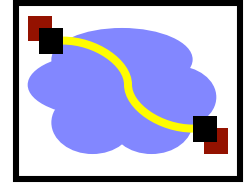


# 15-744 Computer Networking

## Review 2 – Transport Protocols

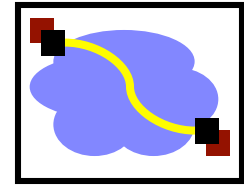
---

# Outline

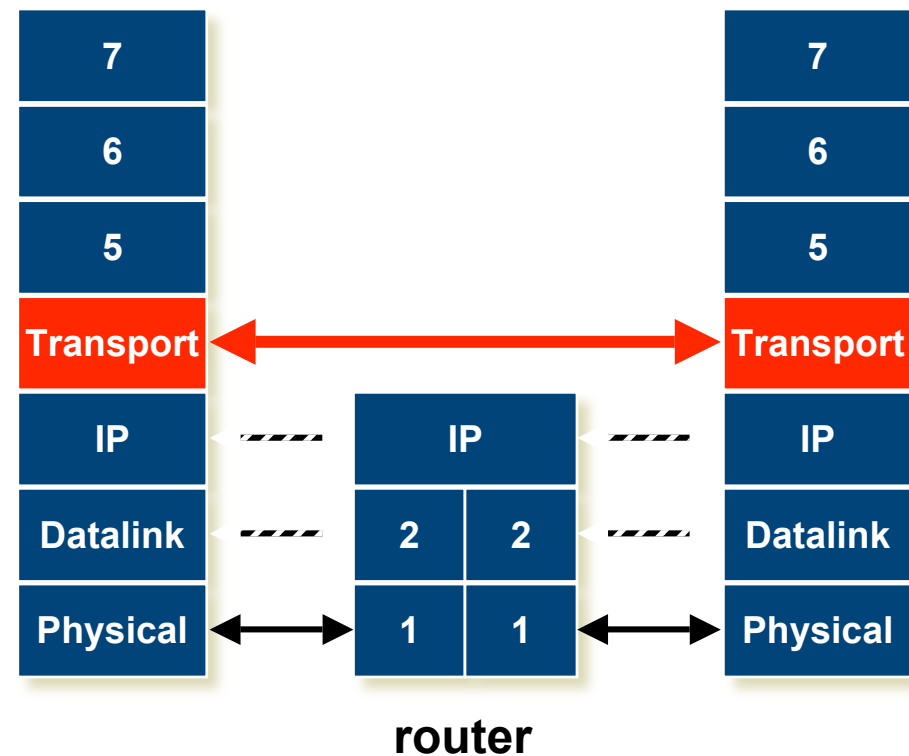


- 
- Transport introduction
  - Error recovery & flow control

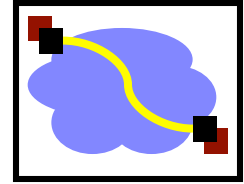
# Transport Protocols



- Lowest level end-to-end protocol.
  - Header generated by sender is interpreted only by the destination
  - Routers view transport header as part of the payload
  - Not always true...
    - Firewalls

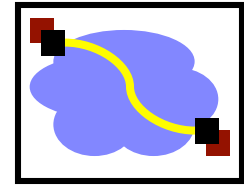


# Functionality Split



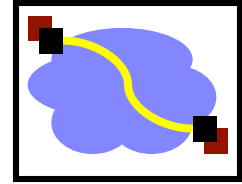
- Network provides best-effort delivery
  - (Hmm, does it anymore? More on this in a few weeks)
- End-systems implement many functions
  - Reliability
  - In-order delivery
  - Demultiplexing
  - Message boundaries
  - Connection abstraction
  - Congestion control
  - ...

# Transport Protocols



- UDP provides just integrity and demux
- TCP adds...
  - Connection-oriented
  - Reliable
  - Ordered
  - Byte-stream
  - Full duplex
  - Flow and congestion controlled
- DCCP, RTP, SCTP -- not widely used.

# UDP: User Datagram Protocol [RFC 768]

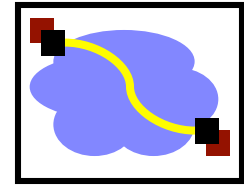


- “No frills,” “bare bones” Internet transport protocol
- “Best effort” service, UDP segments may be:
  - Lost
  - Delivered out of order to app
- *Connectionless*:
  - No handshaking between UDP sender, receiver
  - Each UDP segment handled independently of others

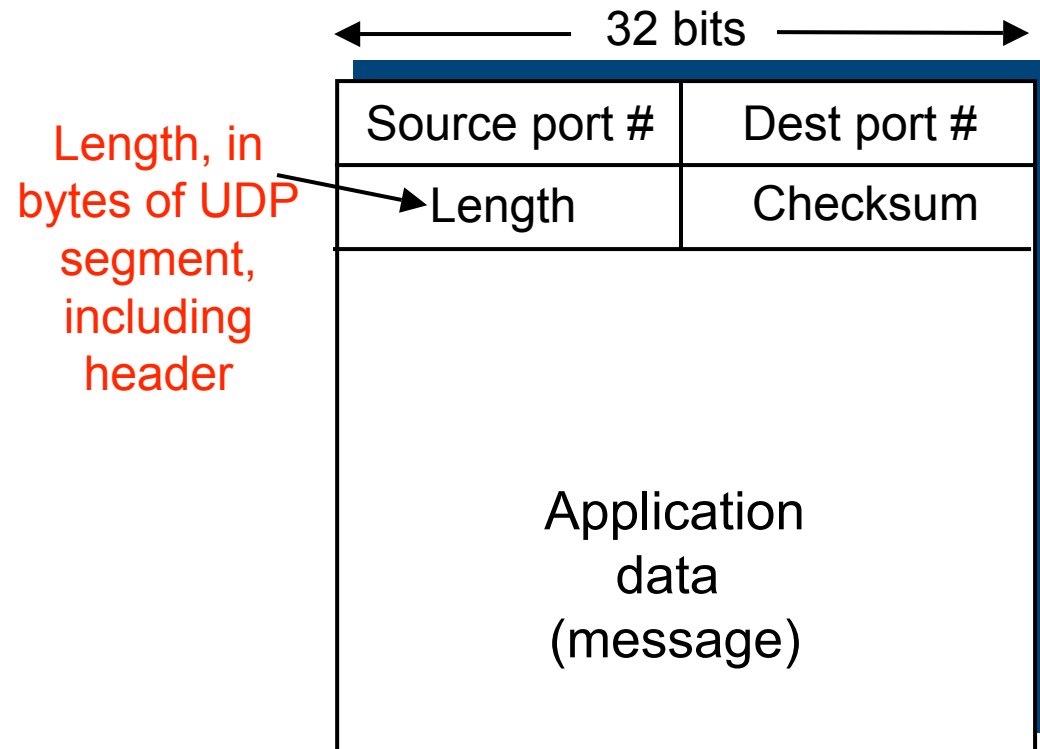
## Why is there a UDP?

- No connection establishment (which can add delay)
- Simple: no connection state at sender, receiver
- Small header
- No congestion control: UDP can blast away as fast as desired

## UDP, cont.

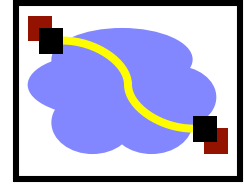


- Often used for streaming multimedia apps
  - Loss tolerant
  - Rate sensitive
- Other UDP uses (why?):
  - DNS
- Reliable transfer over UDP
  - Must be at application layer
  - Application-specific error recovery



UDP segment format

# UDP Checksum



Goal: detect “errors” (e.g., flipped bits) in transmitted segment – optional use!

## Sender:

- Treat segment contents as sequence of 16-bit integers
- Checksum: addition (1’s complement sum) of segment contents
- Sender puts checksum value into UDP checksum field

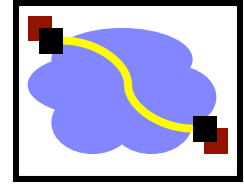
## Receiver:

- Compute checksum of received segment
- Check if computed checksum equals checksum field value:
  - NO - error detected
  - YES - no error detected

*But maybe errors nonetheless?*

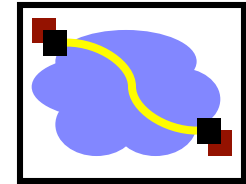


# High-Level TCP Characteristics

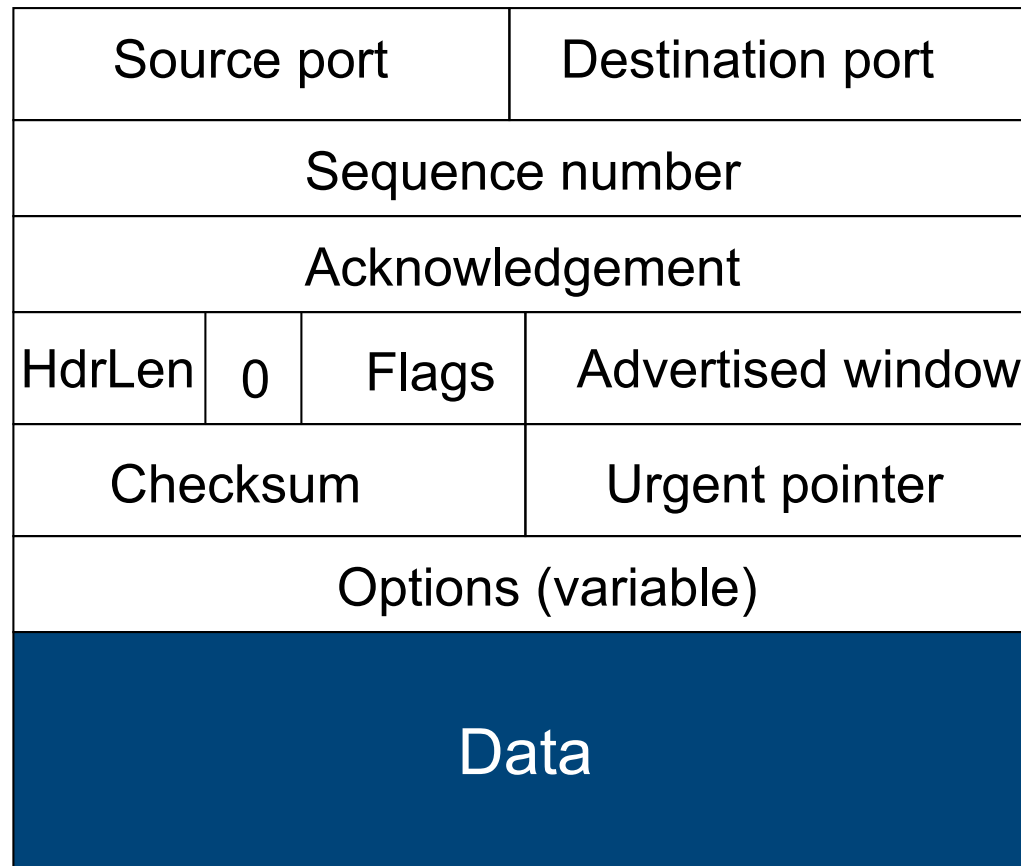


- Protocol implemented entirely at the ends
  - Fate sharing (on IP)
- Protocol has evolved over time and will continue to do so
  - Nearly impossible to change the header
  - Use options to add information to the header
  - Change processing at endpoints
  - Backward compatibility is what makes it TCP

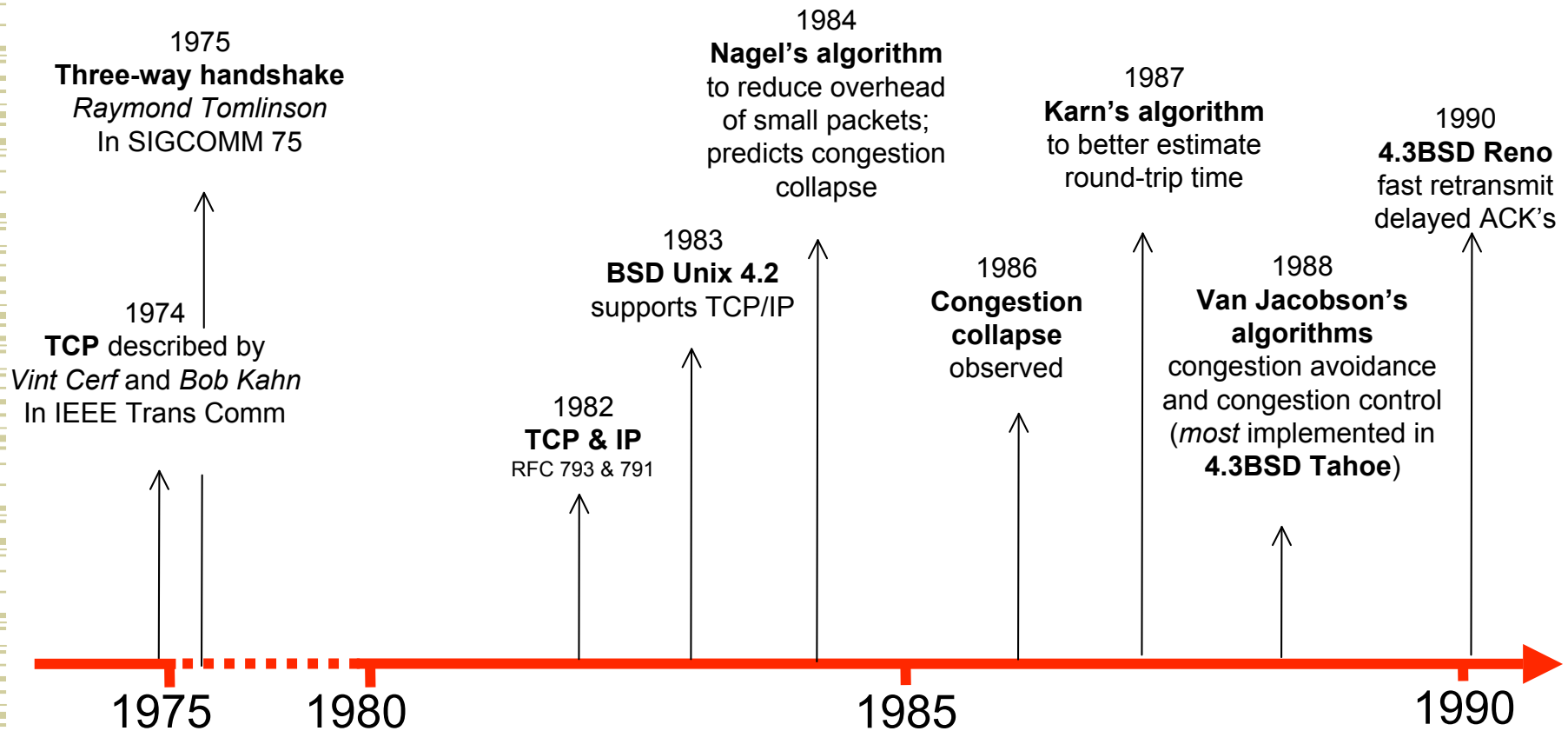
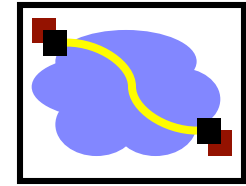
# TCP Header



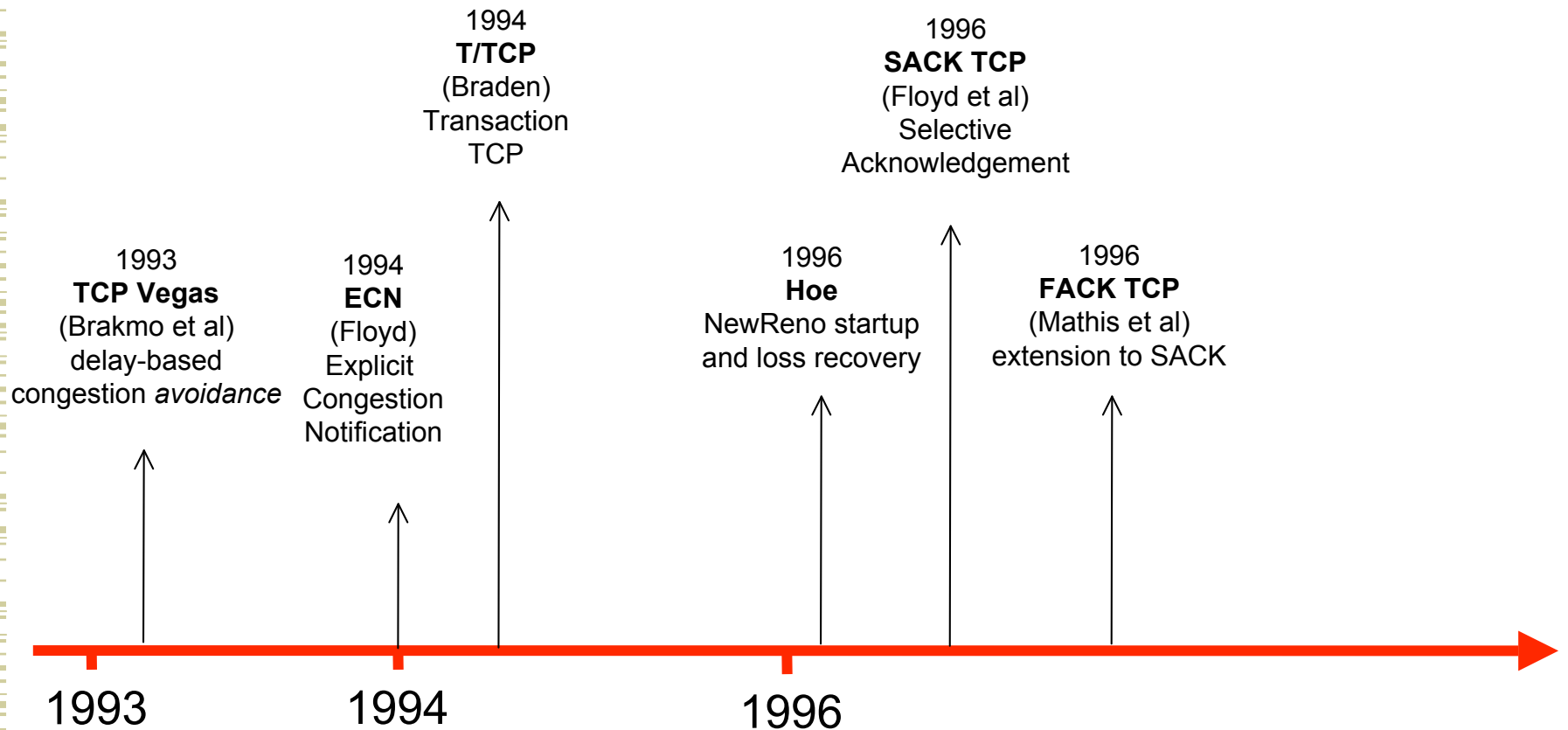
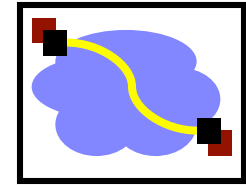
Flags: SYN  
FIN  
RESET  
PUSH  
URG  
ACK



# Evolution of TCP

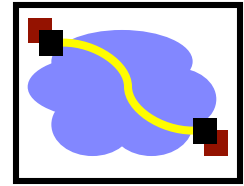


# TCP Through the 1990s



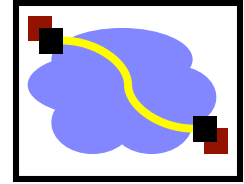
---

# Outline

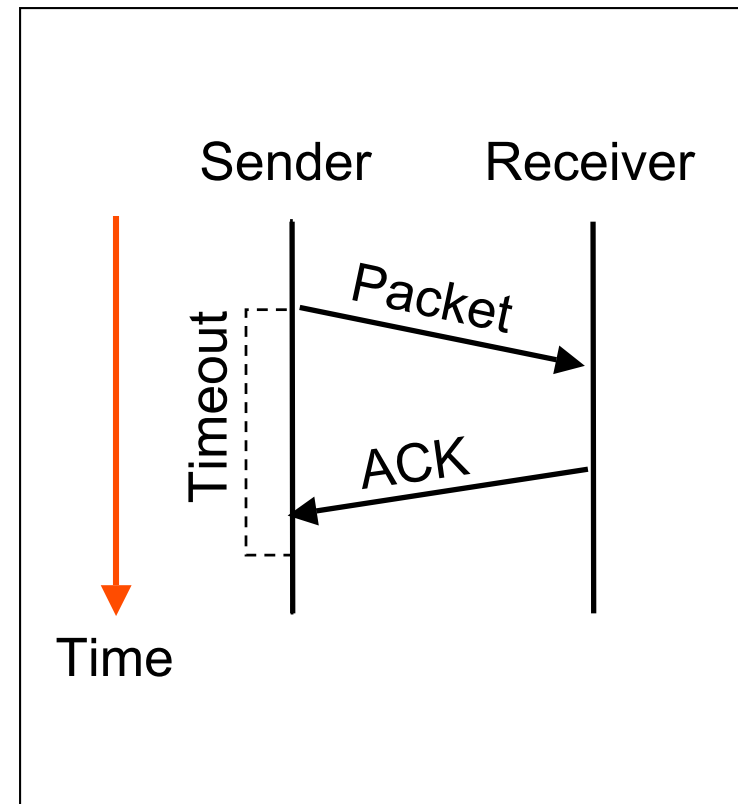


- 
- Transport introduction
  - Error recovery & flow control

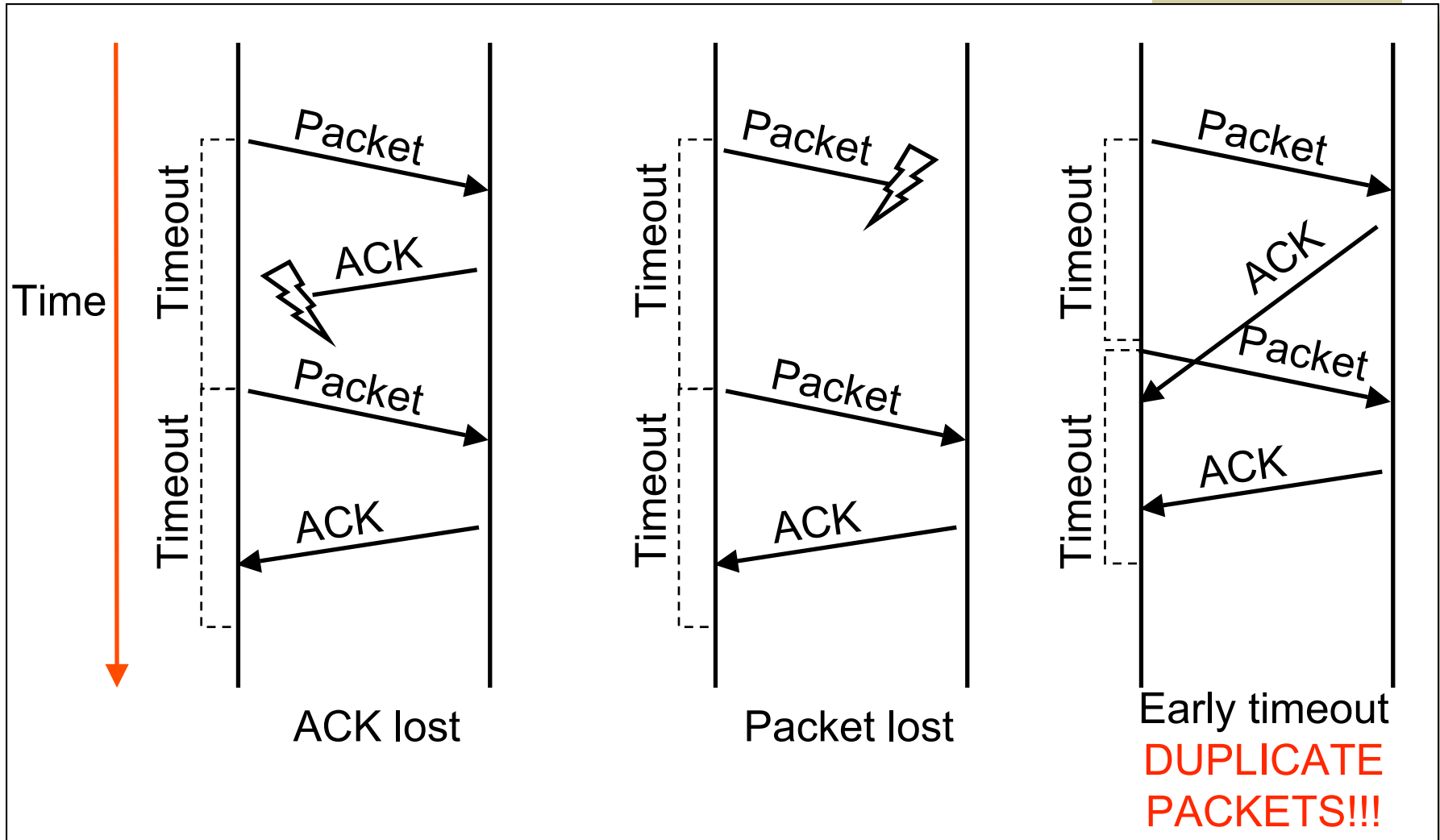
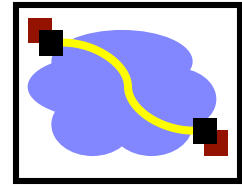
# Stop and Wait



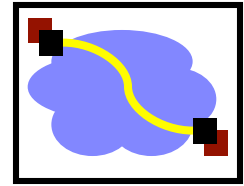
- ARQ
  - Receiver sends acknowledgement (ACK) when it receives packet
  - Sender waits for ACK and timeouts if it does not arrive within some time period
- Simplest ARQ protocol
- Send a packet, stop and wait until ACK arrives



# Recovering from Error



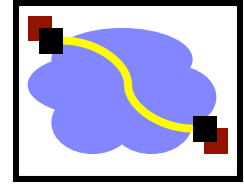
# Problems with Stop and Wait



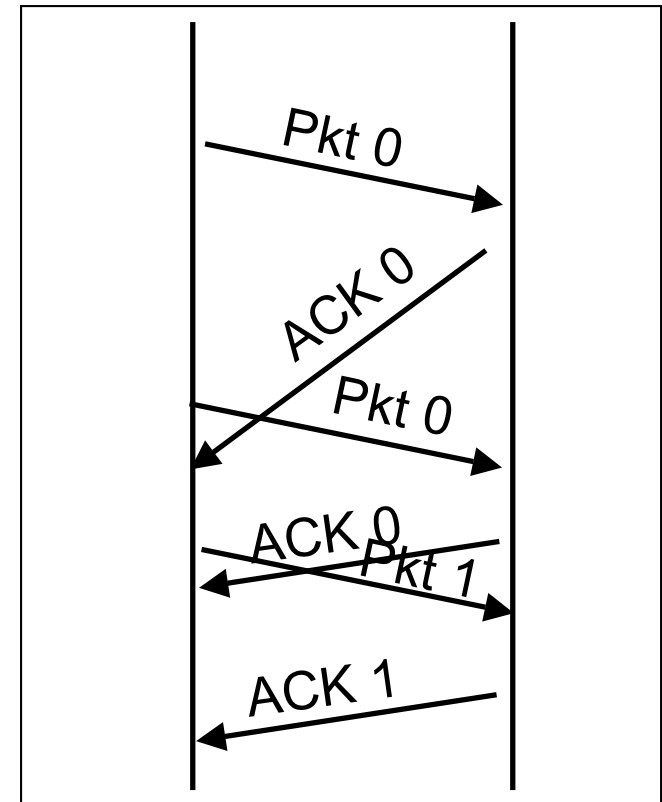
- How to recognize a duplicate
- Performance
  - Can only send one packet per round trip



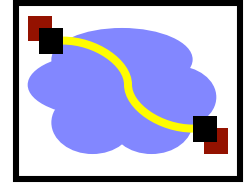
# How to Recognize Resends?



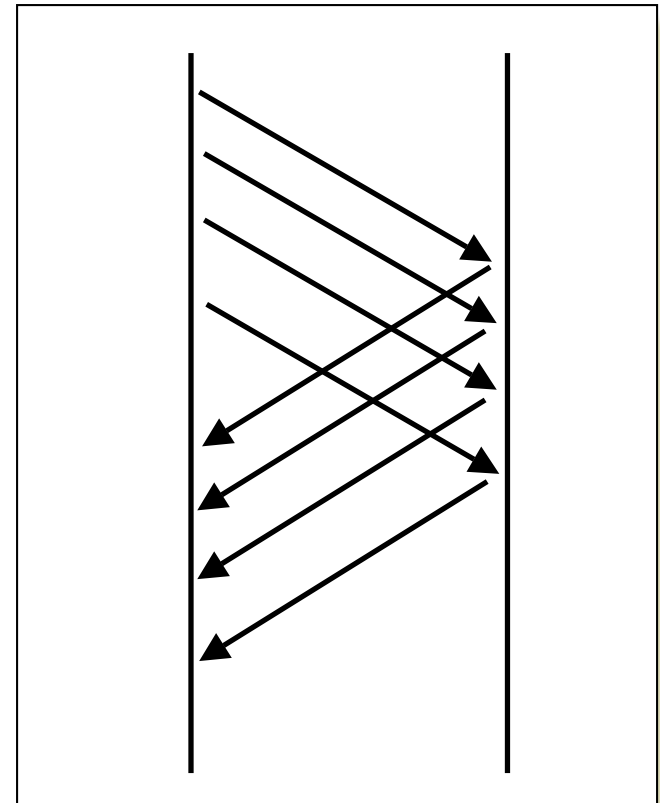
- Use sequence numbers
  - both packets and acks
- Sequence # in packet is finite  
→ How big should it be?
  - For stop and wait?
- One bit – won't send seq #1 until received ACK for seq #0



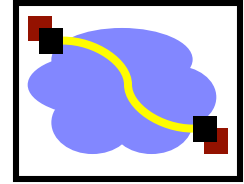
# How to Keep the Pipe Full?



- Send multiple packets without waiting for first to be acked
  - Number of pkts in flight = window:  
Flow control
- Reliable, unordered delivery
  - Several parallel stop & waits
  - Send new packet after each ack
  - Sender keeps list of unack'ed packets; resends after timeout
  - Receiver same as stop & wait
- How large a window is needed?
  - Suppose 10Mbps link, 4ms delay, 500byte pkts
    - 1? 10? 20?
  - Round trip delay \* bandwidth =  
capacity of pipe

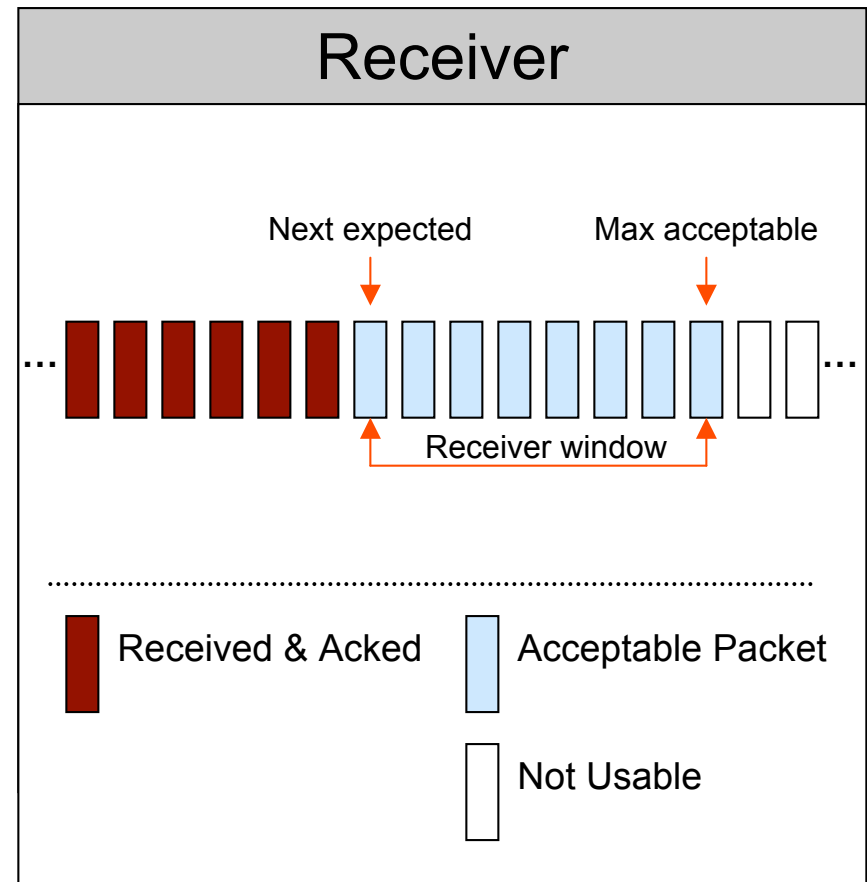
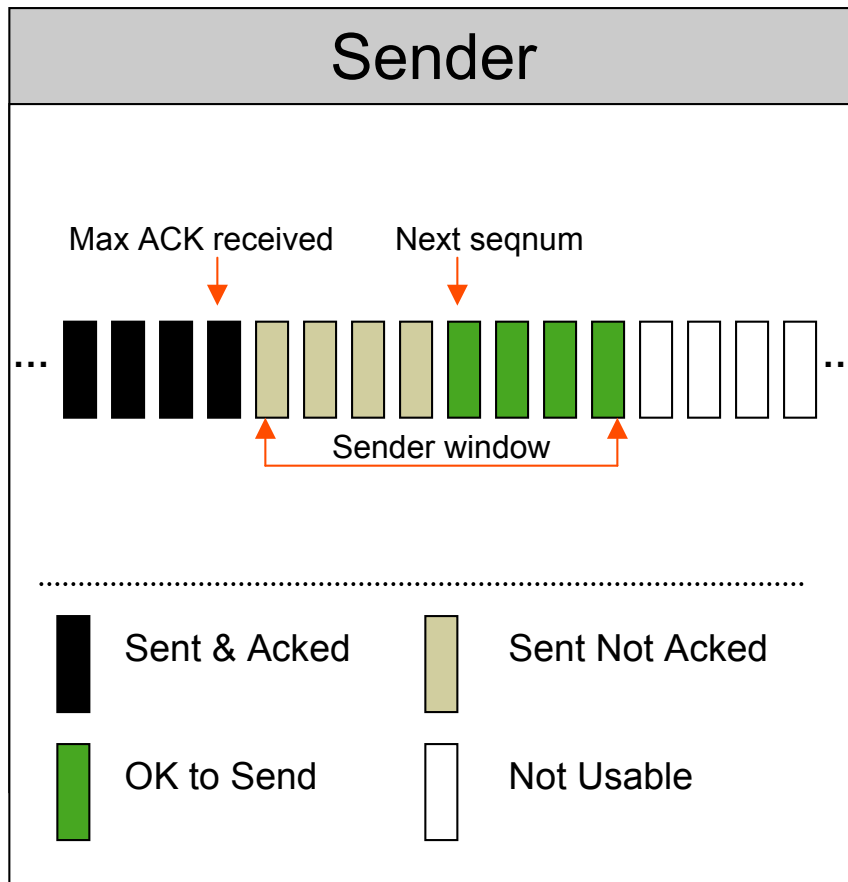
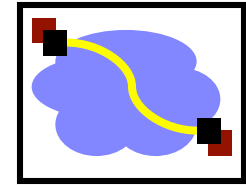


# Sliding Window

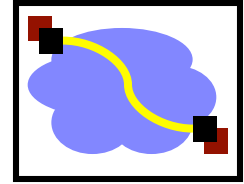


- Reliable, ordered delivery
- Receiver has to hold onto a packet until all prior packets have arrived
  - Why might this be difficult for just parallel stop & wait?
  - Sender must prevent buffer overflow at receiver
- Circular buffer at sender and receiver
  - Packets in transit  $\leq$  buffer size
  - Advance when sender and receiver agree packets at beginning have been received

# Sender/Receiver State

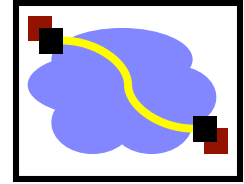


# Sequence Numbers



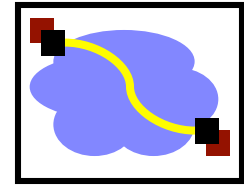
- How large do sequence numbers need to be?
  - Must be able to detect wrap-around
  - Depends on sender/receiver window size
- E.g.
  - Max seq = 7, send win=recv win=7
  - If pkts 0..6 are sent successfully and all acks lost
    - Receiver expects 7,0..5, sender retransmits old 0..6!!!
- Max sequence must be  $\geq$  send window + recv window

## Window Sliding – Common Case



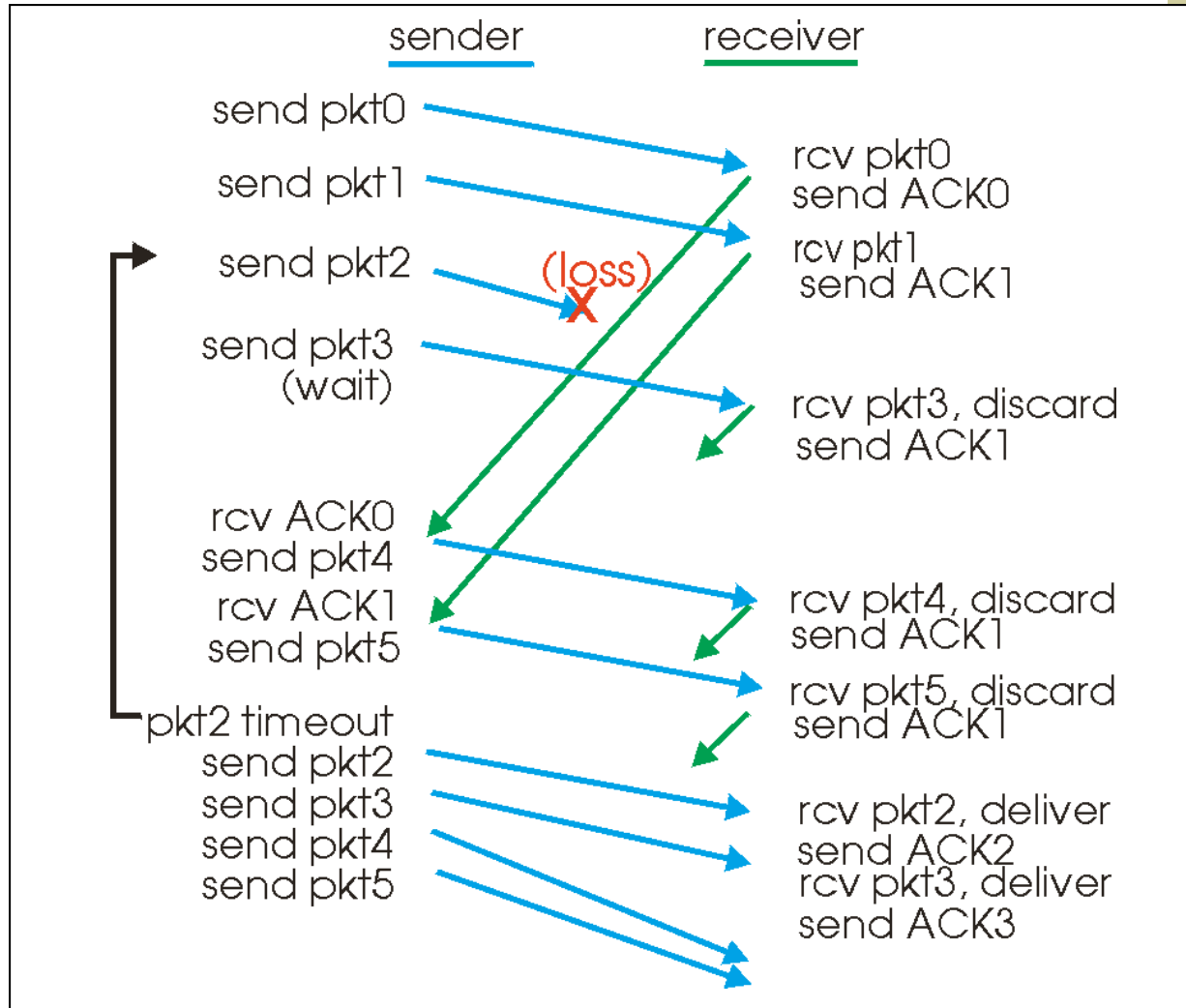
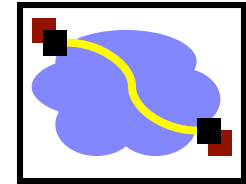
- On reception of new ACK (i.e. ACK for something that was not acked earlier)
  - Increase sequence of max ACK received
  - Send next packet
- On reception of new in-order data packet (next expected)
  - Hand packet to application
  - Send **cumulative ACK** – acknowledges reception of all packets up to sequence number
  - Increase sequence of max acceptable packet

# Loss Recovery



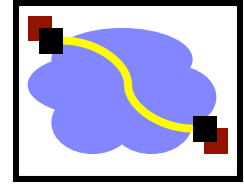
- On reception of out-of-order packet
  - Send nothing (wait for source to timeout)
  - Cumulative ACK (helps source identify loss)
- Timeout (Go-Back-N recovery)
  - Set timer upon transmission of packet
  - Retransmit all unacknowledged packets
- Performance during loss recovery
  - No longer have an entire window in transit
  - Can have much more clever loss recovery

# Go-Back-N in Action



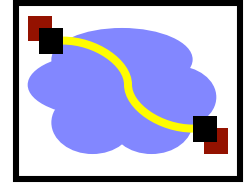


# Important Lessons



- Transport service
  - UDP → mostly just IP service
  - TCP → congestion controlled, reliable, byte stream
- Types of ARQ protocols
  - Stop-and-wait → slow, simple
  - Go-back-n → can keep link utilized (except w/ losses)
  - Selective repeat → efficient loss recovery -- used in SACK
- Sliding window flow control
  - Addresses buffering issues and keeps link utilized

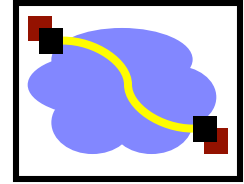
## Good Ideas So Far...



- Flow control
  - Stop & wait
  - Parallel stop & wait
  - Sliding window
- Loss recovery
  - Timeouts
  - Acknowledgement-driven recovery (selective repeat or cumulative acknowledgement)

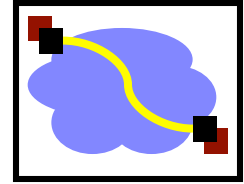
---

## Outline



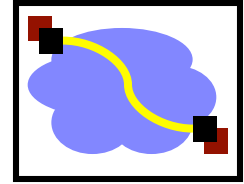
- 
- TCP flow control
  - Congestion sources and collapse
  - Congestion control basics

# More on Sequence Numbers



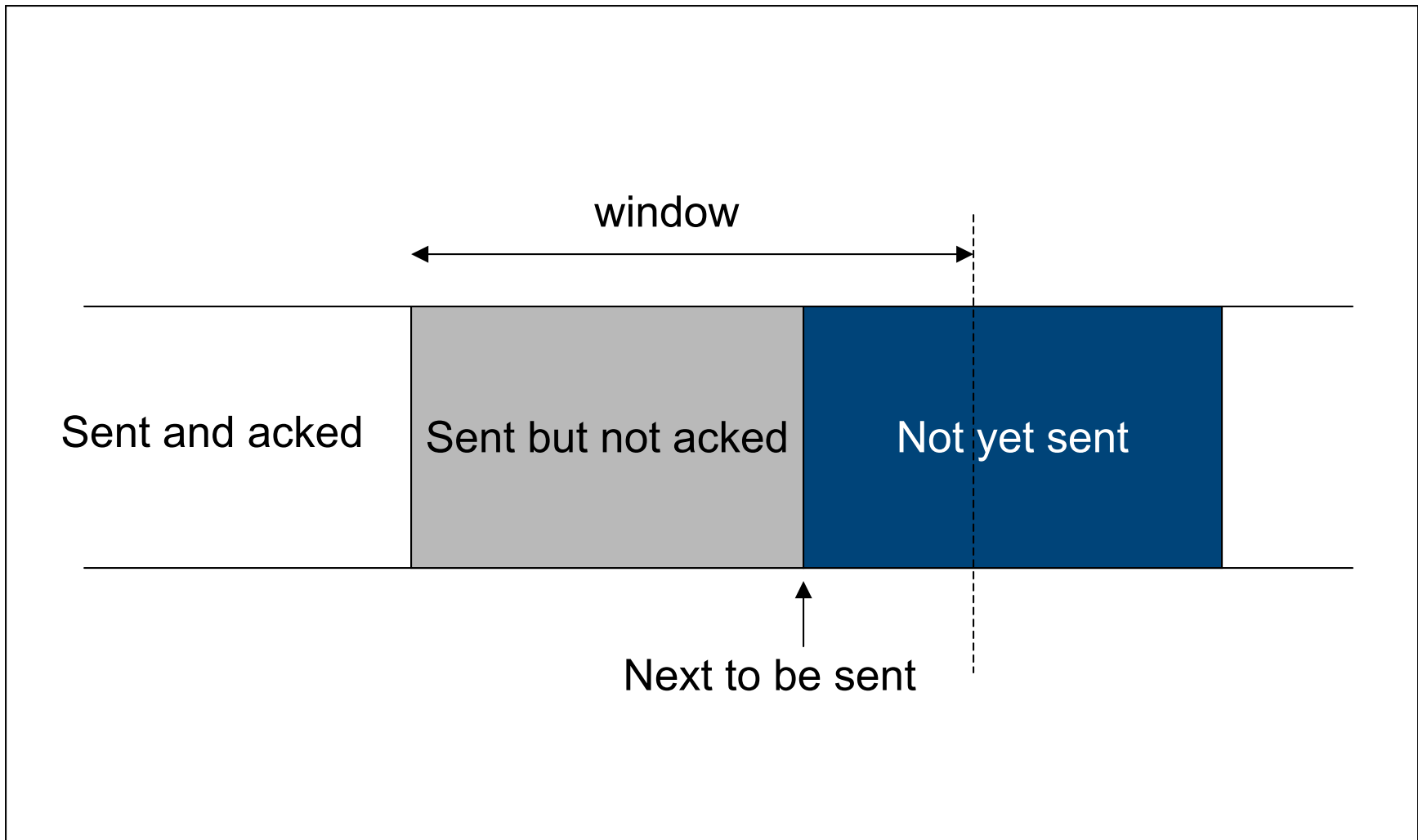
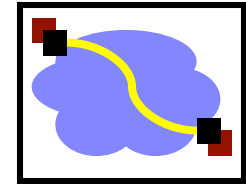
- 32 Bits, Unsigned → for bytes not packets!
- Why So Big?
  - For sliding window, must have  $|\text{Sequence Space}| > |\text{Sending Window}| + |\text{Receiving Window}|$ 
    - No problem
  - Also, want to guard against stray packets
    - With IP, packets have maximum lifetime of 120s
    - Sequence number would wrap around in this time at 286MB/s

# TCP Flow Control

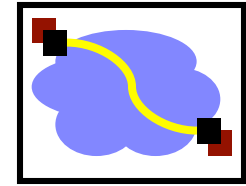


- TCP is a sliding window protocol
  - For window size  $n$ , can send up to  $n$  bytes without receiving an acknowledgement
  - When the data is acknowledged then the window slides forward
- Each packet advertises a window size
  - Indicates number of bytes the receiver has space for
- Original TCP always sent entire window
  - Congestion control now limits this

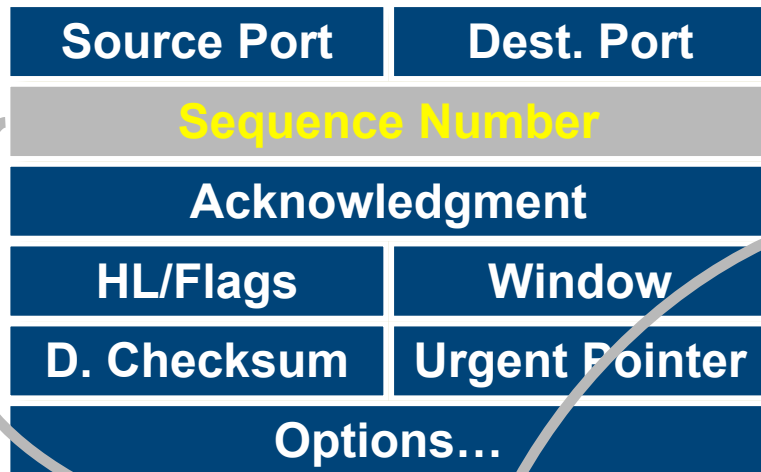
# Window Flow Control: Send Side



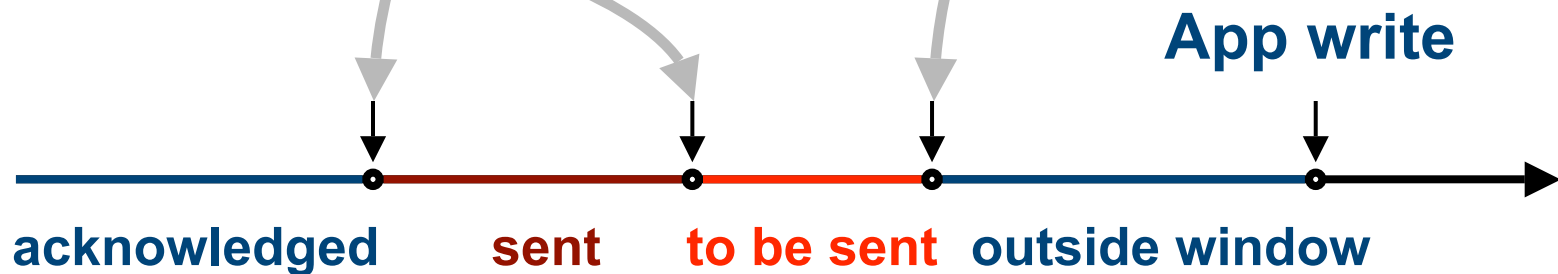
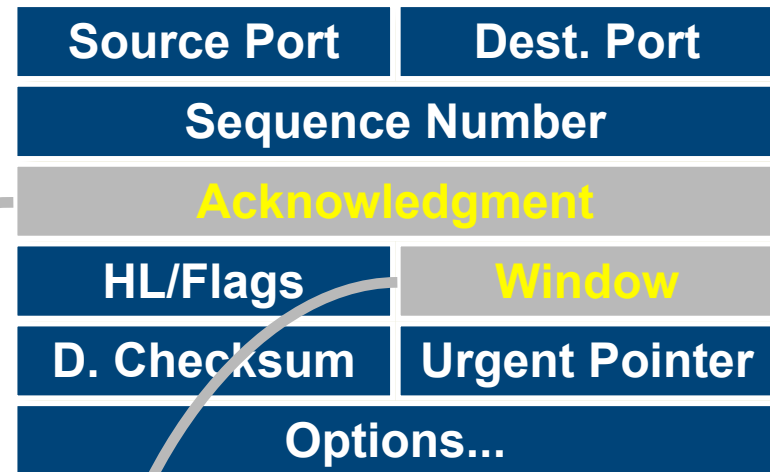
# Window Flow Control: Send Side



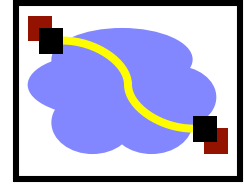
## Packet Sent



## Packet Received



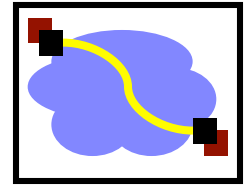
# Performance Considerations



- The window size can be controlled by receiving application
  - Can change the socket buffer size from a default (e.g. 8Kbytes) to a maximum value (e.g. 64 Kbytes)
- The window size field in the TCP header limits the window that the receiver can advertise
  - 16 bits  $\rightarrow$  64 KBytes
  - 10 msec RTT  $\rightarrow$  51 Mbit/second
  - 100 msec RTT  $\rightarrow$  5 Mbit/second
  - TCP options to get around 64KB limit  $\rightarrow$  increases above limit

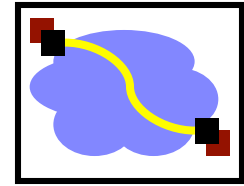


# Outline

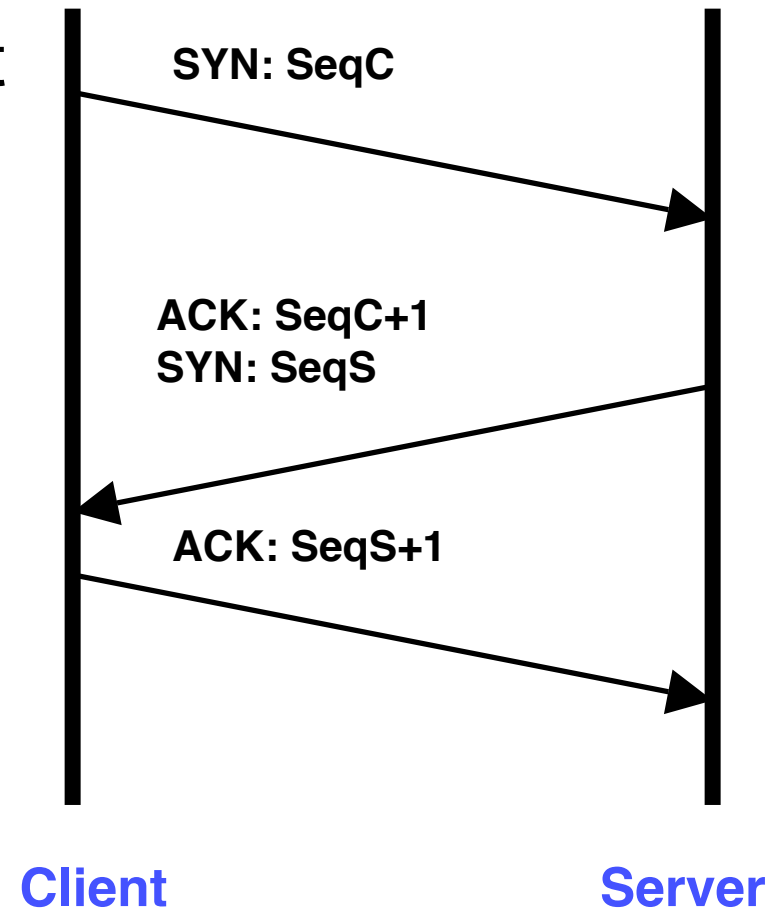


- TCP connection setup/data transfer
- TCP reliability
  - How to recover from lost packets
- TCP congestion avoidance
  - Paper for Monday

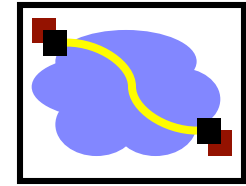
# Establishing Connection: Three-Way handshake



- Each side notifies other of starting sequence number it will use for sending
  - Why not simply chose 0?
    - Must avoid overlap with earlier incarnation
    - Security issues
- Each side acknowledges other's sequence number
  - SYN-ACK: Acknowledge sequence number + 1
- Can combine second SYN with first ACK

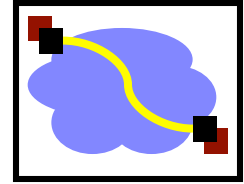


# Outline



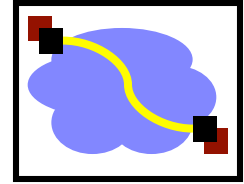
- TCP connection setup/data transfer
- TCP reliability

# Reliability Challenges



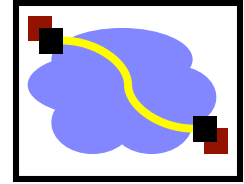
- Congestion related losses
- Variable packet delays
  - What should the timeout be?
- Reordering of packets
  - How to tell the difference between a delayed packet and a lost one?

# TCP = Go-Back-N Variant



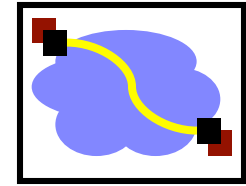
- Sliding window with cumulative acks
  - Receiver can only return a single “ack” sequence number to the sender.
  - Acknowledges all bytes with a lower sequence number
  - Starting point for retransmission
  - Duplicate acks sent when out-of-order packet received
- But: sender only retransmits a single packet.
  - Reason???
  - Only one that it knows is lost
  - Network is congested → shouldn't overload it
- Error control is based on byte sequences, not packets.
  - Retransmitted packet can be different from the original lost packet – Why?

# Round-trip Time Estimation

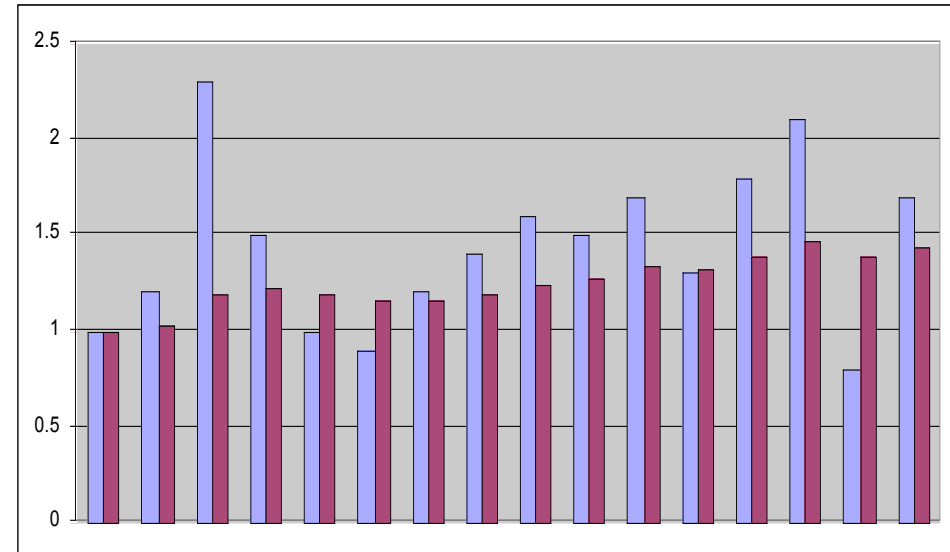


- Wait at least one RTT before retransmitting
- Importance of accurate RTT estimators:
  - Low RTT estimate
    - unneeded retransmissions
  - High RTT estimate
    - poor throughput
- RTT estimator must adapt to change in RTT
  - But not too fast, or too slow!
- Spurious timeouts
  - “Conservation of packets” principle – never more than a window worth of packets in flight

# Original TCP Round-trip Estimator

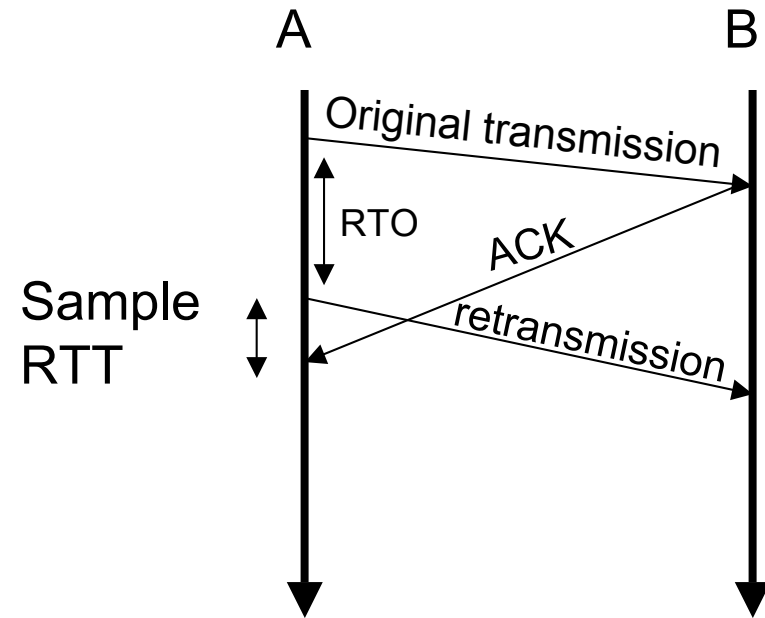
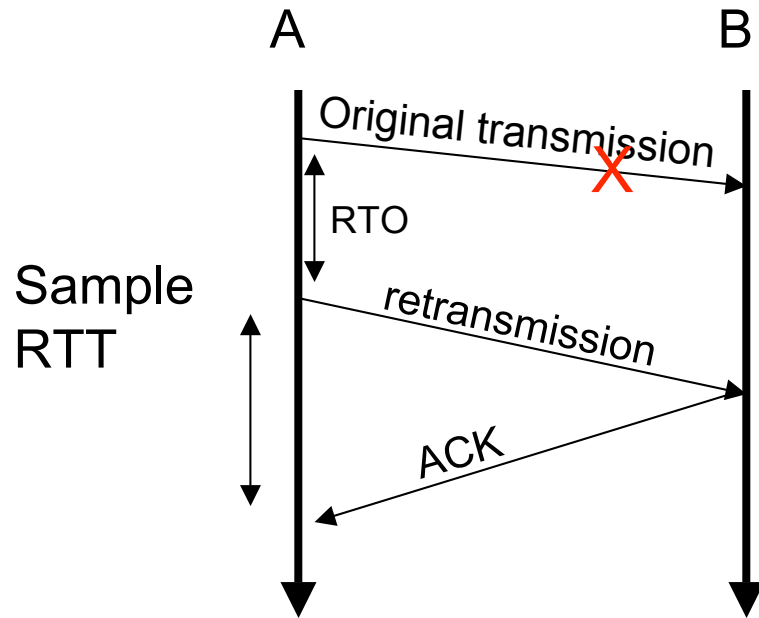
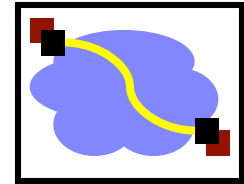


- Round trip times exponentially averaged:
  - $\text{New RTT} = \alpha (\text{old RTT}) + (1 - \alpha) (\text{new sample})$
  - Recommended value for  $\alpha$ : 0.8 - 0.9
    - 0.875 for most TCP's



- Retransmit timer set to  $(b * \text{RTT})$ , where  $b = 2$ 
  - Every time timer expires, RTO exponentially backed-off
- Not good at preventing premature timeouts
  - Why?

# RTT Sample Ambiguity

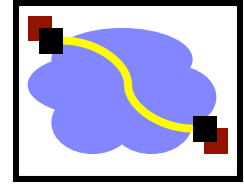


- Karn's RTT Estimator

- If a segment has been retransmitted:
  - Don't count RTT sample on ACKs for this segment
  - Keep backed off time-out for next packet
  - Reuse RTT estimate only after one successful transmission

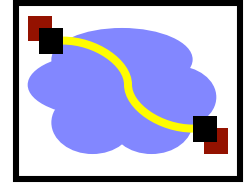


# Jacobson's Retransmission Timeout



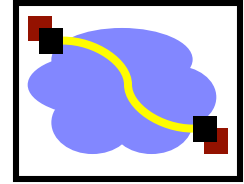
- Key observation:
  - At high loads round trip variance is high
- Solution:
  - Base RTO on RTT and standard deviation
    - $RTO = RTT + 4 * rttvar$
  - $new\_rttvar = \beta * dev + (1 - \beta) old\_rttvar$ 
    - Dev = linear deviation
    - Inappropriately named – actually smoothed linear deviation

## Timestamp Extension



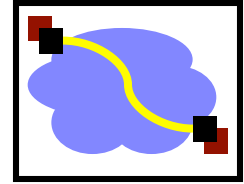
- Used to improve timeout mechanism by more accurate measurement of RTT
- When sending a packet, insert current time into option
  - 4 bytes for time, 4 bytes for echo a received timestamp
- Receiver echoes timestamp in ACK
  - Actually will echo whatever is in timestamp
- Removes retransmission ambiguity
  - Can get RTT sample on any packet

# Timer Granularity



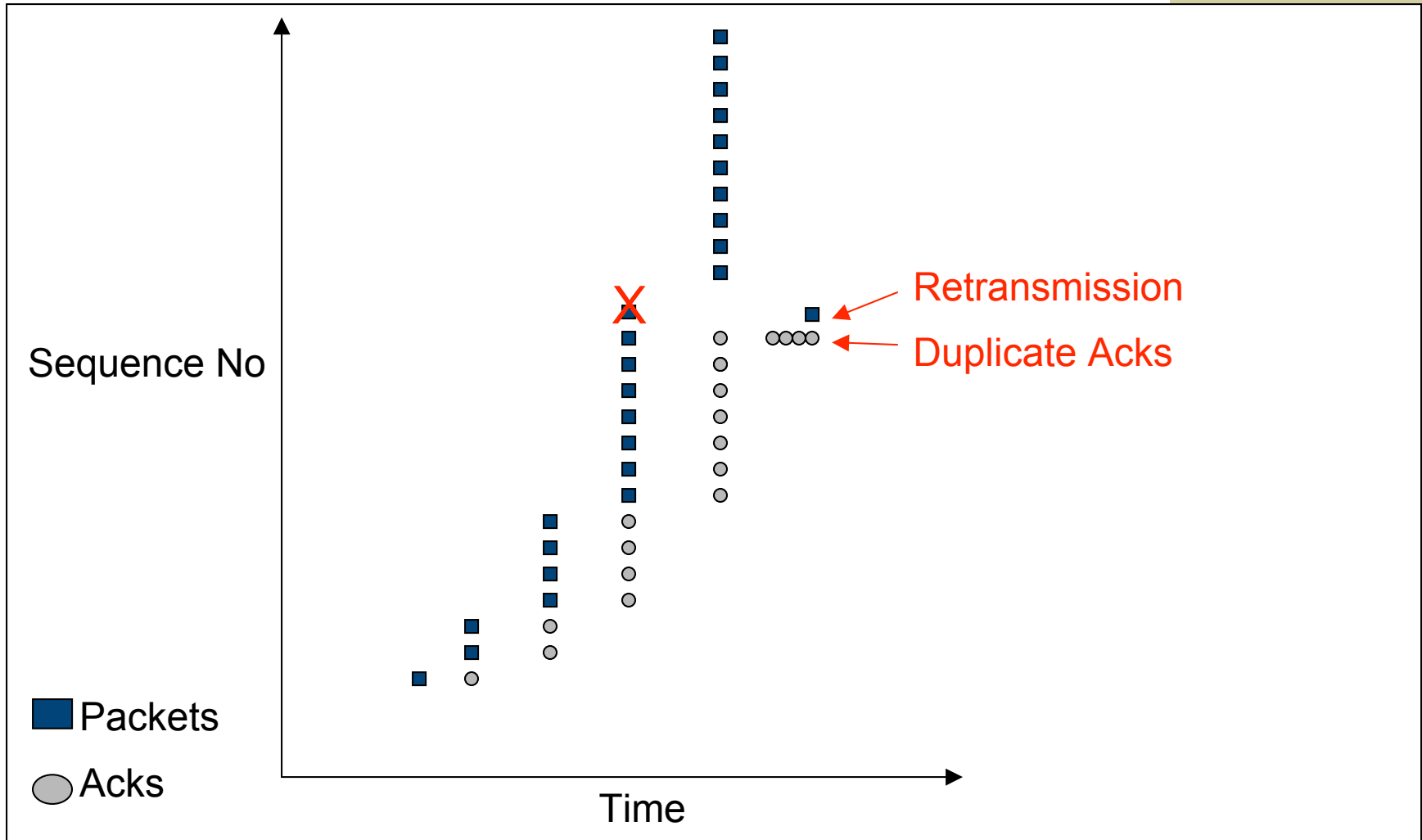
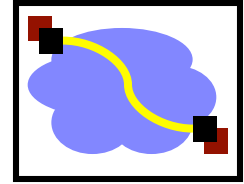
- Many TCP implementations set RTO in multiples of 200,500,1000ms
- Why?
  - Avoid spurious timeouts – RTTs can vary quickly due to cross traffic
- What happens for the first couple of packets?
  - Pick a very conservative value (seconds)

# Fast Retransmit -- Avoiding Timeouts

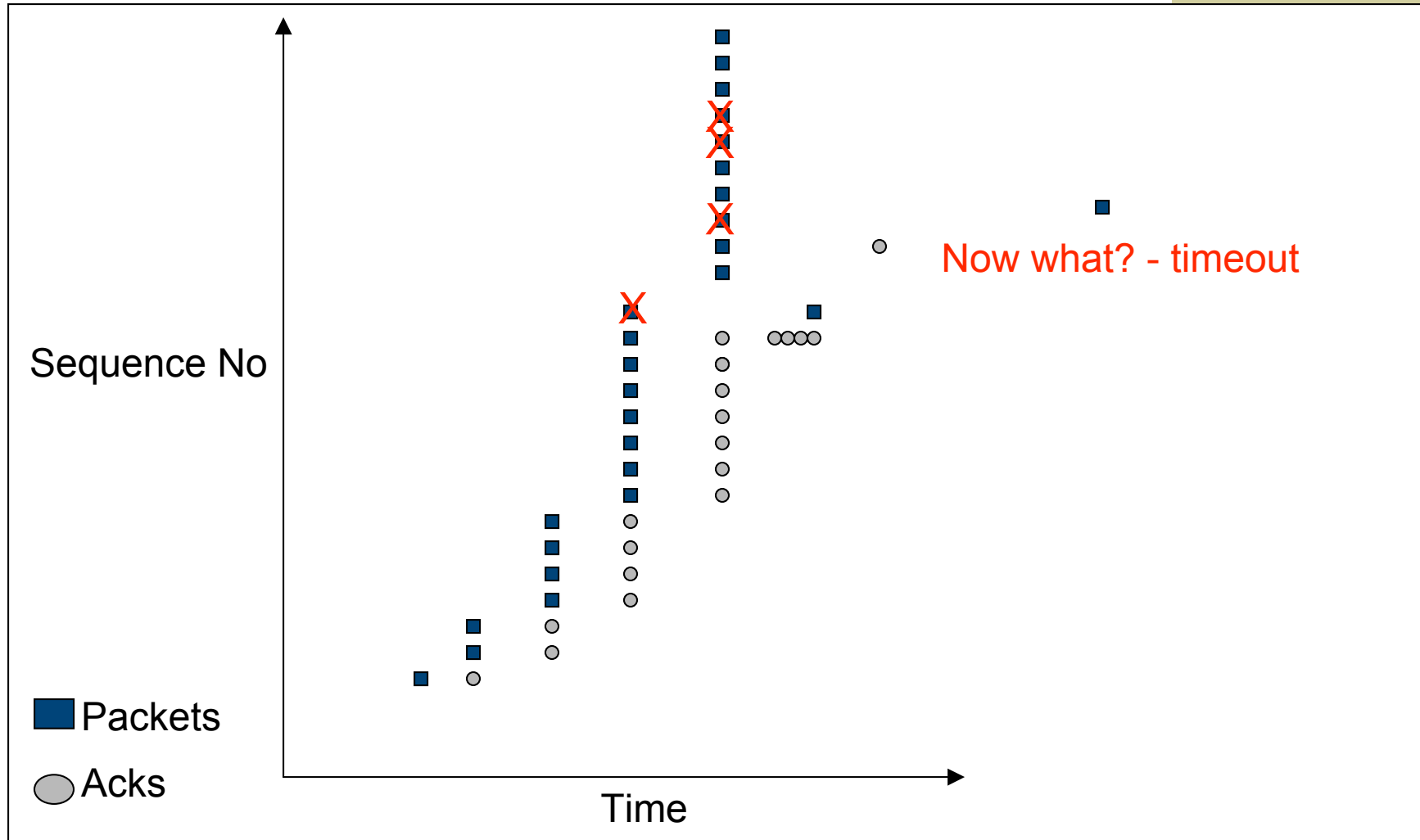
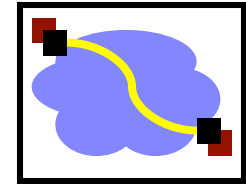


- What are duplicate acks (dupacks)?
  - Repeated acks for the same sequence
- When can duplicate acks occur?
  - Loss
  - Packet re-ordering
  - Window update – advertisement of new flow control window
- Assume re-ordering is infrequent and not of large magnitude
  - Use receipt of 3 or more duplicate acks as indication of loss
  - Don't wait for timeout to retransmit packet

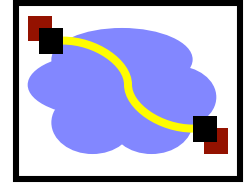
# Fast Retransmit



# TCP (Reno variant)

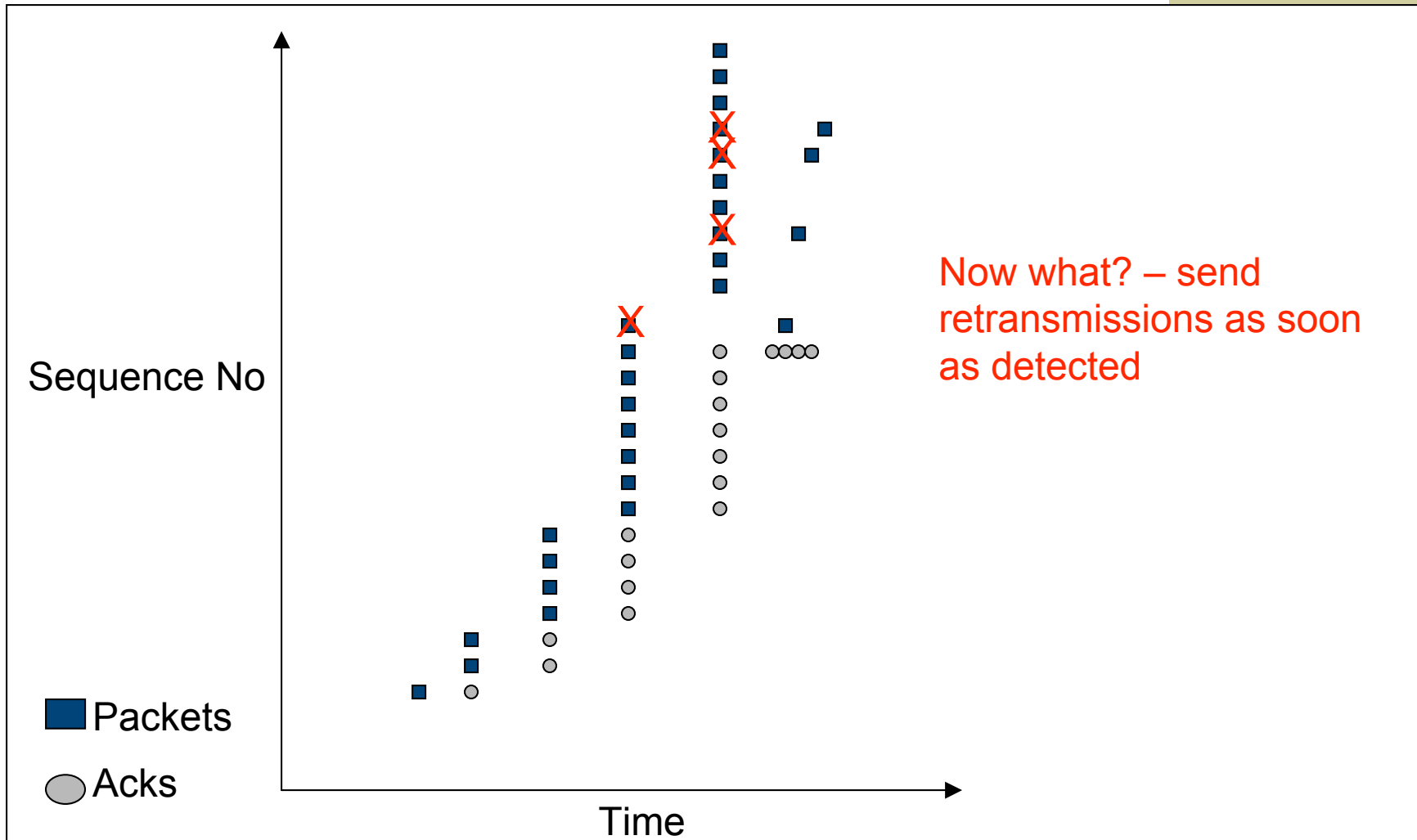
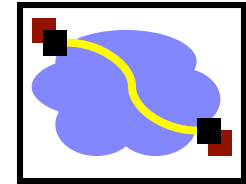


# SACK



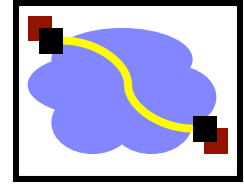
- Basic problem is that cumulative acks provide little information
- Selective acknowledgement (SACK) essentially adds a bitmask of packets received
  - Implemented as a TCP option
  - Encoded as a set of received byte ranges (max of 4 ranges/often max of 3)
- When to retransmit?
  - Still need to deal with reordering → wait for out of order by 3pkts

# SACK



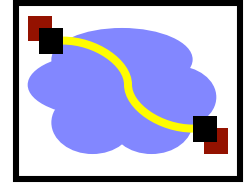


# Performance Issues



- Timeout >> fast retransmit
- Need 3 dupacks/sacks
- Not great for small transfers
  - Don't have 3 packets outstanding
- What are real loss patterns like?

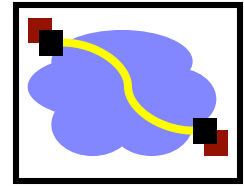
# Important Lessons



- Three-way TCP Handshake
- TCP timeout calculation → how is RTT estimated
- Modern TCP loss recovery
  - Why are timeouts bad?
  - How to avoid them? → e.g. fast retransmit

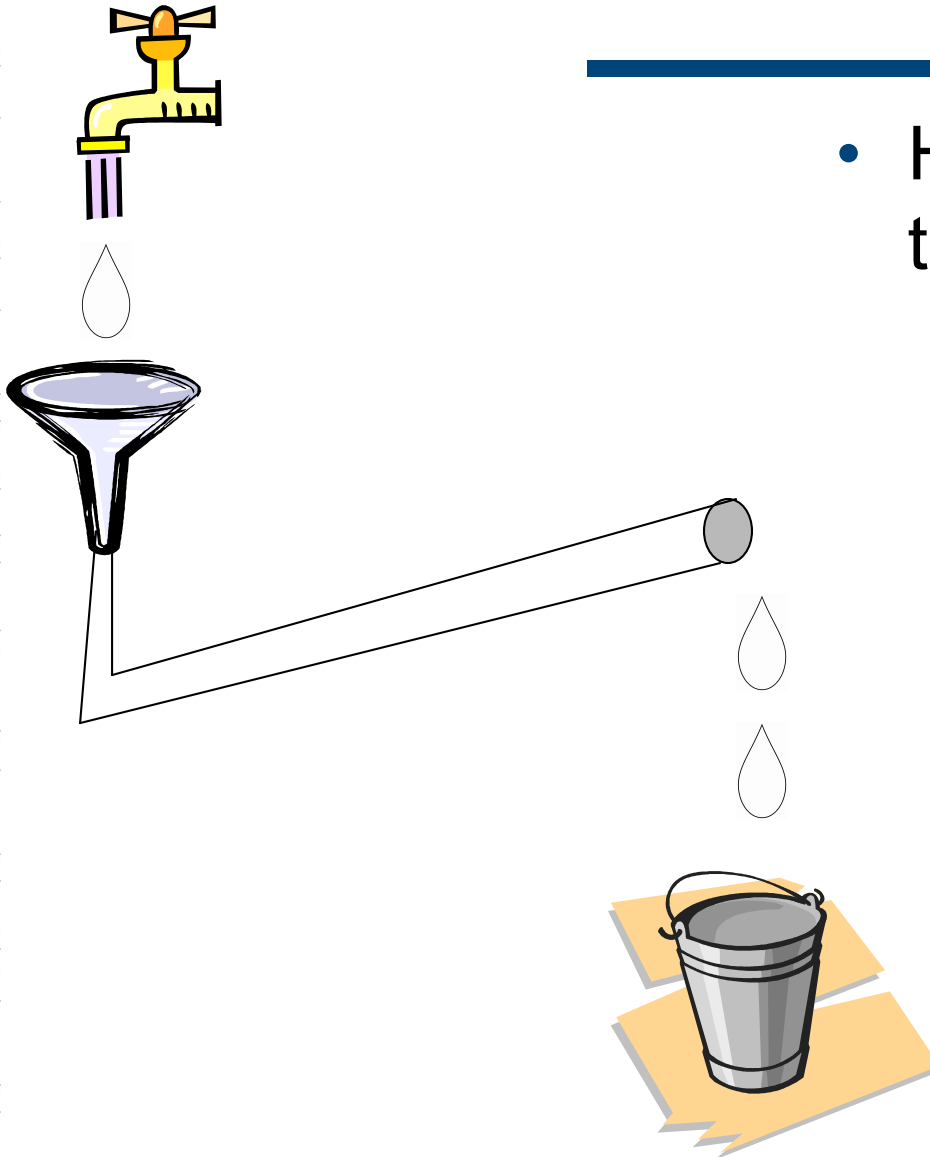
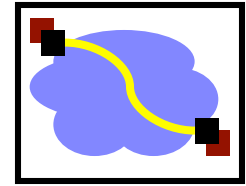
---

# Outline



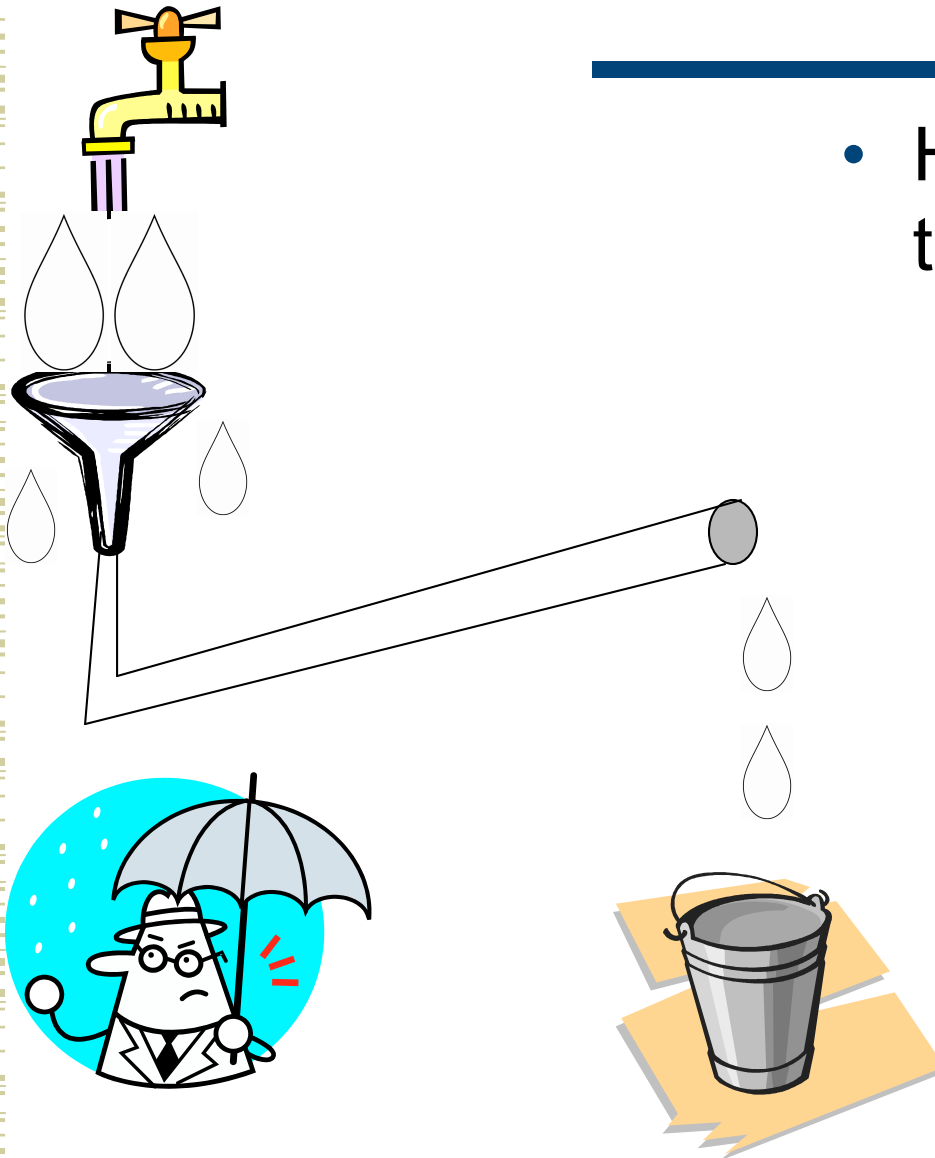
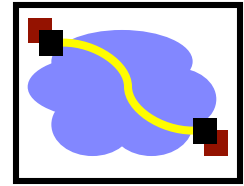
- TCP flow control
- Congestion sources and collapse
- Congestion control basics

# Internet Pipes?



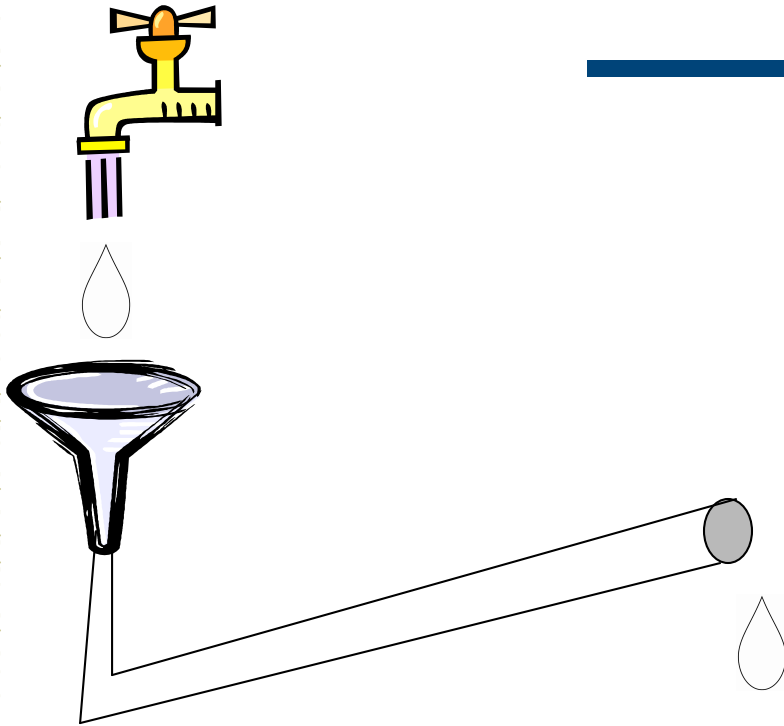
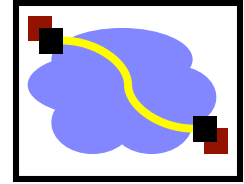
- How should you control the faucet?

# Internet Pipes?



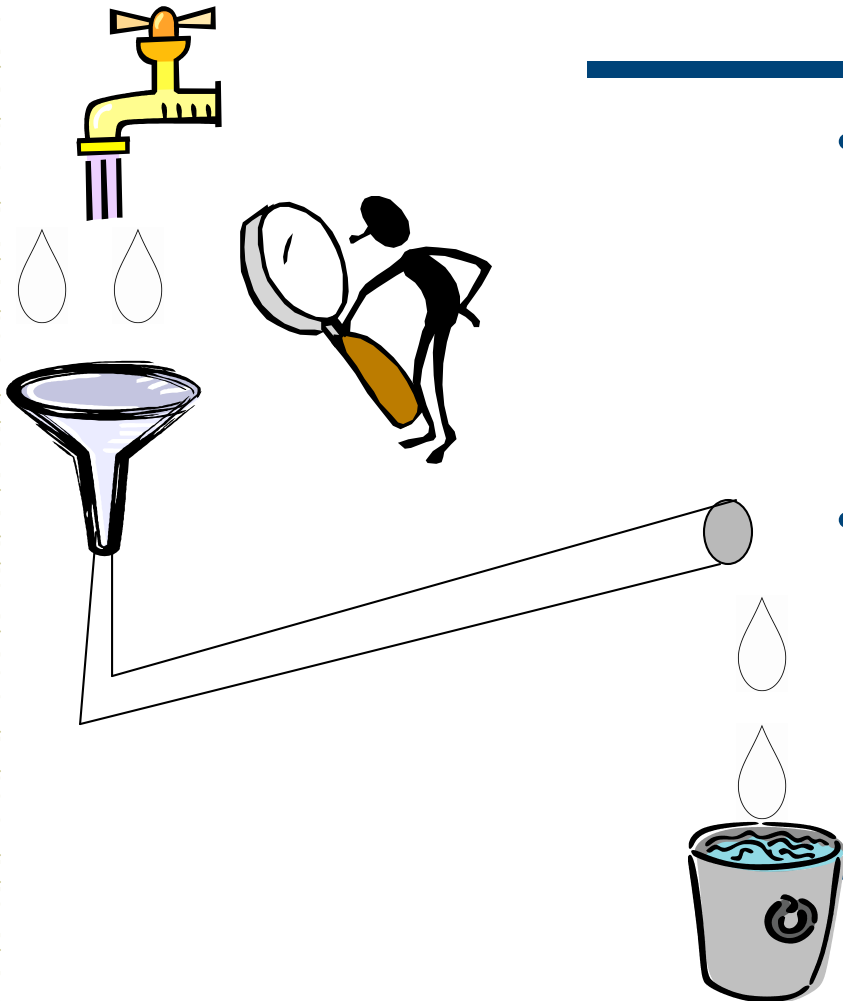
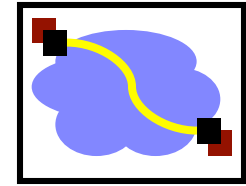
- How should you control the faucet?
  - Too fast – sink overflows!

# Internet Pipes?



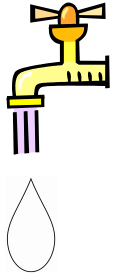
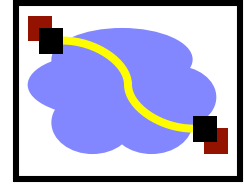
- How should you control the faucet?
  - Too fast – sink overflows!
  - Too slow – what happens?

# Internet Pipes?

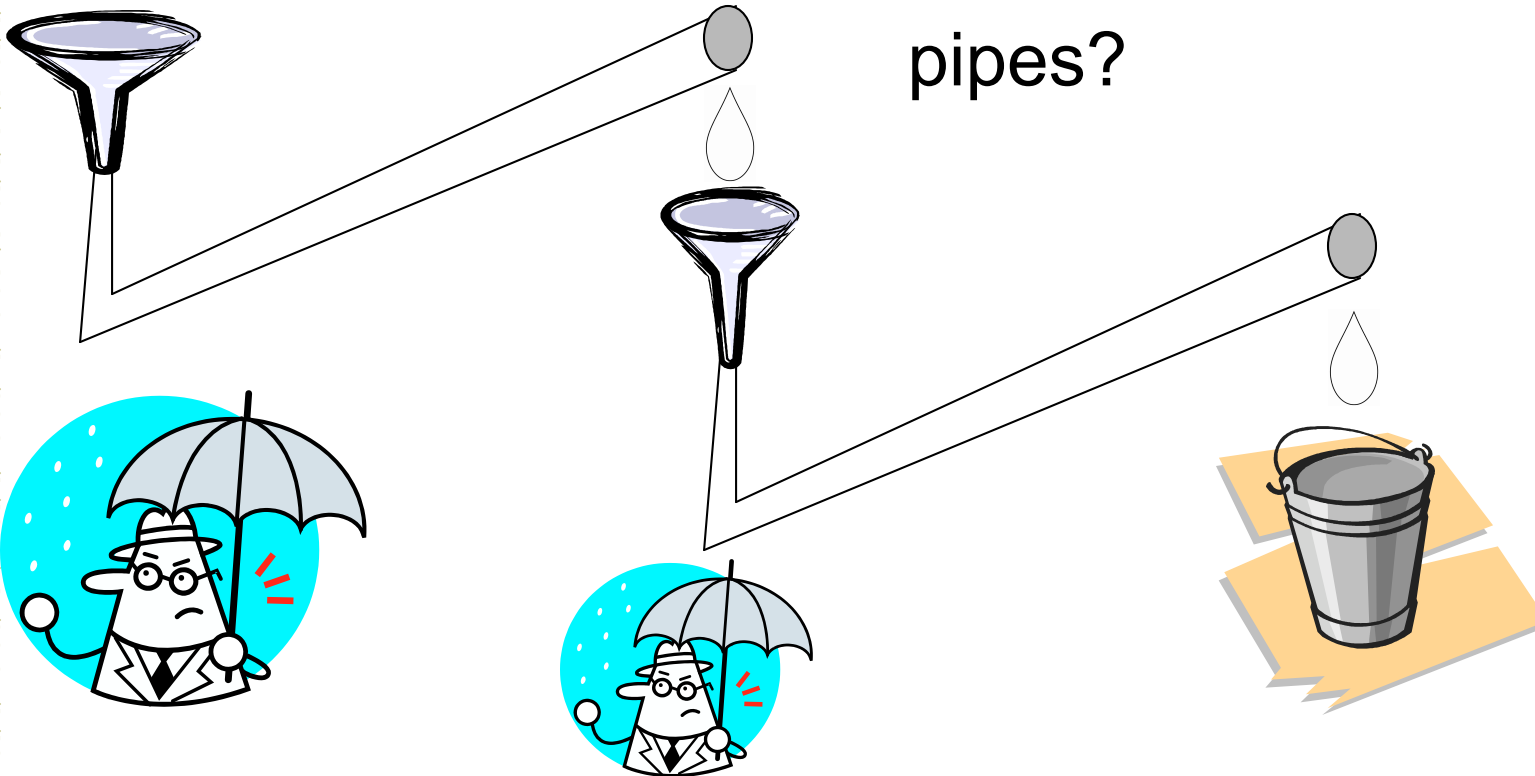


- How should you control the faucet?
  - Too fast – sink overflows
  - Too slow – what happens?
- Goals
  - Fill the bucket as quickly as possible
  - Avoid overflowing the sink
- Solution – watch the sink

# Plumbers Gone Wild!

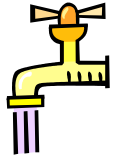
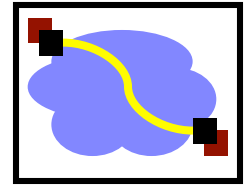


- How do we prevent water loss?
- Know the size of the pipes?

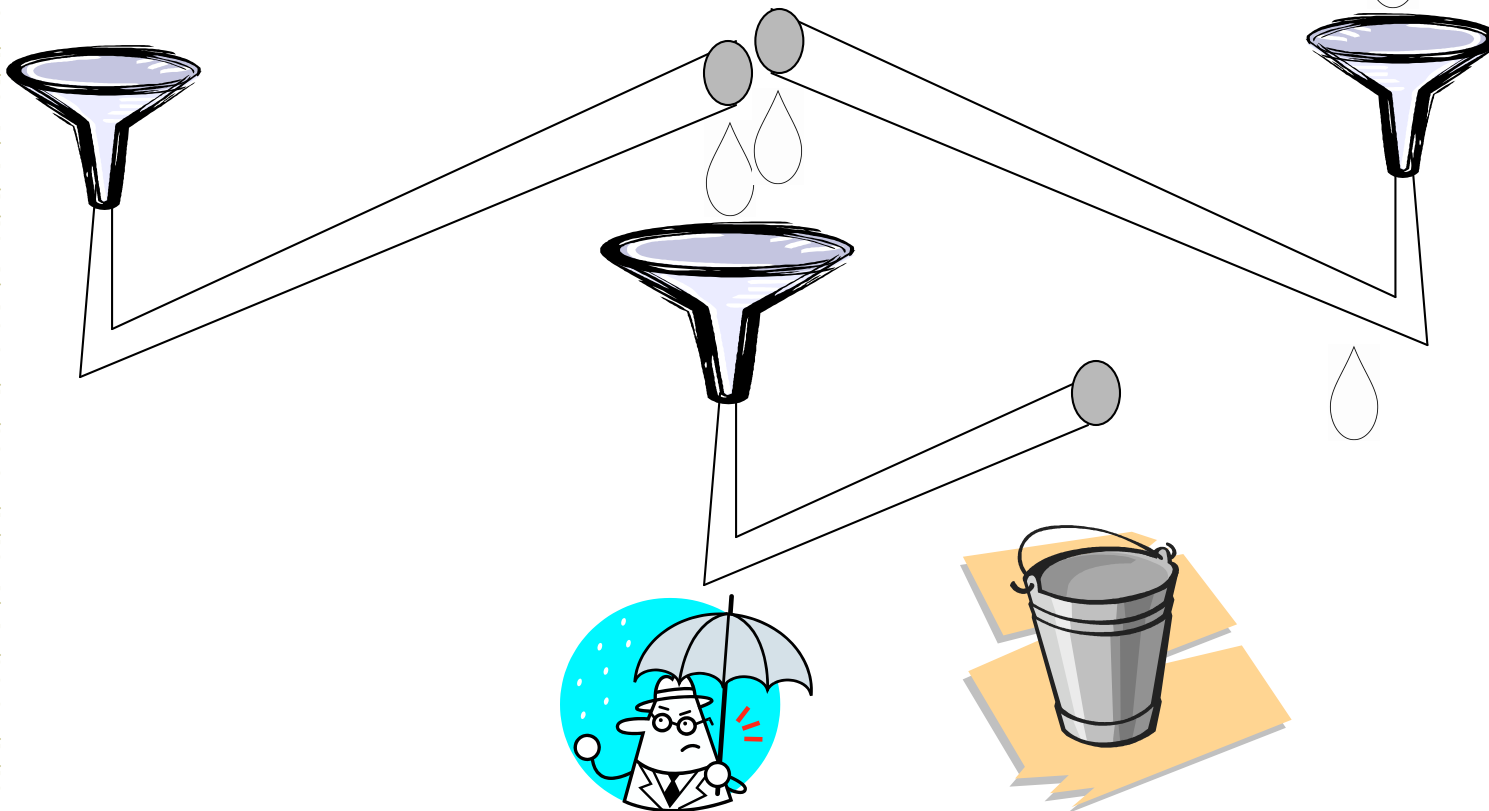




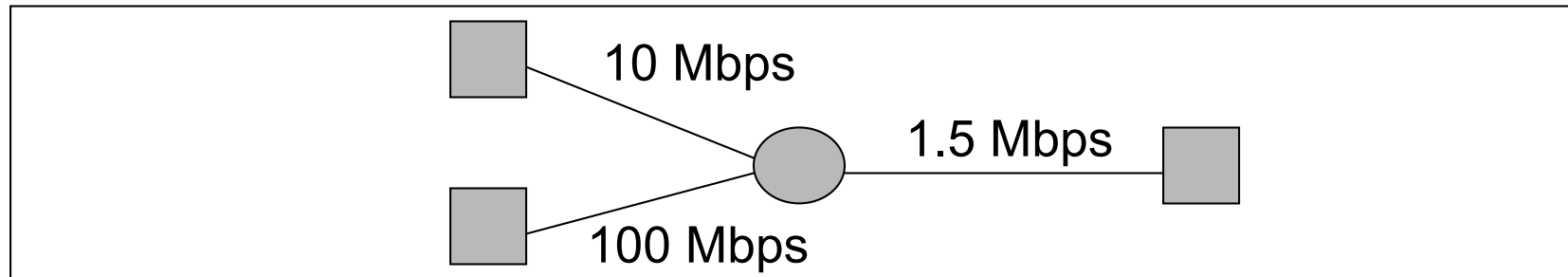
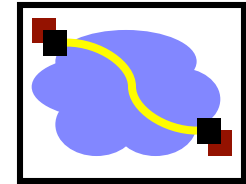
# Plumbers Gone Wild 2!



- Now what?
- Feedback from the bucket or the funnels?

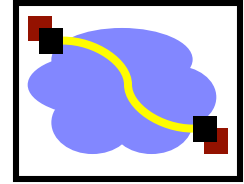


# Congestion



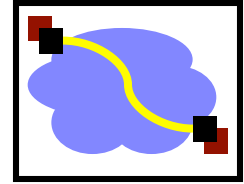
- Different sources compete for resources inside network
- Why is it a problem?
  - Sources are unaware of current state of resource
  - Sources are unaware of each other
- Manifestations:
  - Lost packets (buffer overflow at routers)
  - Long delays (queuing in router buffers)
  - Can result in throughput less than bottleneck link (1.5Mbps for the above topology) → a.k.a. congestion collapse

# Congestion Collapse



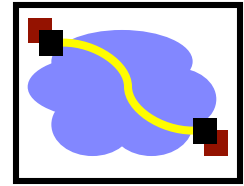
- Definition: *Increase in network load results in decrease of useful work done*
- Many possible causes
  - Spurious retransmissions of packets still in flight
    - Classical congestion collapse
    - How can this happen with packet conservation
    - Solution: better timers and TCP congestion control
  - Undelivered packets
    - Packets consume resources and are dropped elsewhere in network
    - Solution: congestion control for ALL traffic

# Congestion Control and Avoidance



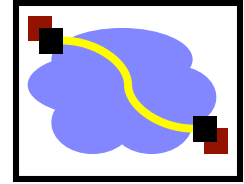
- A mechanism which:
  - Uses network resources efficiently
  - Preserves fair network resource allocation
  - Prevents or avoids collapse
- Congestion collapse is not just a theory
  - Has been frequently observed in many networks

# Approaches Towards Congestion Control



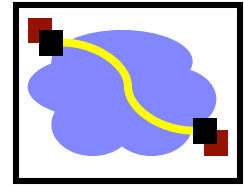
- Two broad approaches towards congestion control:
- **End-end congestion control:**
  - No explicit feedback from network
  - Congestion inferred from end-system observed loss, delay
  - Approach taken by TCP
- **Network-assisted congestion control:**
  - Routers provide feedback to end systems
    - Single bit indicating congestion (SNA, DECbit, TCP/IP ECN, ATM)
    - Explicit rate sender should send at
  - Problem: makes routers complicated

# Example: TCP Congestion Control



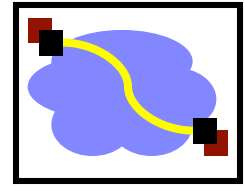
- Very simple mechanisms in network
  - FIFO scheduling with shared buffer pool
  - Feedback through packet drops
- TCP interprets packet drops as signs of congestion and slows down
  - This is an assumption: packet drops are not a sign of congestion in all networks
    - E.g. wireless networks
- Periodically probes the network to check whether more bandwidth has become available.

# Important Lessons



- Transport service
  - UDP → mostly just IP service
  - TCP → congestion controlled, reliable, byte stream
- Types of ARQ protocols
  - Stop-and-wait → slow, simple
  - Go-back-n → can keep link utilized (except w/ losses)
  - Selective repeat → efficient loss recovery
- Sliding window flow control
- TCP flow control
  - Sliding window → mapping to packet headers
  - 32bit sequence numbers (bytes)

## Important Lessons



- Why is congestion control needed?
- Next paper: How to evaluate congestion control algorithms?
  - Why is AIMD the right choice for congestion control?
- Later: Is AIMD always the right choice? (XCP)