a WAN

- WAN connects multiple physical (datalink) layer networks (LANs)
- Datalink layer networks are connected by **routers**
 - Different LANs can use different communication technology (e.g., wireless, cellular, optics, wired)



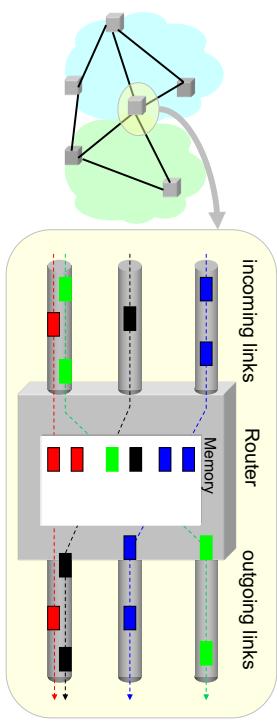
Lec 24.22

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

Lec 24.21

Routers

- Forward each packet received on an **incoming link** to an **outgoing link** based on packet's destination IP address (towards its destination)
- Store & forward: packets are buffered before being forwarded
- Forwarding table: mapping between IP address and the output link



4/21/22

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

Lec 24.22

IP Addresses vs. MAC Addresses

- Why not use MAC addresses for routing?
 - Doesn't scale
- Analogy
 - MAC address → SSN
 - IP address → (unreadable) home address
- MAC address: uniquely associated with device for the entire lifetime of the device
- IP address: changes as the device location changes
 - Your notebook IP address at school is different from home



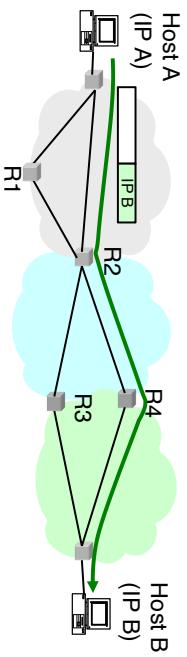
4/21/22

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

Lec 24.24

Packet Forwarding

- Upon receiving a packet, a router
 - read the IP destination address of the packet
 - consults its forwarding table → output port
 - fowards packet to corresponding output port
- Default route (for subnets without explicit entries)
 - Forward to more authoritative router



4/21/22

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

Lec 24.23

IP Addresses vs. MAC Addresses

- Why does packet forwarding using IP addr. scale?
 - Because IP addresses can be aggregated
 - E.g. all IP addresses at UC Berkeley start with **0xA9E5**, i.e., any address of form **0xA9E5******* belongs to Berkeley
 - Thus, a router in NY needs to keep a **single** entry for **all** hosts at Berkeley
 - If we were using MAC addresses the NY router would need to maintain **an entry for every** Berkeley host!!
- Analogy:
 - Give this letter to person with SSN: 123-45-6789 vs.
— Give this letter to "John Smith, 123 First Street, LA, US"



4/21/22

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

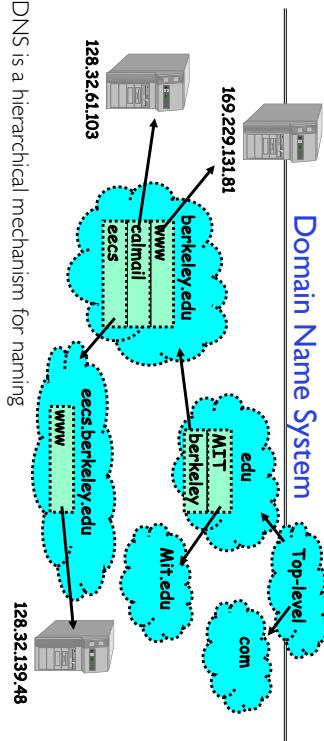
Lec 24.25

Administrivia

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

4/21/22

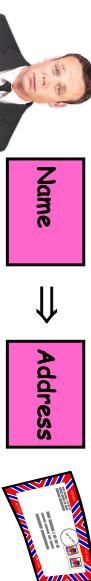
Joseph & Kubiatowicz CS 162 © UCB Spring 2022



Lec 24.27

4/21/22

Joseph & Kubiatowicz CS 162 © UCB Spring 2022



Naming in the Internet

- How to map human-readable names to IP addresses?
 - Eg. www.berkeley.edu ⇒ 128.32.139.48
 - Eg. 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)

Lec 24.27

4/21/22

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

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

Lec 24.29

4/21/22

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

Lec 24.29

4/21/22

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

Lec 24.30

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:

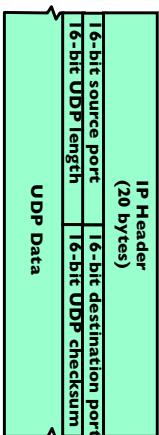
Physical Reality: Packets	Abstraction: Messages
Limited Size (MTU)	Arbitrary Size
Unordered (sometimes)	Ordered
Unreliable	Reliable
Machine-to-machine	Process-to-process
Only on local area net	Routed anywhere
Asynchronous	Synchronous
Insecure	Secure

Lec 24.3 |

4/21/22

Building a messaging service on IP

- 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]
- For example: The Unreliable Datagram Protocol (UDP)
 - Layered on top of basic IP (**IP Protocol 17**)
 - » **Datagram:** an unreliable, unordered, packet sent from source user → dest user (Call it UDP/IP)



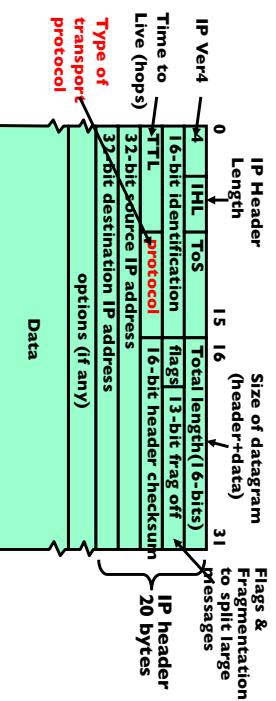
- Important aspect: low overhead!
 - » Often used for high-bandwidth video streams
 - » Many uses of UDP considered “anti-social” – none of the “well-behaved” aspects of (say) TCP/IP

4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.3 |

Recall: IPv4 Packet Format



Lec 24.3 |

4/21/22

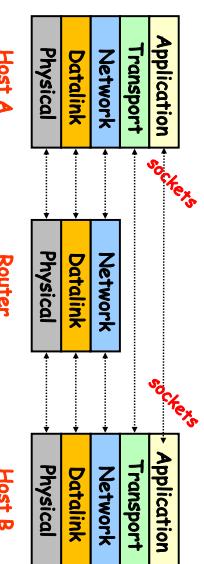
- IP Packet Format:
 - Function of network – deliver datagrams

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.3 |

Internet Architecture: Five Layers

- Lower three layers implemented everywhere
- Top two layers implemented only at hosts

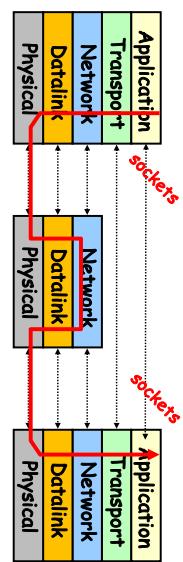


Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.3 |

Internet Architecture: Five Layers

- Communication goes down to physical network
- Then from network peer to peer
- Then up to relevant layer

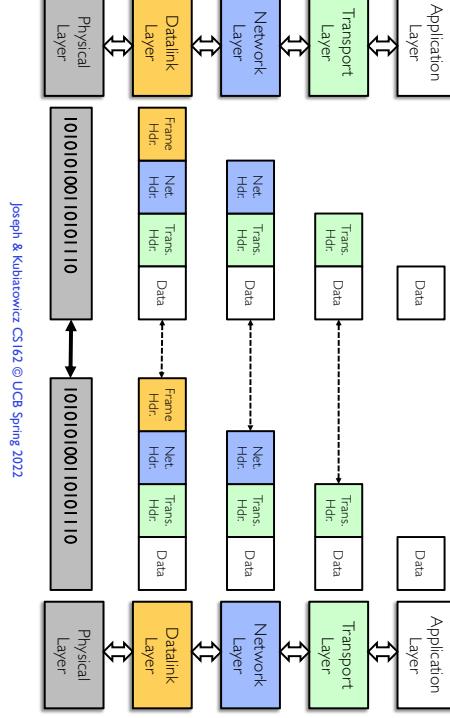


4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.35

4/21/22



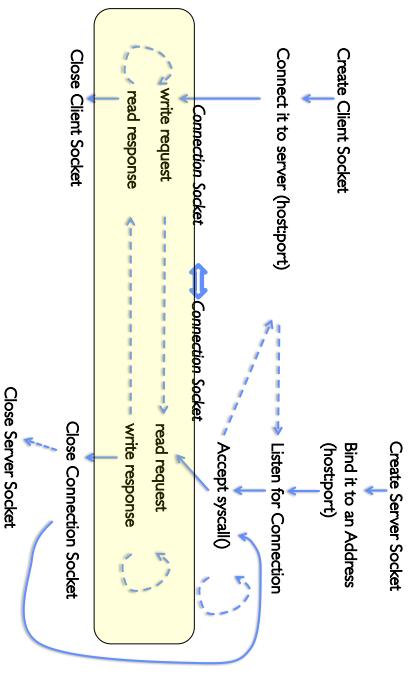
Lec 24.36

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Recall: Sockets in concept

Client

Server



4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

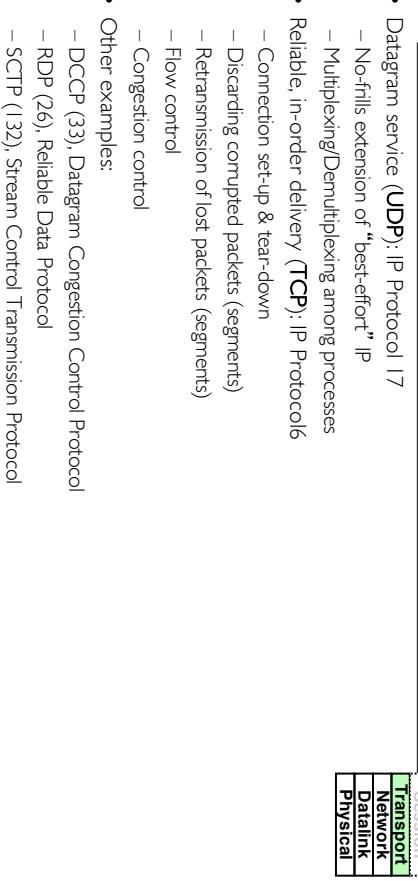
Lec 24.37

4/21/22

Layering Analogy: Packets in Envelopes

Application
Session
Transport
Network
Datalink
Physical

Lec 24.38



- Datagram service (**UDP**): IP Protocol 17
 - No-fills extension of “best-effort” IP
 - Multiplexing/Demultiplexing among processes
- Reliable, in-order delivery (**TCP**): IP Protocol 6
 - Connection set-up & tear-down
 - Discarding corrupted packets (segments)
 - Retransmission of lost packets (segments)
 - Flow control
 - Congestion control
- Other examples:
 - DCCP (33), Datagram Congestion Control Protocol
 - RDP (26), Reliable Data Protocol
 - SCTP (132), Stream Control Transmission Protocol

- Datagram service (**UDP**): IP Protocol 17
 - No-fills extension of “best-effort” IP
 - Multiplexing/Demultiplexing among processes
- Reliable, in-order delivery (**TCP**): IP Protocol 6
 - Connection set-up & tear-down
 - Discarding corrupted packets (segments)
 - Retransmission of lost packets (segments)
 - Flow control
 - Congestion control
- Other examples:
 - DCCP (33), Datagram Congestion Control Protocol
 - RDP (26), Reliable Data Protocol
 - SCTP (132), Stream Control Transmission Protocol

Lec 24.38

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

4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.39

Problem: Dropped Packets

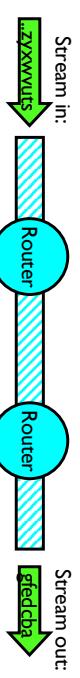
- All physical networks can garble or drop packets
 - Physical hardware problems (bad wire, bad signal)
- Therefore, IP can garble or drop packets
 - It doesn't repair this itself (end-to-end principle!)
- Building reliable message delivery
 - Confirm that packets aren't garbled
 - Confirm that packets arrive **exactly once**

Lec 24.41

4/21/22

Lec 24.39

Joseph & Kubiatowicz CS162 © UCB Spring 2022

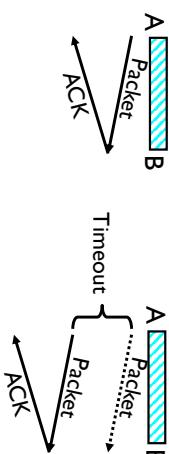


Transmission Control Protocol (TCP)

- 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
 - 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"

Lec 24.40

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, ACK each

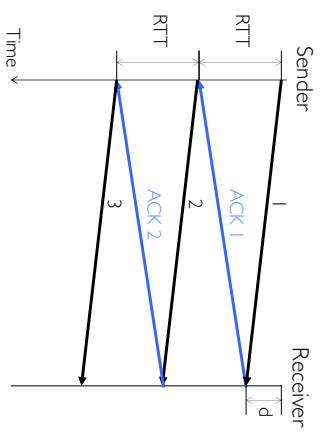
Lec 24.42

4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Stop-and-Wait (No Packet Loss)

- Send; wait for ACK; repeat
- Round Trip Time (RTT): time it takes a packet to travel from sender to receiver and back
 - One-way latency (d): one way delay from sender and receiver
- For symmetric latency,
 $RTT = 2d$



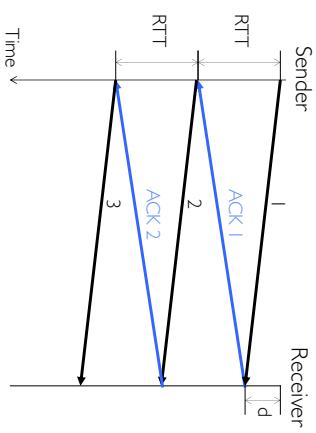
4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.43

Stop-and-Wait (No Packet Loss)

- So bandwidth is 1 packet per RTT
 - Depends only on latency, not network capacity (!)



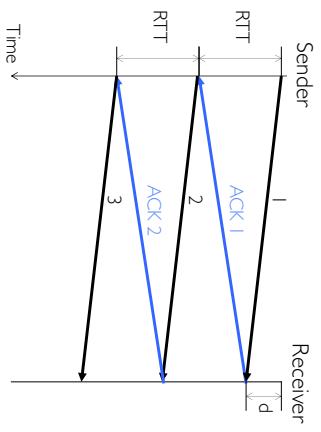
4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.45

Stop-and-Wait (No Packet Loss)

- How fast can you send data?
- Little's Law applied to the network:
 $n = B \cdot RTT$
- For Stop-and-Wait, $n = 1$ packet
- So bandwidth is 1 packet per RTT
 - Depends only on latency, not network capacity (!)



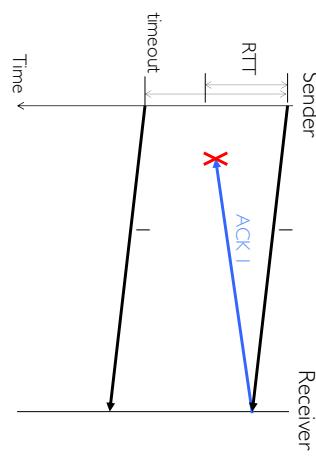
4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.44

Stop-and-Wait with Packet Loss

- Loss recovery relies on timeouts
- How to choose a good timeout?
 - Too short – lots of duplication
 - Too long – packet loss is really disruptive!
- How to deal with duplication?
 - Retransmission certainly opens up the possibility for
- Very inefficient if we have a 100 Mbps link



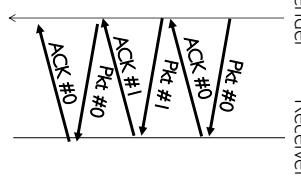
4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.46

How to Deal with Message Duplication?

- Solution: put sequence number in message to identify re-transmitted packets
 - Receiver checks for duplicate number's; Discard if detected
- Requirements:
 - Sender keeps copy of unACK'd messages
 - » Easy: only need to buffer messages
 - » Hard: when ok to forget about received message?
 - 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 number of last message received
 - Pros: simple, small overhead
 - Con: doesn't work if network can delay or duplicate messages arbitrarily



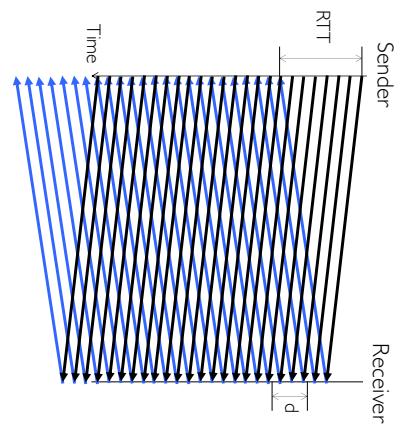
4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.47

Advantages of Moving Away From Stop-and-Wait

- Larger space of acknowledgements
 - Pipelining: don't wait for ACK before sending next packet
- ACKs serve dual purpose:**
 - Reliability: Confirming packet received
 - Ordering: Packets can be reordered at destination
- How much data is in flight now?
 - Bytes in flight: $W_{\text{send}} = \text{RTT} \times B$
 - Here B is in "bytes/second"
 - $W_{\text{send}} \equiv$ Sender's "Window Size"
 - Packets in flight = $(W_{\text{send}} / \text{packet size})$
- How long does the sender have to keep the packets around?
- How long does the receiver have to keep the packets' data?
- What if sender is sending packets faster than the receiver can process the data?

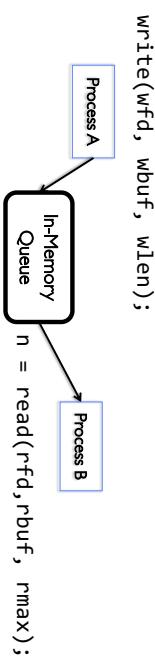


4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.48

Recall: Communication Between Processes



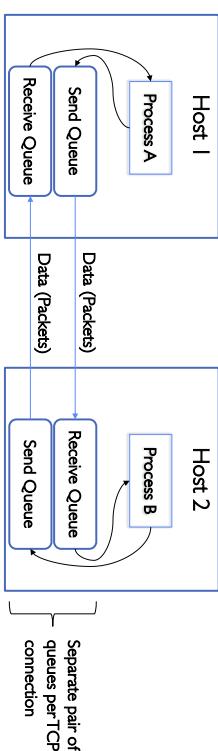
write(wfd, wbuf, wlen);

Process A → In-Memory Queue → Process B

$n = \text{read}(rfd, rbuf, rmax);$

- Data written by A is held in memory until B reads it
- Queue has a fixed capacity
 - Writing to the queue blocks if the queue is full
 - Reading from the queue blocks if the queue is empty
- POSIX provides this abstraction in the form of **pipes**

Buffering in a TCP Connection



Separate pair of queues per TCP connection

- A single TCP connection needs **four** in-memory queues:

- Send buffer: add data on write syscall, remove data when ACK received
- Receive buffer: add data when packets received, remove data on read syscall

4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

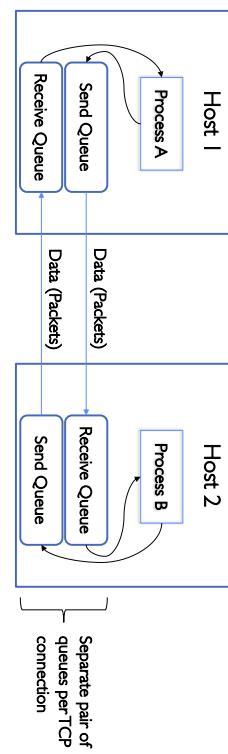
Lec 24.49

4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.50

Window Size: Space in Receive Queue



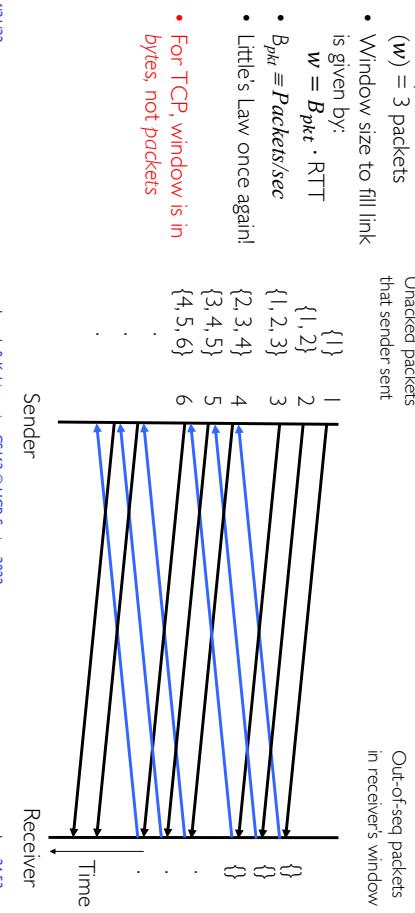
- A host's window size for a TCP connection is how much remaining space it has in its receive queue
- △ host advertises its window size in every TCP packet it sends!
- Sender never sends more than receiver's advertised window size**

4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.51

Sliding Window (No Packet Loss)



4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.53

Sliding Window Protocol

- TCP sender knows receiver's window size, and aims never to exceed it
- But packets that it previously send may arrive, filling the window size!

Rule: TCP sender ensures that:
Number of Sent but UnACKed Bytes < Receiver's Advertised Window Size

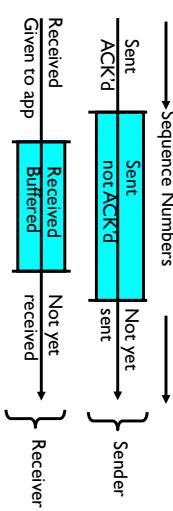
- Can send new packets as long as sent-but-unacked packets haven't already filled the advertised window size

4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.52

TCP Windows and Sequence Numbers: PER BYTE!



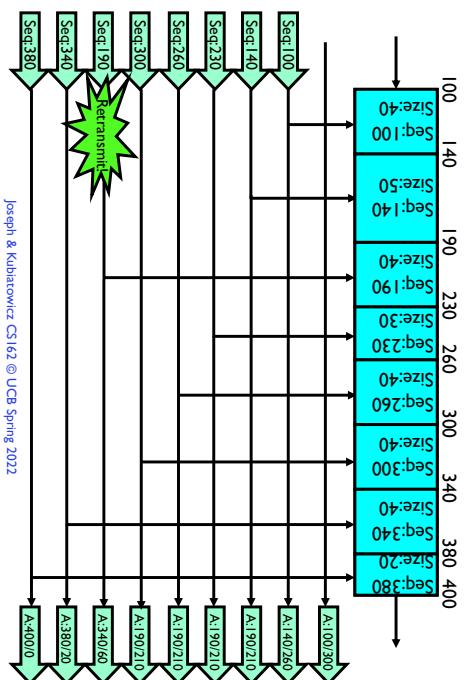
- Sender has three regions:
 - Sequence regions
 - » sent and ACK'd
 - » sent and not ACK'd
 - » not yet sent
- Receiver has three regions:
 - Sequence regions
 - » received and ACK'd (given to application)
 - » received and buffered
 - » received and not yet received (or discarded because out of order)

4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.54

Window-Based Acknowledgements (TCP)



Lec 24.55

4/21/22

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

Congestion Avoidance

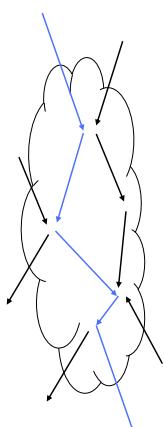
- Congestion
 - How long should timeout be for re-sending messages?
 - » Too long → wastes time if message lost
 - » Too short → retransmit even though ACK will arrive shortly
 - Stability problem: more congestion \Rightarrow ACK is delayed \Rightarrow unnecessary timeout \Rightarrow more traffic
 - More congestion \Rightarrow more retransmission
 - 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 window until acknowledgements start being delayed/lost
 - TCP solution: "slow start" (start sending slowly)
 - If no timeout, slowly increase window size (throughput) by 1 for each ACK received
 - Timeout \Rightarrow congestion, so cut window size in half
 - "Additive Increase, Multiplicative Decrease"

4/21/22

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

Congestion

- Too much data trying to flow through some part of the network
 - » Network diagram showing multiple paths from source to destination. One path is heavily congested, indicated by a large number of overlapping arrows.
- IP's solution: Drop packets
- What happens to TCP connection?
 - Lots of retransmission – wasted work and wasted bandwidth (when bandwidth is scarce)

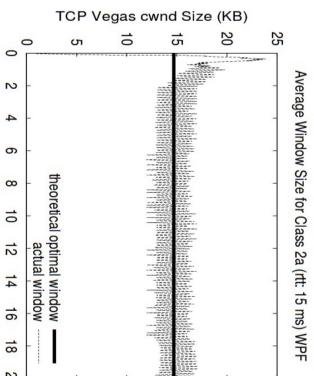


Lec 24.55

4/21/22

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

Congestion Management



From Low, Peterson, and Wang, "Understanding vegas: Duality Model", J. ACM, March 2002.

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

Lec 24.55

- TCP artificially restricts the window size if it sees packet loss

- Careful control loop to make sure:
 1. We don't send too fast and overwhelm the network
 2. We utilize most of the bandwidth the network has available
- In general, these are conflicting goals!

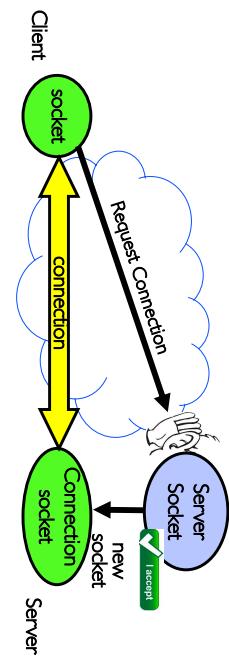
4/21/22

Lec 24.55

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

Lec 24.55

Recall: Connection Setup over TCP/IP



- 5-Tuple identifies each connection:

1. Source IP Address
2. Destination IP Address
3. Source Port Number
4. Destination Port Number
 - Protocol (always TCP here)
5. Well-known ports from 0—1023

4/21/22

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

Lec 24.59

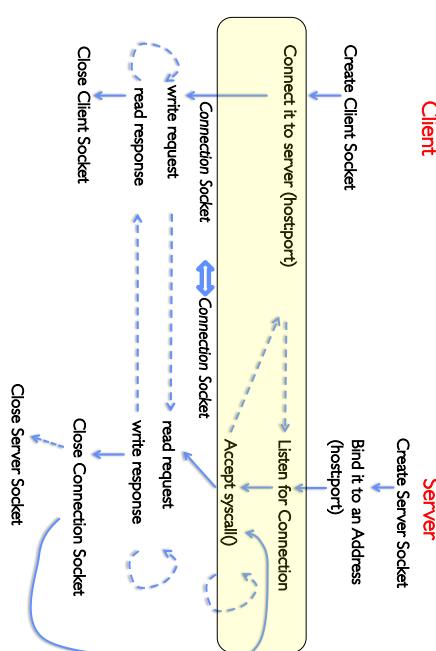
1. Open connection: 3-way handshaking
2. Reliable byte stream transfer from (IPa, TCP_Port1) to (IPb, TCP_Port2)
 - Indication if connection fails: Reset
3. Close (tear-down) connection

4/21/22

Joseph & Kubiatowicz CS 162 © UCB Spring 2022

Lec 24.60

Sockets in concept

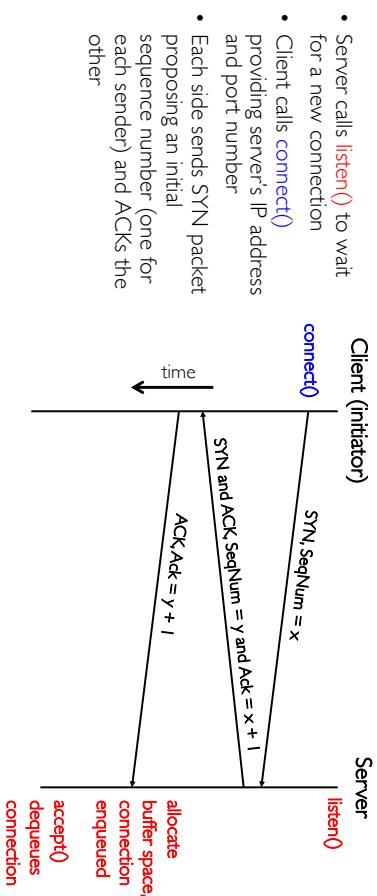


Lec 24.51

4/21/22

Establishing TCP Service

Open Connection: 3-Way Handshake

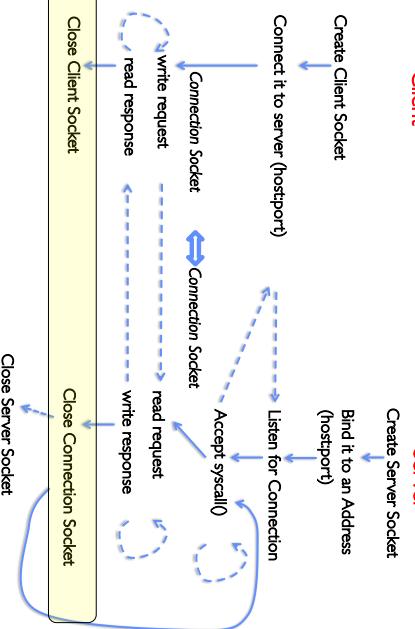


Joseph & Kubiatowicz CS 162 © UCB Spring 2022

Lec 24.62

4/21/22

Sockets in concept

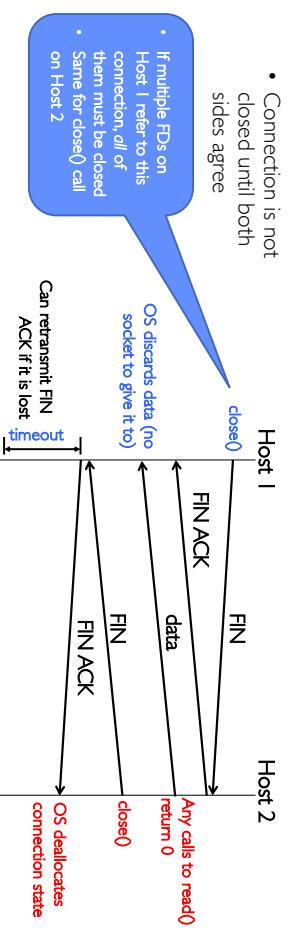


4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.63

Close Connection: 4-Way Teardown



4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.64

Recall: Distributed Applications Build With Messages

- How do you actually program a distributed application?
 - Need to synchronize multiple threads, running on different machines
 - » No shared memory, so cannot use test&set
-
- The diagram shows a central 'Network' cloud connected to two computer icons labeled 'Host 1' and 'Host 2'. Arrows indicate data flow from Host 1 to Host 2 and vice versa through the network.
- One Abstraction: send/receive messages
 - » Already atomic: no receiver gets portion of a message and two receivers cannot get same message
 - Interface:
 - Mailbox (`inbox`) 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

4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.65

Question: Data Representation

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

4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.66

Simple Data Types

- Suppose I want to write a **x** to a file
 - First, open the file: **FILE* f = fopen("foo.txt", "w");**
 - Then, I have two choices:
 1. **fprintf(f, "%lu", x);**
 2. **fwrite(&x, sizeof(uint32_t), 1, f);**
 - » Or equivalently, **write(fd, &x, sizeof(uint32_t));** (perhaps with a loop to be safe)
- Neither one is “wrong” but sender and receiver should be consistent!

4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.67

- For a byte-address machine, which end of a machine-recognized object (e.g., int) does its byte-address refer to?
- Big Endian: address is the most-significant bits
- Little Endian: address is the least-significant bits

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

```
int main(int argc, char *argv[])
{
    (base). CulterMac19:code00 culles .endian
    val = 12345678
    int val = 0x12345678;
    int i;
    printf("val = %X\n", val);
    for (i = 0; i < sizeof(val); i++) {
        printf("val[%d] = %X\n", i, ((uint8_t *) &val)[i]);
    }
}
```

4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.69

Endianness

Machine Representation

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

4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.68

What Endian is the Internet?

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

NAME	SYNOPSIS	DESCRIPTION
sys/types.h	#include <sys/types.h>	The <code>in_port_t</code> and <code>in_addr_t</code> types shall be defined as described in <code><sys/types.h></code> . The <code>in_addr</code> structure shall be defined as described in <code><sys/types.h></code> . The <code>NETINET_PTON</code> and <code>NETINET_ATON</code> macros shall be defined as described in <code><sys/types.h></code> . The following shall either be declared as functions, defined as macros, or both. If functions are declared, function prototypes shall be provided: <code>uint16_t htons(uint16_t);</code> <code>uint16_t ntohs(uint16_t);</code> <code>uint32_t htonl(uint32_t);</code> <code>uint32_t ntohl(uint32_t);</code> <code>uint64_t htonl64(uint64_t);</code> <code>uint64_t ntohl64(uint64_t);</code>

The `uint32_t` and `int16_t` types shall be defined as described in `<sys/types.h>`.
The following shall be declared as functions and may also be defined as macros. Function prototypes shall be provided:
`int addrlen_t *inet_addr(const char *);`
`int addrlen_t *inet_ntoa(struct sockaddr_in const *restrict, char *restrict);`
`*inet_ntop(int family, const void *restrict, char *restrict, size_t n);`
`int inet_pton(int family, const char *restrict, void *restrict);`

Inclusion of the `<sys/types.h>` header may also make visible all symbols from `<sys/types.h>` and `<sys/socket.h>`.

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.70

Remote Procedure Call (RPC)

- Raw messaging is a bit too low-level for programming
 - Must wrap up information into message at source
 - Must decide what to do with message at destination
 - May need to sit and wait for multiple messages to arrive
 - **And must deal with machine representation by hand**

- Another option: Remote Procedure Call (RPC)
 - Calls a procedure on a remote machine
 - Idea: Make communication look like an ordinary function call
 - Automate all of the complexity of translating between representations
 - Client calls:
`r = f(v1, v2);`
 - Translated automatically into call on server:
`fileSys->Read("rutabaga");`

4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

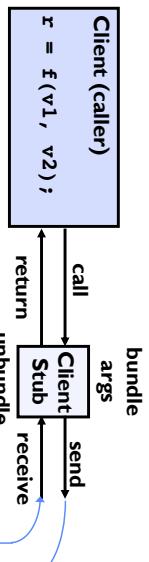
Lec 24.75

4/21/22

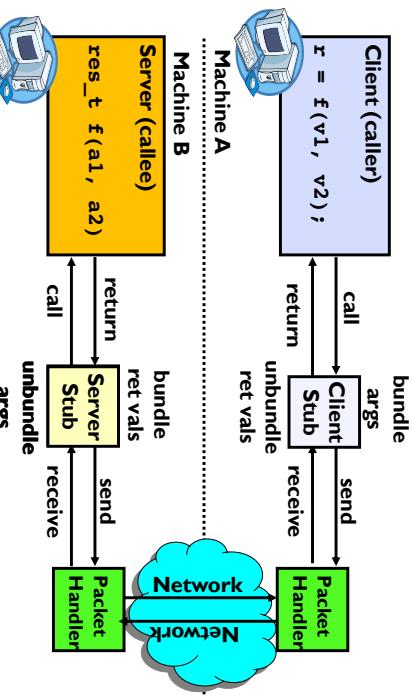
Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.76

RPC Concept



RPC Information Flow



RPC Implementation

- Request-response message passing (under covers!)

- "Stub" provides glue on client/server
 - Client stub is responsible for "marshalling" arguments and "unmarshalling" the return values
 - Server-side stub is responsible for "unmarshalling" arguments and "marshalling" the return values

- **Marshalling** involves (depending on system)
 - Converting values to a canonical form, serializing objects, copying arguments passed by reference, etc.

4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.77

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.78

RPC Details (1/3)

- Equivalence with regular procedure call
 - Parameters \Leftrightarrow Request Message
 - Result \Leftrightarrow Reply message
 - Name of Procedure: Passed in request message
 - Return Address: mbox2 (client return mail box)
- Stub generator: Compiler that generates stubs
 - Input: interface definitions in an “interface definition language (IDL)”
 - » Contains, among other things, types of arguments/return
 - Output: stub code in the appropriate source language
 - » Code for client to pack message, send it off, wait for result, unpack result and return to caller
 - » Code for server to unpack message, call procedure, pack results, send them off

4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

RPC Details (2/3)

- Cross-platform issues:
 - What if client/server machines are different architectures/ languages?
 - » Convert everything to/from some canonical form
 - » Tag every item with an indication of how it is encoded (avoids unnecessary conversions)
- How does client know which mbox (destination queue) to send to?
 - Need to translate name of remote service into network endpoint (Remote machine, port, possibly other info)
 - **Binding**: the process of converting a user-visible name into a network endpoint
 - » This is another word for ‘naming’ at network level
 - » Static: fixed at compile time
 - » Dynamic: performed at runtime

Lec 24.79

4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

RPC Details (3/3)

- Dynamic Binding
 - Most RPC systems use dynamic binding via name service
 - » Name service provides dynamic translation of service \rightarrow mbox
 - Why dynamic binding?
 - » Access control: check who is permitted to access service
 - » Fail-over: If server fails, use a different one
- What if there are multiple servers?
 - Could give flexibility at binding time
 - » Choose unloaded server for each new client
 - Could provide same mbox (router level redirect)
 - » Choose unloaded server for each new request
 - » Only works if no state carried from one call to next
- What if multiple clients?
 - Pass pointer to client-specific return mbox in request

Lec 24.81

4/21/22

Problems with RPC: Non-Atomic Failures

- Different failure modes in dist. system than on a single machine
- Consider many different types of failures
 - User-level bug causes address space to crash
 - Machine failure, kernel bug causes all processes on same machine to fail
 - Some machine is compromised by malicious party
- Before RPC: whole system would crash/die
- After RPC: One machine crashes/compromised while others keep working
- Can easily result in inconsistent view of the world
 - Did my cached data get written back or not?
 - Did server do what I requested or not?
- Answer? Distributed transactions/Byzantine Commit

Lec 24.82

Joseph & Kubiatowicz CS162 © UCB Spring 2022

4/21/22

Problems with RPC: Performance

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

4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.83

Cross-Domain Communication/Location Transparency

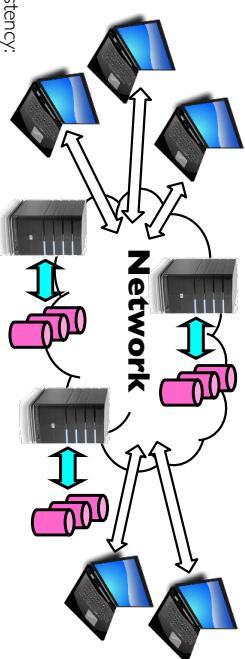
- How do address spaces communicate with one another?
 - Shared Memory with Semaphores, monitors, etc...
 - File System
 - Pipes (1-way communication)
 - “Remote” procedure call (2-way communication)
 - RPC’s can be used to communicate between address spaces on different machines or the same machine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

4/21/22

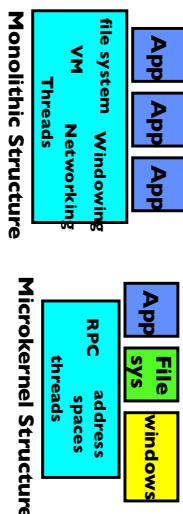
Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.84

Network-Attached Storage and the CAP Theorem



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



- Why split the OS into separate domains?

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

4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.85

4/21/22

Joseph & Kubiatowicz CS162 © UCB Spring 2022

Lec 24.86

Summary

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