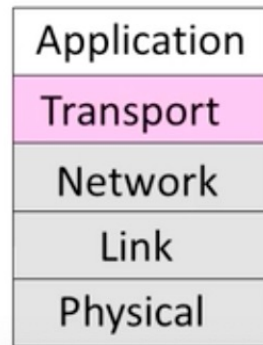# Computer Network Design
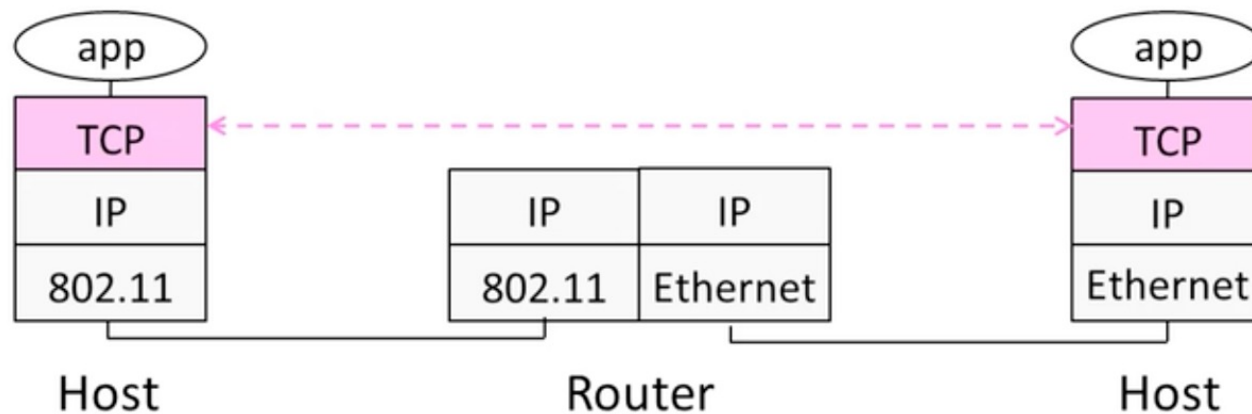## Transport Layer I

Yalda Edalat – Spring 23

# Where We are in the Course

- Starting the Transport layer
  - Builds on the network layer to deliver data across networks for applications with the desired reliability or quality
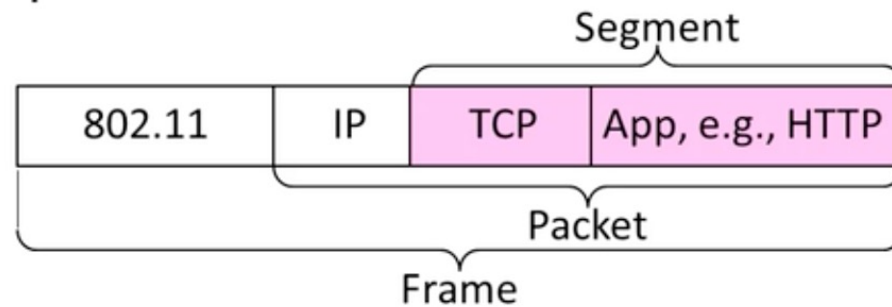
| Application |
|:---:|
| Transport |
| Network |
| Link |
| Physical |

# Recall

- Transport layer provides end-to-end connectivity across the network

# Recall (2)

- Segments carry application data across the network
- Segments are carried within packets within frames

# Transport Layer Services

- Provide different kinds of data delivery across the network to applications

|  | Unreliable | Reliable |
|---|---|---|
| Messages | Datagrams (UDP) | |
| Bytestream | | Streams (TCP) |

# Comparison of Internet Transports

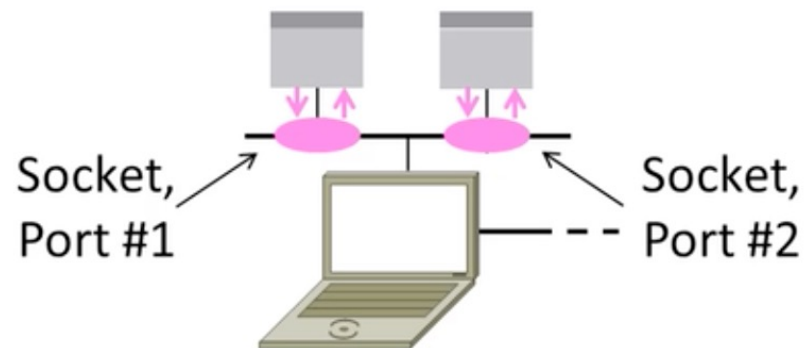- TCP is full-featured, UDP is a glorified packet

| TCP (Streams) | UDP (Datagrams) |
|---|---|
| Connections | Datagrams |
| Bytes are delivered once, reliably, and in order | Messages may be lost, reordered, duplicated |
| Arbitrary length content | Limited message size |
| Flow control matches sender to receiver | Can send regardless of receiver state |
| Congestion control matches sender to network | Can send regardless of network state |

# Socket API

- Simple abstraction to use the network
  - The "network" API (really transport service) used to write all Internet apps
  - Part of all major Oses and languages;  originally Berkeley (Unix)

- Supports both Internet transport services (Stream and Datagrams)

# Socket API (2)

- Sockets let apps attach to the local network at different ports

# Socket API (3)

- Some API used for Streams and Datagrams

| Primitive | Meaning |
|---|---|
| SOCKET | Create a new communication endpoint |
| BIND | Associate a local address (port) with a socket |
| LISTEN | Announce willingness to accept connections |
| ACCEPT | Passively establish an incoming connection |
| CONNECT | Actively attempt to establish a connection |
| SEND(TO) | Send some data over the socket |
| RECEIVE(FROM) | Receive some data over the socket |
| CLOSE | Release the socket |

Only needed for Streams { LISTEN, ACCEPT, CONNECT

To/From forms for Datagrams { SEND(TO), RECEIVE(FROM)

# Ports

- Application process is identified by the tuple IP address, protocol and port
  - Ports are 16-bit integers representing local "mailboxes" that a process leases

- Servers often bind to "well-known ports"
  - <1024, require administrative privileges
- Clients often assigned "ephemeral" ports
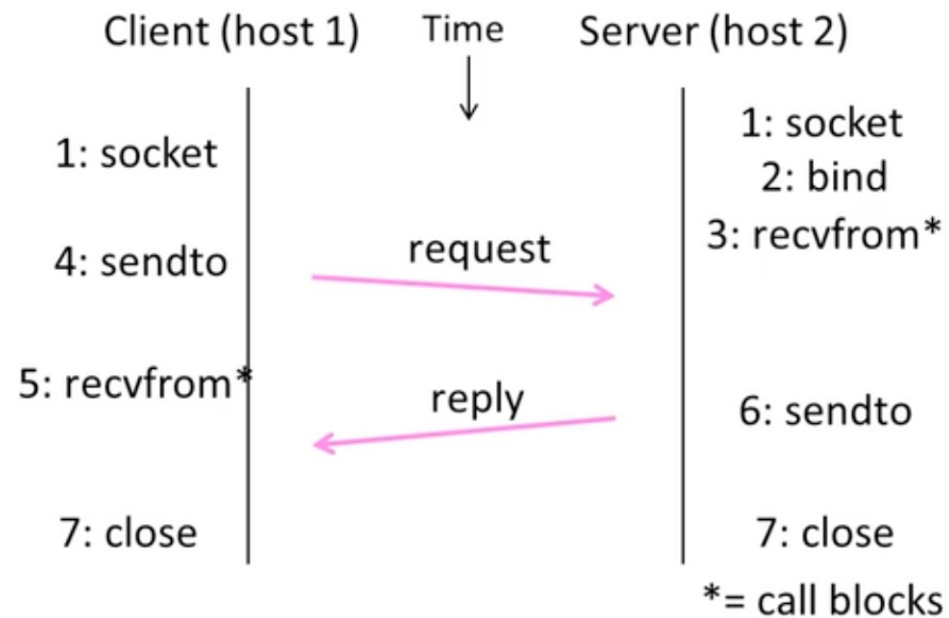  - Chosen by OS, used temporarily

# Some Well-Known Ports

| Port | Protocol | Use |
|---|---|---|
| 20, 21 | FTP | File transfer |
| 22 | SSH | Remote login, replacement for Telnet |
| 25 | SMTP | Email |
| 80 | HTTP | World Wide Web |
| 110 | POP-3 | Remote email access |
| 143 | IMAP | Remote email access |
| 443 | HTTPS | Secure Web (HTTP over SSL/TLS) |
| 543 | RTSP | Media player control |
| 631 | IPP | Printer sharing |

# User datagram Protocol (UDP)

- Used by apps that don't want reliability or bytestreams
  - Voice-over-IP (unreliable)
  - DNS, RPC (message-oriented)
  - DHCP (bootstrapping)

- If application wants reliability and messages then it has work to do!
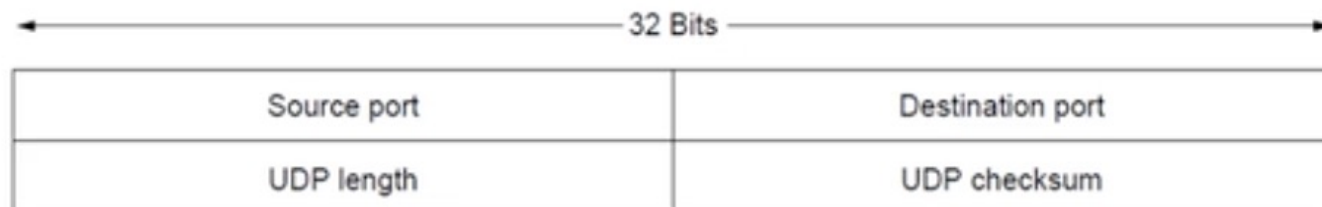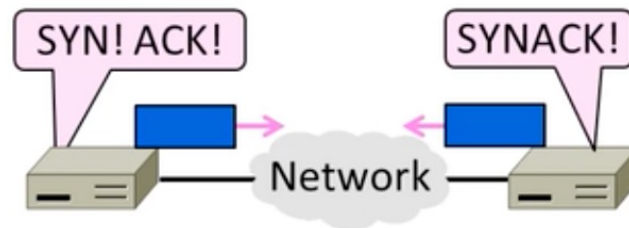
# Datagram Sockets

# UDP Buffering

# UDP Header

- Uses ports to identify sending and receiving application processes
- Datagram length up to 64K
- Checksum (16 bits) for reliability

| 32 Bits | |
|---|---|
| Source port | Destination port |
| UDP length | UDP checksum |

# Connection Establishment

- TCP implement a connection oriented stream service

- How to set up connections?
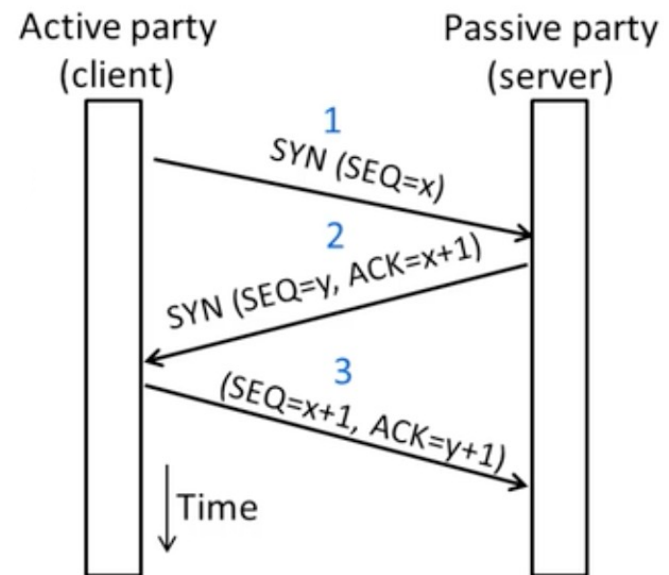  - We will see how TCP does

# Connection Establishment (2)

- Both sender and receiver must be ready before we start the transfer of data
  - Need to agree on a set of parameters
  - E.g., the Maximum Segment Size (MSS)

- This is signaling
  - It sets up state at the endpoints
  - Like "dialing" for a telephone call

# Three-Way Handshake

- Used in TCP; opens connection for data in both directions
- Each side probes the other with a fresh Initial Sequence Number (ISN)
  - Sends on a SYNchronize segment
  - Echo on an ACKnowledge segment

# Three-Way Handshake (2)

- Three steps:
  - Client sends SYN (x)
  - Server replies with SYN(y)ACK(x+1)
  - Client replies with ACK(y+1)
  - SYNs are retransmitted if lost

- Sequence and ack numbers carried on further segments

# Connection Release

- Orderly release by both parties when done
  - Delivers all pending data and "hangs up"


- Key problem is to provide reliability while releasing
  - TCP uses a "symmetric" close in which both sides shutdown independently

# TCP Connection Release

- Two steps:
  - Active sends FIN(x), passive ACKs
  - Passive sends FIN(y), active ACKs
  - FINs are retransmitted if lost

- Each FIN/ACK closes one direction of data transfer

Active party          Passive party

FIN (SEQ=x)
1
(SEQ=y, ACK=x+1)
FIN (SEQ=y, ACK=x+1)
2
(SEQ=x+1, ACK=y+1)

# Recall

- Stop-and-wait

# Limitation of Stop-and-Wait

- It allows only a single message to be outstanding from the sender:
  - Fine for LAN (only one frame fit)
  - Not efficient for network paths with BD>>1 packet

# Limitation of Stop-and-Wait (2)

- Example: R = 1 Mbps, D = 50 ms, packet size = 1250 bytes
  - RTT (Round Trip Time) = 2D = 100 ms
  - How many packets/sec?



  - What if R = 10 Mbps?

# Sliding Window

- Generalization of stop-and-wait
  - Allow W packets to be outstanding
  - Can send W packets per RTT (= 2D)

  - Pipelining improves performance
  - Need W = 2BD to fill network path

# Sliding Window (2)

- What W will use the network capacity?
- Ex: R = 1 Mbps, D = 50 ms




- Ex: What if R = 10 Mbps?

# Sliding Window (3)

- Ex: R = 1 Mbps, D = 50 ms
  - 2BD = $10^6$ b/sec x $100.10^{-3}$ sec = 100 kbit
  - W = 2BD = 10 packets of 1250 bytes


- Ex: What if R = 10 Mbps?
  - 2BD = 1000 kbit
  - W = 2BD = 100 packets of 1250 bytes

# Sliding Window Protocol

- Many variations, depending on how buffers, acknowledgements, and retransmissions are handled


- Go-Back-N
  - Simplest version, can be inefficient
- Selective Repeat
  - More complex, better performance

# Sliding Window - Sender

- Sender buffers up to W segments until they are acknowledged
  - LFS = LAST FRAME SENT
  - LAR = LAST ACK RECEIVED
  - Sends while LFS-LAR <= W

# Sliding Window – Sender (2)

- Transport accepts another segment of data from the application …
  - Transport sends it (as LFS-LAR -> 5)

# Sliding Window – Sender (3)

- Next higher ACK arrives from peer…
  - Window advances, buffer is feed
  - LFS-LAR -> 4 (can send one more)

# Sliding Window – Go-Back-N

- Receiver keeps only a single packet buffer for the next segment
  - State variable, LAS = LAST ACK SENT

- On receive:
  - If seq. number is LAS+1, accept and pass it to app, update LAS, send ACK
  - Otherwise discard (as out of order)

# Sliding Window – Selective Repeat

- Receiver passes data to app in order, and buffers out-of-order segments to reduce retransmissions

- ACK conveys highest in-order segment, plus hints about out-of-order segments

- TCP uses a selective repeat design; we will see the details later

# Sliding Window – Selective Repeat (2)

- Buffers W segments, keep state variable, LAS = LAST ACK SENT

- On receive:
  - Buffer segments [LAS+1, LAS+W]
  - Pass up to app in-order segments from LAS+1, and update LAS
  - Send ACK for LAS

# Sliding Window - Retransmission

- Go-Back-N sender uses single timer to detect losses
  - On timeout, resends buffered packets starting at LAR+1

- Selective repeat sender uses a timer per unacked segment to detect losses
  - On timeout for segment, resend it
  - Hope to resend fewer segments

# Sequence Number

- Need more than 0/1 for Stop-and-Wait …
    - But how many?

- For selective repeat, need W numbers for packets, plus W for acks of earlier packets
    - 2W sequence numbers
    - Fewer for Go-Back-N (W+1)

- Typically implement seq. number with an N-bit counter that wraps around at $2^N-1$
    - E.g., N = 8:  …, 253, 254, 255, 0, 1, 2, …

# Sequence Time Plot

# Sequence Time Plot (2)



Go-Back-N scenario

# Sequence Time Plot (3)

# Problem

- Sliding window uses pipelining to keep the network busy
  - What if the receiver is overloaded?



- Solution: Adding **flow control** to the sliding window algorithm
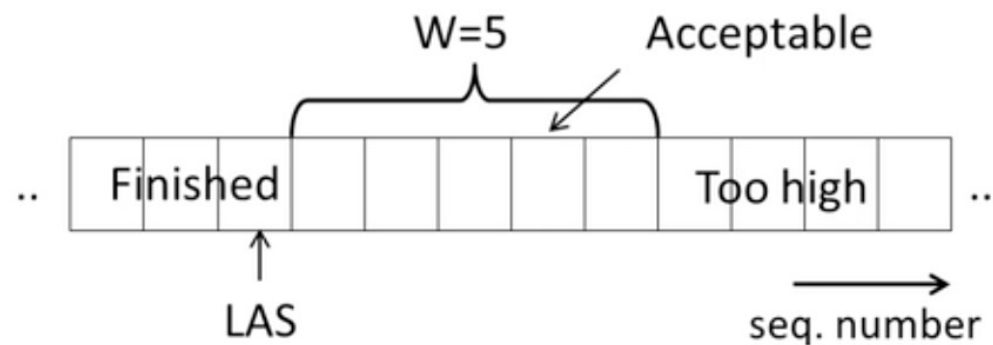  - To slow the over-enthusiastic sender

# Sliding Window - Receiver

- Consider receiver with W buffers
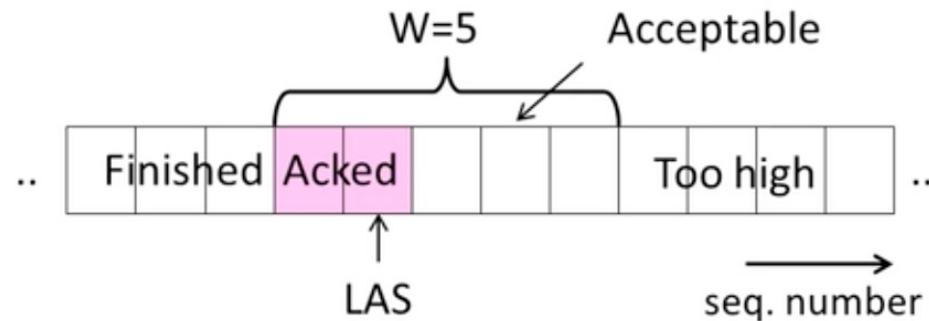  - LAS = LAST ACK SENT, app pulls in-order data from buffer with recv() call

# Sliding Window – Receiver (2)
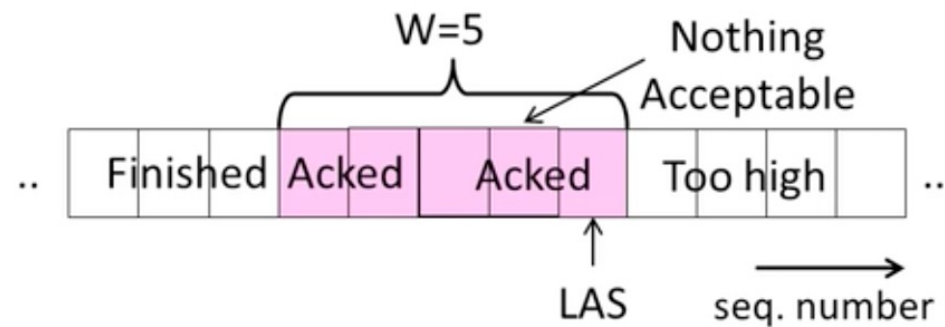
- Suppose the next two segments arrive but app does not call recv()

# Sliding Window – Receiver (3)

- Suppose the next two segments arrive but app does not call recv()
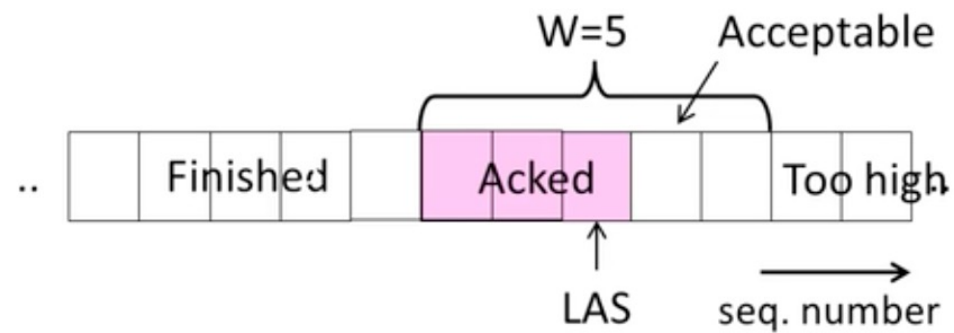  - LAS rises, but we can't slide window!

# Sliding Window – Receiver (4)

- If further segments arrive (even in order) we can fill the buffer
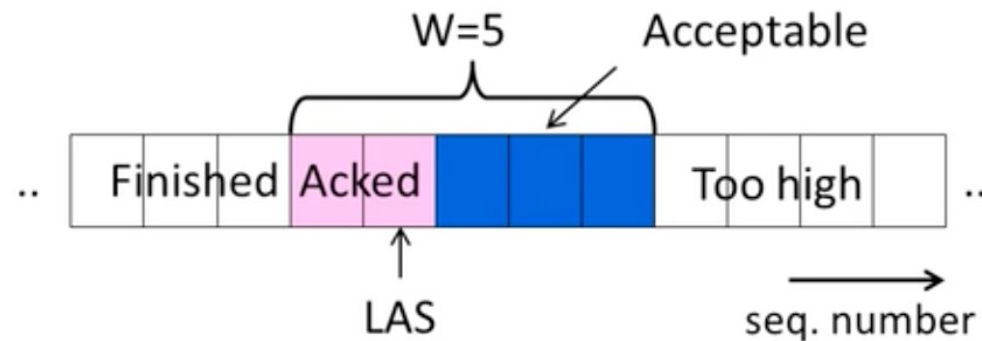  - Must drop segments until app receives!

# Sliding Window – Receiver (5)
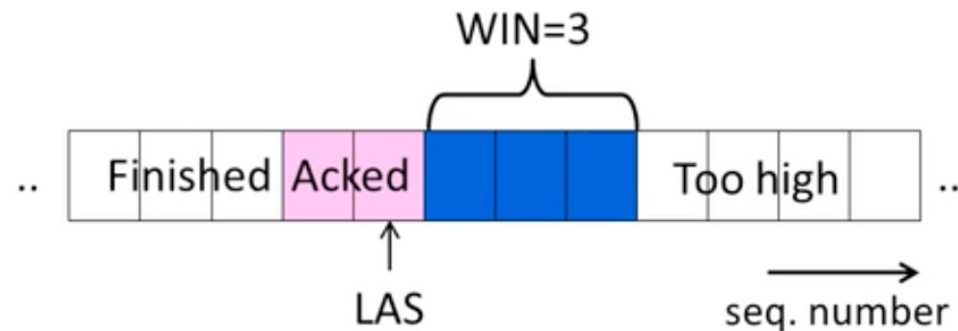
- App recv() takes two segments

# Flow Control

- Avoid loss at receiver by telling sender the available buffer space
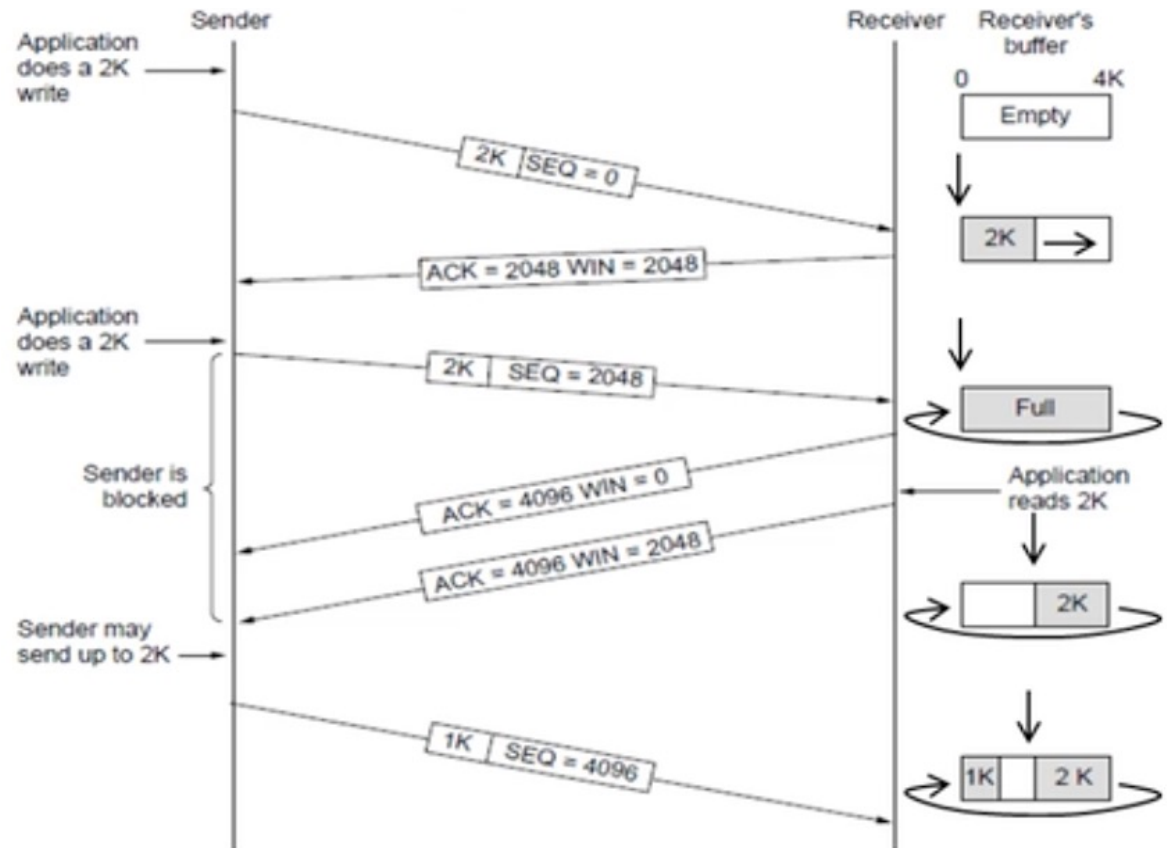  - WIN = #acceptable, not W (from LAS)

# Flow Control (2)

- Sender uses the lower of the sliding window and flow control window (WIN) as the effective window size
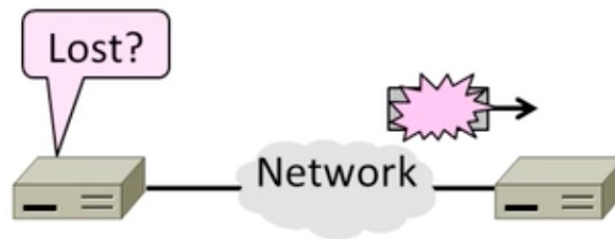
# Flow Control (3)

- TCP-style example
  - SEQ/ACK sliding window
  - Flow control with WIN
  - SEQ + length < ACK + WIN
  - 4KB buffer at receiver
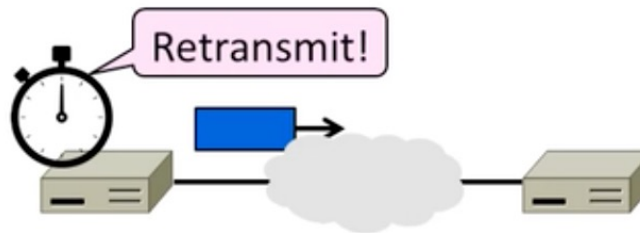  - Circular buffer of bytes

# Retransmission Timer

- How to set the timeout for sending a retransmission?
  - Adapting to the network path
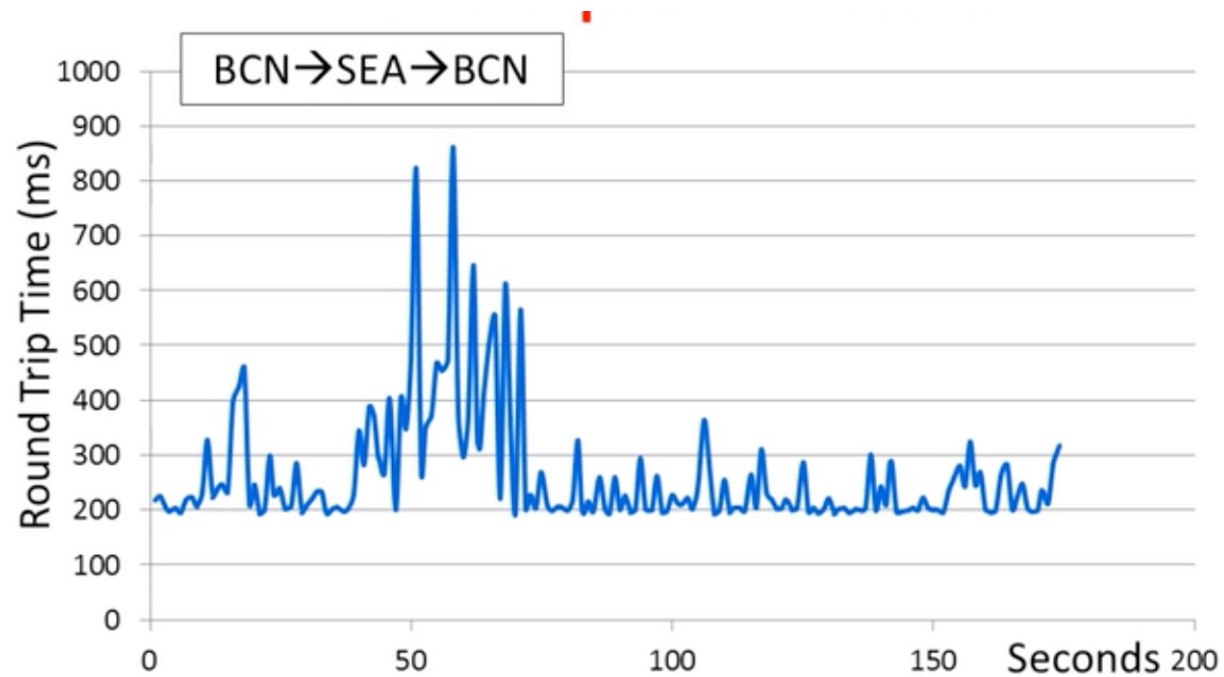
# Retransmissions

- With sliding window, the strategy for detecting loss is the timeout
    - Set timer when a segment is sent
    - Cancel timer when ack is received
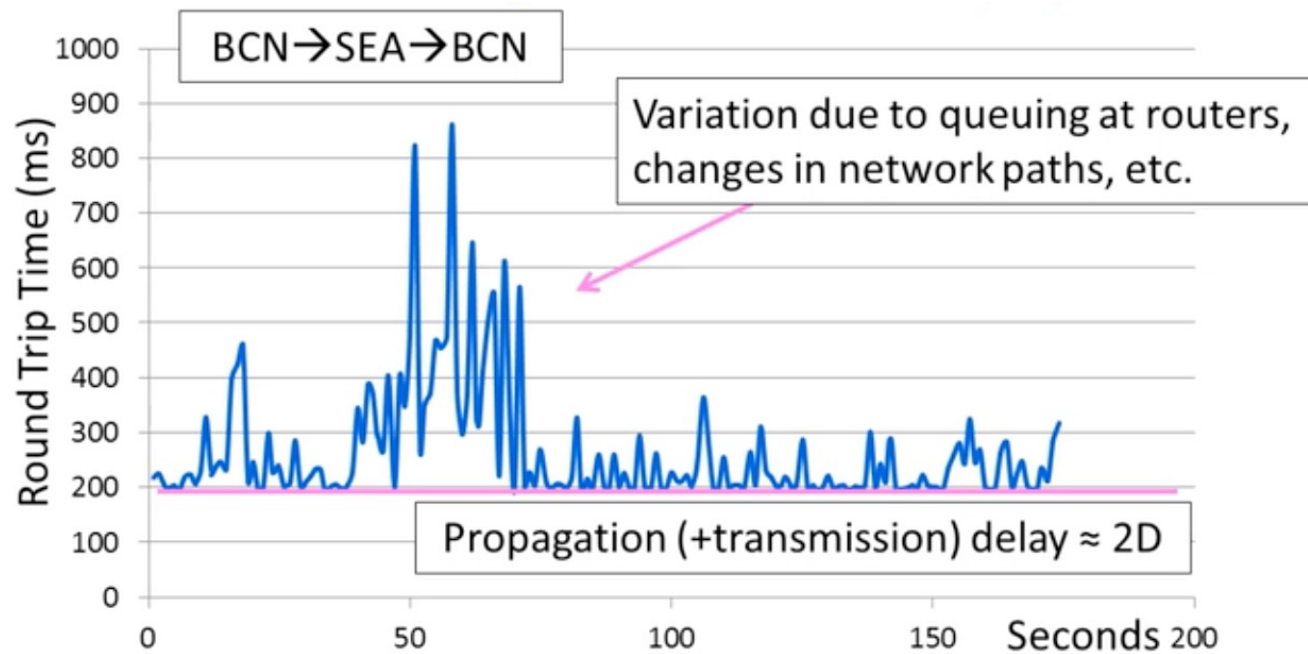    - If timer fires, retransmission data as lost

# Timeout Problem

- Timeout should be "just right"
  - Too long wastes network capacity
  - Too short leads to spurious resends
  - But what is "just right"?

- Easy to set on a LAN (link)
  - Short, fixed, predictable RTT

- Hard on the Internet (Transport)
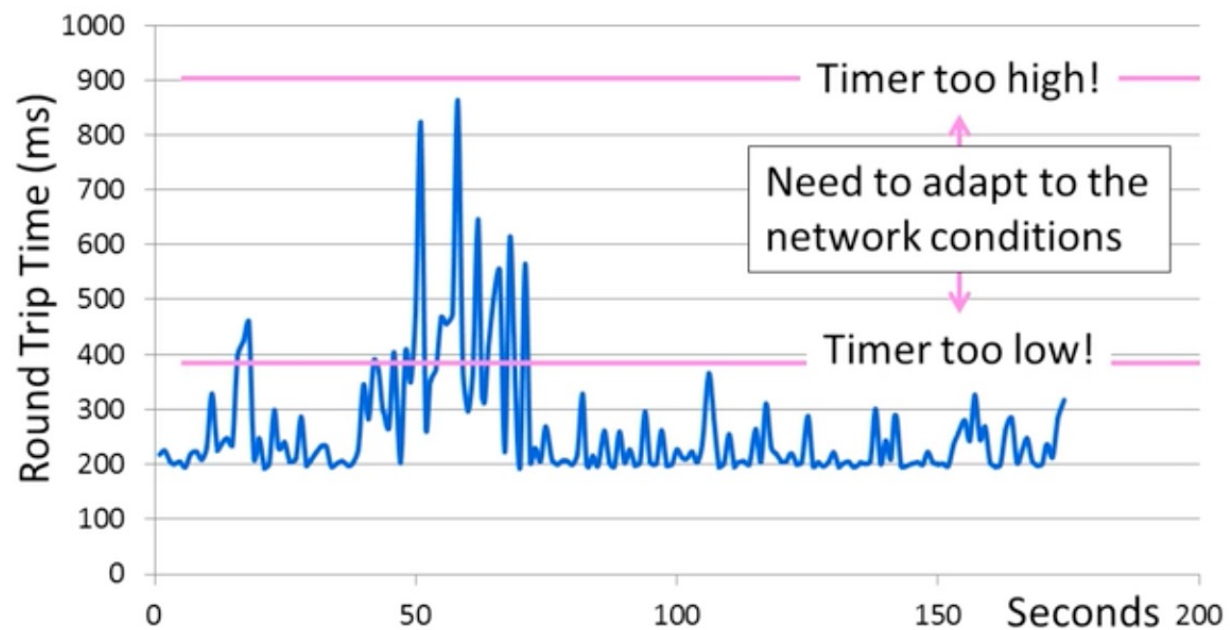  - Wide range, variable RTT
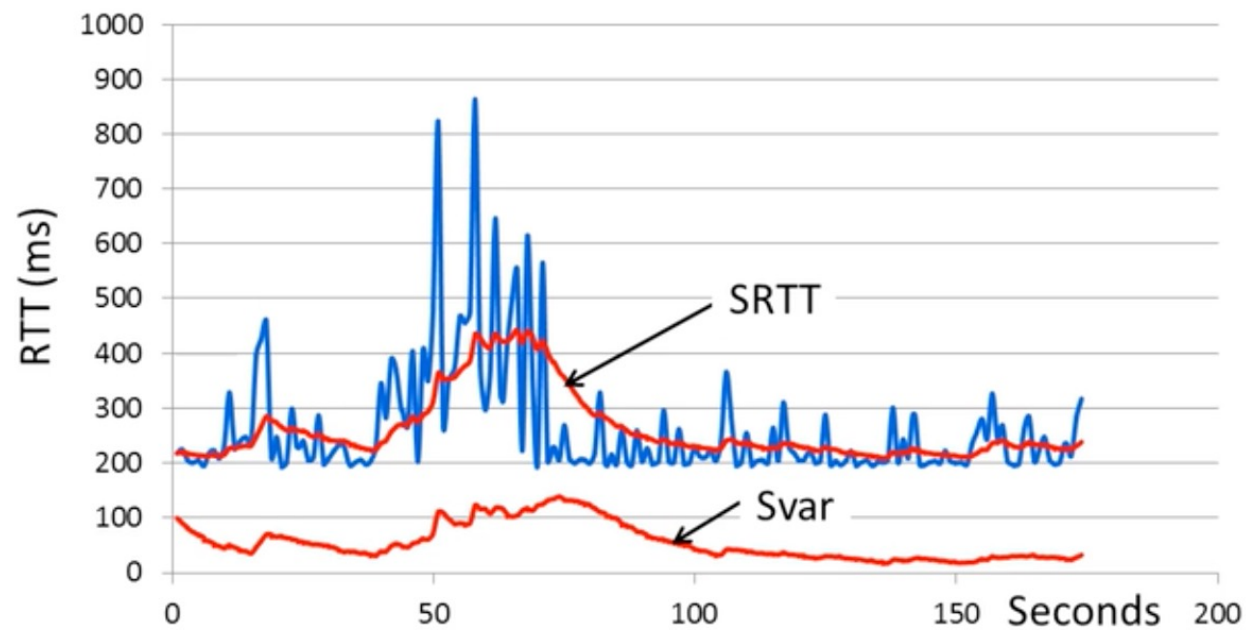
# Example of RTTs
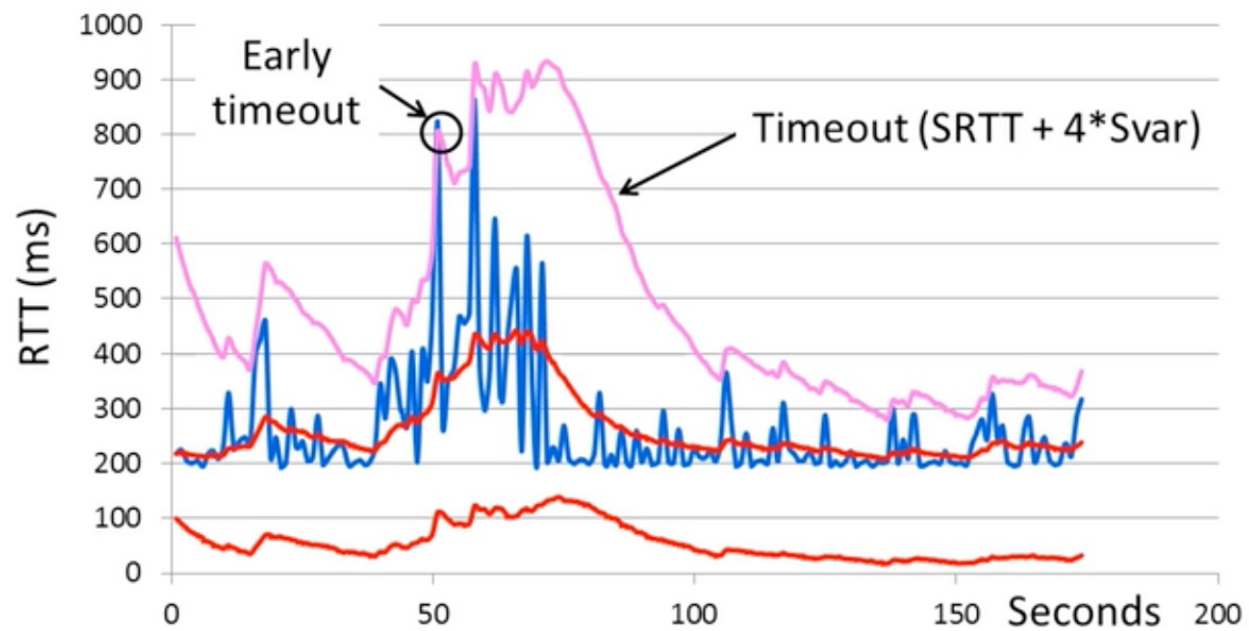
# Example of RTTs (2)

# Example of RTTs (3)

# Adaptive Timeout

- Keep smoothed estimates of the RTT (1) and variance in RTT (2)
  - Update estimates with a moving average
  1. $SRTT_{N+1} = 0.9*SRTT_N + 0.1*RTT_{N+1}$
  2. $Svar_{N+1} = 0.9*Svar_N + 0.1*|RTT_{N+1}-SRTT_{N+1}|$

- Set timeout to a multiple of estimates
  - To estimate the upper RTT in practice
  - TCP timeout$_N$ = $SRTT_N+4*Svar_N$

# Example of Adaptive Timeout
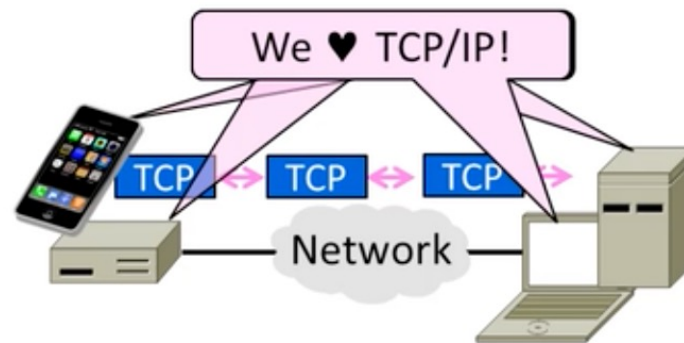
# Example of Adaptive Timeout (2)

# Adaptive Timeout (2)

- Simple to compute, does a good job of tracking actual RTT
  - Little "headroom" to lower
  - Yet very few early timeouts

- Turns out to be important for good performance and robustness

# Transmission Control Protocol (TCP)

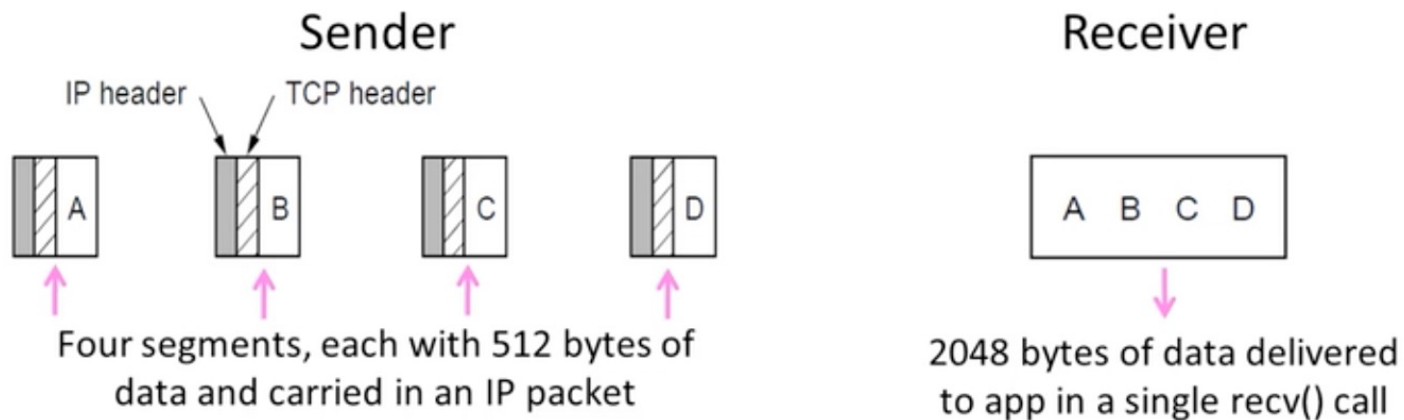• The transport protocol used for most content on the Internet

# TCP Features

- A reliable bytestream service

- Based on connections

- Sliding window for reliability
  - With adaptive timeout

- Flow control for slow receivers

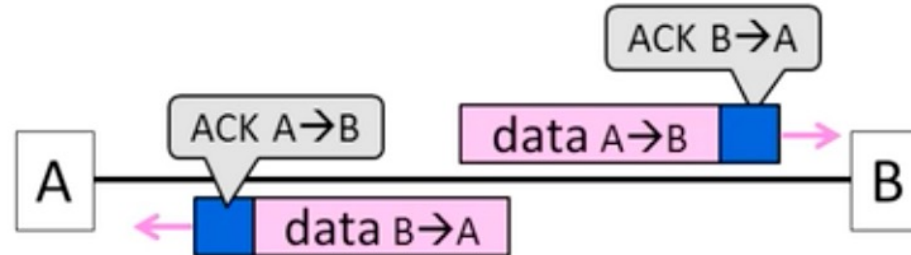- Congestion control to allocate network bandwidth

# Reliable Bytestream

- Message boundaries not preserved from send() to recv()
  - But reliable and ordered (receive bytes in same order as sent)

Sender

IP header    TCP header

A    B    C    D

Four segments, each with 512 bytes of
data and carried in an IP packet

Receiver

A  B  C  D

2048 bytes of data delivered
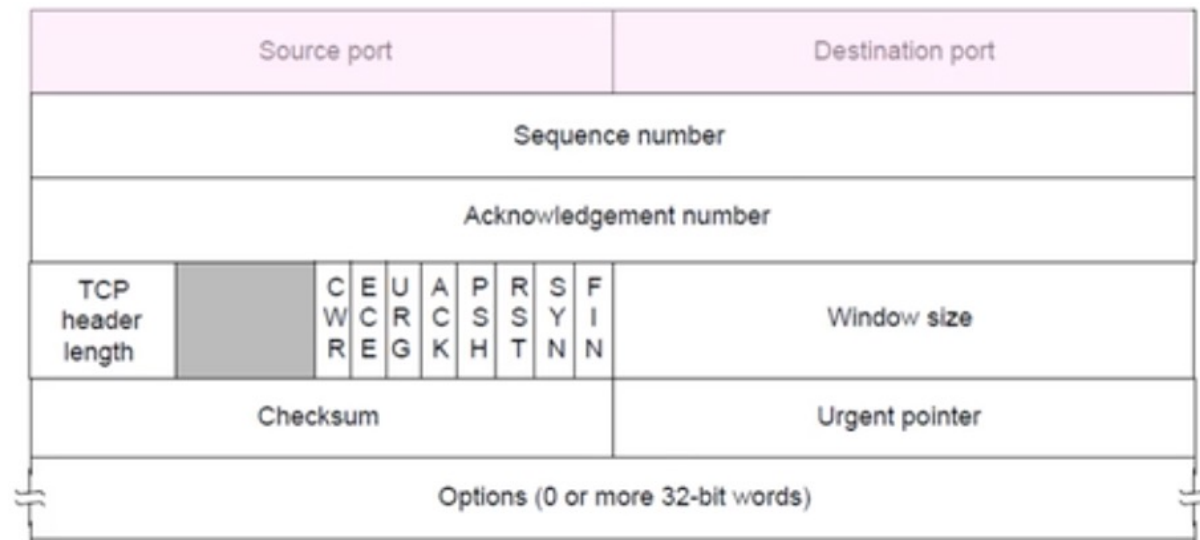to app in a single recv() call

# Reliable Bytestream (2)

- Bidirectional data transfer
  - Control information (e.g., ACK) piggybacks on data segments in reverse direction
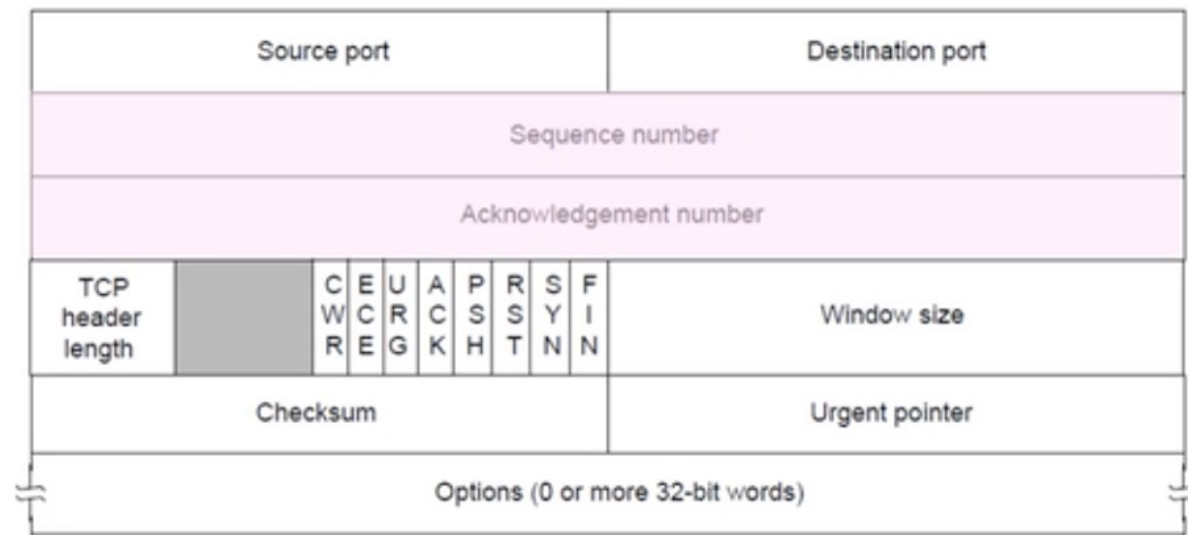
# TCP Header

- Ports identify apps (socket API)
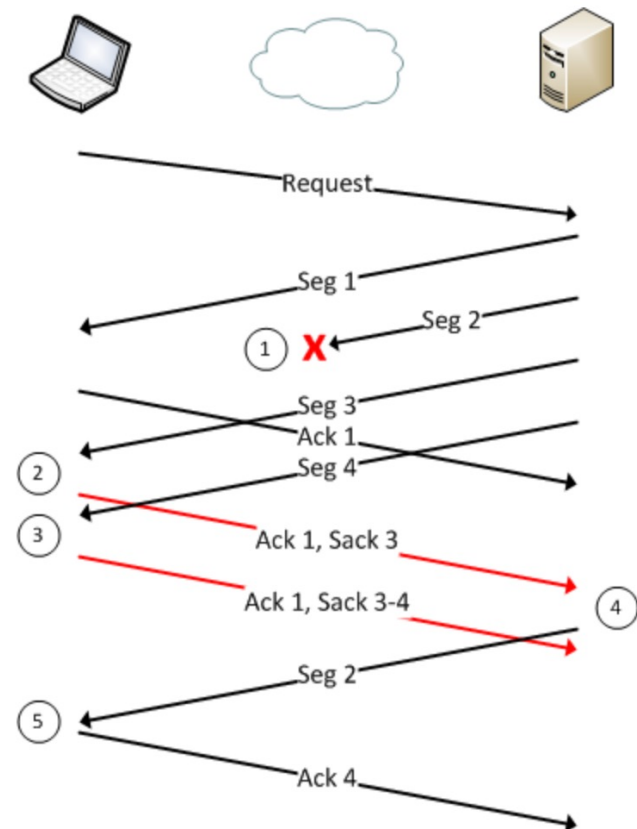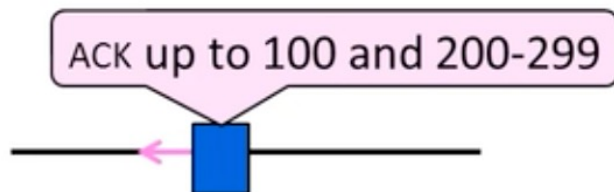  - 16-bit identifiers

# TCP Headers (2)

- SEQ/ACK used for sliding window
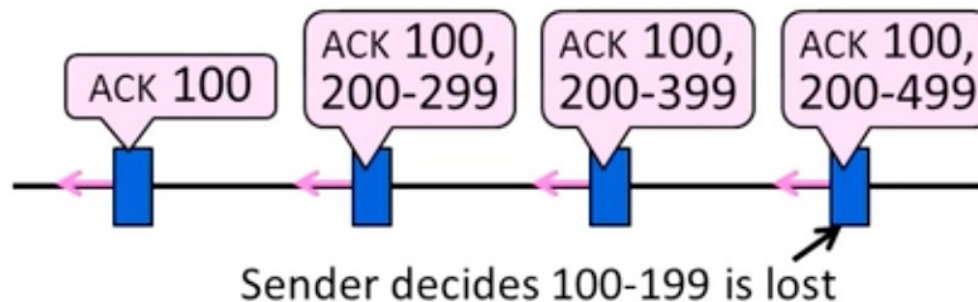  - Selective repeat, with byte positions

# TCP Sliding Window - Receiver

- Cumulative ACK tells next expected byte sequence number ("LAS+1")

- Optionally, selective ACKs (SACK) give hints for receiver buffer state
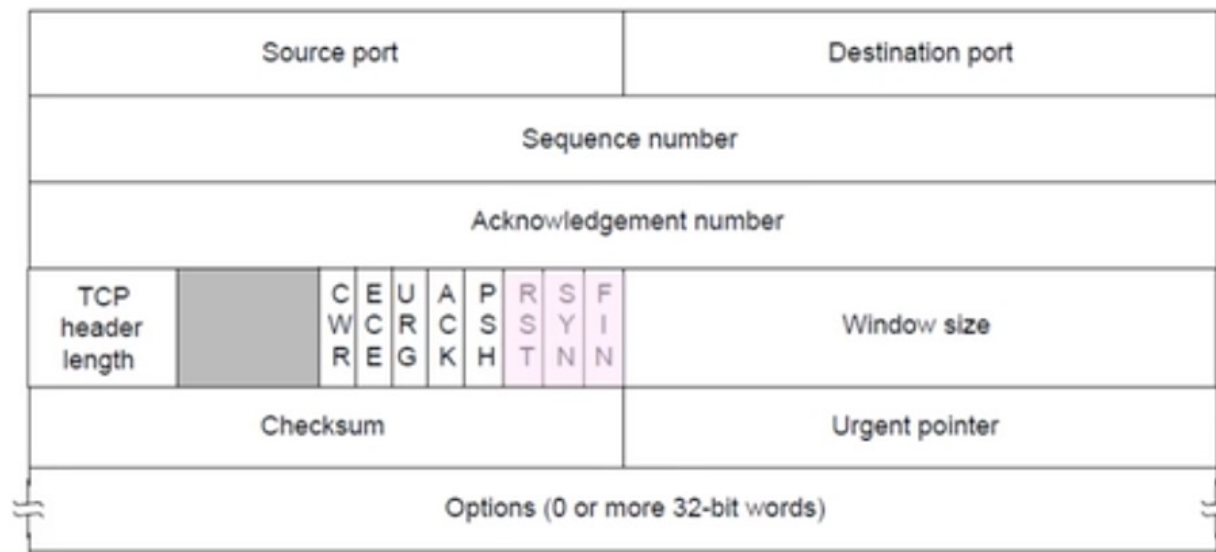  - List up to 3 ranges of received bytes

ACK up to 100 and 200-299

# TCP Sliding Window - Sender

- Uses an adaptive retransmission timeout to resend data from LAS+1
- Uses heuristics to infer loss quickly and resend to avoid timeouts
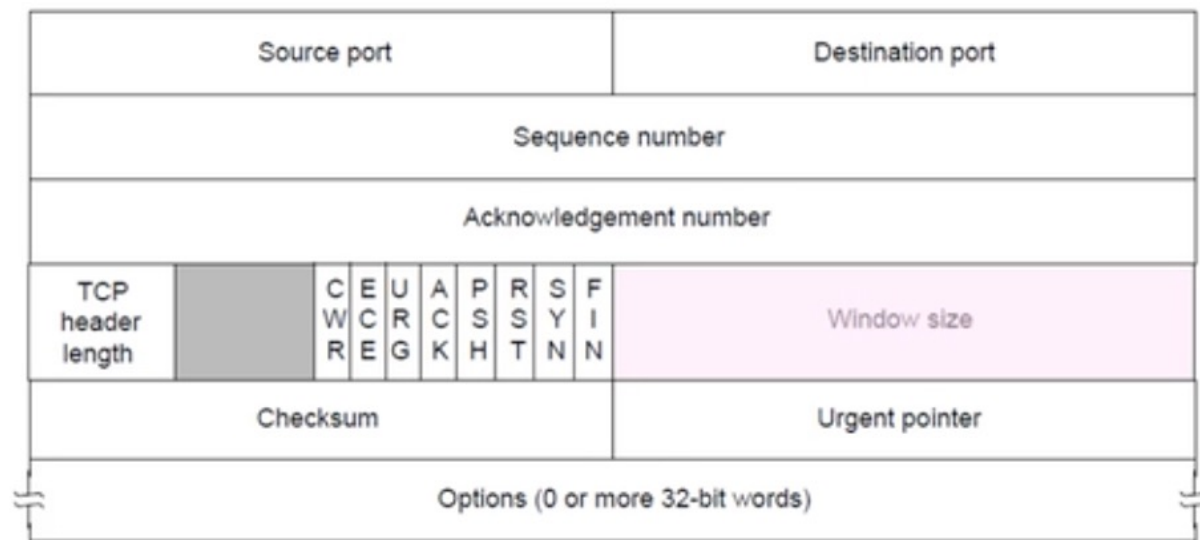  - "Three duplicate ACKs" treated as loss



Sender decides 100-199 is lost

# TCP Header (3)

- SYN/FIN/RST flags for connections
  - Flag indicates segment is a SYN etc.

| Source port | | | | | | | | | | Destination port |
|---|---|---|---|---|---|---|---|---|---|---|
| Sequence number | | | | | | | | | | |
| Acknowledgement number | | | | | | | | | | |
| TCP header length | | C W R | E C E | U R G | A C K | P S H | R S T | S Y N | F I N | Window size |
| Checksum | | | | | | | | | | Urgent pointer |
| Options (0 or more 32-bit words) | | | | | | | | | | |

# TCP Header (4)

- Window size for flow control
  - Relative to ACK, and in bytes

# Other TCP Details

- Many, many quirks you can learn about its operation
  - But they are the details

- Biggest remaining mystery is the working of congestion control

# To-do

- Quiz next week
- Work on your research
- Lab2 due date is April 27th