

Congestion Control

Arvind Krishnamurthy

Material based on courses at Stanford, Princeton, and MIT

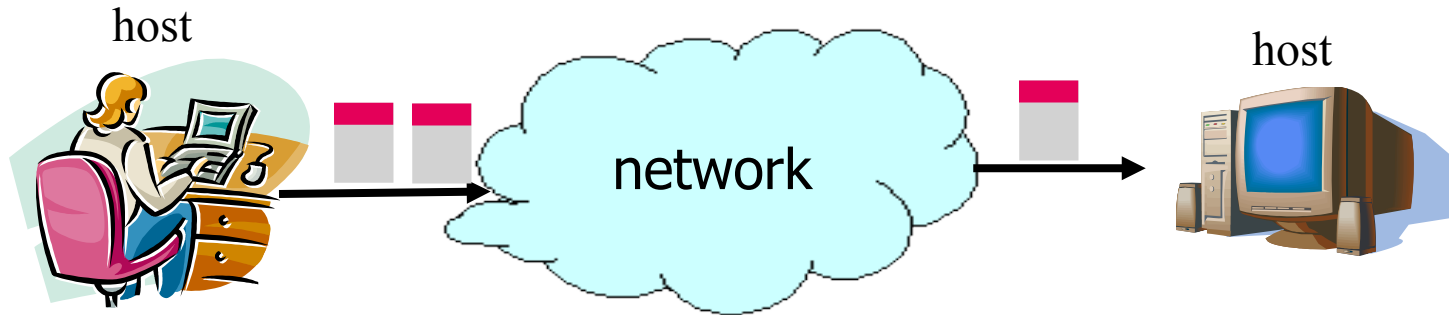
Recap

- Intro to networks
- Two interfaces: socket & traceroute
- Protocols and layering

Best-Effort Packet-Delivery Service

Host-Network Division of Labor

- Packet switching
 - Divide messages into a sequence of packets
 - Headers with source and destination address
- Best-effort delivery
 - Packets may be lost
 - Packets may be corrupted
 - Packets may be delivered out of order



- What is an alternative to packet switching?
- What are the advantages of packet switching over other approaches?

Host-Network Interface: Why Packets?

- Data traffic is bursty
- Don't want to waste bandwidth
 - No traffic exchanged during idle periods
- Better to allow multiplexing
 - Different transfers share access to same links
- Packets can be delivered by most anything
 - RFC 1149: IP Datagrams over Avian Carriers

Host-Network Interface: Why Best-Effort?

- Never having to say you're sorry...
 - Don't reserve bandwidth and memory
 - Don't do error detection & correction
 - Don't remember from one packet to next
- Easier to survive failures
 - Transient disruptions are okay during failover
- Can run on nearly any link technology
 - Greater interoperability and evolution

Intermediate Transport Layer

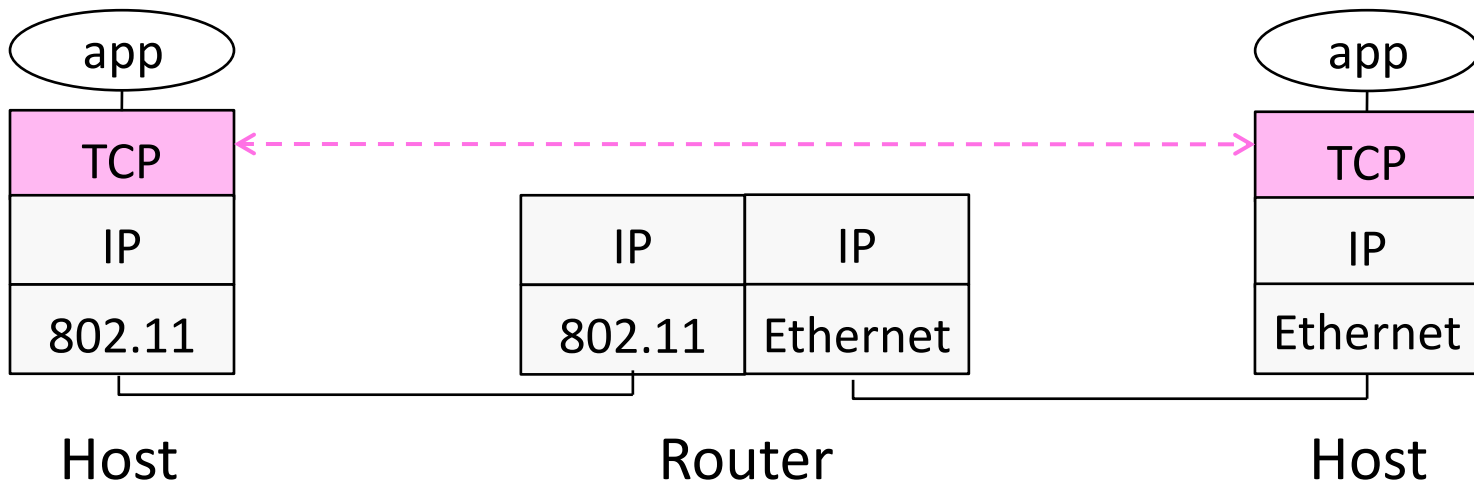
- But, *applications* want efficient, accurate transfer of data in order, in a timely fashion
 - Let the end hosts handle all of that
 - (An example of the “end-to-end argument”)
- Transport layer can optionally...
 - Detect and retransmit lost packets
 - Put out-of-order packets back in order
 - Detect and handle corrupted packets
 - Avoid overloading the receiver
 - <insert your requirement here>

Computer Networks

Transport Layer Overview

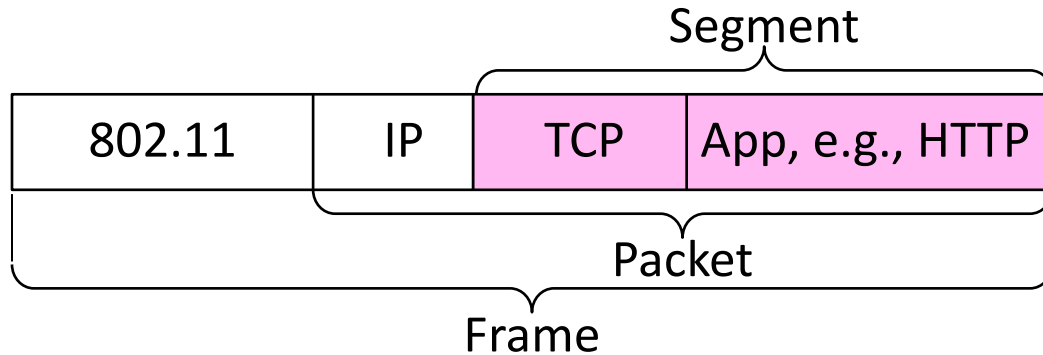
Recall

- 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

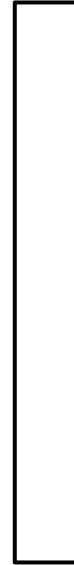
Connection Establishment

- 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
- Chosen to be robust even against delayed duplicates

Active party
(client)

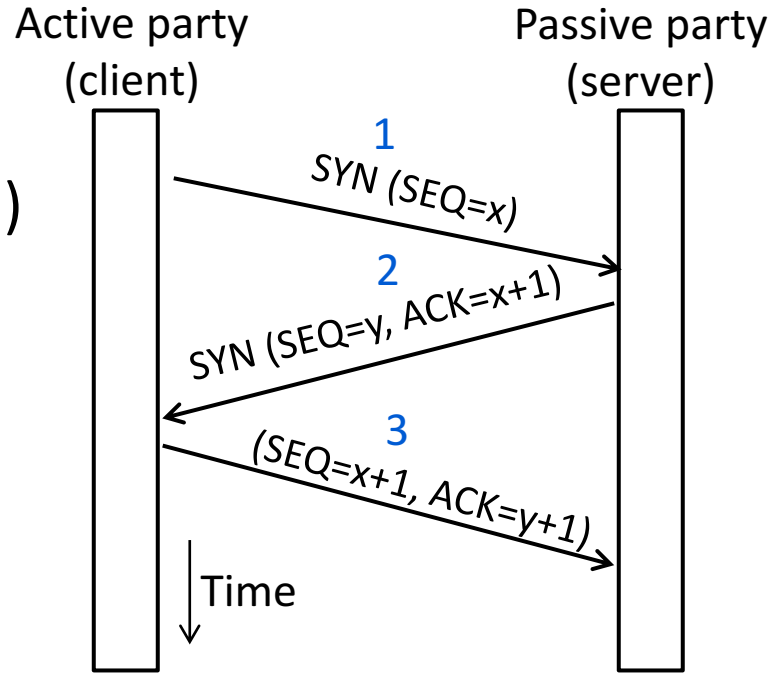


Passive party
(server)



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



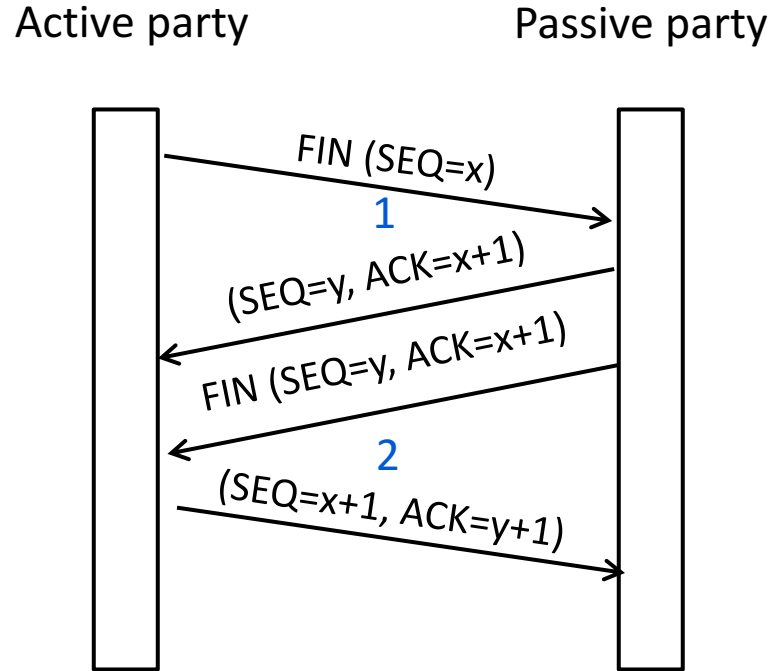
- What are the implications of using sequence numbers?
 - What happens when packets are delayed?
 - What are the security issues here?

Connection Release

- Orderly release by both parties when done
 - Delivers all pending data and “hangs up”
 - Cleans up state in sender and receiver
- 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), ACKs
 - Passive sends FIN(y), ACKs
 - FINs are retransmitted if lost
- Each FIN/ACK closes one direction of data transfer

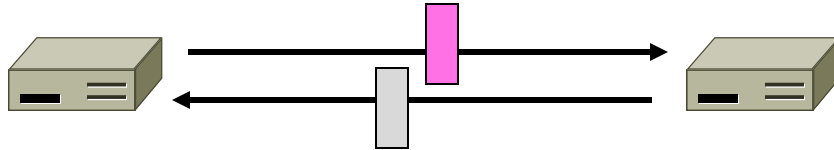


Computer Networks

Sliding Windows

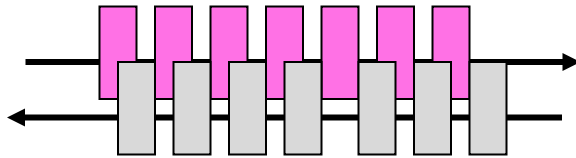
Limitation of Stop-and-Wait

- 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 $BD \gg 1$ packet



Sliding Window

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



- Pipelining improves performance
- Need $W=2BD/\text{pkt_size}$ to fill network path

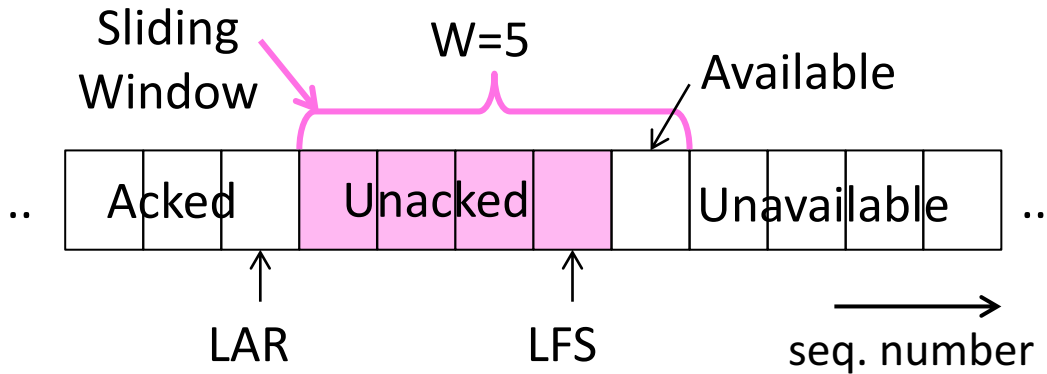
- Bandwidth-delay product (BDP) is an useful concept
- How would a datacenter differ from a wide-area Internet in terms of BDP?
 - What are the implications?

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

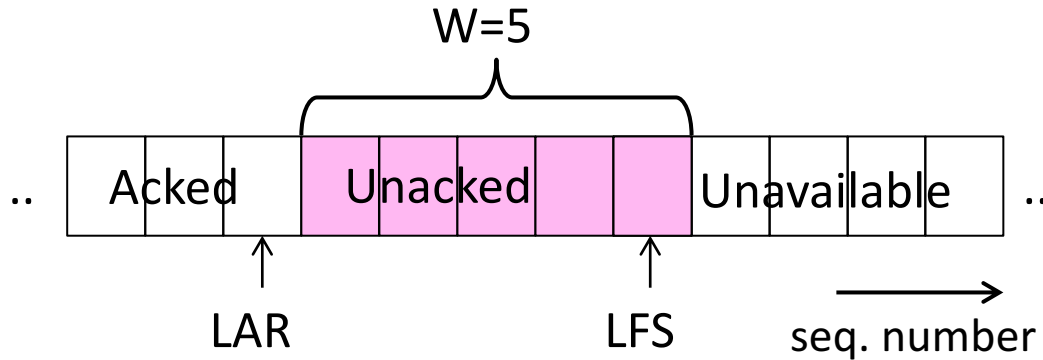
Sliding Window – Sender

- Sender buffers up to W segments until they are acknowledged
 - LFS=LAST FRAME SENT, LAR=LAST ACK REC'D
 - Sends while $LFS - LAR \leq W$



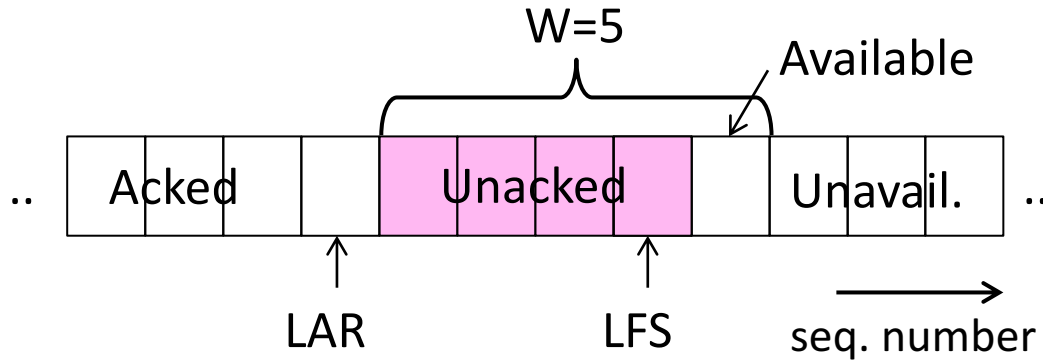
Sliding Window – Sender (2)

- Transport accepts another segment of data from the Application ...
 - Transport sends it (as LFS–LAR \rightarrow 5)



Sliding Window – Sender (3)

- Next higher ACK arrives from peer...
 - Window advances, buffer is freed
 - LFS–LAR \rightarrow 4 (can send one more)



Sliding Window – Go-Back-N

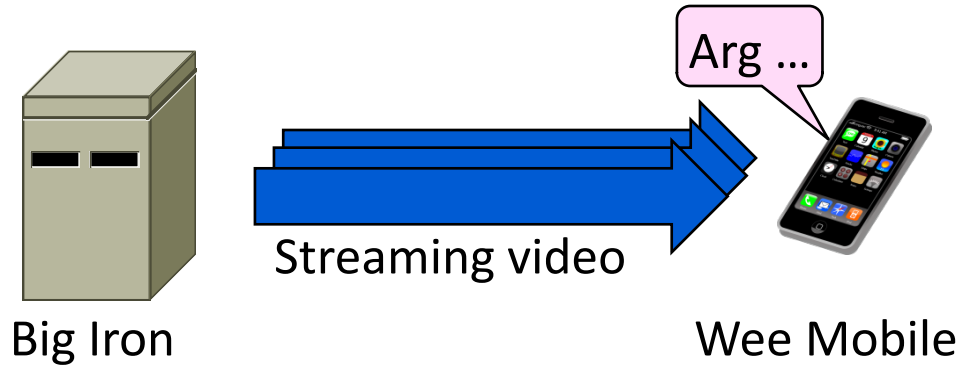
- Receiver keeps only a single packet buffer for the next segment
 - State variable, $LAS = \text{LAST ACK SENT}$
- On receive:
 - If seq. number is $LAS+1$, accept and pass it to app, update LAS , send ACK
 - Otherwise discard (as out of order)

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'll see the details later

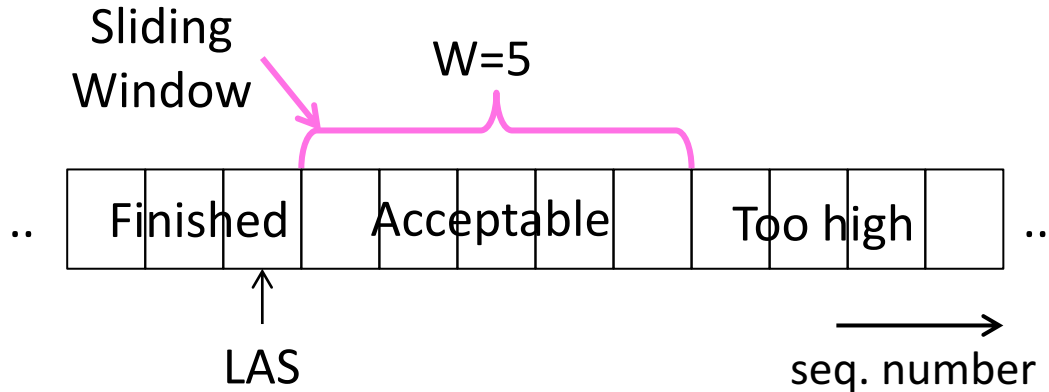
Problem

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



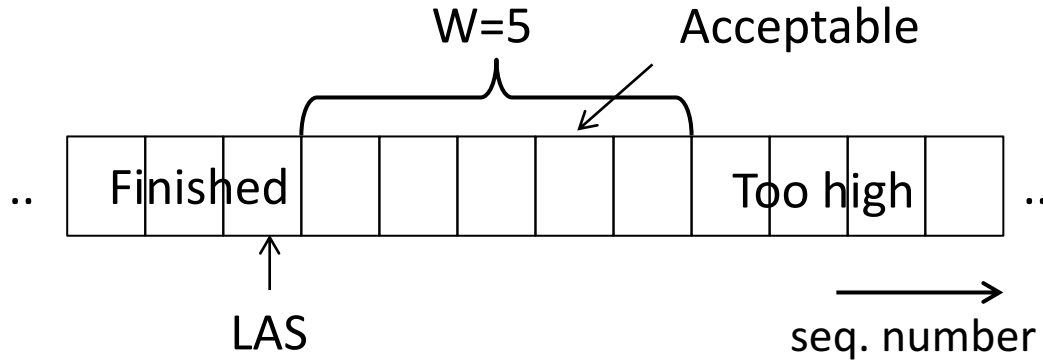
Sliding Window – Receiver

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



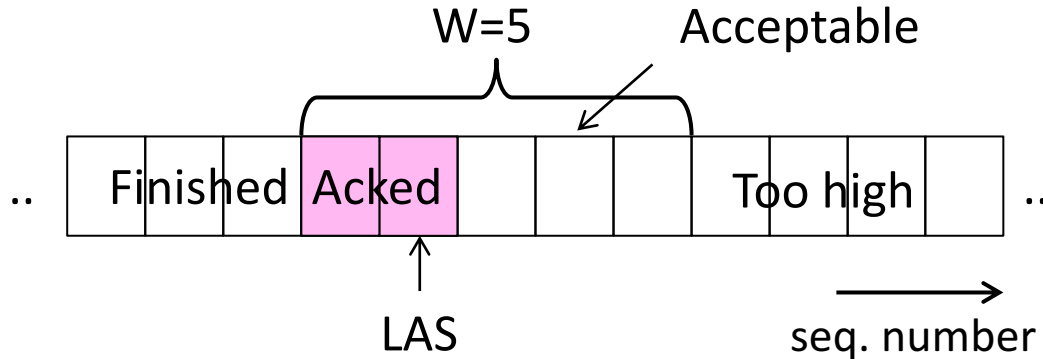
Sliding Window – Receiver (2)

- Suppose the next two segments arrive but app does not call `recv()`



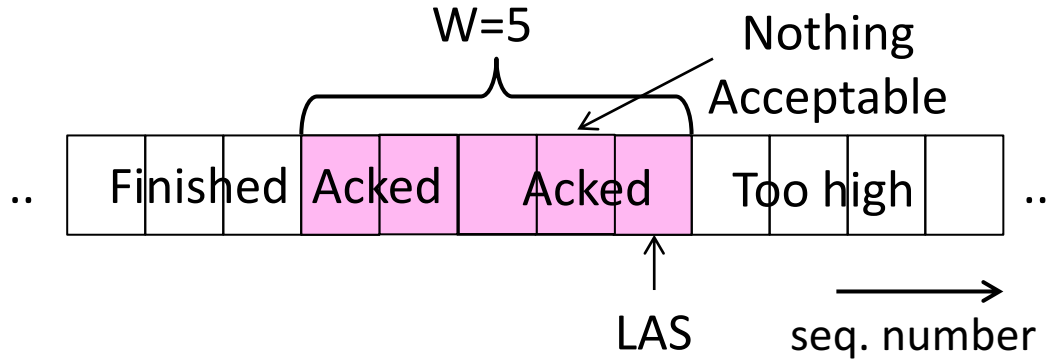
Sliding Window – Receiver (3)

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



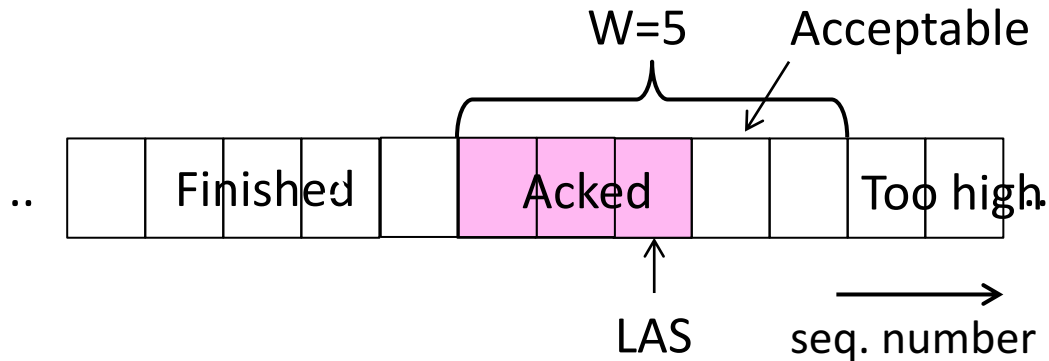
Sliding Window – Receiver (4)

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



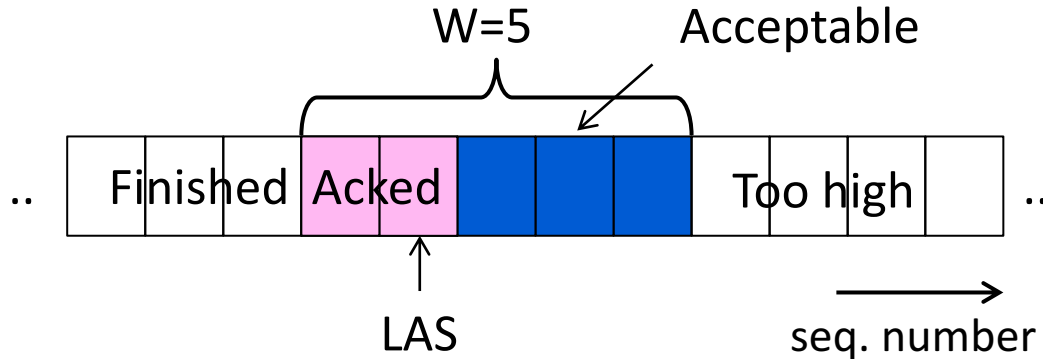
Sliding Window – Receiver (5)

- App recv() takes two segments
 - Window slides



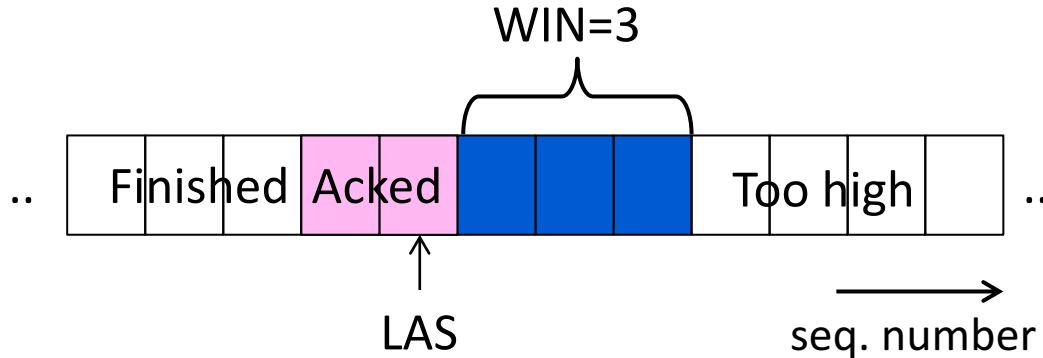
Flow Control

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



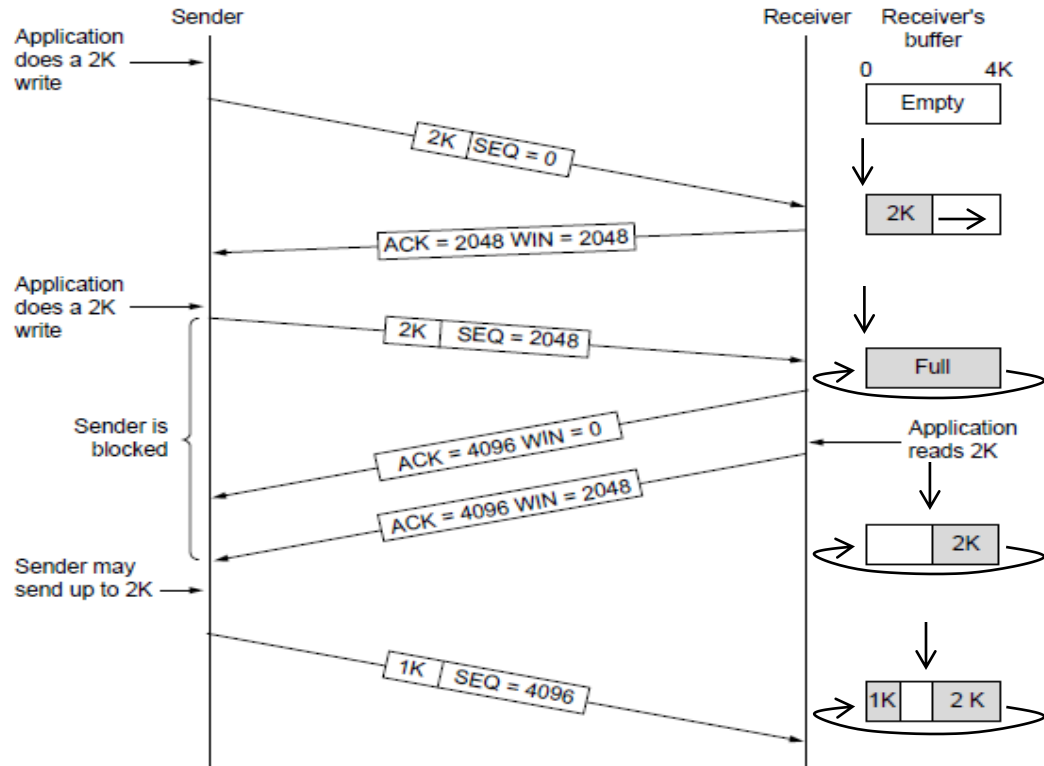
Flow Control (2)

- 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
 - $SEQ + \text{length} < ACK + WIN$
 - 4KB buffer at receiver
 - Circular buffer of bytes

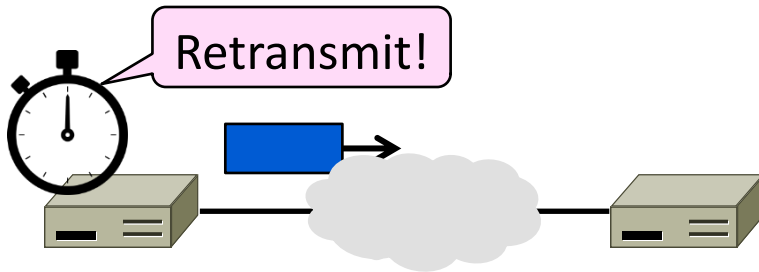


Computer Networks

Retransmission Timeouts

Retransmissions

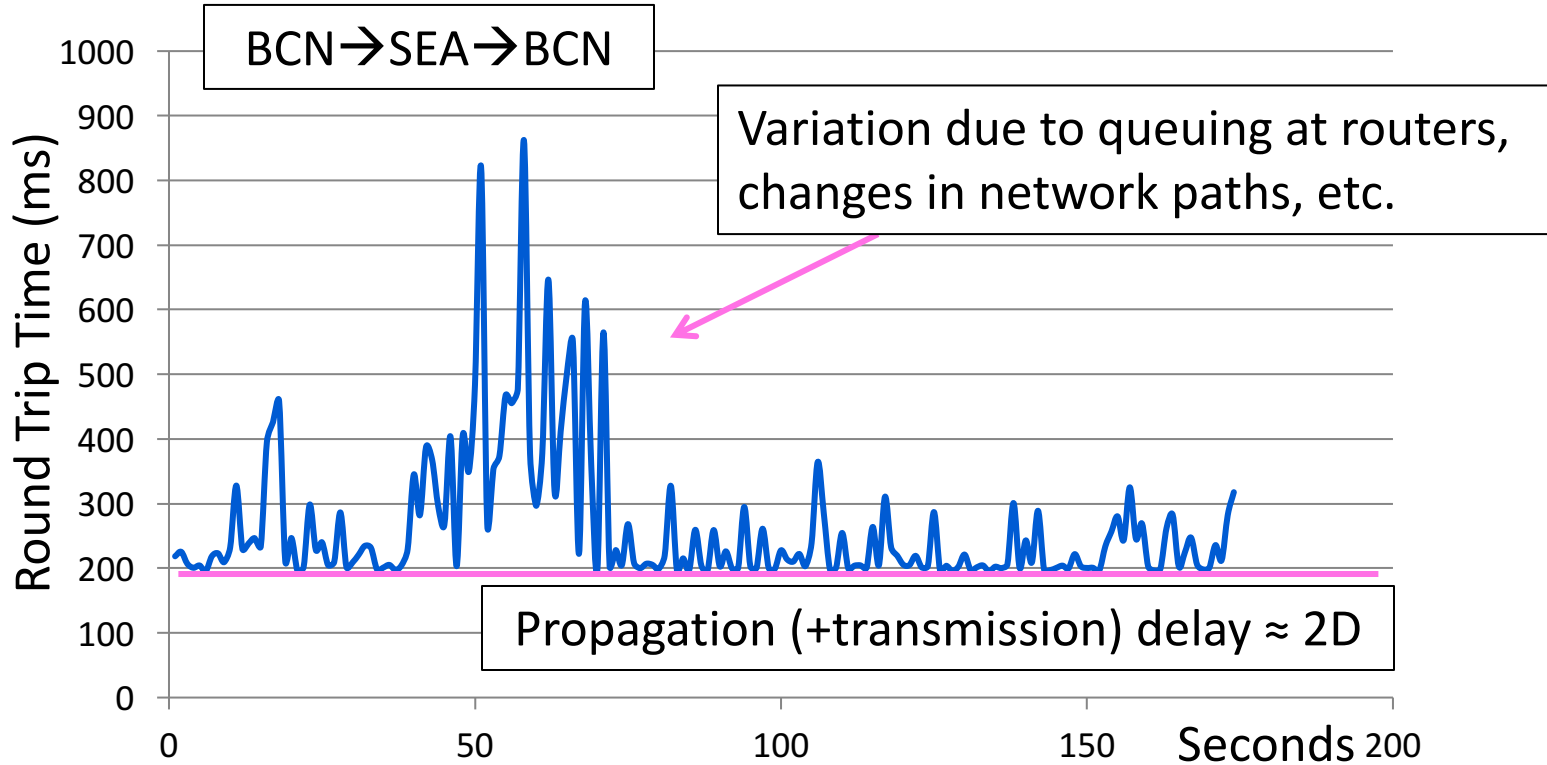
- With sliding window, the strategy for detecting loss is the timeout
 - Set timer when a segment is sent
 - Cancel timer when ack is received
 - If timer fires, retransmit data as lost



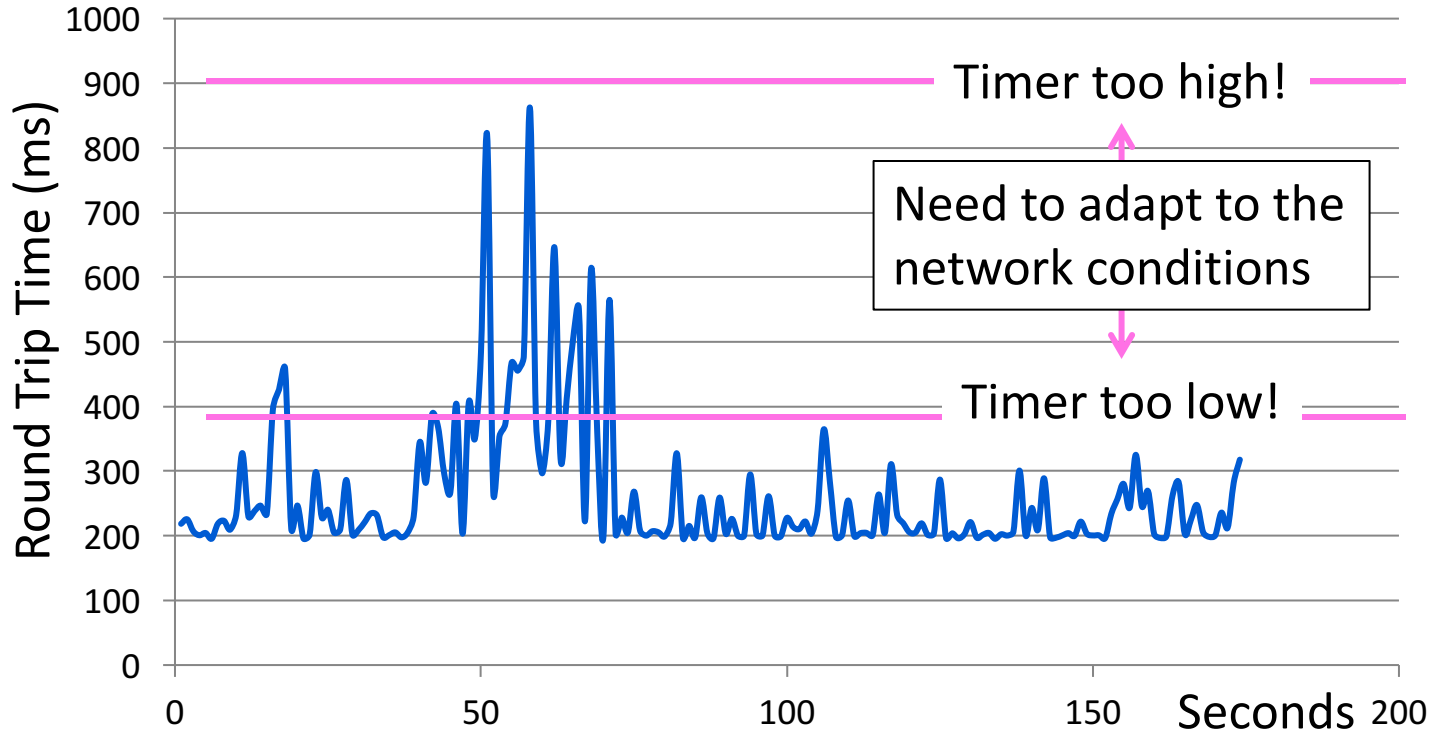
Timeout Problem

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

Example of RTTs



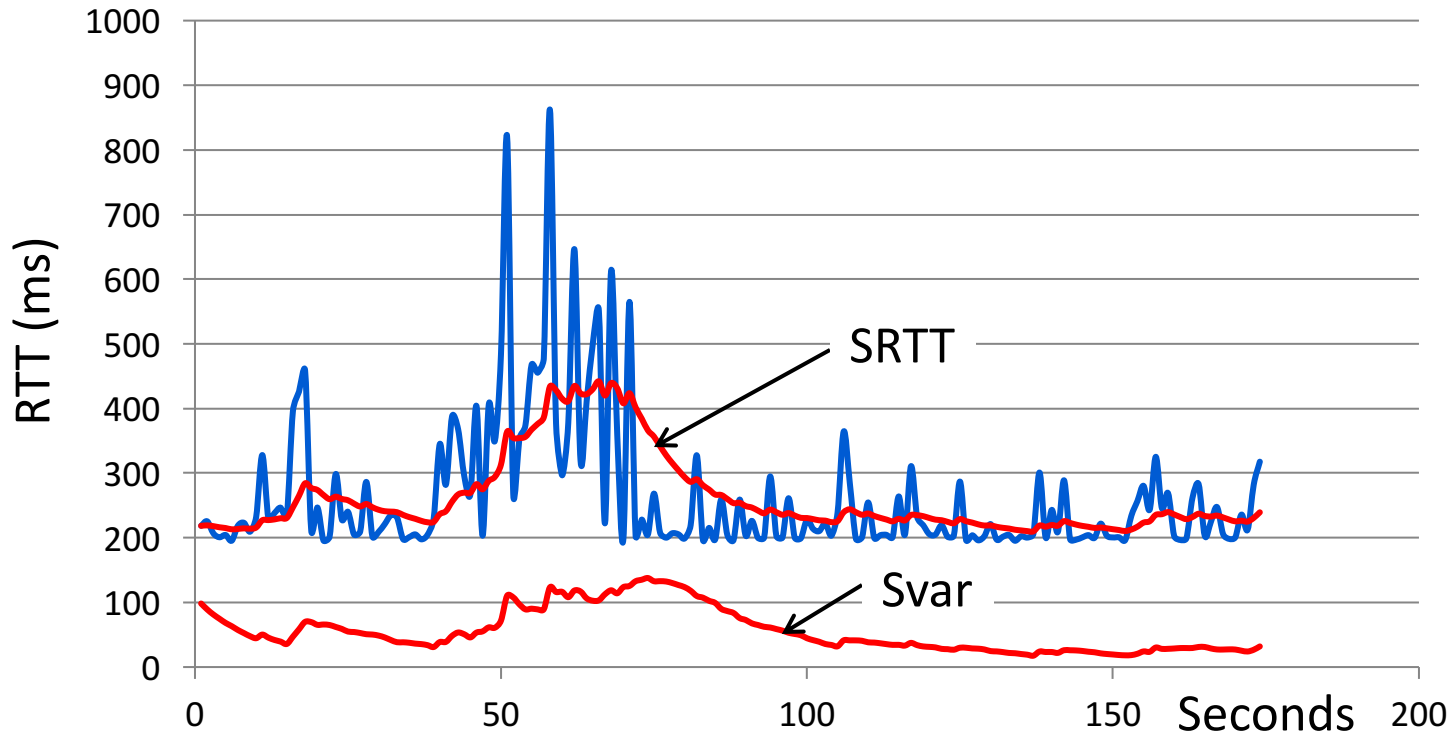
Example of RTTs



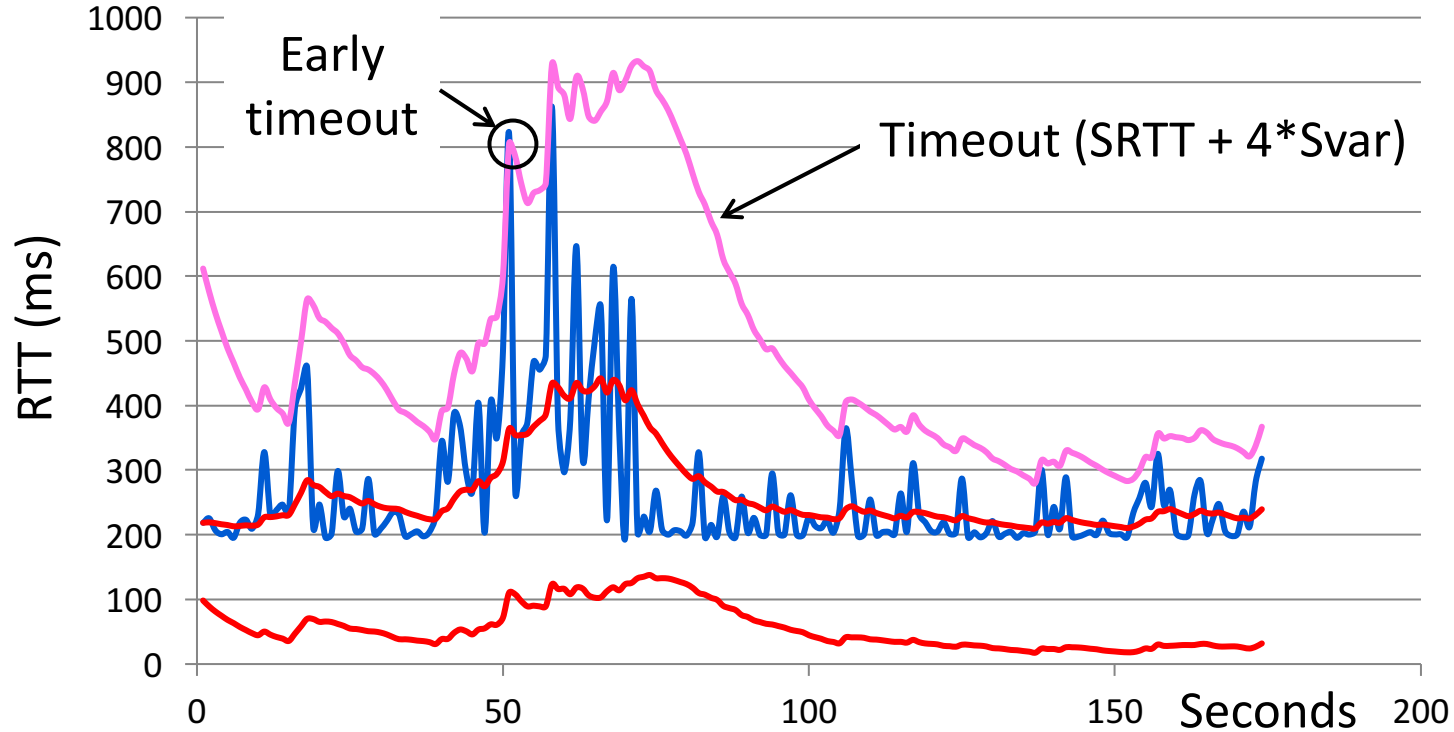
Adaptive Timeout

- 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)

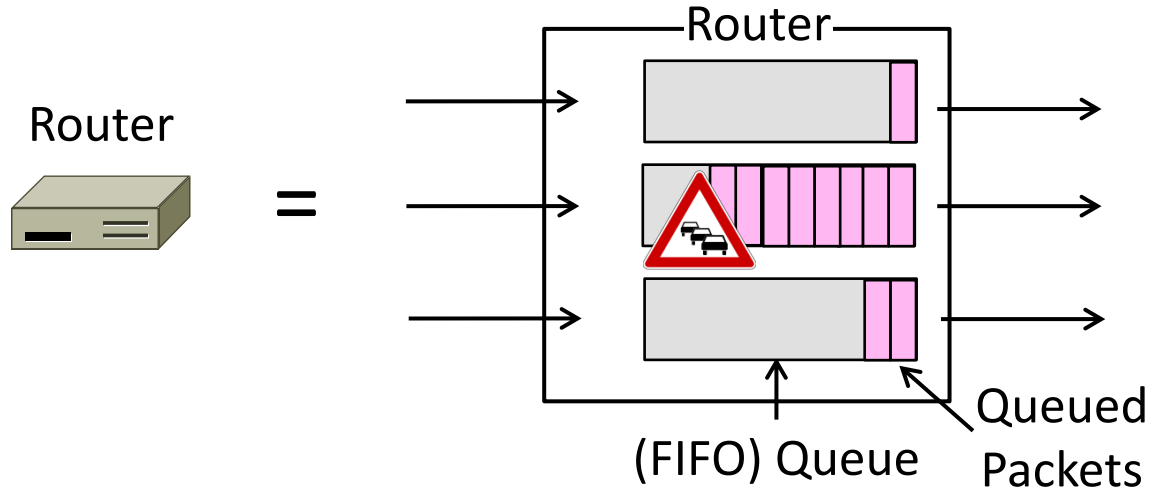


Computer Networks

Congestion Overview

Nature of Congestion

- Simplified view of per port output queues
 - Typically FIFO (First In First Out), discard when full

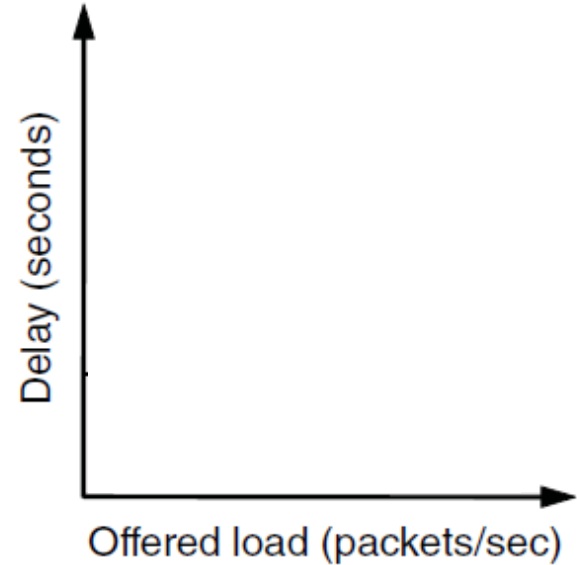
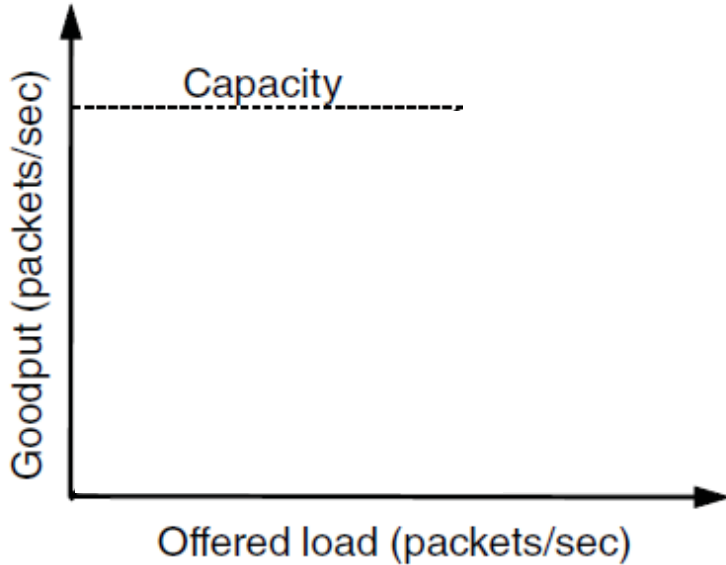


Nature of Congestion (2)

- Queues help by absorbing bursts when input $>$ output rate
- But if input $>$ output rate persistently, queue will overflow
 - This is congestion
- Congestion is a function of the traffic patterns – can occur even if every link have the same capacity

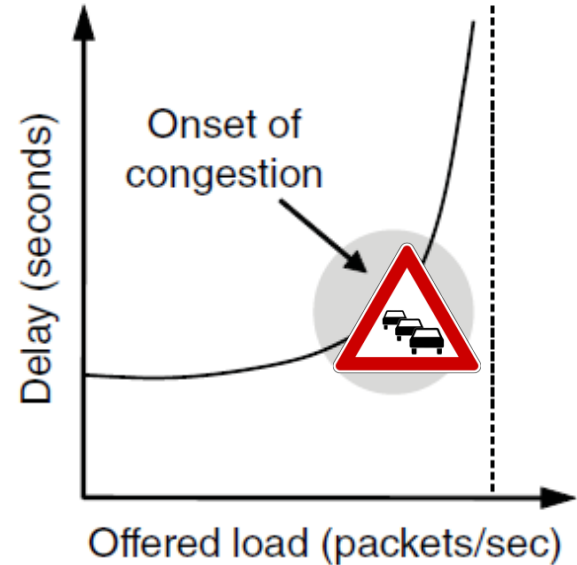
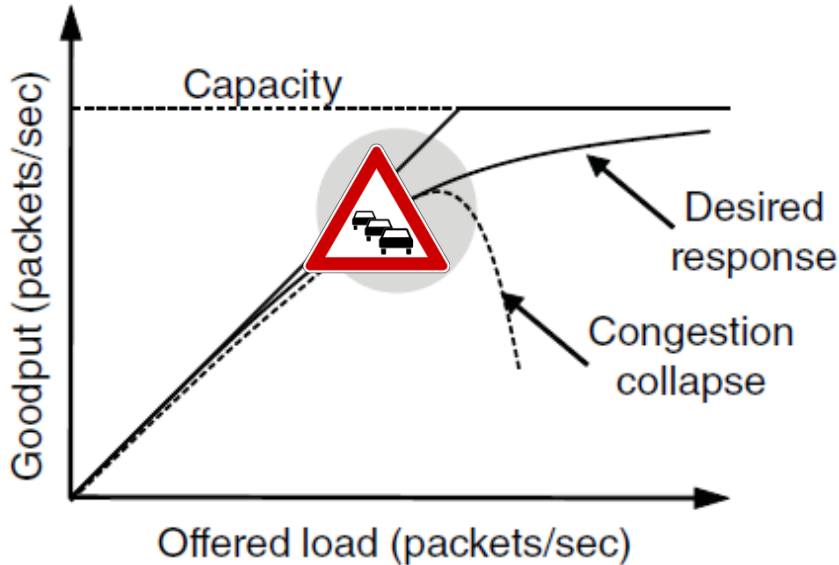
Effects of Congestion

- What happens to performance as we increase the load?



Effects of Congestion (2)

- What happens to performance as we increase the load?



Effects of Congestion

- As offered load rises, congestion occurs as queues begin to fill:
 - Delay and loss rise sharply with more load
 - Throughput falls below load (due to loss)
 - Goodput may fall below throughput (due to spurious retransmissions)
- None of the above is good!
 - Want to operate network just before the onset of congestion



Bandwidth Allocation

- Important task for network is to allocate its capacity to senders
 - Good allocation is efficient and fair
- Efficient means most capacity is used but there is no congestion
- Fair means every sender gets a reasonable share the network

Bandwidth Allocation (2)

- Why is it hard? (Just split equally!)
 - Number of senders and their offered load is constantly changing
 - Senders may lack capacity in different parts of the network
 - Network is distributed; no single party has an overall picture of its state

Bandwidth Allocation (3)

- Key observation:
 - In an effective solution, Transport and Network layers must work together
- Network layer witnesses congestion
 - Only it can provide direct feedback
- Transport layer causes congestion
 - Only it can reduce offered load

Bandwidth Allocation (4)

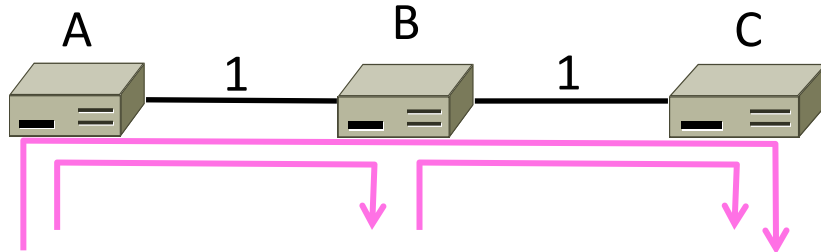
- Solution context:
 - Senders adapt concurrently based on their own view of the network
 - Design this adaptation so the network usage as a whole is efficient and fair
 - Adaptation is continuous since offered loads continue to change over time

Computer Networks

Fairness of Bandwidth
Allocation

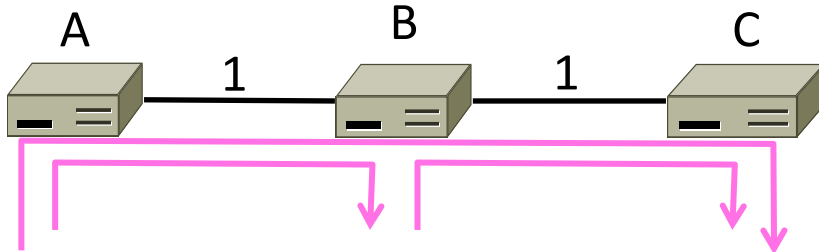
Efficiency vs. Fairness

- Cannot always have both!
 - Example network with traffic $A \rightarrow B$, $B \rightarrow C$ and $A \rightarrow C$
 - How much traffic can we carry?



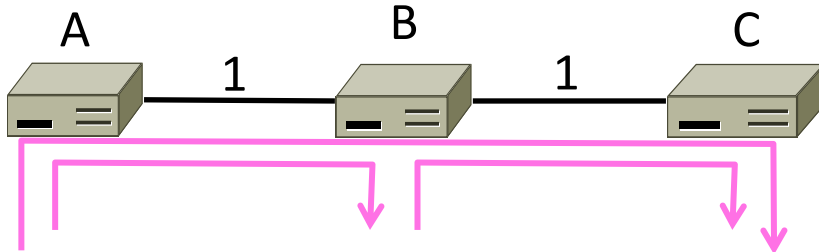
Efficiency vs. Fairness (2)

- If we care about fairness:
 - Give equal bandwidth to each flow
 - $A \rightarrow B$: $\frac{1}{2}$ unit, $B \rightarrow C$: $\frac{1}{2}$, and $A \rightarrow C$, $\frac{1}{2}$
 - Total traffic carried is $1 \frac{1}{2}$ units



Efficiency vs. Fairness (3)

- If we care about efficiency:
 - Maximize total traffic in network
 - $A \rightarrow B$: 1 unit, $B \rightarrow C$: 1, and $A \rightarrow C$, 0
 - Total traffic rises to 2 units!

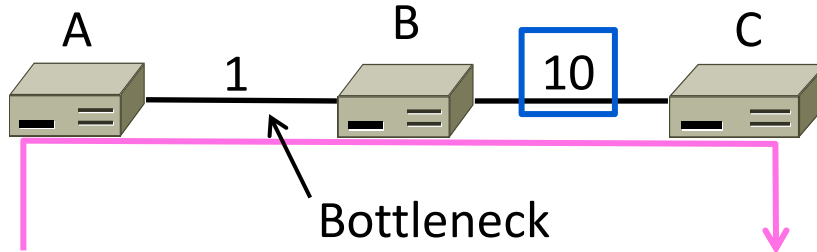


The Slippery Notion of Fairness

- Why is “equal per flow” fair anyway?
 - $A \rightarrow C$ uses more network resources (two links) than $A \rightarrow B$ or $B \rightarrow C$
 - Host A sends two flows, B sends one
- Not productive to seek exact fairness
 - More important to avoid starvation
 - “Equal per flow” is good enough

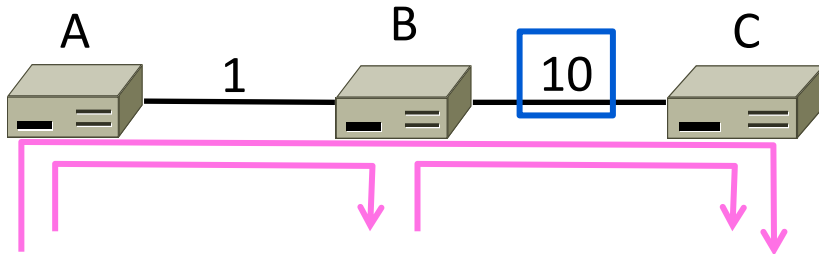
Generalizing “Equal per Flow”

- Bottleneck for a flow of traffic is the link that limits its bandwidth
 - Where congestion occurs for the flow
 - For $A \rightarrow C$, link A–B is the bottleneck



Generalizing “Equal per Flow” (2)

- Flows may have different bottlenecks
 - For $A \rightarrow C$, link $A-B$ is the bottleneck
 - For $B \rightarrow C$, link $B-C$ is the bottleneck
 - Can no longer divide links equally ...



Max-Min Fairness

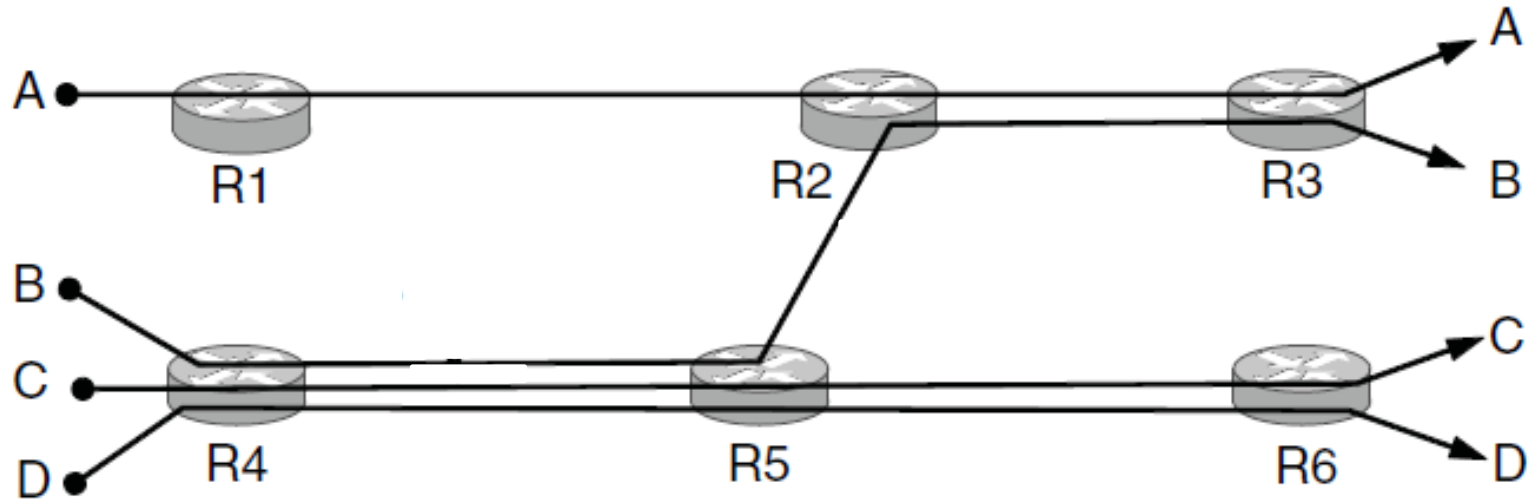
- Intuitively, flows bottlenecked on a link get an equal share of that link
- Max-min fair allocation is one that:
 - Increasing the rate of one flow will decrease the rate of a smaller flow
 - This “maximizes the minimum” flow

Max-Min Fairness (2)

- To find it given a network, imagine “pouring water into the network”
 1. Start with all flows at rate 0
 2. Increase the flows until there is a new bottleneck in the network
 3. Hold fixed the rate of the flows that are bottlenecked
 4. Go to step 2 for any remaining flows

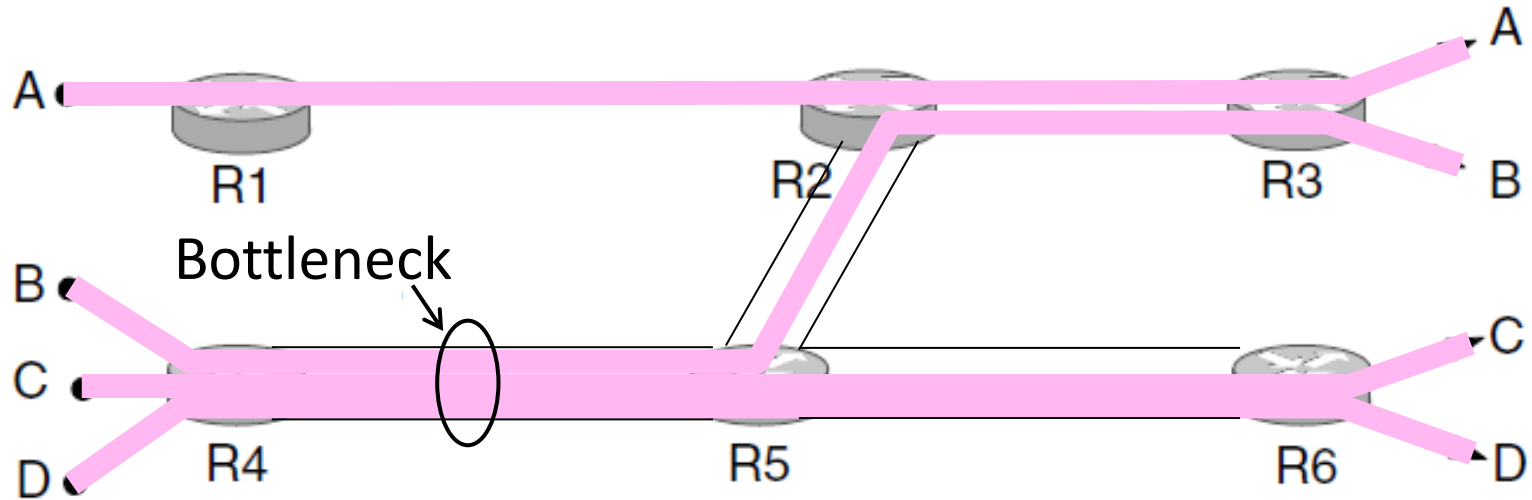
Max-Min Example

- Example: network with 4 flows, links equal bandwidth
 - What is the max-min fair allocation?



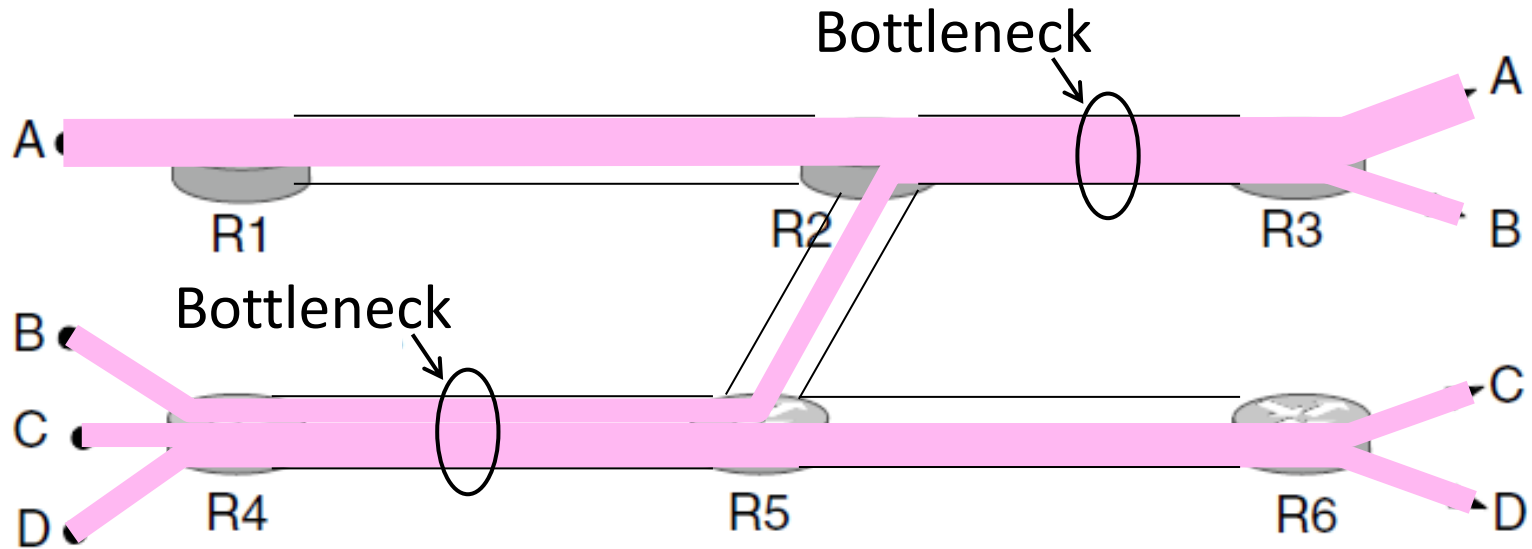
Max-Min Example (2)

- When rate=1/3, flows B, C, and D bottleneck R4—R5
 - Fix B, C, and D, continue to increase A



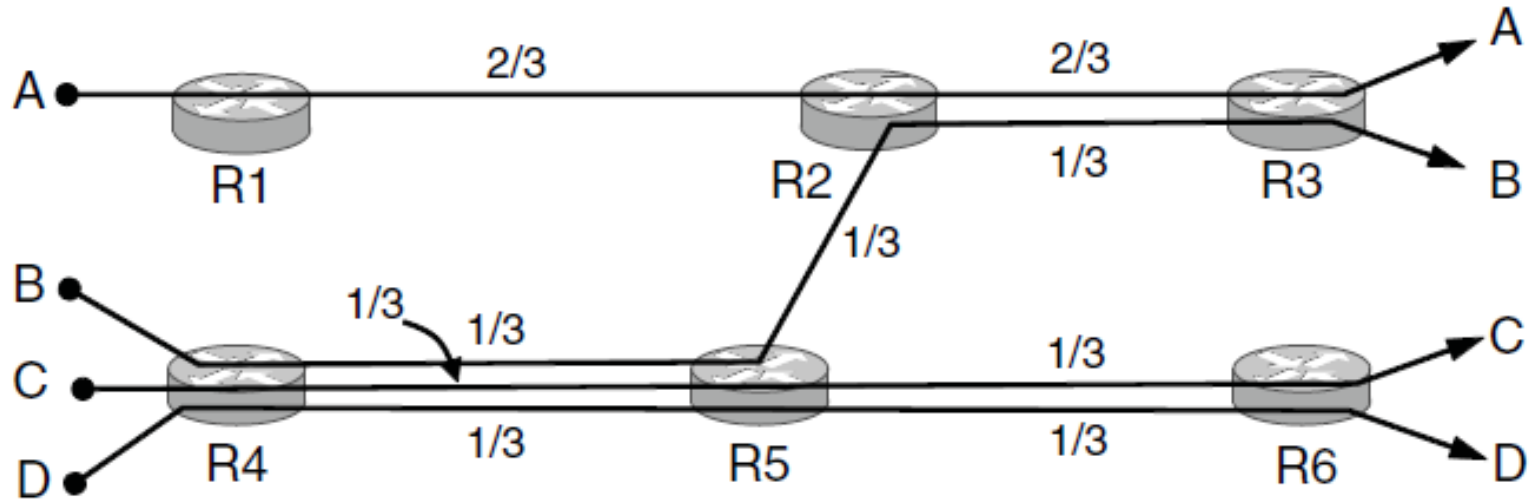
Max-Min Example (3)

- When rate=2/3, flow A bottlenecks R2—R3. Done.



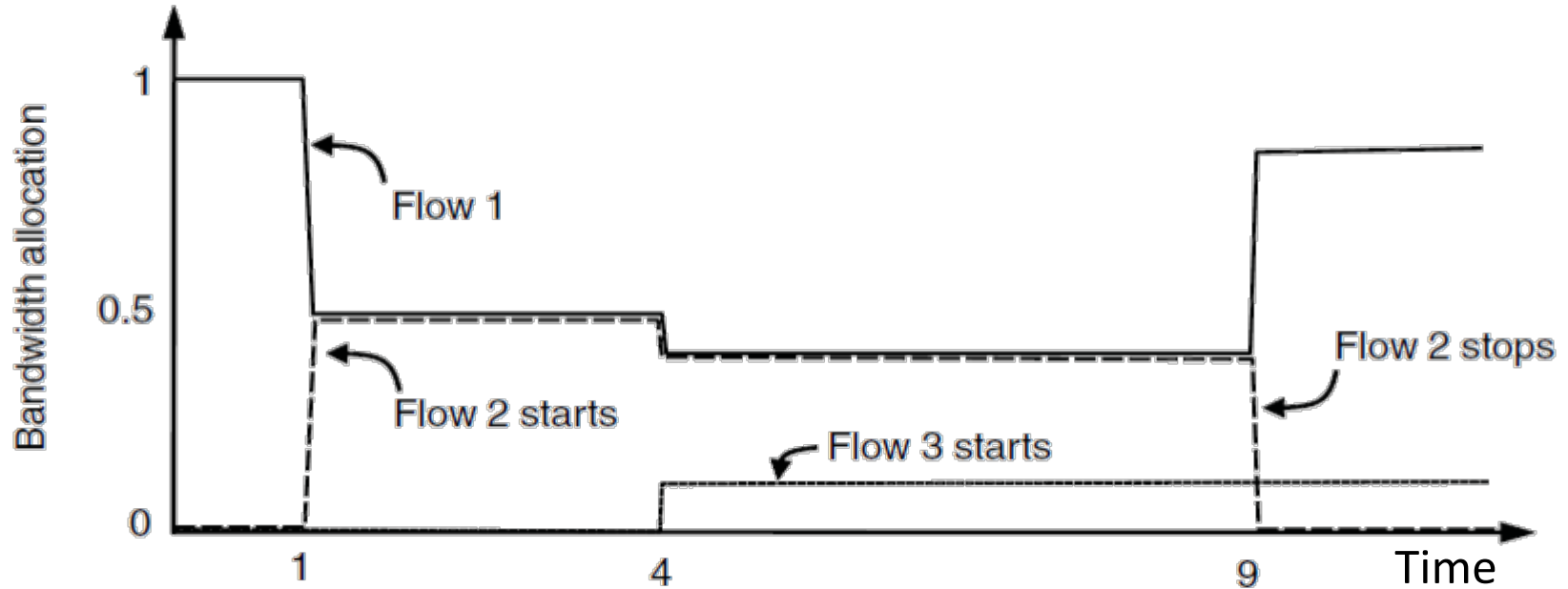
Max-Min Example (4)

- End with $A=2/3$, $B, C, D=1/3$, and $R2-R3$, $R4-R5$ full
 - Other links have extra capacity that can't be used



Adapting over Time

- Allocation changes as flows start and stop



Computer Networks

Additive Increase Multiplicative
Decrease (AIMD)

Bandwidth Allocation Models

- Open loop versus closed loop
 - Open: reserve bandwidth before use
 - Closed: use feedback to adjust rates
- Host versus Network support
 - Who sets/enforces allocations?
- Window versus Rate based
 - How is allocation expressed?

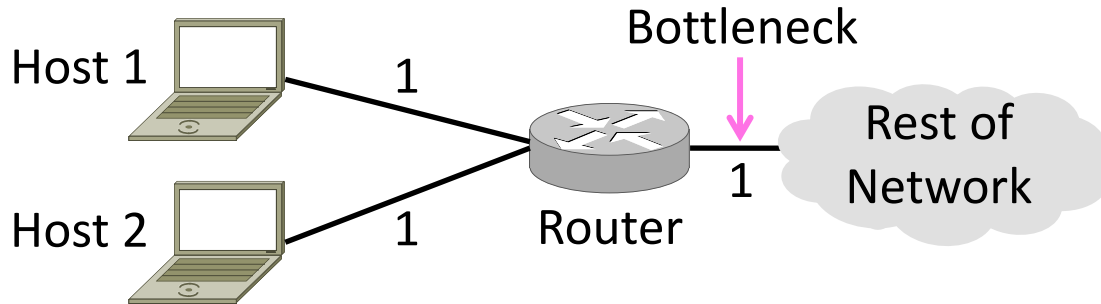
TCP is a closed loop, host-driven, and window-based

Additive Increase Multiplicative Decrease

- AIMD is a control law hosts can use to reach a good allocation
 - Hosts additively increase rate while network is not congested
 - Hosts multiplicatively decrease rate when congestion occurs
 - Used by TCP
- Let's explore the AIMD game ...

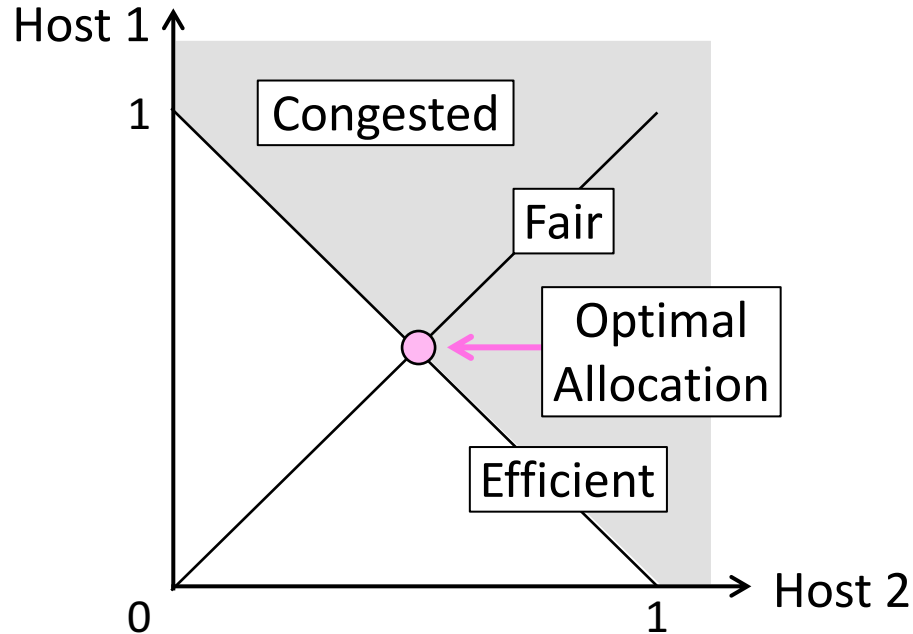
AIMD Game

- Hosts 1 and 2 share a bottleneck
 - But do not talk to each other directly
- Router provides binary feedback
 - Tells hosts if network is congested



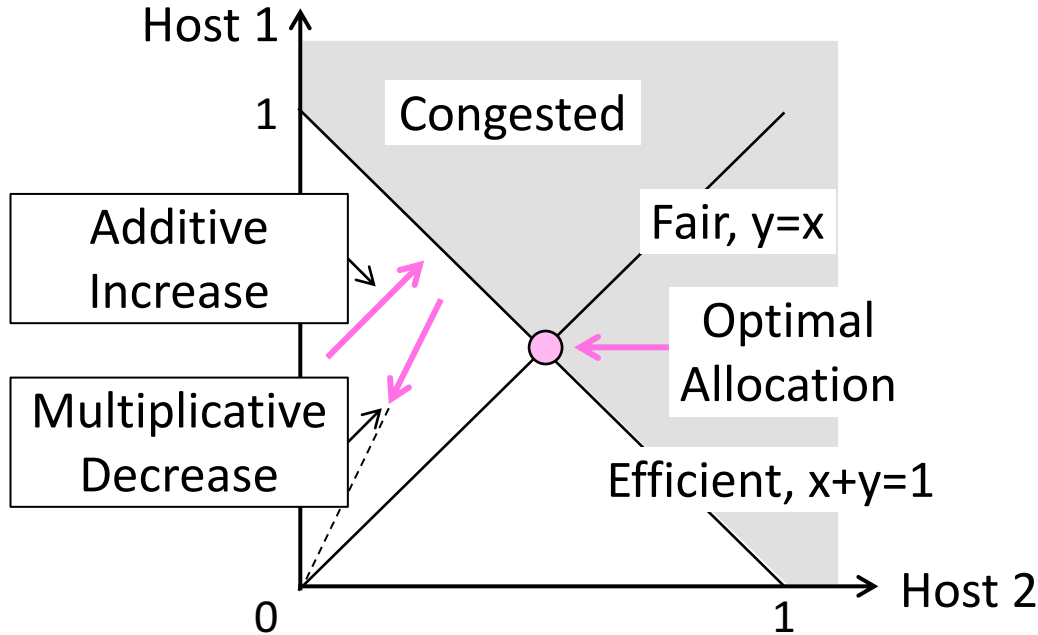
AIMD Game (2)

- Each point is a possible allocation



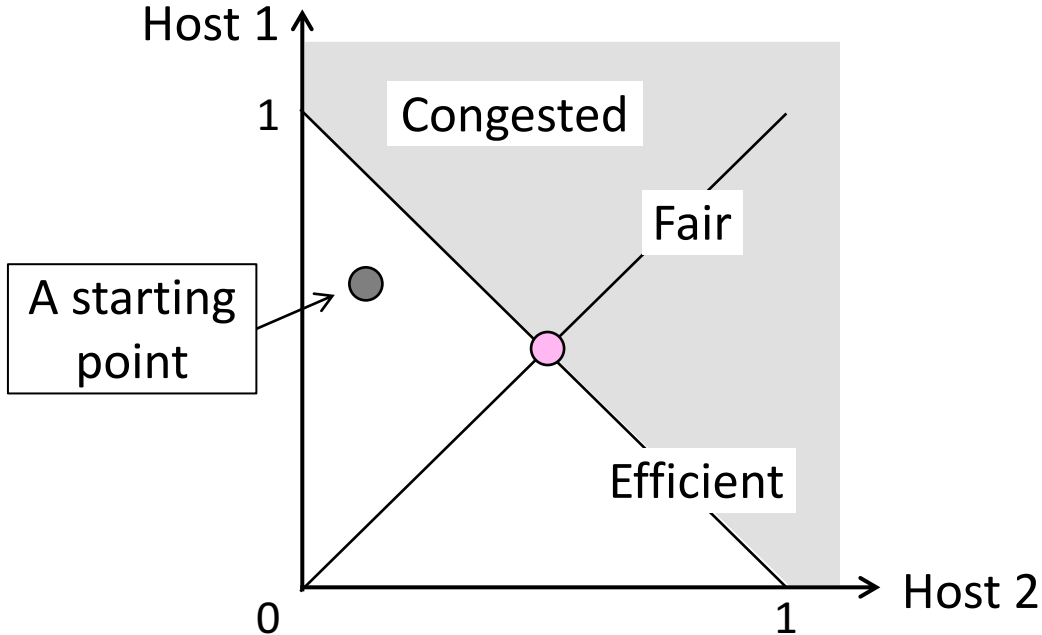
AIMD Game (3)

- AI and MD move the allocation



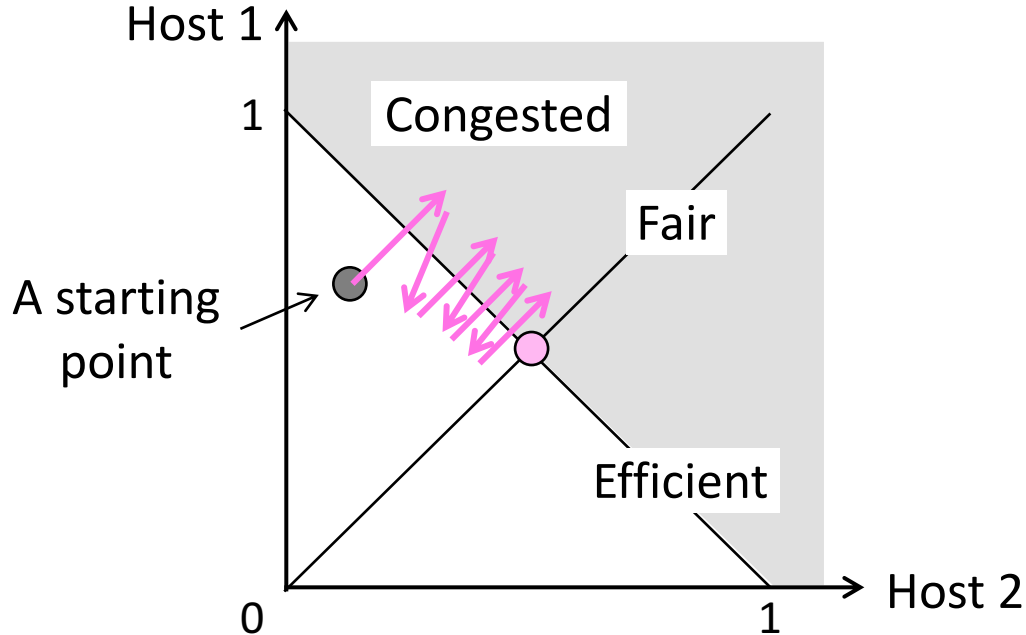
AIMD Game (4)

- Play the game!



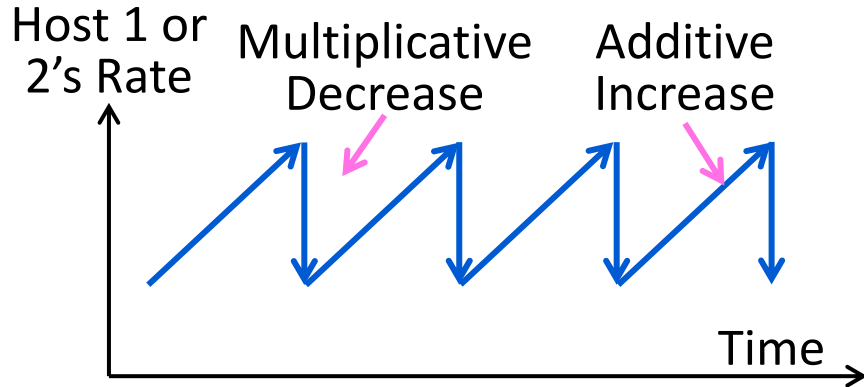
AIMD Game (5)

- Always converge to good allocation!



AIMD Sawtooth

- Produces a “sawtooth” pattern over time for rate of each host
 - This is the TCP sawtooth (later)



- What are alternatives to AIMD?
- What are their strengths/weaknesses?

Feedback Signals

- Several possible signals, with different pros/cons
 - We'll look at classic TCP that uses packet loss as a signal

Signal	Example Protocol	Pros / Cons
Packet loss	TCP NewReno Cubic TCP (Linux)	Hard to get wrong Hear about congestion late
Packet delay	Compound TCP (Windows)	Hear about congestion early Need to infer congestion
Router indication	TCPs with Explicit Congestion Notification	Hear about congestion early Require router support

TCP Tahoe/Reno

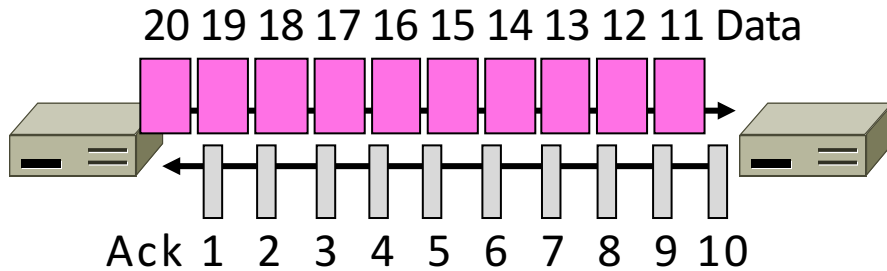
- Avoid congestion collapse without changing routers (or even receivers)
- Idea is to fix timeouts and introduce a congestion window (cwnd) over the sliding window to limit queues/loss
- TCP Tahoe/Reno implements AIMD by adapting cwnd using packet loss as the network feedback signal

TCP Tahoe/Reno (2)

- TCP behaviors we will study:
 - ACK clocking
 - Adaptive timeout (mean and variance)
 - Slow-start
 - Fast Retransmission
 - Fast Recovery
- Together, they implement AIMD

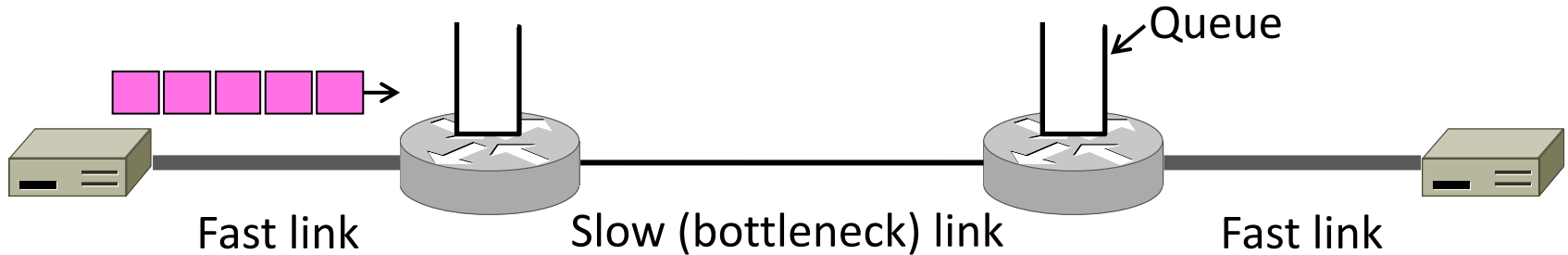
Sliding Window ACK Clock

- Each in-order ACK advances the sliding window and lets a new segment enter the network
 - ACKs “clock” data segments



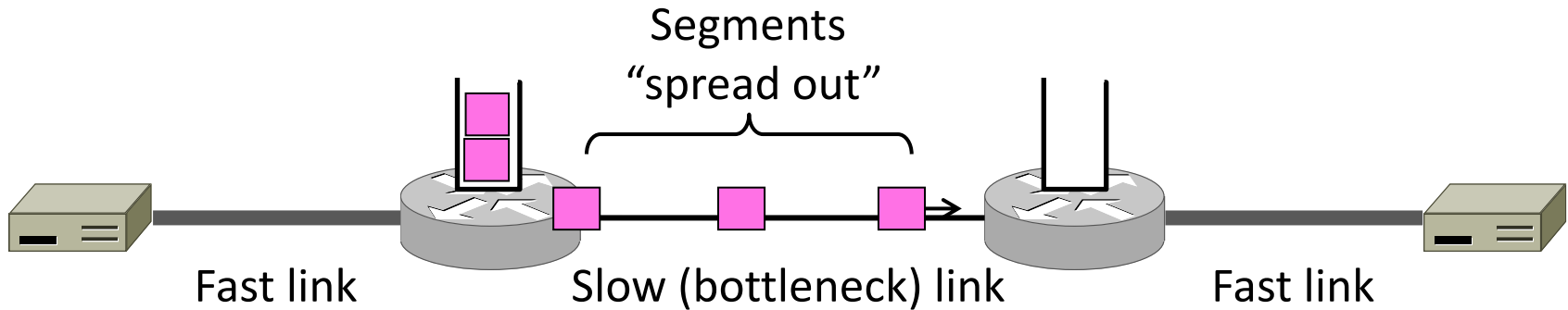
Benefit of ACK Clocking

- Consider what happens when sender injects a burst of segments into the network



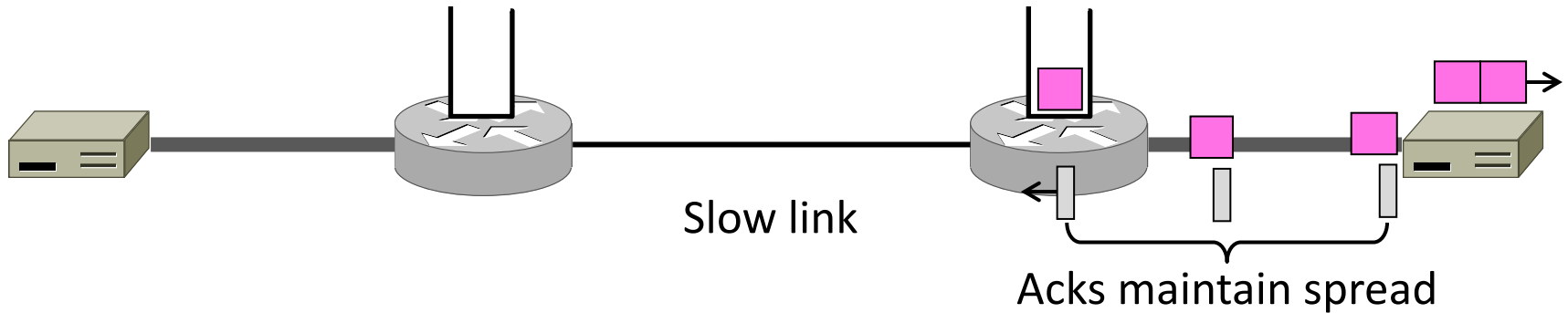
Benefit of ACK Clocking (2)

- Segments are buffered and spread out on slow link



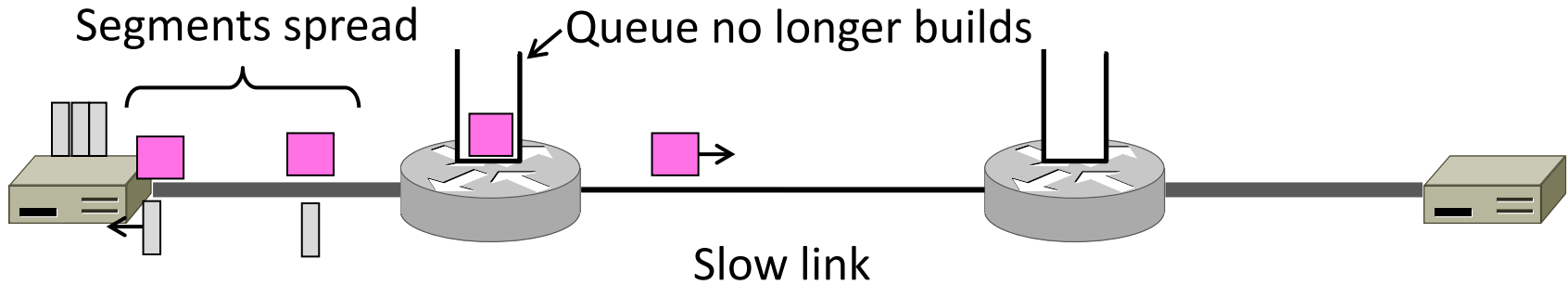
Benefit of ACK Clocking (3)

- ACKs maintain the spread back to the original sender



Benefit of ACK Clocking (4)

- Sender clocks new segments with the spread
 - Now sending at the bottleneck link without queuing!



Benefit of ACK Clocking (4)

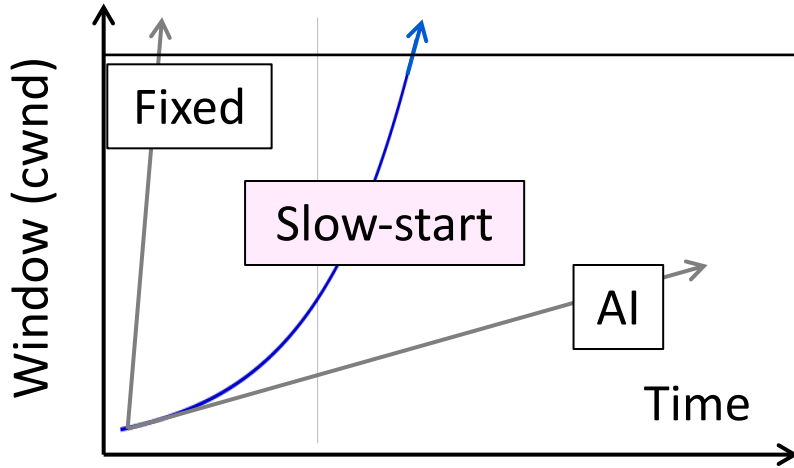
- Helps the network run with low levels of loss and delay!
- The network has smoothed out the burst of data segments
- ACK clock transfers this smooth timing back to the sender
- Subsequent data segments are not sent in bursts so do not queue up in the network

TCP Startup Problem

- We want to quickly reach the right rate, $\text{cwnd}_{\text{IDEAL}}$, but it varies greatly
 - Fixed sliding window doesn't adapt and is rough on the network (loss!)
 - AI with small bursts adapts cwnd gently to the network, but might take a long time to become efficient

Slow-Start Solution

- Start by doubling cwnd every RTT
 - Exponential growth (1, 2, 4, 8, 16, ...)
 - Start slow, quickly reach large values

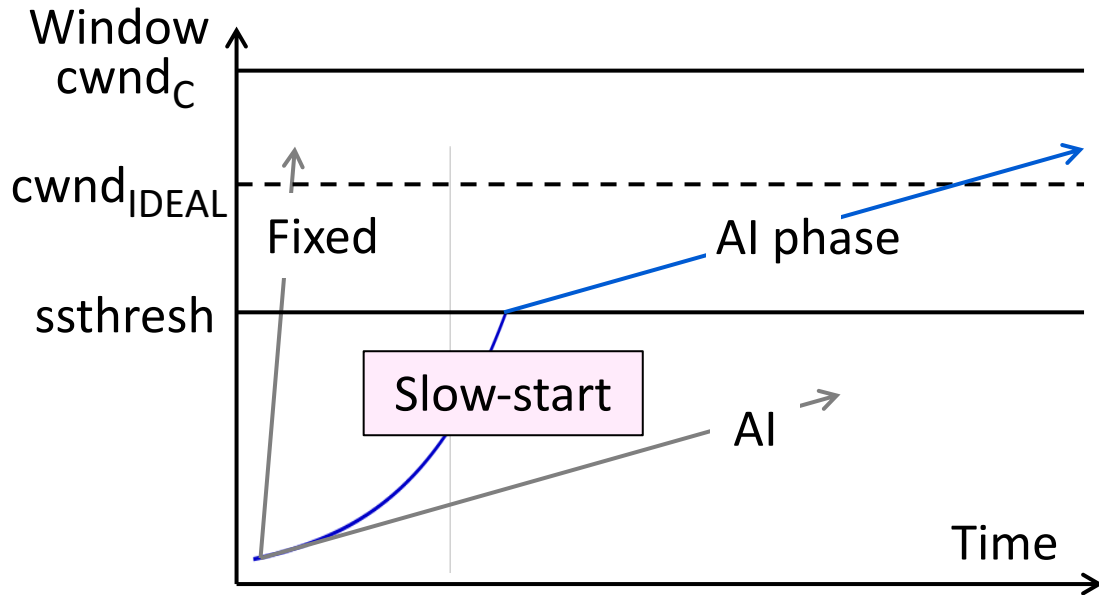


Slow-Start Solution (2)

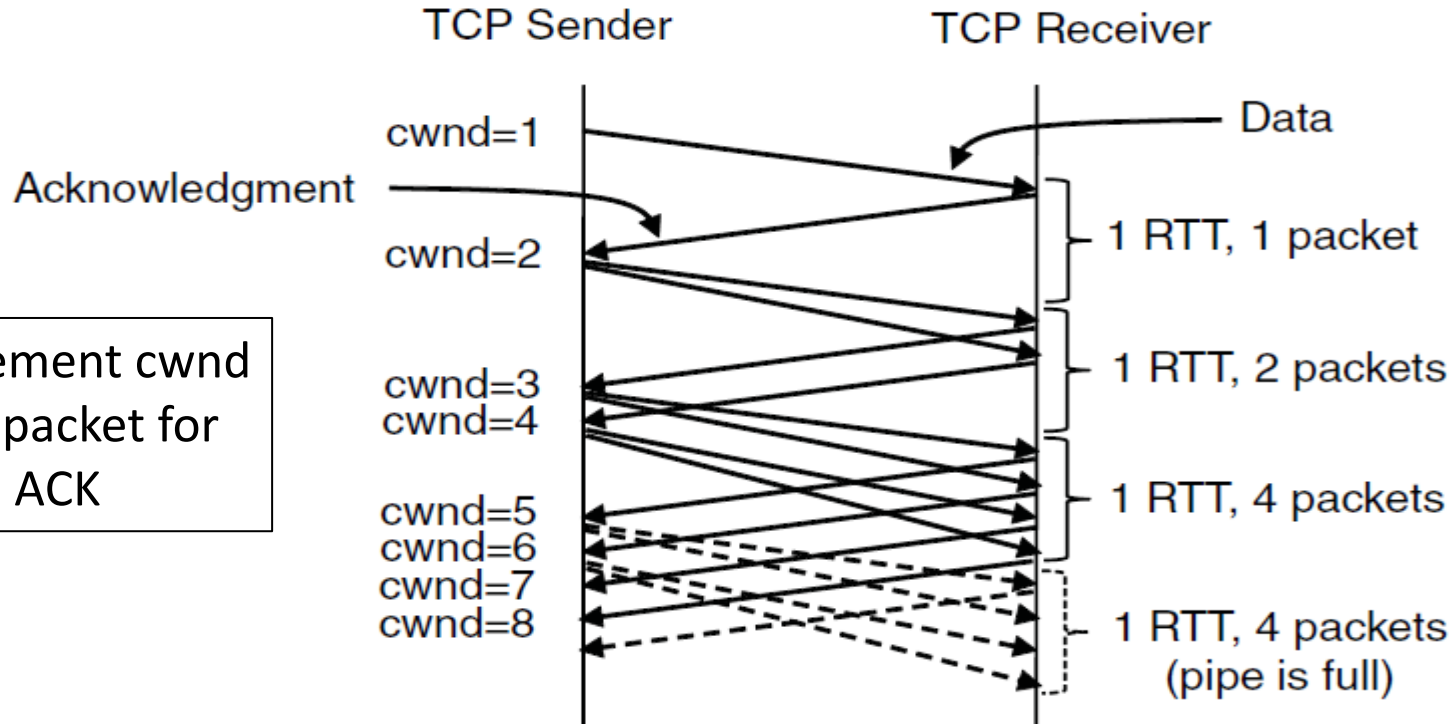
- Eventually packet loss will occur when the network is congested
 - Loss timeout tells us cwnd is too large
 - Next time, switch to AI beforehand
 - Slowly adapt cwnd near right value
- In terms of cwnd:
 - Expect loss for $\text{cwnd}_c \approx 2BD + \text{queue}$
 - Use $\text{ssthresh} = \text{cwnd}_c / 2$ to switch to AI

Slow-Start Solution (3)

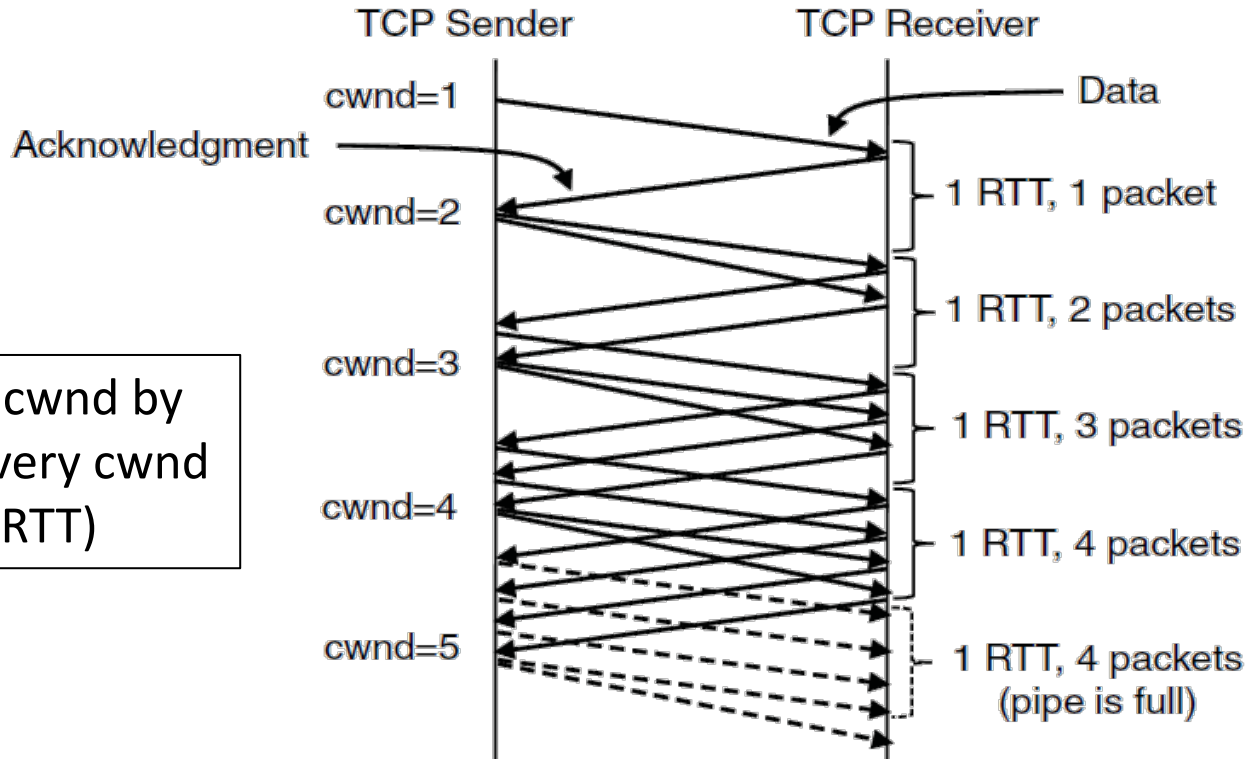
- Combined behavior, after first time
 - Most time spent near right value



Slow-Start (Doubling) Timeline



Additive Increase Timeline



Increment cwnd by
1 packet every cwnd
ACKs (or 1 RTT)

TCP Tahoe (Implementation)

- Initial slow-start (doubling) phase
 - Start with $\text{cwnd} = 1$ (or small value)
 - $\text{cwnd} += 1$ packet per ACK
- Later Additive Increase phase
 - $\text{cwnd} += 1/\text{cwnd}$ packets per ACK
 - Roughly adds 1 packet per RTT
- Switching threshold (initially infinity)
 - Switch to AI when $\text{cwnd} > \text{ssthresh}$
 - Set $\text{ssthresh} = \text{cwnd}/2$ after loss
 - Begin with slow-start after timeout

Timeout Misfortunes

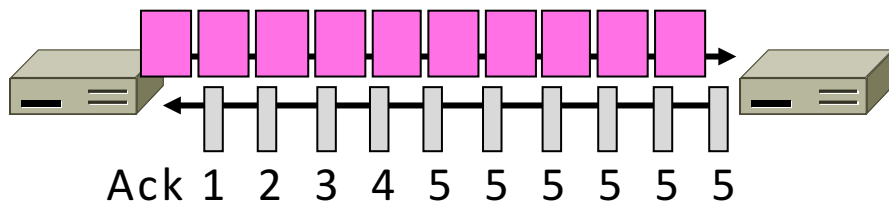
- Why do a slow-start after timeout?
 - Instead of MD cwnd (for AIMD)
- Timeouts are sufficiently long that the ACK clock will have run down
 - Slow-start ramps up the ACK clock
- We need to detect loss before a timeout to get to full AIMD
 - Done in TCP Reno

Inferring Loss from ACKs

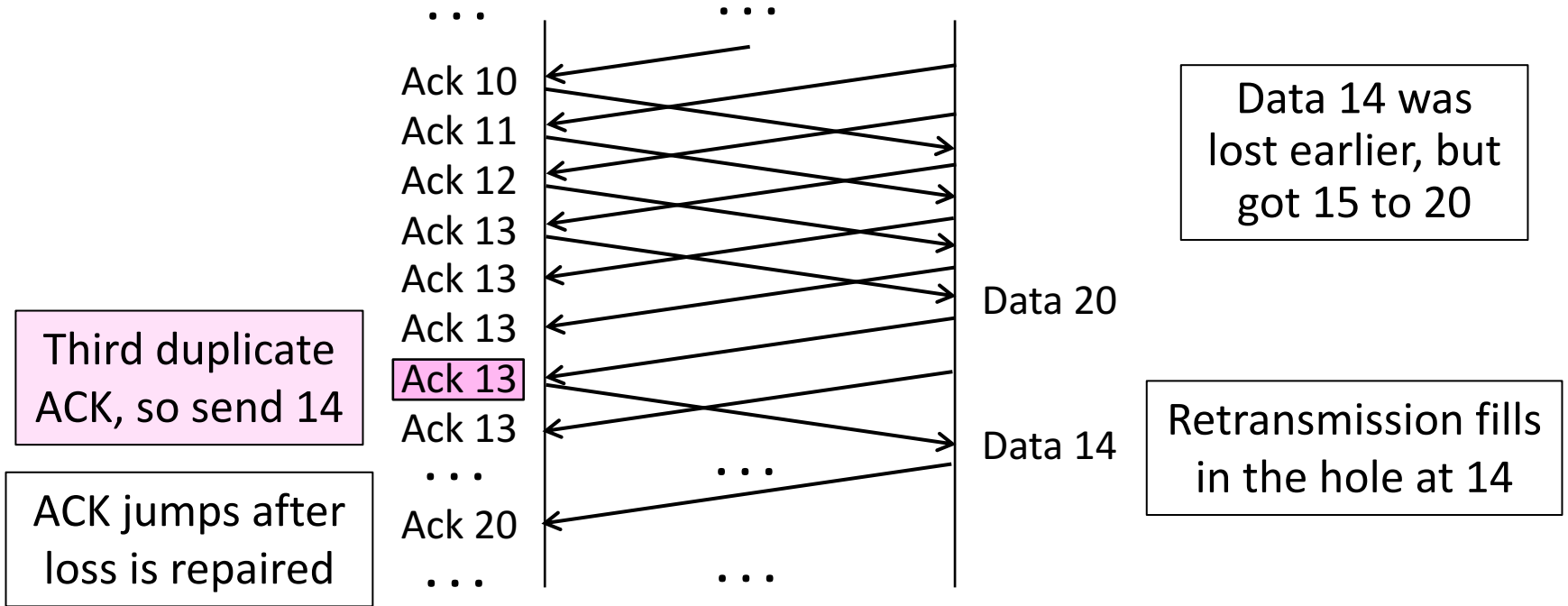
- TCP uses a cumulative ACK
 - Carries highest in-order seq. number
 - Normally a steady advance
- Duplicate ACKs give us hints about what data hasn't arrived
 - Tell us some new data did arrive, but it was not next segment
 - Thus the next segment may be lost

Fast Retransmit

- Treat three duplicate ACKs as a loss
 - Retransmit next expected segment
 - Some repetition allows for reordering, but still detects loss quickly



Fast Retransmit (2)



Fast Retransmit (3)

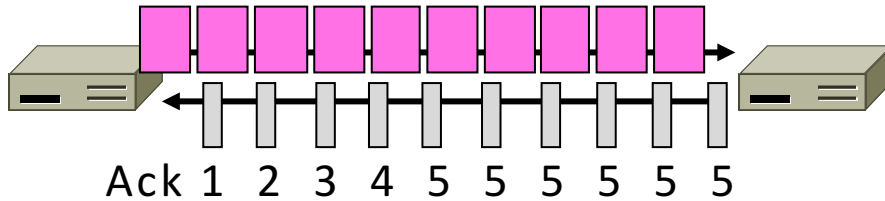
- It can repair single segment loss quickly, typically before a timeout
- However, we have quiet time at the sender/receiver while waiting for the ACK to jump
- And we still need to MD cwnd ...

Inferring Non-Loss from ACKs

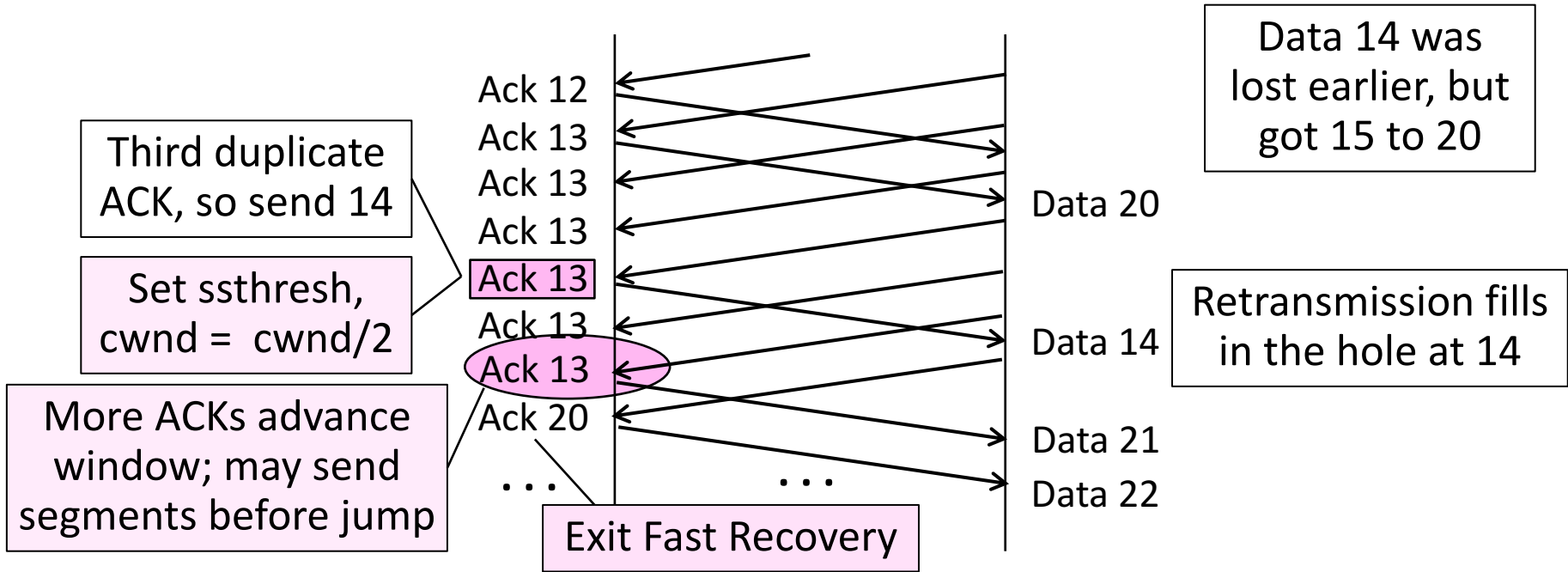
- Duplicate ACKs also give us hints about what data has arrived
 - Each new duplicate ACK means that some new segment has arrived
 - It will be the segments after the loss
 - Thus advancing the sliding window will not increase the number of segments stored in the network

Fast Recovery

- First fast retransmit, and MD cwnd
- Then pretend further duplicate ACKs are the expected ACKs
 - Lets new segments be sent for ACKs
 - Reconcile views when the ACK jumps



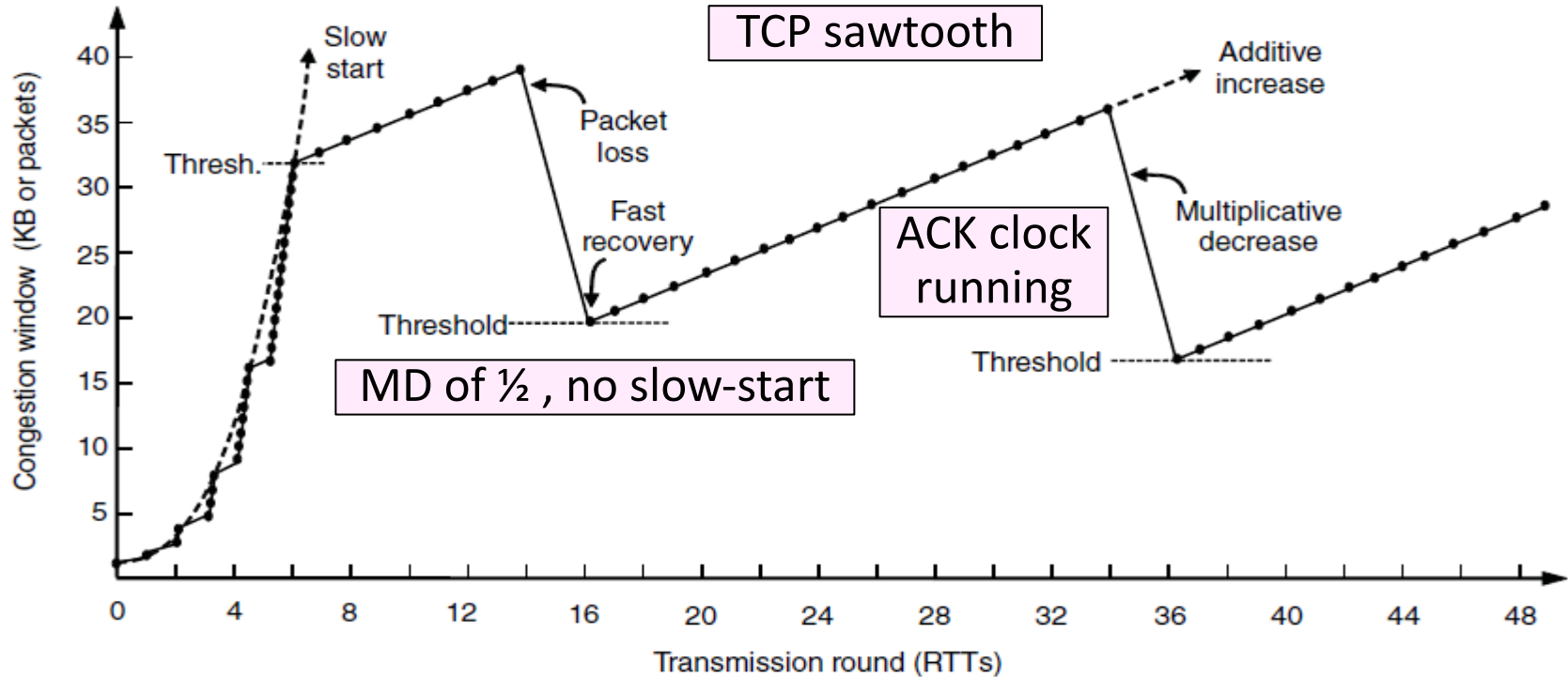
Fast Recovery (2)



Fast Recovery (3)

- With fast retransmit, it repairs a single segment loss quickly and keeps the ACK clock running
- This allows us to realize AIMD
 - No timeouts or slow-start after loss, just continue with a smaller cwnd
- TCP Reno combines slow-start, fast retransmit and fast recovery
 - Multiplicative Decrease is $\frac{1}{2}$

TCP Reno



TCP Reno, NewReno, and SACK

- Reno can repair one loss per RTT
 - Multiple losses cause a timeout
- NewReno further refines ACK heuristics
 - Repairs multiple losses without timeout
- SACK is a better idea
 - Receiver sends ACK ranges so sender can retransmit without guesswork

Computer Networks

Explicit Congestion Notification

Congestion Avoidance vs. Control

- Classic TCP drives the network into congestion and then recovers
 - Needs to see loss to slow down
- Would be better to use the network but avoid congestion altogether!
 - Reduces loss and delay
- But how can we do this?

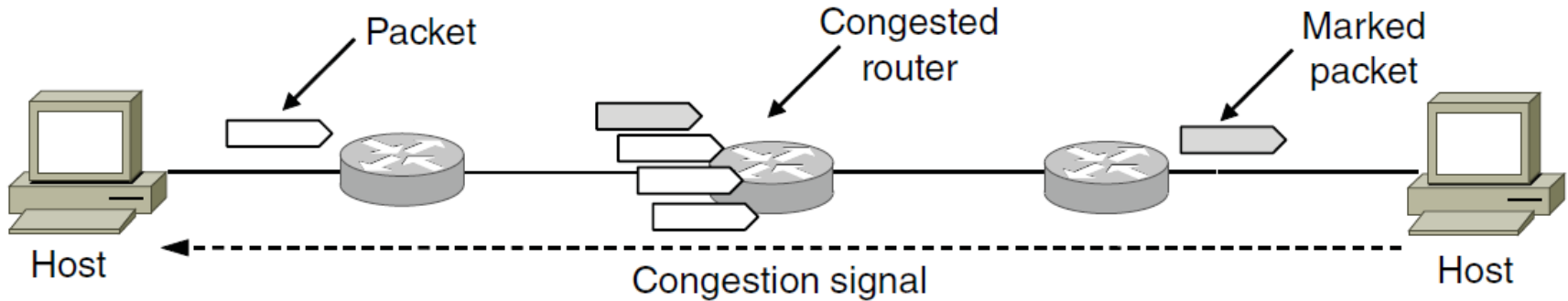
Feedback Signals

- Delay and router signals can let us avoid congestion

Signal	Example Protocol	Pros / Cons
Packet loss	Classic TCP Cubic TCP (Linux)	Hard to get wrong Hear about congestion late
Packet delay	Compound TCP (Windows)	Hear about congestion early Need to infer congestion
Router indication	TCPs with Explicit Congestion Notification	Hear about congestion early Require router support

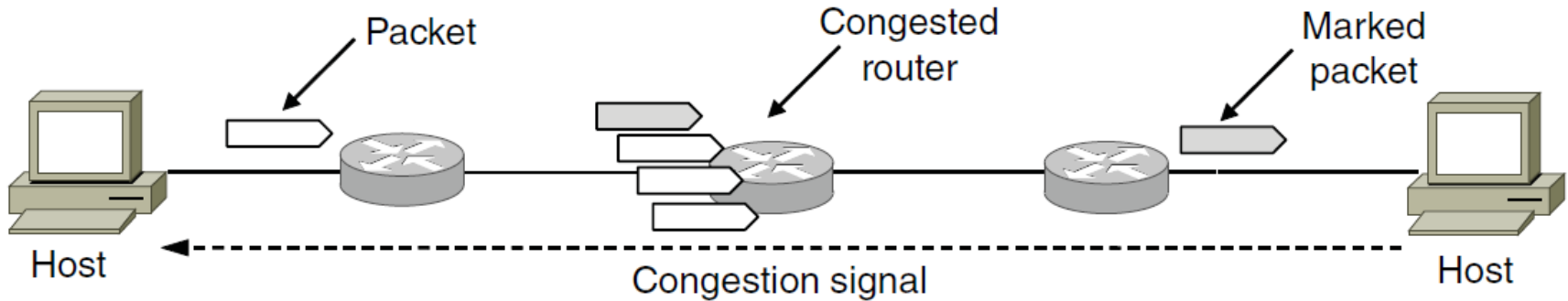
ECN (Explicit Congestion Notification)

- Router detects the onset of congestion via its queue
 - When congested, it marks affected packets (IP header)



ECN (2)

- Marked packets arrive at receiver; treated as loss
 - TCP receiver reliably informs TCP sender of the congestion



ECN (3)

- Advantages:
 - Routers deliver clear signal to hosts
 - Congestion is detected early, no loss
 - No extra packets need to be sent
- Disadvantages:
 - Routers and hosts must be upgraded