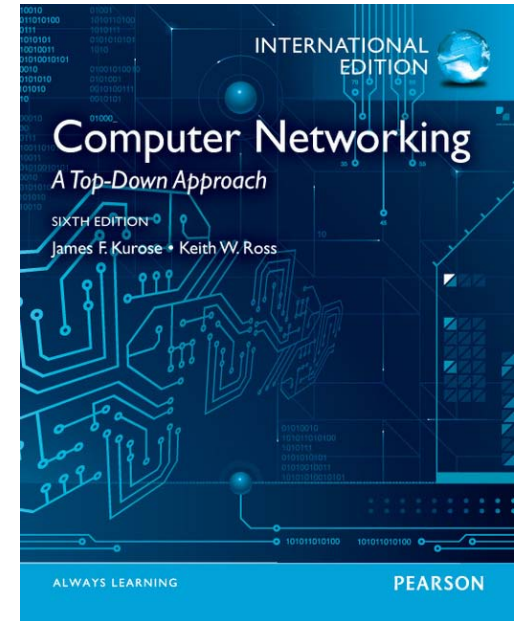


# Chapter 3

## Transport Layer



# Chapter 3: Transport Layer

## Our goals:

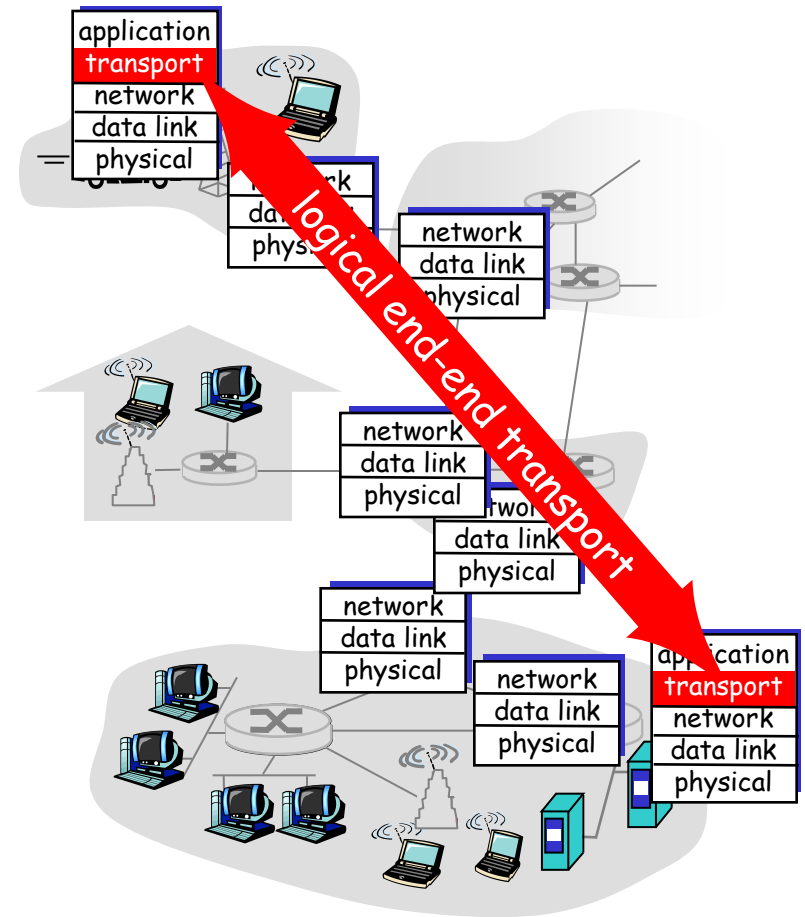
- r understand principles behind transport layer services:
  - m multiplexing/demultiplexing
  - m reliable data transfer
  - m flow control
  - m congestion control
- r learn about transport layer protocols in the Internet:
  - m UDP: connectionless transport
  - m TCP: connection-oriented transport
  - m TCP congestion control

# Chapter 3 outline

- r 3.1 Transport-layer services
- r 3.2 Multiplexing and demultiplexing
- r 3.3 Connectionless transport: UDP
- r 3.4 Principles of reliable data transfer
- r 3.5 Connection-oriented transport: TCP
  - m segment structure
  - m reliable data transfer
  - m flow control
  - m connection management
- r 3.6 Principles of congestion control
- r 3.7 TCP congestion control

# Internet transport-layer protocols

- r reliable, in-order delivery (TCP)
  - m congestion control
  - m flow control
  - m connection setup
- r unreliable, unordered delivery: UDP
  - m no-frills extension of “best-effort” IP
- r services not available:
  - m delay guarantees
  - m bandwidth guarantees



# Chapter 3 outline

- r 3.1 Transport-layer services
- r 3.2 Multiplexing and demultiplexing
- r 3.3 Connectionless transport: UDP
- r 3.4 Principles of reliable data transfer
- r 3.5 Connection-oriented transport: TCP
  - m segment structure
  - m reliable data transfer
  - m flow control
  - m connection management
- r 3.6 Principles of congestion control
- r 3.7 TCP congestion control

# Multiplexing/demultiplexing

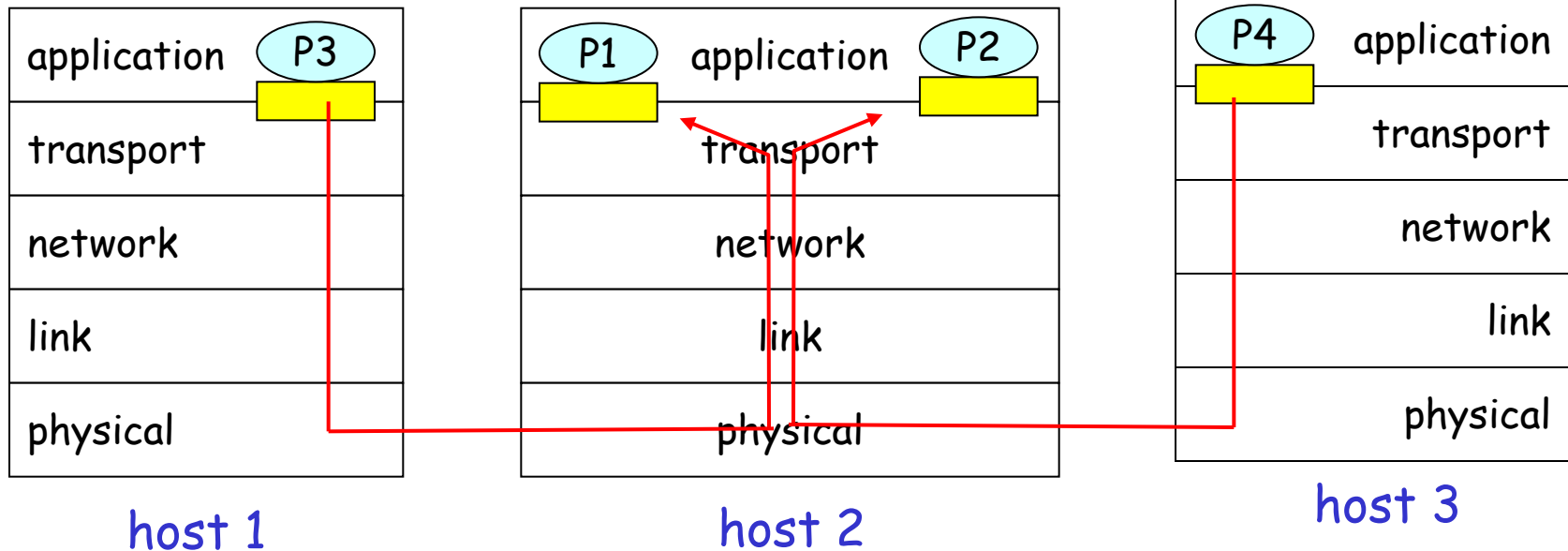
## Demultiplexing at rcv host:

delivering received segments to correct socket

## Multiplexing at send host:

gathering data from multiple sockets, enveloping data with header (later used for demultiplexing)

■ = socket      ○ = process



# Connectionless demultiplexing

- r Create sockets with port numbers:

```
DatagramSocket mySocket1 = new  
    DatagramSocket(12534);
```

```
DatagramSocket mySocket2 = new  
    DatagramSocket(12535);
```

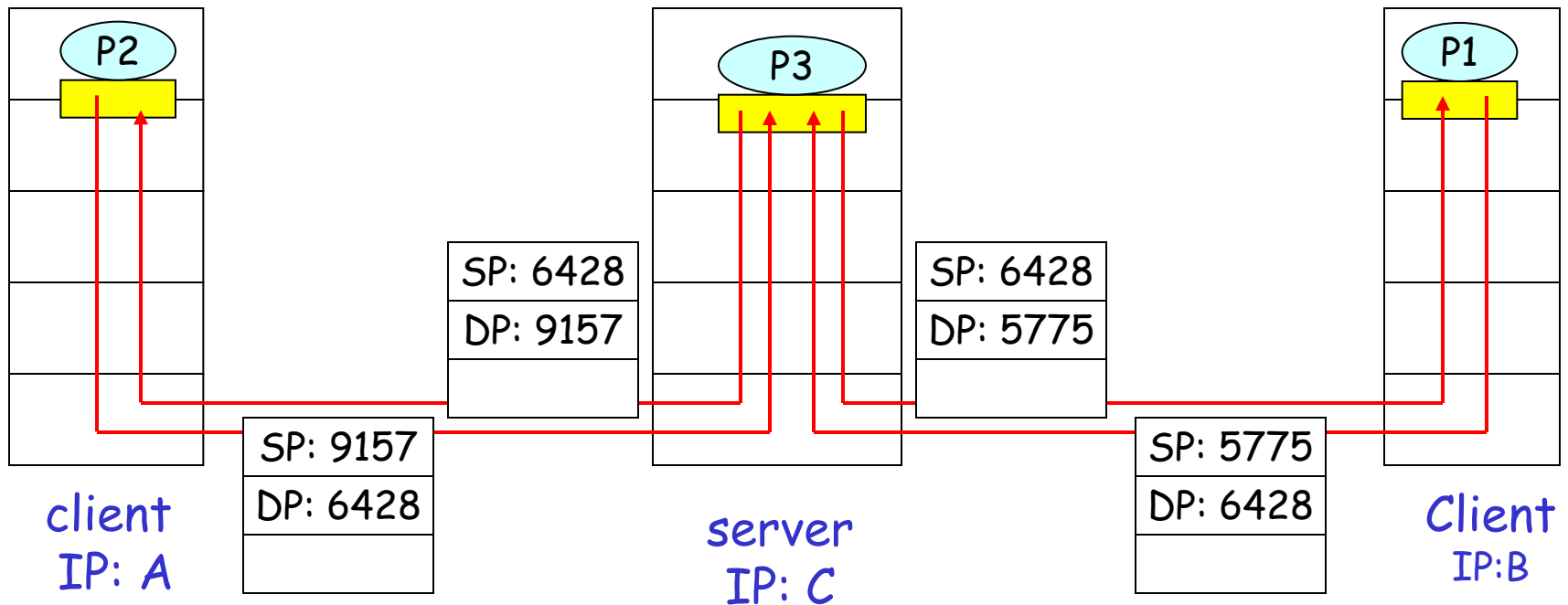
- r UDP socket identified by two-tuple:

(dest IP address, dest port number)

- r When host receives UDP segment:
  - m checks destination port number in segment
  - m directs UDP segment to socket with that port number
- r IP datagrams with different source IP addresses and/or source port numbers directed to same socket

# Connectionless demux (cont)

```
DatagramSocket serverSocket = new DatagramSocket(6428);
```



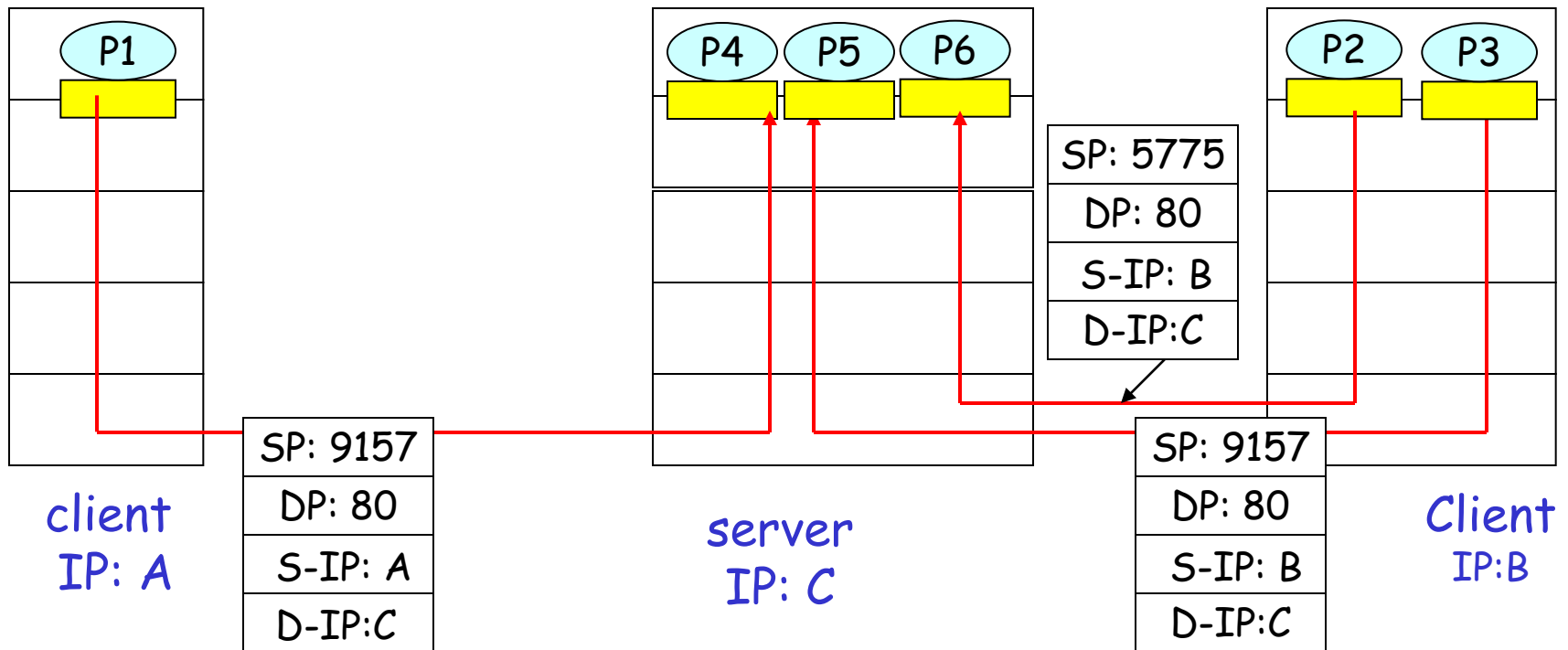
SP provides "return address"



# Connection-oriented demux

- r TCP socket identified by 4-tuple:
  - m source IP address
  - m source port number
  - m dest IP address
  - m dest port number
- r receiving host uses all four values to direct segment to appropriate socket
- r Server host may support many simultaneous TCP sockets:
  - m each socket identified by its own 4-tuple
- r Web servers have different sockets for each connecting client
  - m non-persistent HTTP will have different socket for each request

# Connection-oriented demux (cont)



# Chapter 3 outline

- r 3.1 Transport-layer services
- r 3.2 Multiplexing and demultiplexing
- r 3.3 Connectionless transport: UDP
- r 3.4 Principles of reliable data transfer
- r 3.5 Connection-oriented transport: TCP
  - m segment structure
  - m reliable data transfer
  - m flow control
  - m connection management
- r 3.6 Principles of congestion control
- r 3.7 TCP congestion control

# UDP: User Datagram Protocol [RFC 768]

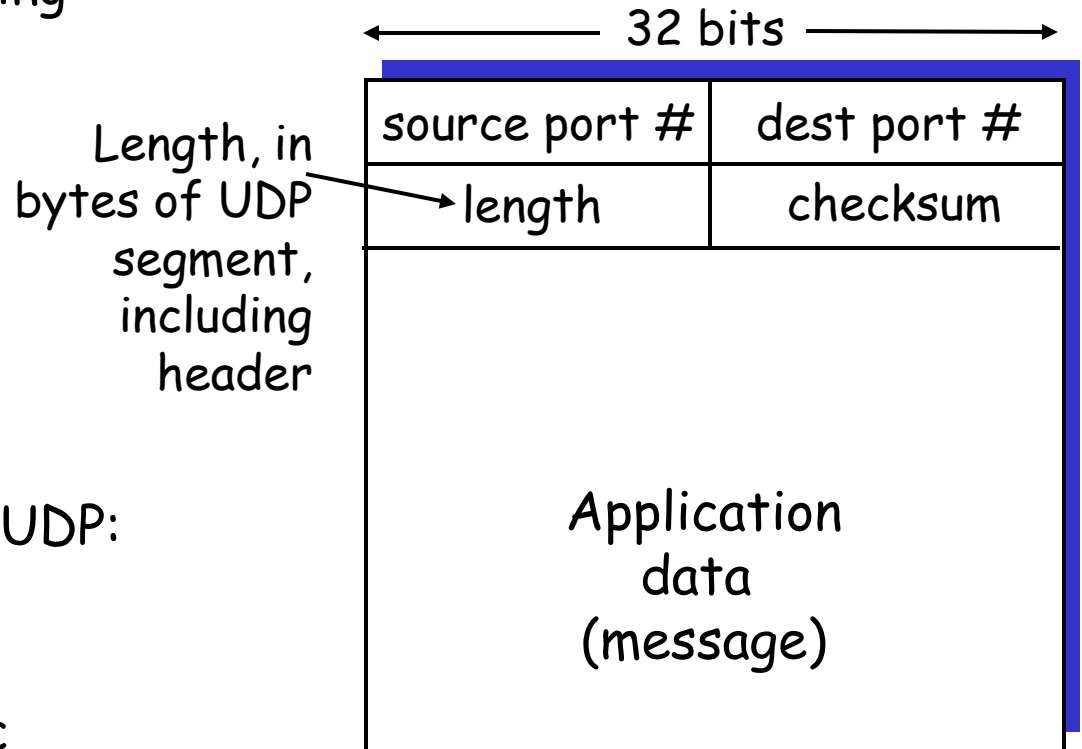
- r “no frills,” “bare bones” Internet transport protocol
- r “best effort” service, UDP segments may be:
  - m lost
  - m delivered out of order to app
- r *connectionless*:
  - m no handshaking between UDP sender, receiver
  - m each UDP segment handled independently of others

## Why is there a UDP?

- r no connection establishment (which can add delay)
- r simple: no connection state at sender, receiver
- r small segment header
- r no congestion control: UDP can blast away as fast as desired

# UDP: more

- r often used for streaming multimedia apps
  - m loss tolerant
  - m rate sensitive
- r other UDP uses
  - m DNS
  - m SNMP
- r reliable transfer over UDP:  
add reliability at application layer
  - m application-specific error recovery!



UDP segment format

# UDP checksum

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

## Sender:

- r treat segment contents as sequence of 16-bit integers
- r checksum: addition (1's complement sum) of segment contents
- r sender puts checksum value into UDP checksum field

## Receiver:

- r compute checksum of received segment
- r check if computed checksum equals checksum field value:
  - m NO - error detected
  - m YES - no error detected.  
*But maybe errors nonetheless? More later*
- ....

# Internet Checksum Example

r Note

m When adding numbers, a carryout from the most significant bit needs to be added to the result

r Example: add two 16-bit integers

	1	1	1	0	0	1	1	0	0	1	1	0	0	1	1	0	
	1	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1	
<hr/>																	
wraparound	1	1	0	1	1	1	0	1	1	1	0	1	1	1	0	1	1
<hr/>																	
sum	1	0	1	1	1	0	1	1	1	0	1	1	1	1	0	0	
checksum	0	1	0	0	0	1	0	0	0	1	0	0	0	0	1	1	

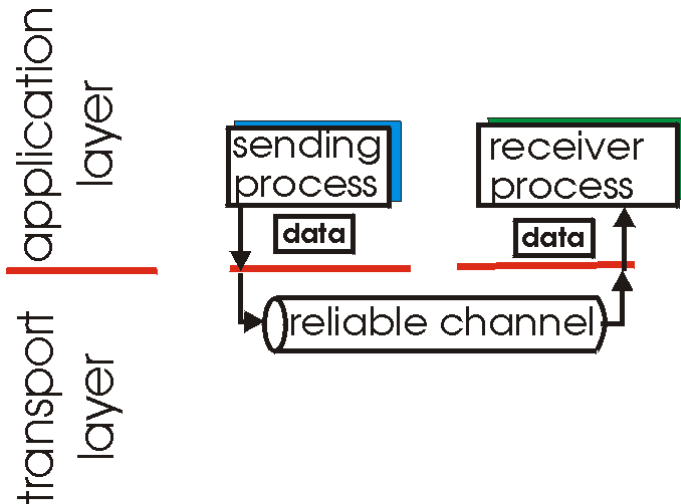
# Chapter 3 outline

- r 3.1 Transport-layer services
- r 3.2 Multiplexing and demultiplexing
- r 3.3 Connectionless transport: UDP
- r 3.4 Principles of reliable data transfer
- r 3.5 Connection-oriented transport: TCP
  - m segment structure
  - m reliable data transfer
  - m flow control
  - m connection management
- r 3.6 Principles of congestion control
- r 3.7 TCP congestion control



# Principles of Reliable data transfer

- r important in app., transport, link layers
- r top-10 list of important networking topics!

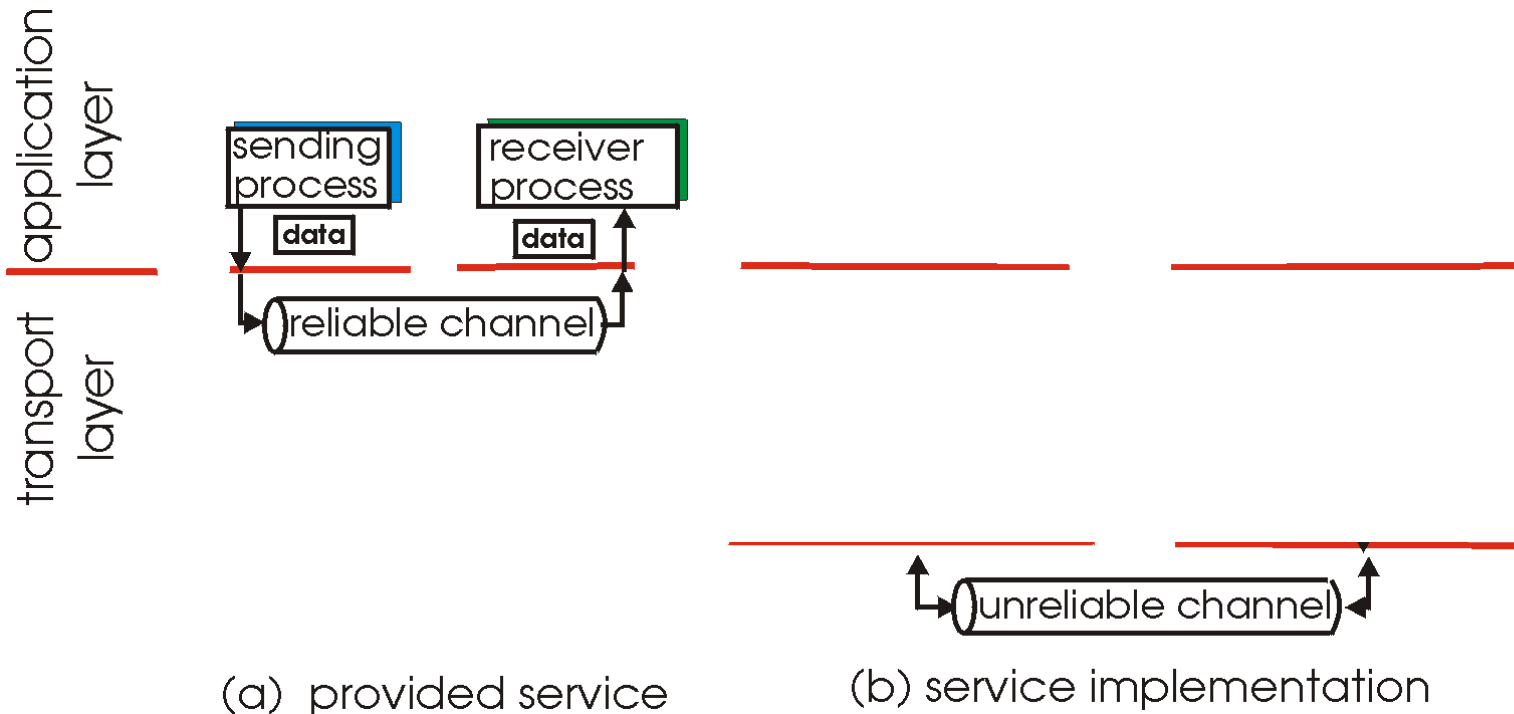


(a) provided service

- r characteristics of unreliable channel will determine complexity of reliable data transfer protocol (rdt)

# Principles of Reliable data transfer

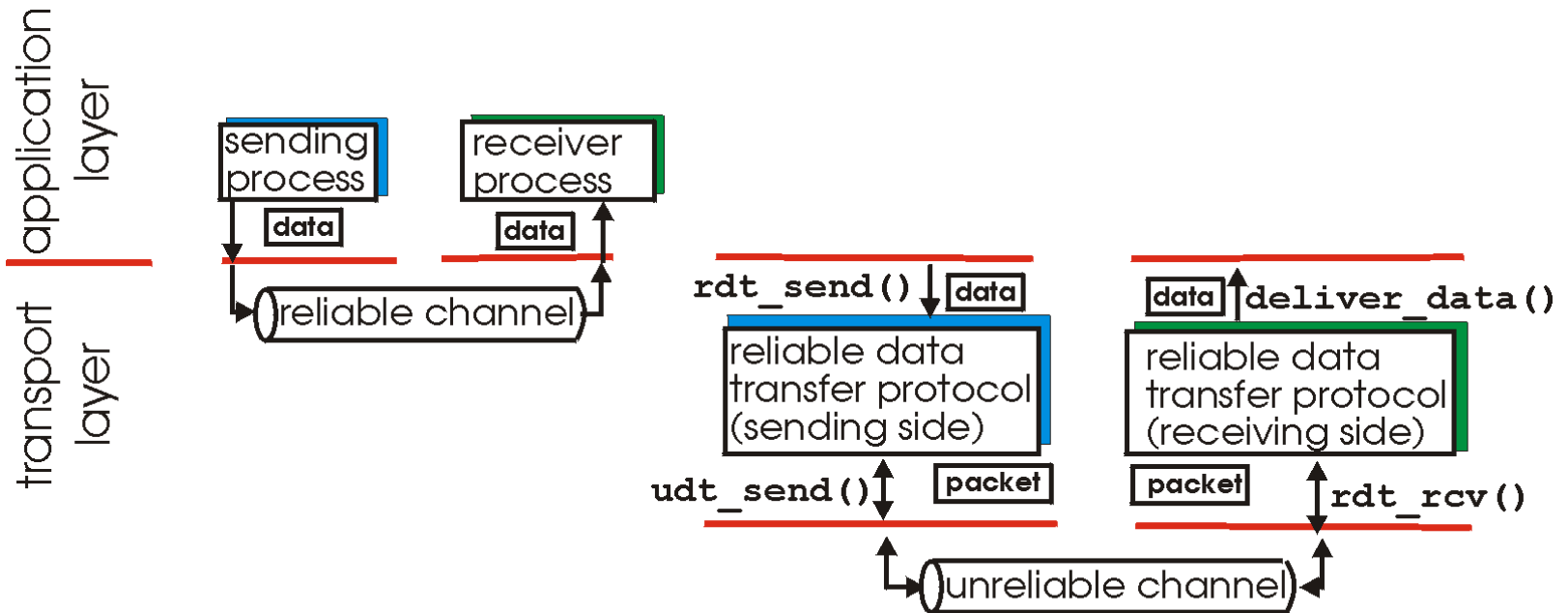
- r important in app., transport, link layers
- r top-10 list of important networking topics!



- r characteristics of unreliable channel will determine complexity of reliable data transfer protocol (rdt)

# Principles of Reliable data transfer

- r important in app., transport, link layers
- r top-10 list of important networking topics!



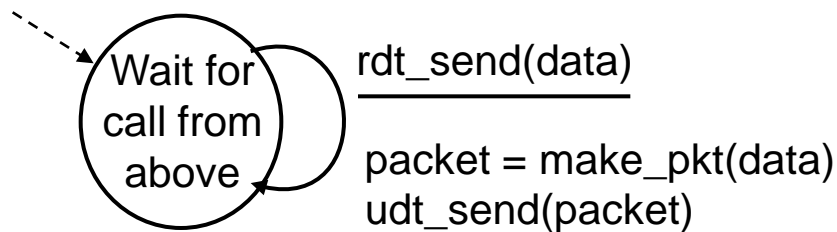
(a) provided service

(b) service implementation

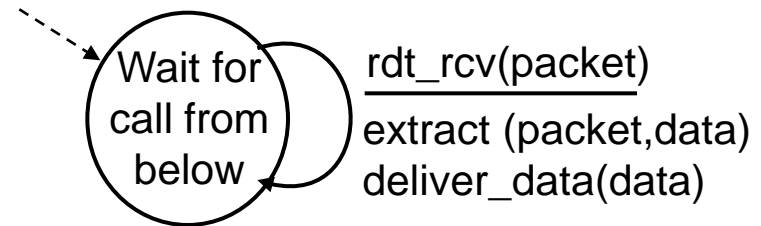
- r characteristics of unreliable channel will determine complexity of reliable data transfer protocol (rdt)

# Rdt1.0: reliable transfer over a reliable channel

- r underlying channel perfectly reliable
  - m no bit errors
  - m no loss of packets
- r separate FSMs for sender, receiver:
  - m sender sends data into underlying channel
  - m receiver read data from underlying channel



sender



receiver

## Rdt2.0: channel with bit errors

- r underlying channel may flip bits in packet
  - m checksum to detect bit errors
- r *the question: how to recover from errors:*
  - m *acknowledgements (ACKs):* receiver explicitly tells sender that pkt received OK
  - m *negative acknowledgements (NAKs):* receiver explicitly tells sender that pkt had errors
  - m sender retransmits pkt on receipt of NAK
- r new mechanisms in rdt2.0 (beyond rdt1.0):
  - m error detection
  - m receiver feedback: control msgs (ACK,NAK) rcvr->sender

# rdt2.0 has a fatal flaw!

## What happens if ACK/NAK corrupted?

- r sender doesn't know what happened at receiver!
- r can't just retransmit: possible duplicate

## Handling duplicates:

- r sender retransmits current pkt if ACK/NAK garbled
- r sender adds *sequence number* to each pkt
- r receiver discards (doesn't deliver up) duplicate pkt

### stop and wait

Sender sends one packet, then waits for receiver response

# rdt3.0: channels with errors *and* loss

## New assumption:

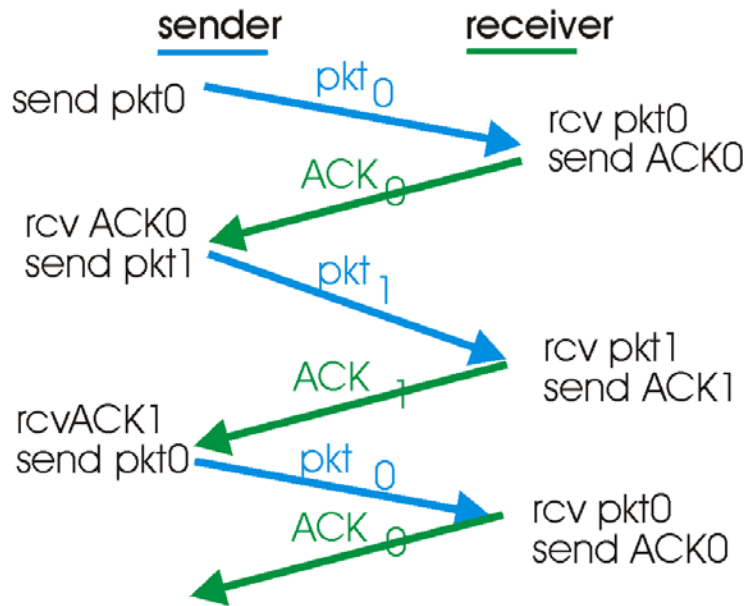
underlying channel can also lose packets (data or ACKs)

- m checksum, seq. #, ACKs, retransmissions will be of help, but not enough

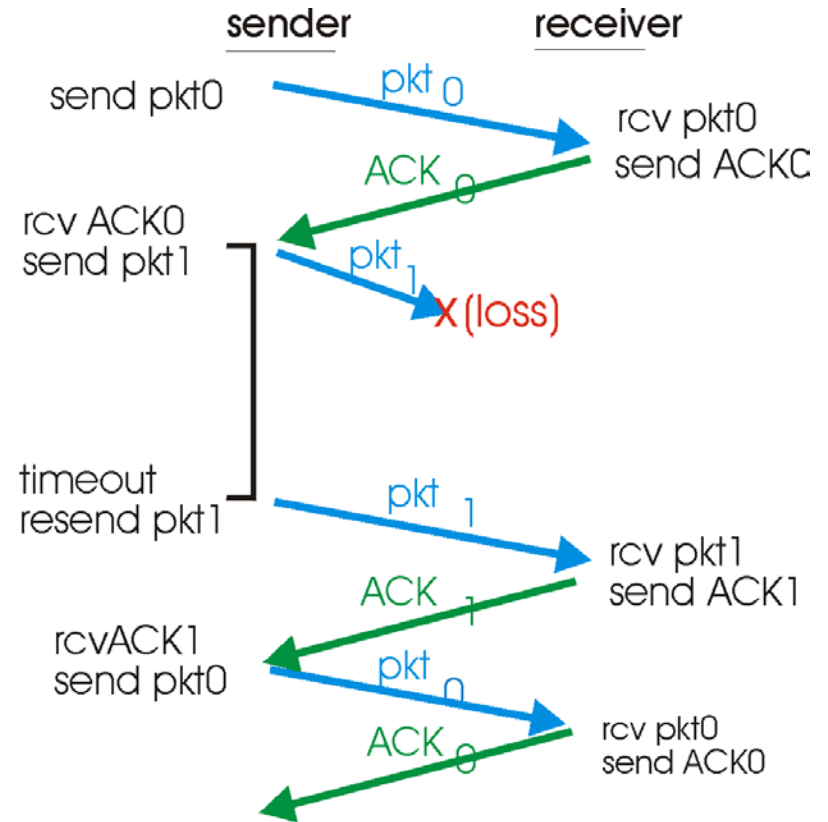
Approach: sender waits “reasonable” amount of time for ACK

- r retransmits if no ACK received in this time
- r if pkt (or ACK) just delayed (not lost):
  - m retransmission will be duplicate, but use of seq. #'s already handles this
  - m receiver must specify seq # of pkt being ACKed
- r requires countdown timer

# rdt3.0 in action



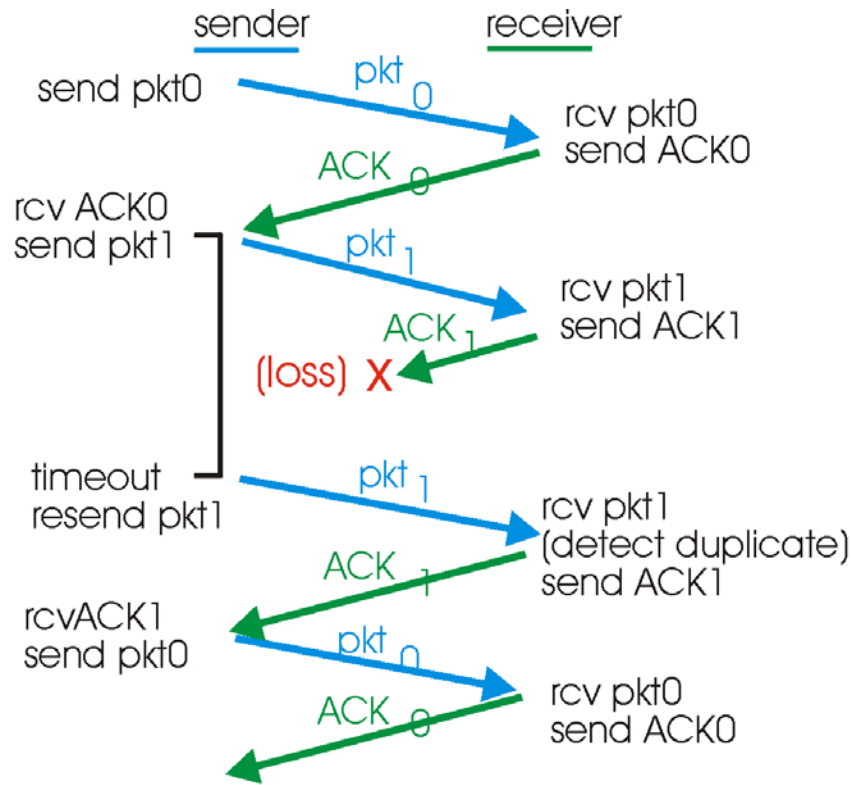
(a) operation with no loss



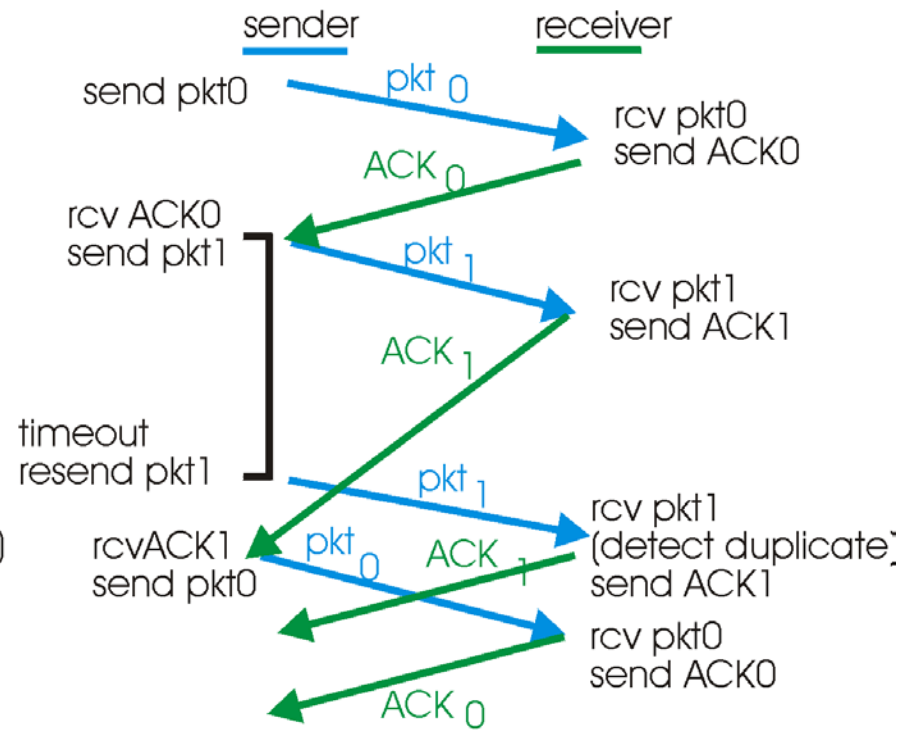
(b) lost packet



# rdt3.0 in action



(c) lost ACK



(d) premature timeout

# Performance of rdt3.0

- r rdt3.0 works, but performance stinks
- r ex: 1 Gbps link, 15 ms prop. delay, 8000 bit packet:

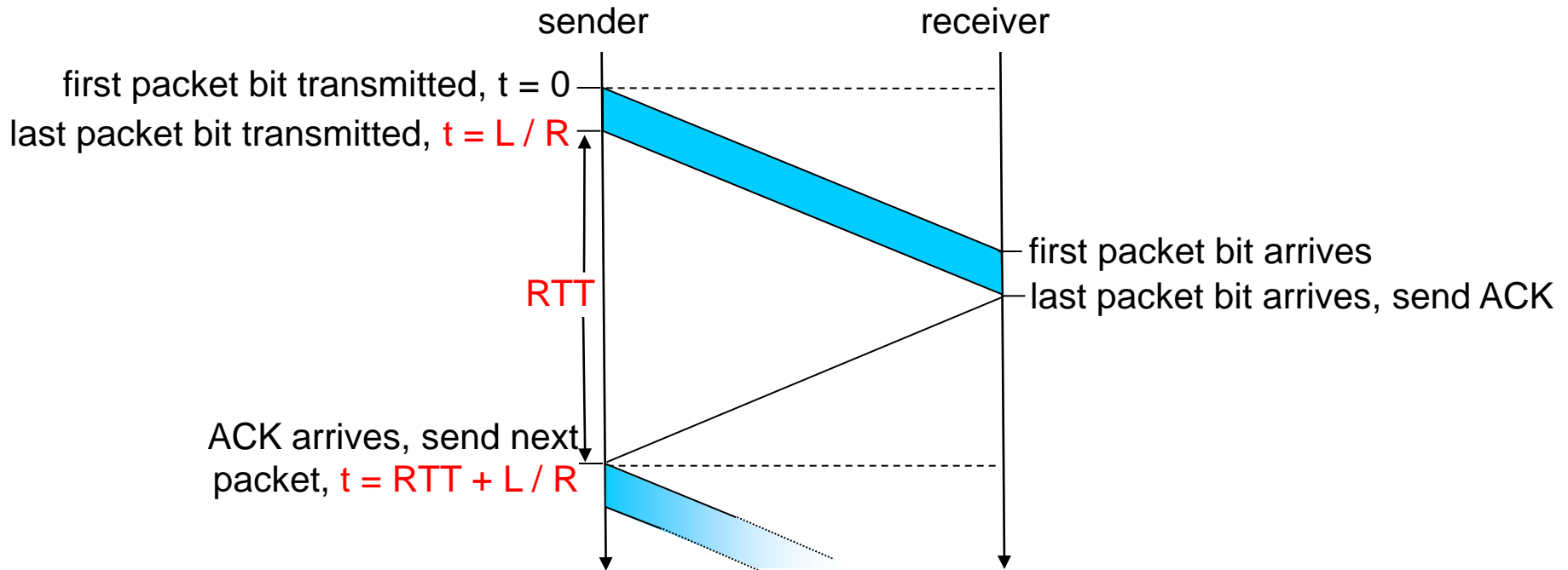
$$d_{trans} = \frac{L}{R} = \frac{8000\text{bits}}{10^9 \text{bps}} = 8\text{microseconds}$$

- m  $U_{\text{sender}}$ : **utilization** - fraction of time sender busy sending

U

- m 1KB pkt every 30 msec -> 33kB/sec thrupt over 1 Gbps link
- m network protocol limits use of physical resources!

# rdt3.0: stop-and-wait operation



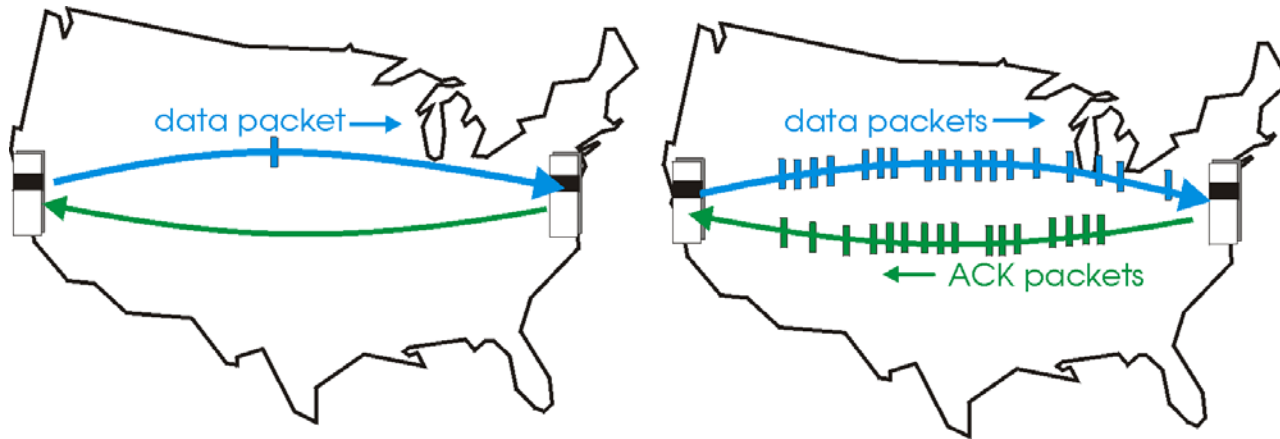
U

# Pipelined protocols

**Pipelining:** sender allows multiple, “in-flight”, yet-to-be-acknowledged pkts

m range of sequence numbers must be increased

m buffering at sender and/or receiver

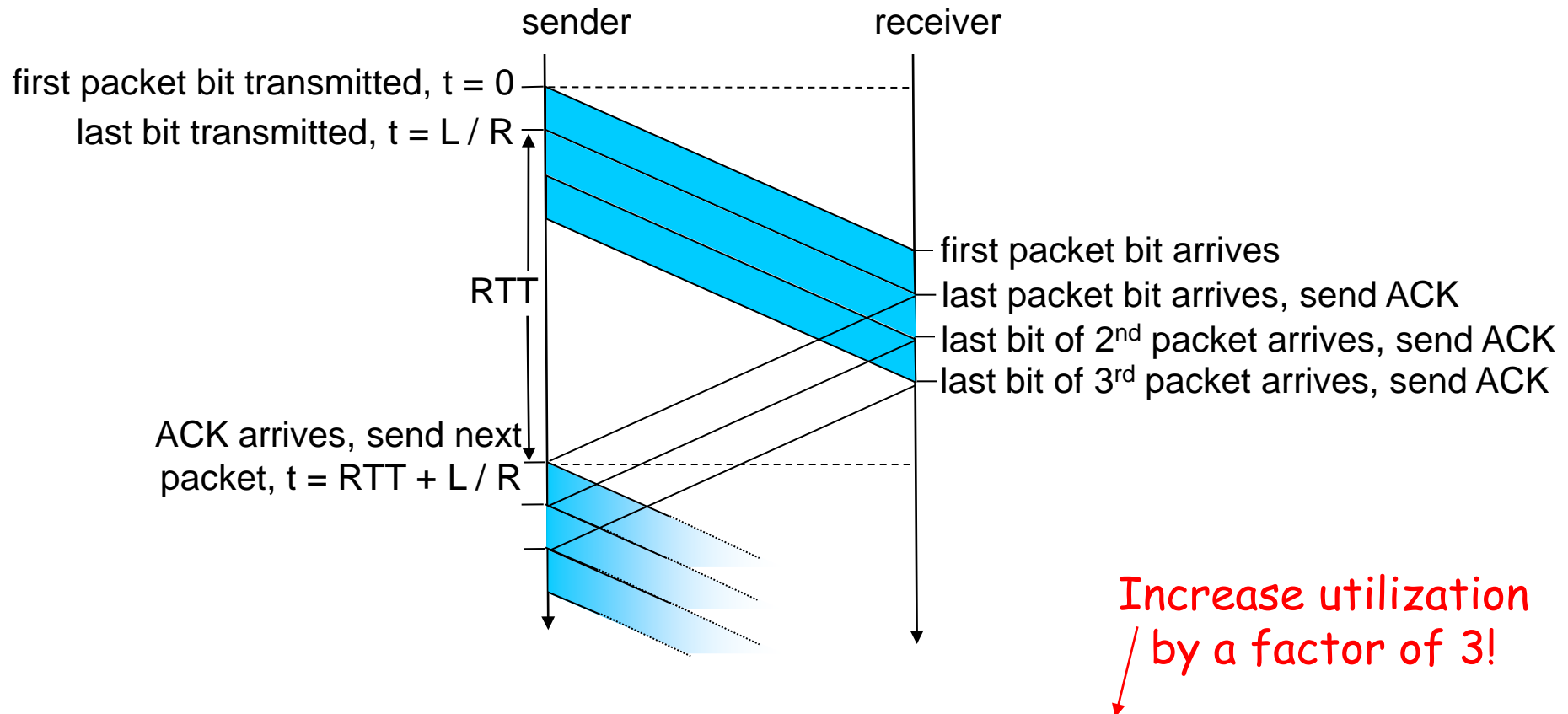


(a) a stop-and-wait protocol in operation

(b) a pipelined protocol in operation

r Two generic forms of pipelined protocols: *go-Back-N*, *selective repeat*

# Pipelining: increased utilization



U

# Pipelining Protocols

## Go-back-N: overview

- r *sender*: up to N unACKed pkts in pipeline
- r *receiver*: only sends cumulative ACKs
  - m doesn't ACK pkt if there's a gap
- r *sender*: has timer for oldest unACKed pkt
  - m if timer expires: retransmit all unACKed packets

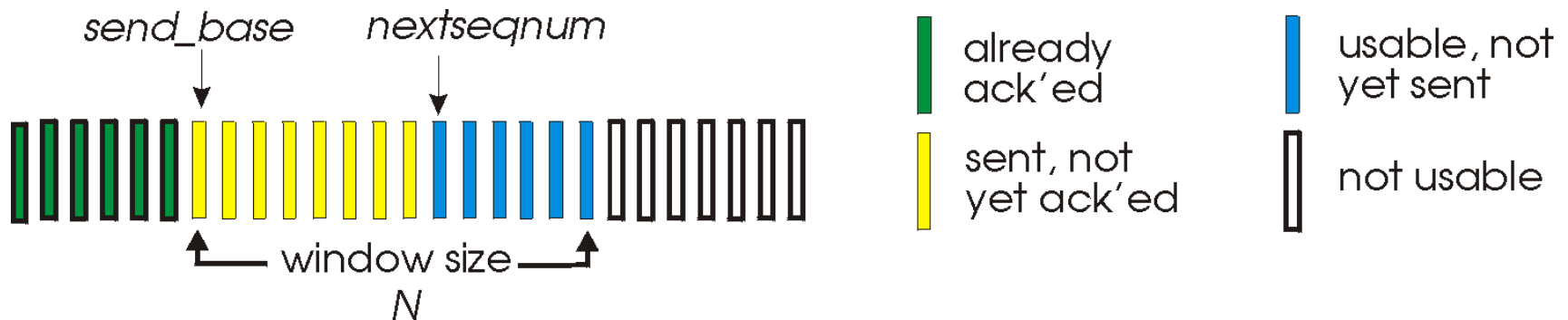
## Selective Repeat: overview

- r *sender*: up to N unACKed packets in pipeline
- r *receiver*: ACKs individual pkts
- r *sender*: maintains timer for each unACKed pkt
  - m if timer expires: retransmit only unACKed packet

# Go-Back-N

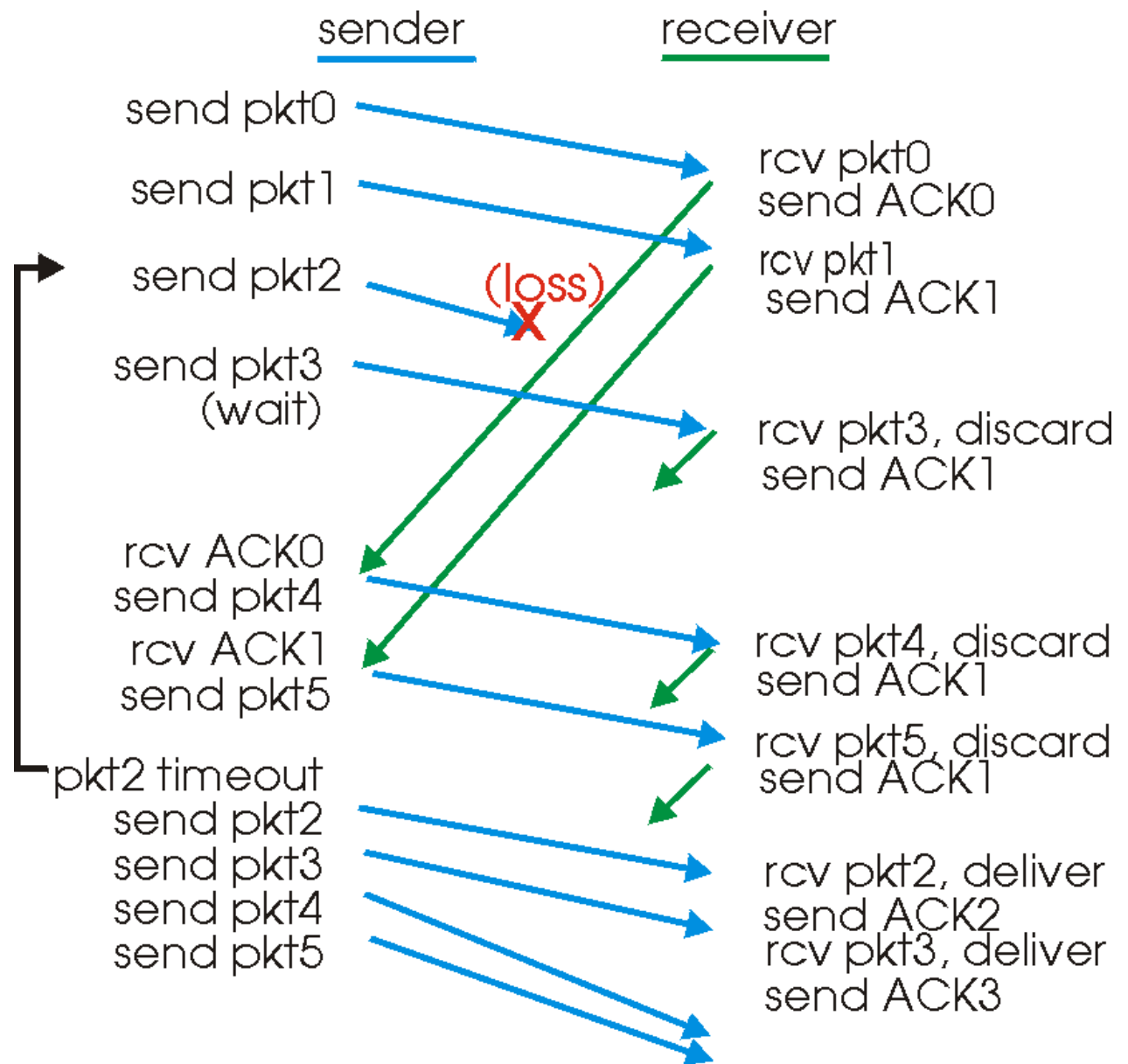
## Sender:

- r k-bit seq # in pkt header
- r “window” of up to  $N$ , consecutive unACKed pkts allowed



- r ACK( $n$ ): ACKs all pkts up to, including seq #  $n$  - “cumulative ACK”
  - m may receive duplicate ACKs (see receiver)
- r timer for each in-flight pkt
- r *timeout*( $n$ ): retransmit pkt  $n$  and all higher seq # pkts in window

# GBN in action

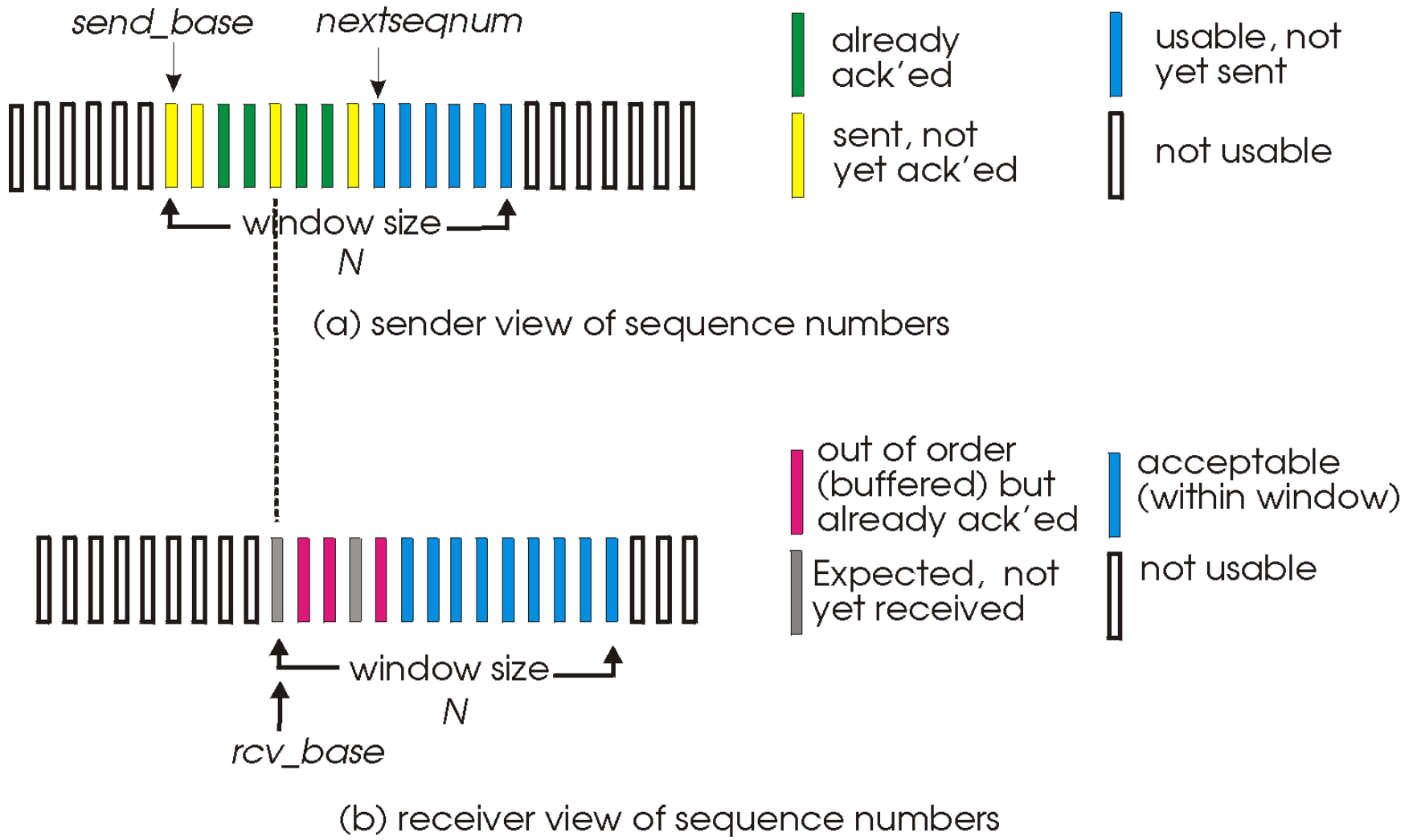




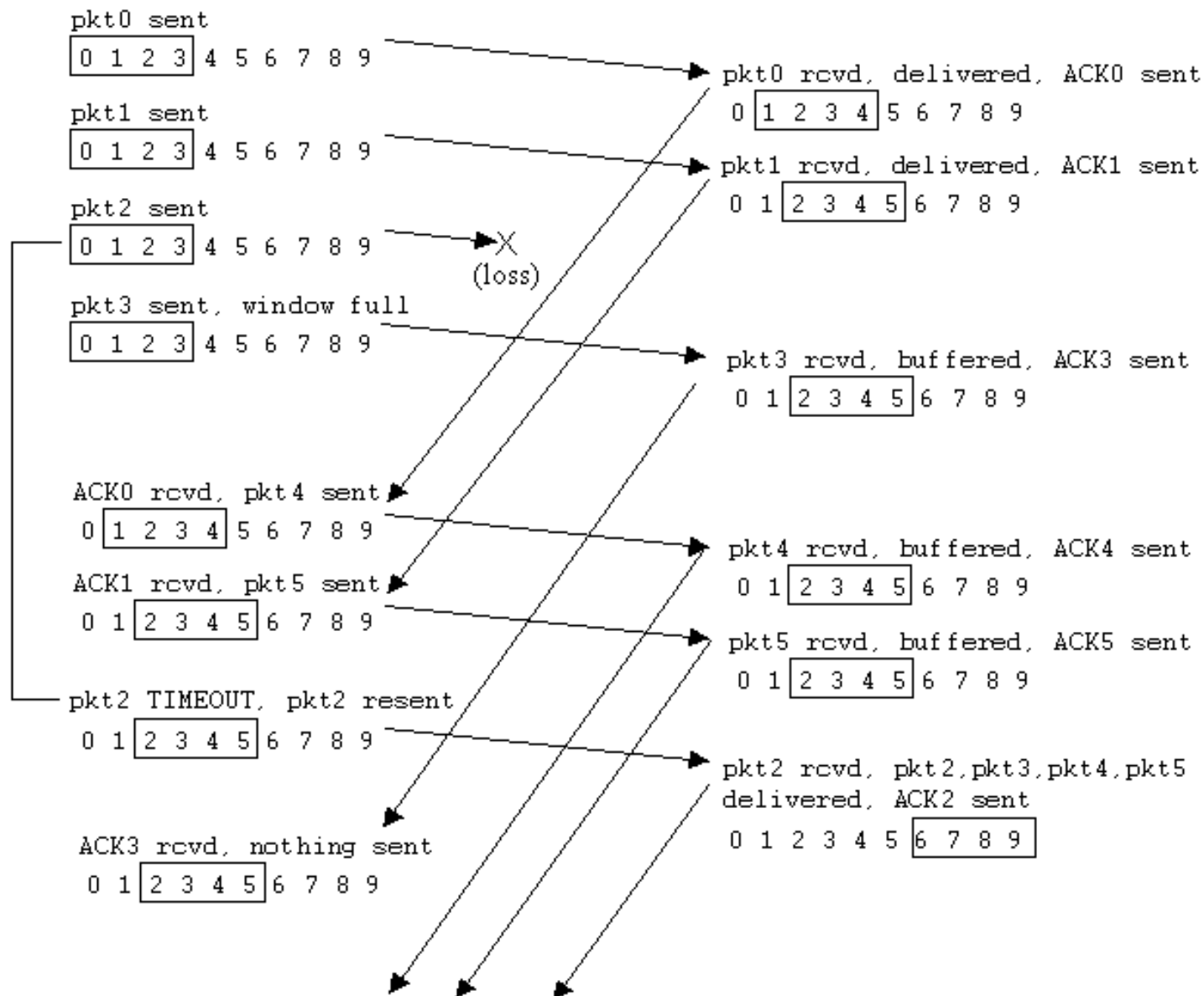
# Selective Repeat

- r receiver *individually* acknowledges all correctly received pkts
  - m buffers pkts, as needed, for eventual in-order delivery to upper layer
- r sender only resends pkts for which ACK not received
  - m sender timer for each unACKed pkt
- r sender window
  - m N consecutive seq #'s
  - m again limits seq #'s of sent, unACKed pkts

# Selective repeat: sender, receiver windows



# Selective repeat in action



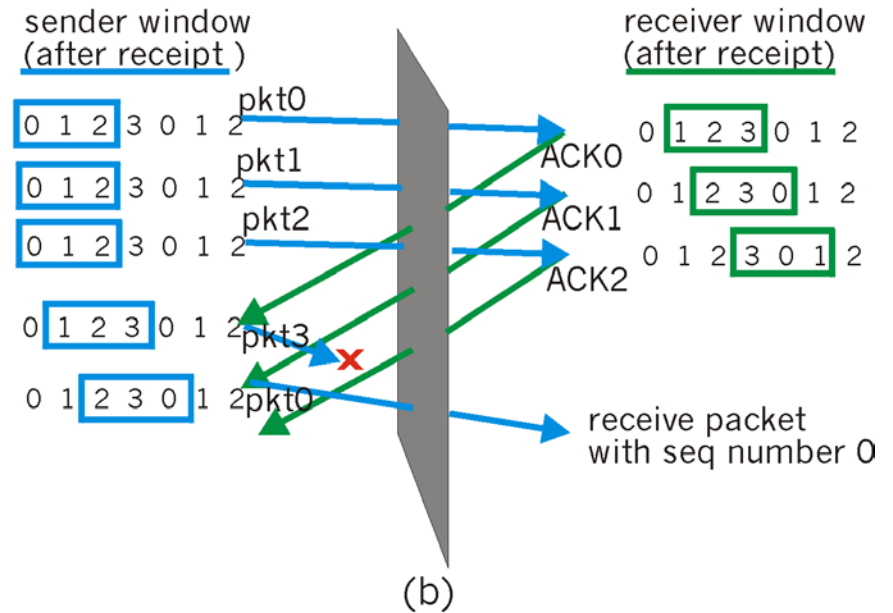
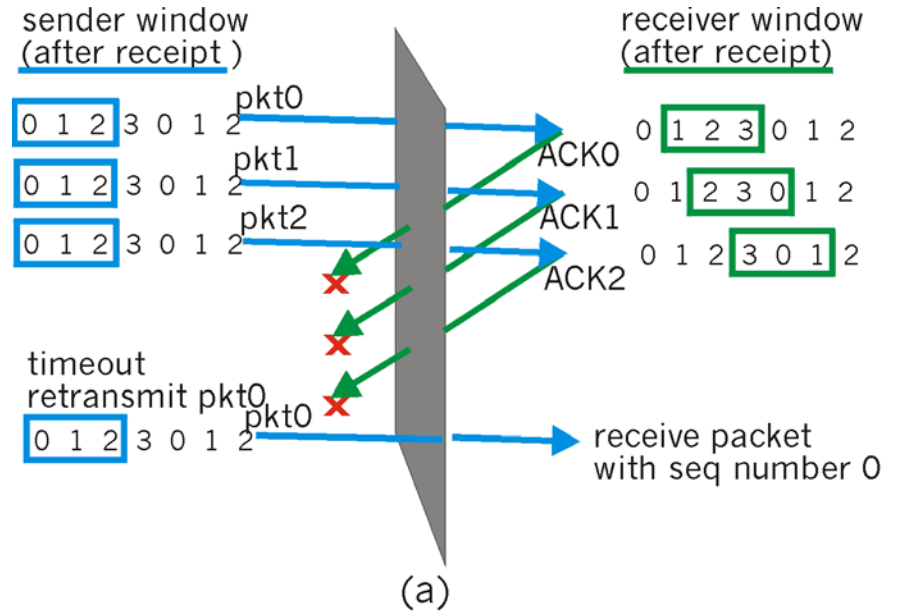
# Selective repeat: dilemma

Example:

- r seq #'s: 0, 1, 2, 3
- r window size=3

- r receiver sees no difference in two scenarios!
- r incorrectly passes duplicate data as new in (a)

Q: what relationship between seq # size and window size?



# Chapter 3 outline

- r 3.1 Transport-layer services
- r 3.2 Multiplexing and demultiplexing
- r 3.3 Connectionless transport: UDP
- r 3.4 Principles of reliable data transfer
- r 3.5 Connection-oriented transport: TCP
  - m segment structure
  - m reliable data transfer
  - m flow control
  - m connection management
- r 3.6 Principles of congestion control
- r 3.7 TCP congestion control

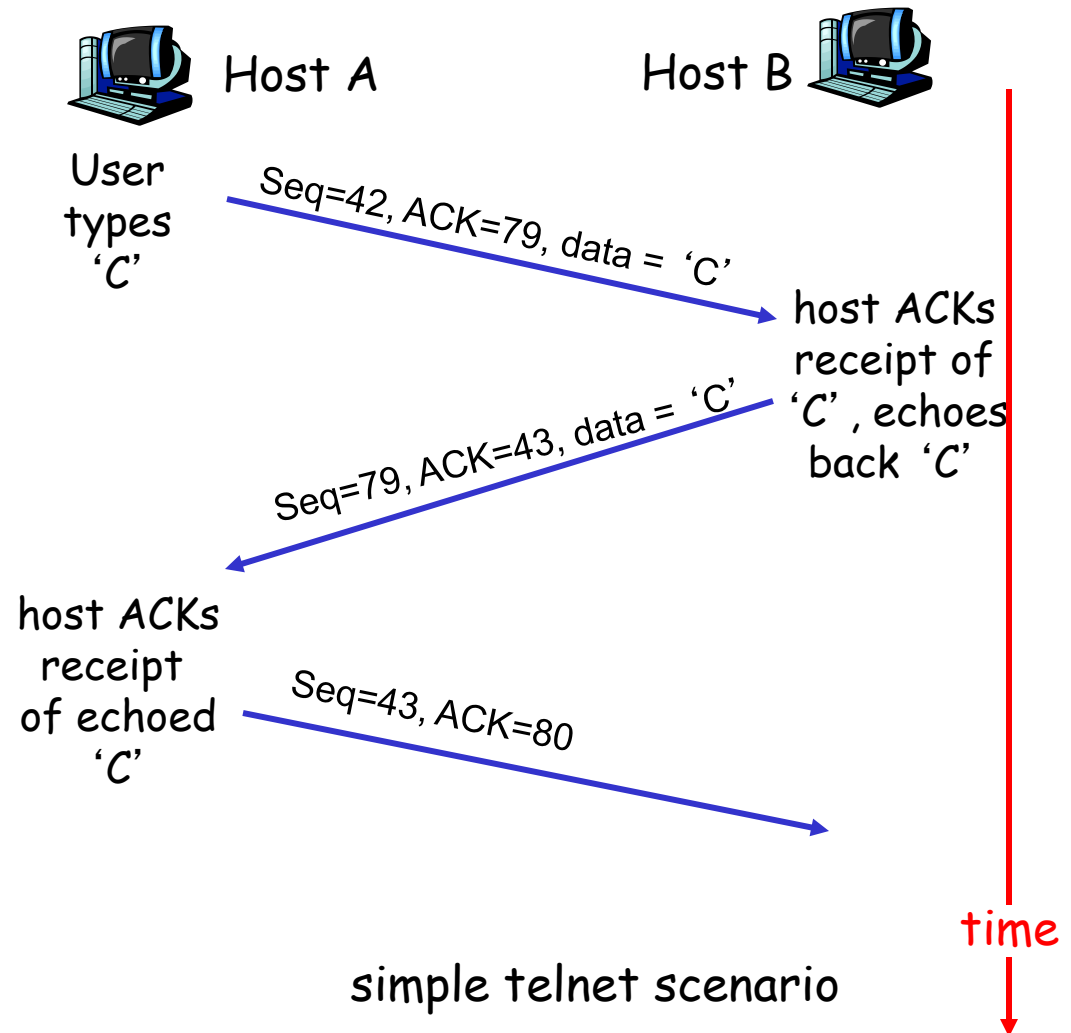
# TCP seq. #'s and ACKs

## Seq. #'s:

- m byte stream
- “number” of first byte in segment's data

## ACKs:

- m seq # of next byte expected from other side
- m cumulative ACK



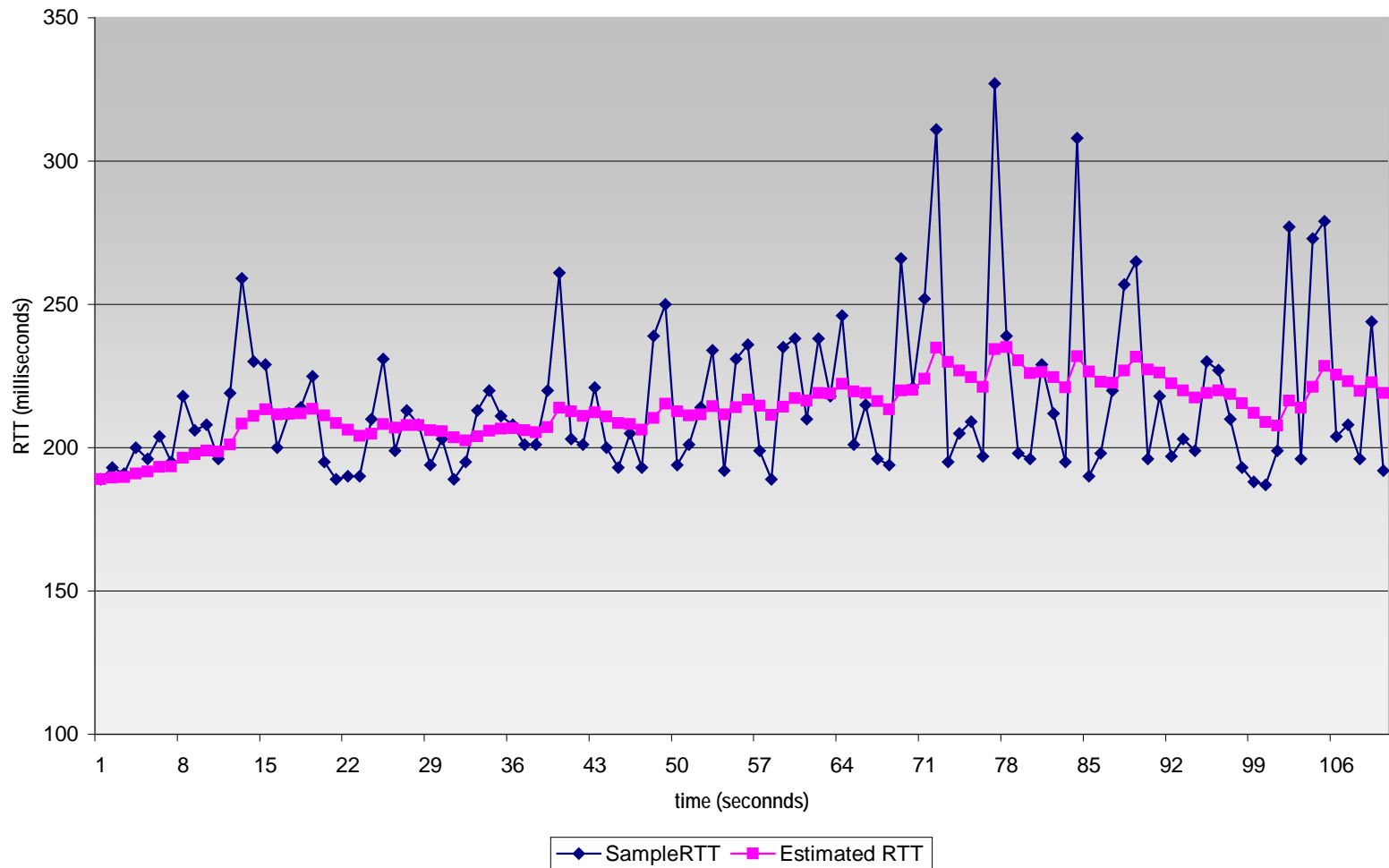
# TCP Round Trip Time and Timeout

$$\text{EstimatedRTT} = (1 - \alpha) * \text{EstimatedRTT} + \alpha * \text{SampleRTT}$$

- r Exponential weighted moving average
- r influence of past sample decreases exponentially fast
- r typical value:  $\alpha = 0.125$

# Example RTT estimation:

RTT: gaia.cs.umass.edu to fantasia.eurecom.fr

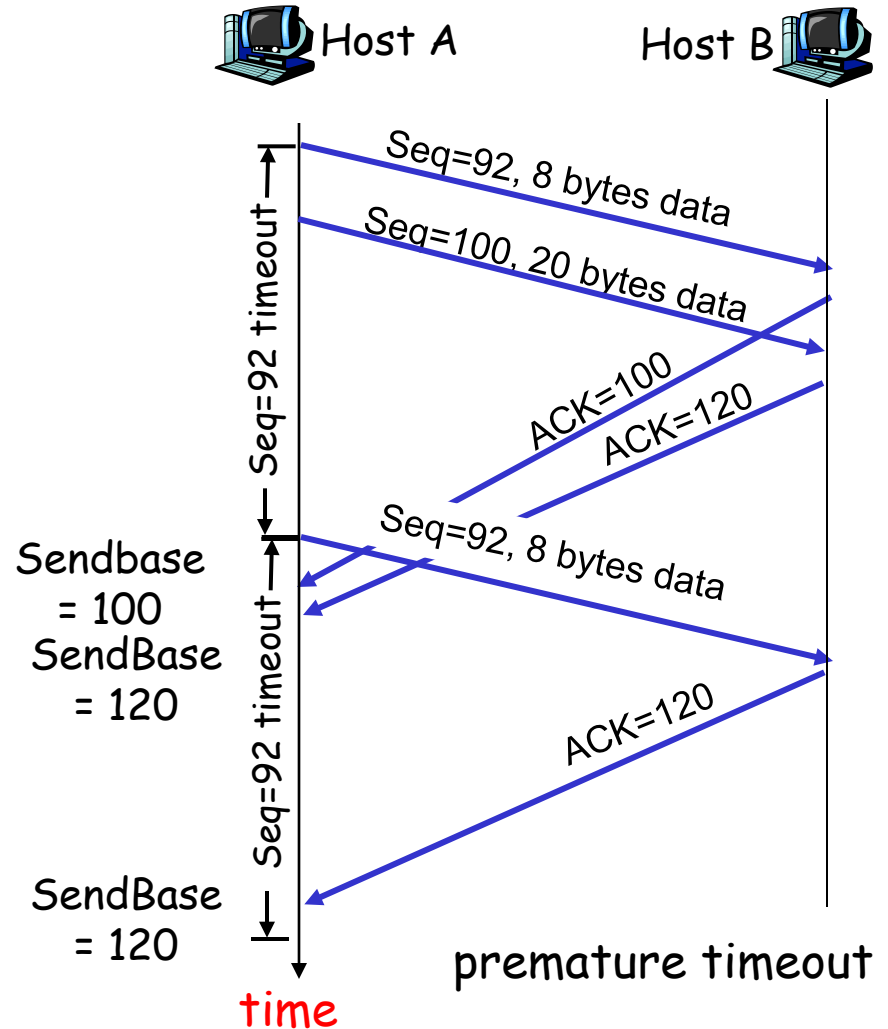
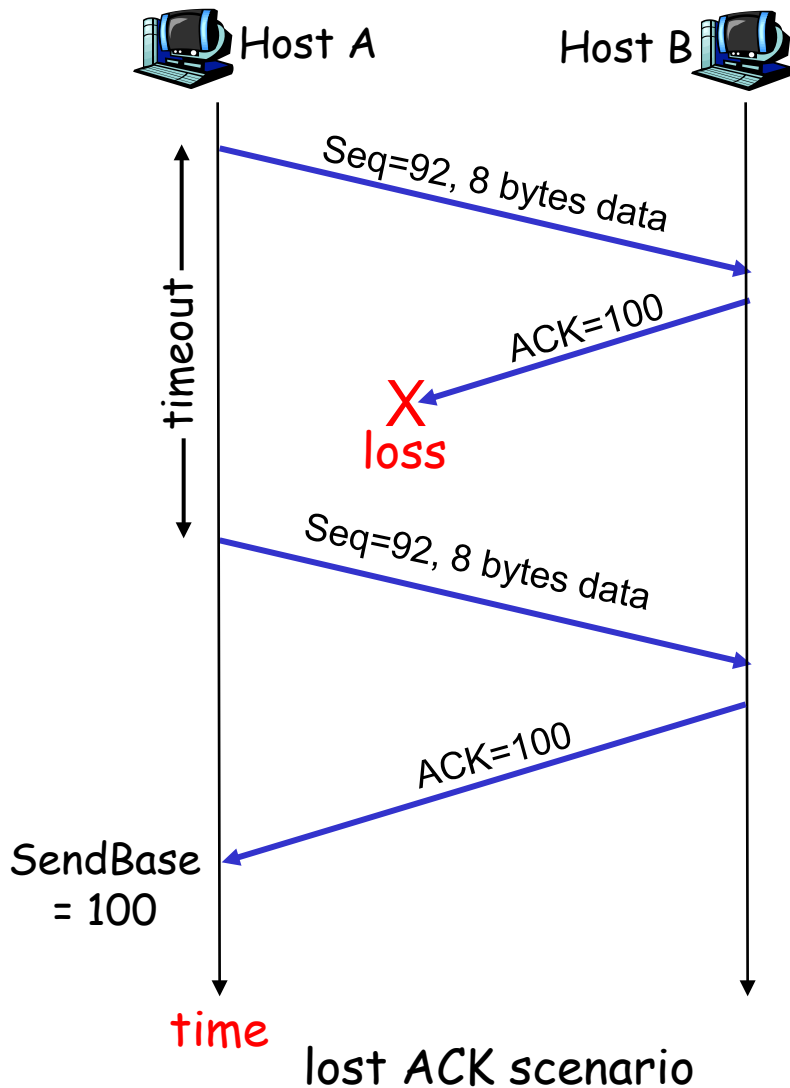




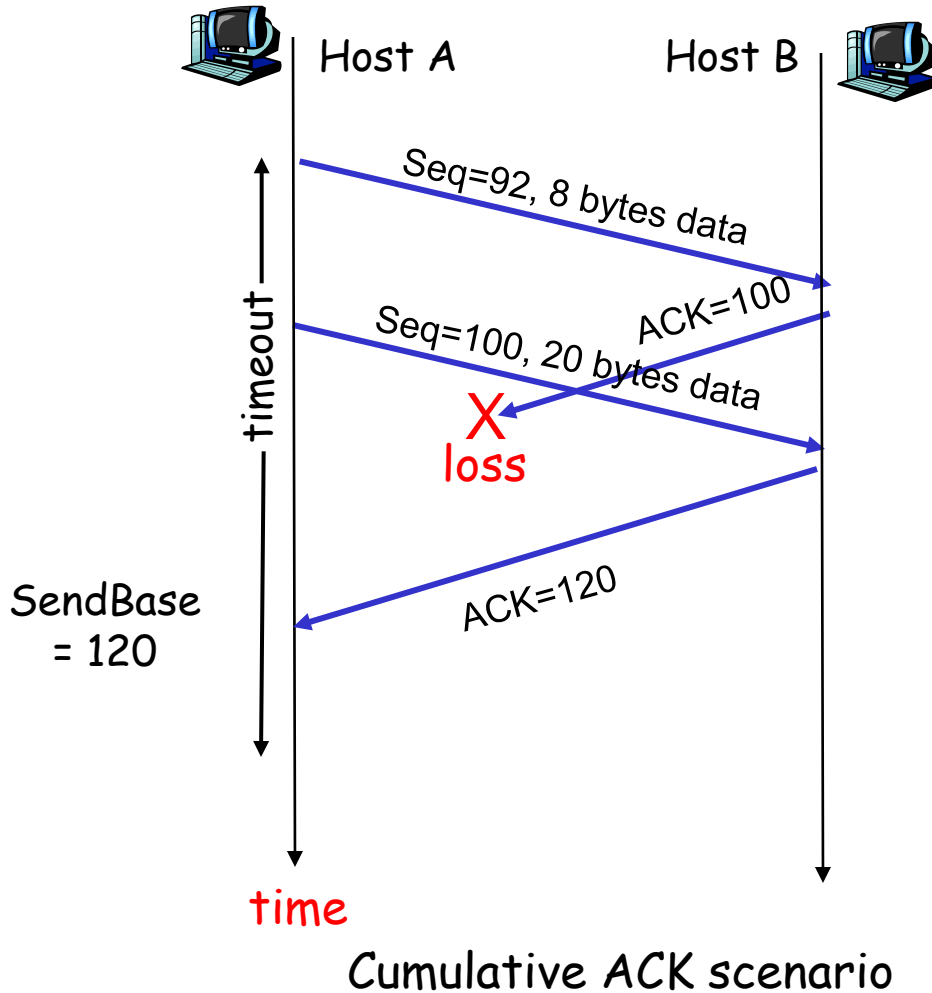
# Chapter 3 outline

- r 3.1 Transport-layer services
- r 3.2 Multiplexing and demultiplexing
- r 3.3 Connectionless transport: UDP
- r 3.4 Principles of reliable data transfer
- r 3.5 Connection-oriented transport: TCP
  - m segment structure
  - m **reliable data transfer**
  - m flow control
  - m connection management
- r 3.6 Principles of congestion control
- r 3.7 TCP congestion control

# TCP: retransmission scenarios

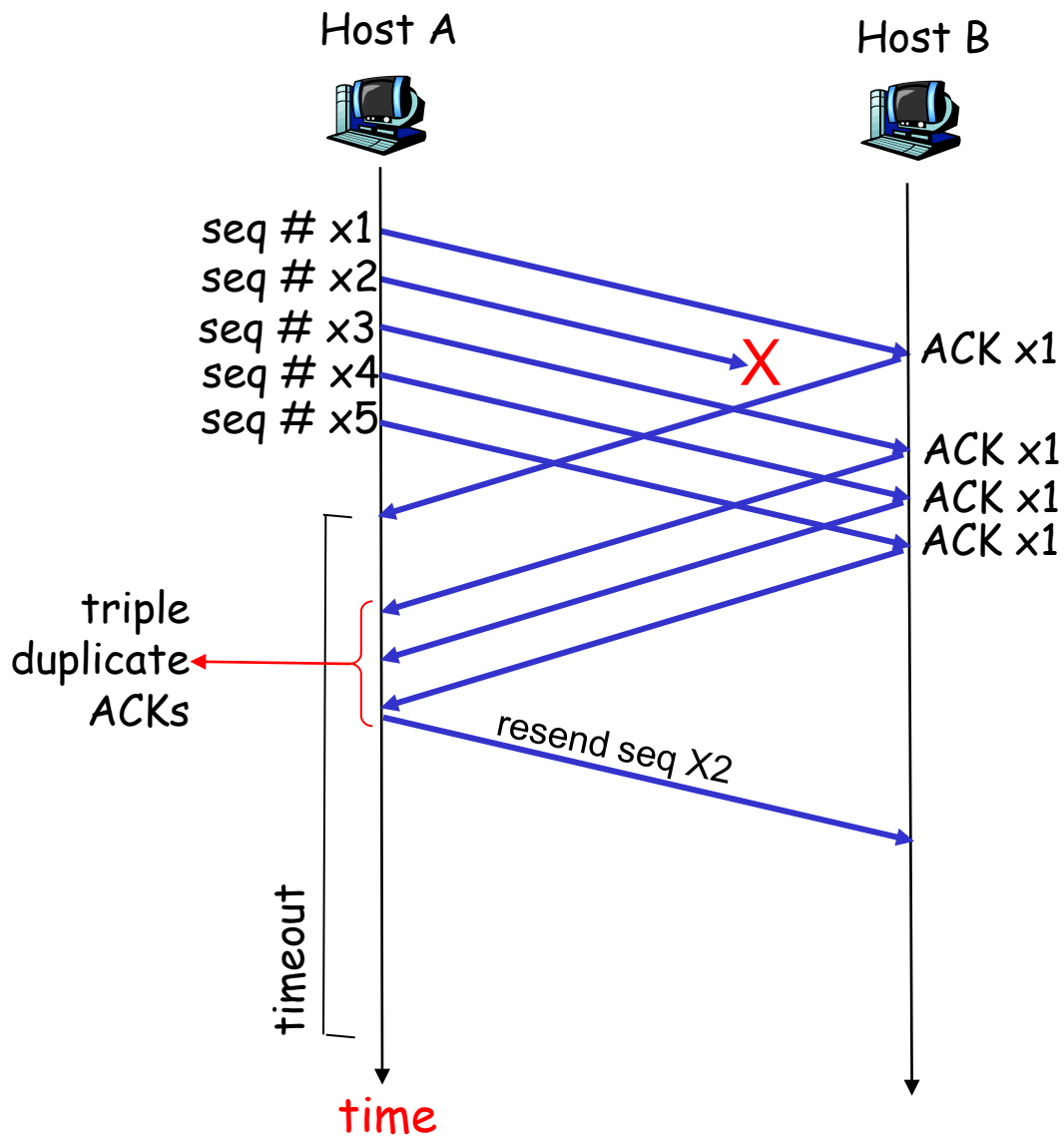


# TCP retransmission scenarios (more)



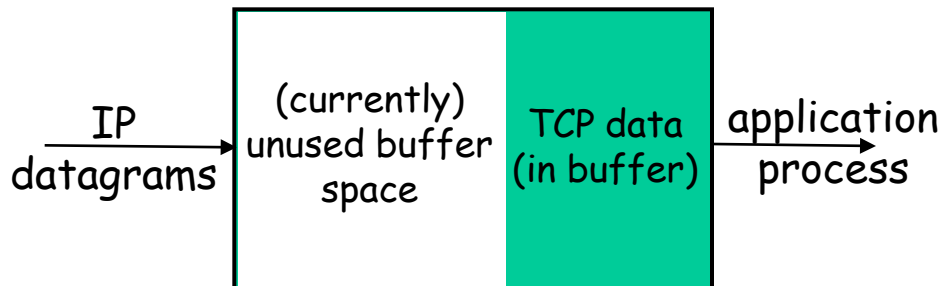
# Fast Retransmit

- r time-out period often relatively long:
  - m long delay before resending lost packet
- r detect lost segments via duplicate ACKs.
  - m sender often sends many segments back-to-back
  - m if segment is lost, there will likely be many duplicate ACKs for that segment
- r If sender receives 3 ACKs for same data, it assumes that segment after ACKed data was lost:
  - m fast retransmit: resend segment before timer expires



# TCP Flow Control

- r receive side of TCP connection has a receive buffer:



- r app process may be slow at reading from buffer

## flow control

sender won't overflow receiver's buffer by transmitting too much, too fast

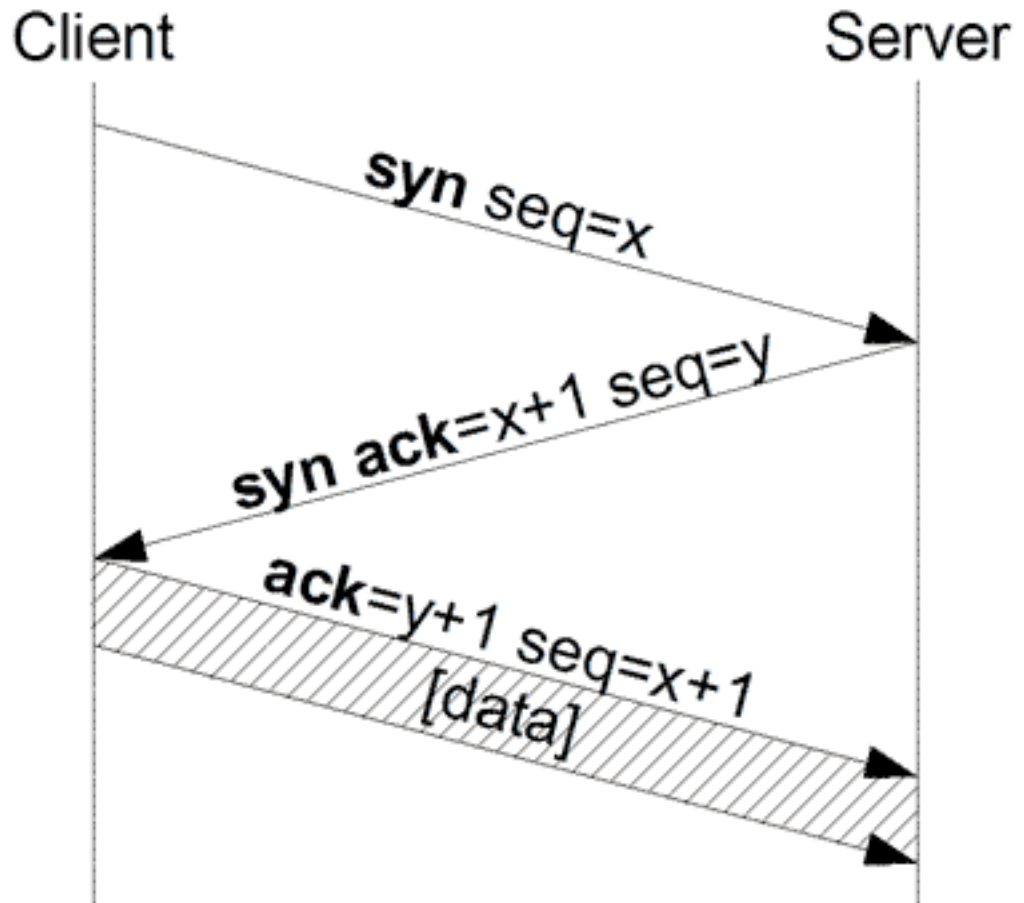
- r *speed-matching service*: matching send rate to receiving application's drain rate

# TCP Three way handshake

Step 1: client host sends TCP SYN segment to server

Step 2: server host receives SYN, replies with SYNACK segment

Step 3: client receives SYNACK, replies with ACK segment, which may contain data



# Chapter 3 outline

- r 3.1 Transport-layer services
- r 3.2 Multiplexing and demultiplexing
- r 3.3 Connectionless transport: UDP
- r 3.4 Principles of reliable data transfer
- r 3.5 Connection-oriented transport: TCP
  - m segment structure
  - m reliable data transfer
  - m flow control
  - m connection management
- r 3.6 Principles of congestion control
- r 3.7 TCP congestion control



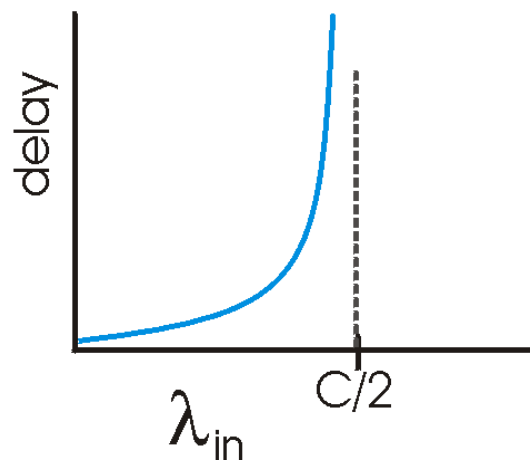
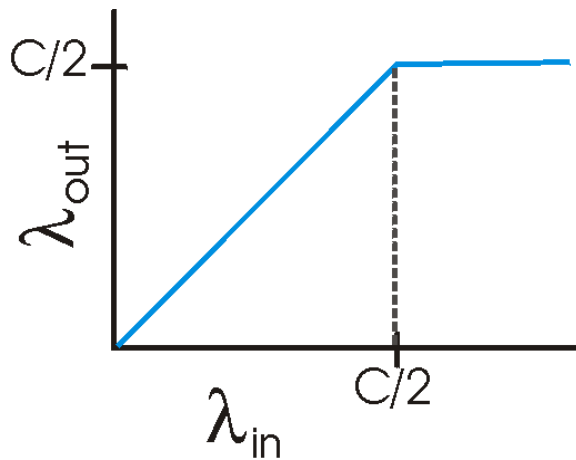
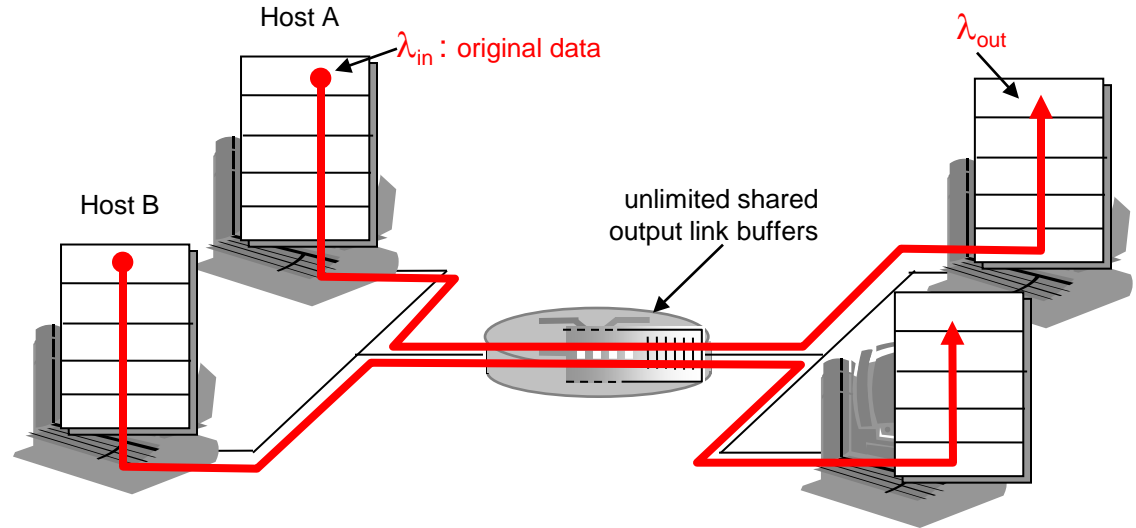
# Principles of Congestion Control

## Congestion:

- r informally: “too many sources sending too much data too fast for *network* to handle”
- r different from flow control!
- r manifestations:
  - m lost packets (buffer overflow at routers)
  - m long delays (queueing in router buffers)
- r a top-10 problem!

# Causes/costs of congestion: scenario 1

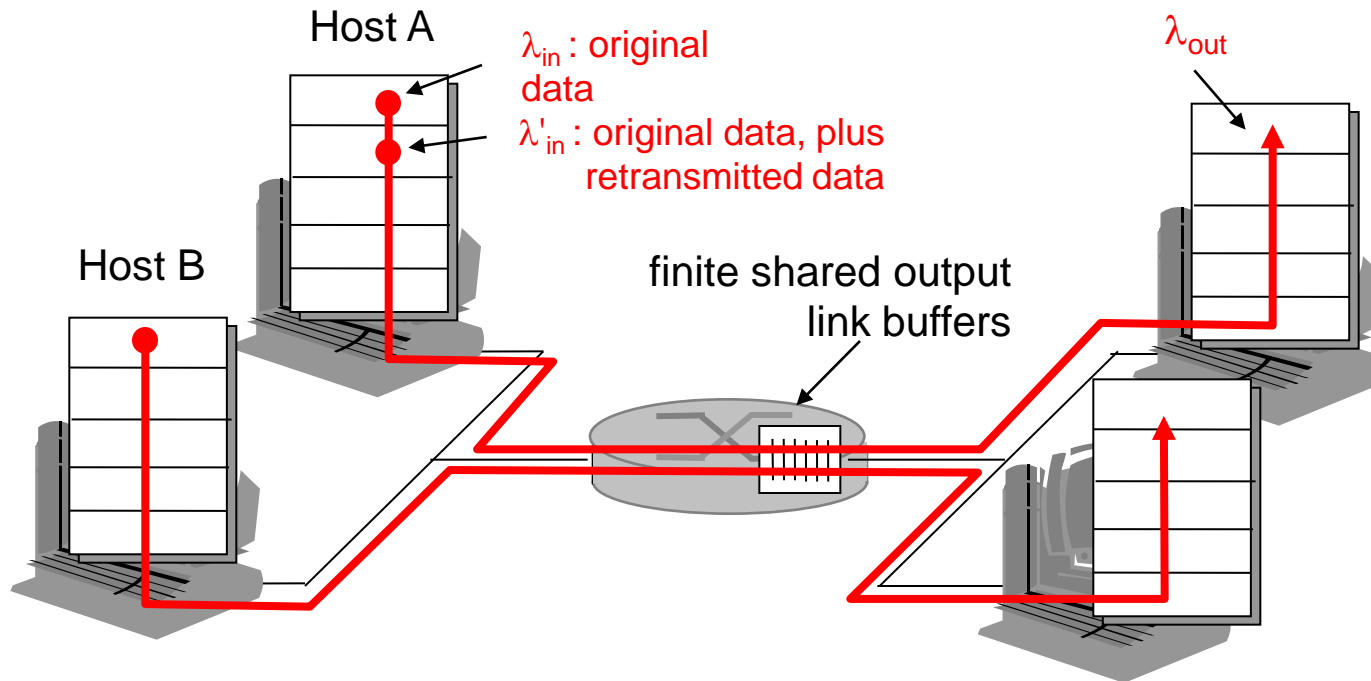
- r two senders, two receivers
- r one router, infinite buffers
- r no retransmission



- r large delays when congested
- r maximum achievable throughput

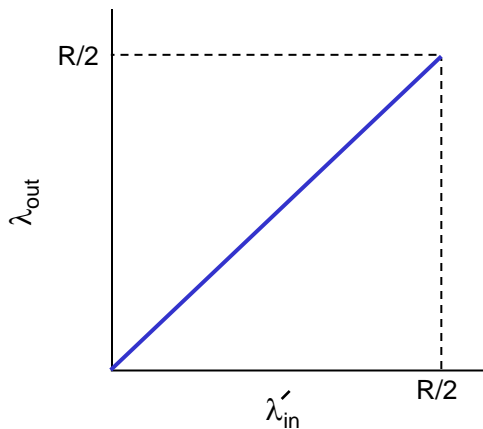
# Causes/costs of congestion: scenario 2

- r one router, *finite* buffers
- r sender retransmission of lost packet

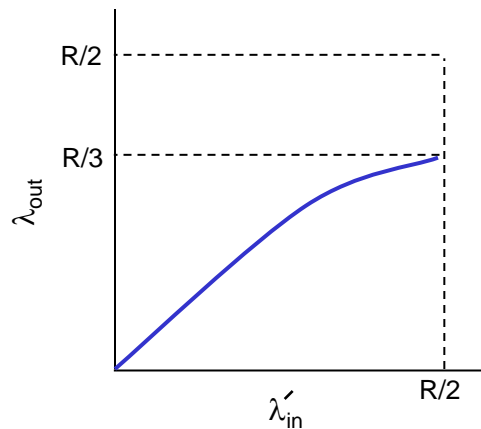


# Causes/costs of congestion: scenario 2

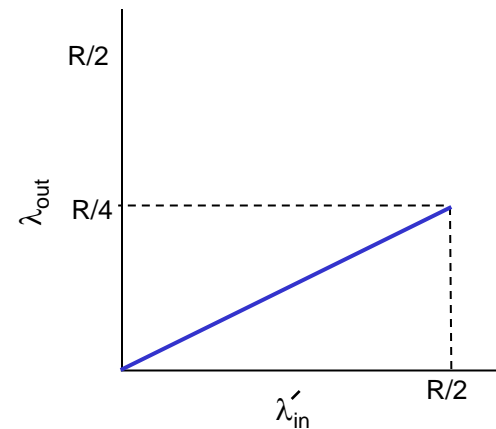
- r always:  $\lambda_{in} = \lambda_{out}$  (goodput)
- r “perfect” retransmission only when loss:  $\lambda'_{in} > \lambda_{out}$
- r retransmission of delayed (not lost) packet makes  $\lambda'_{in}$  larger (than perfect case) for same  $\lambda_{out}$



a.



b.



c.

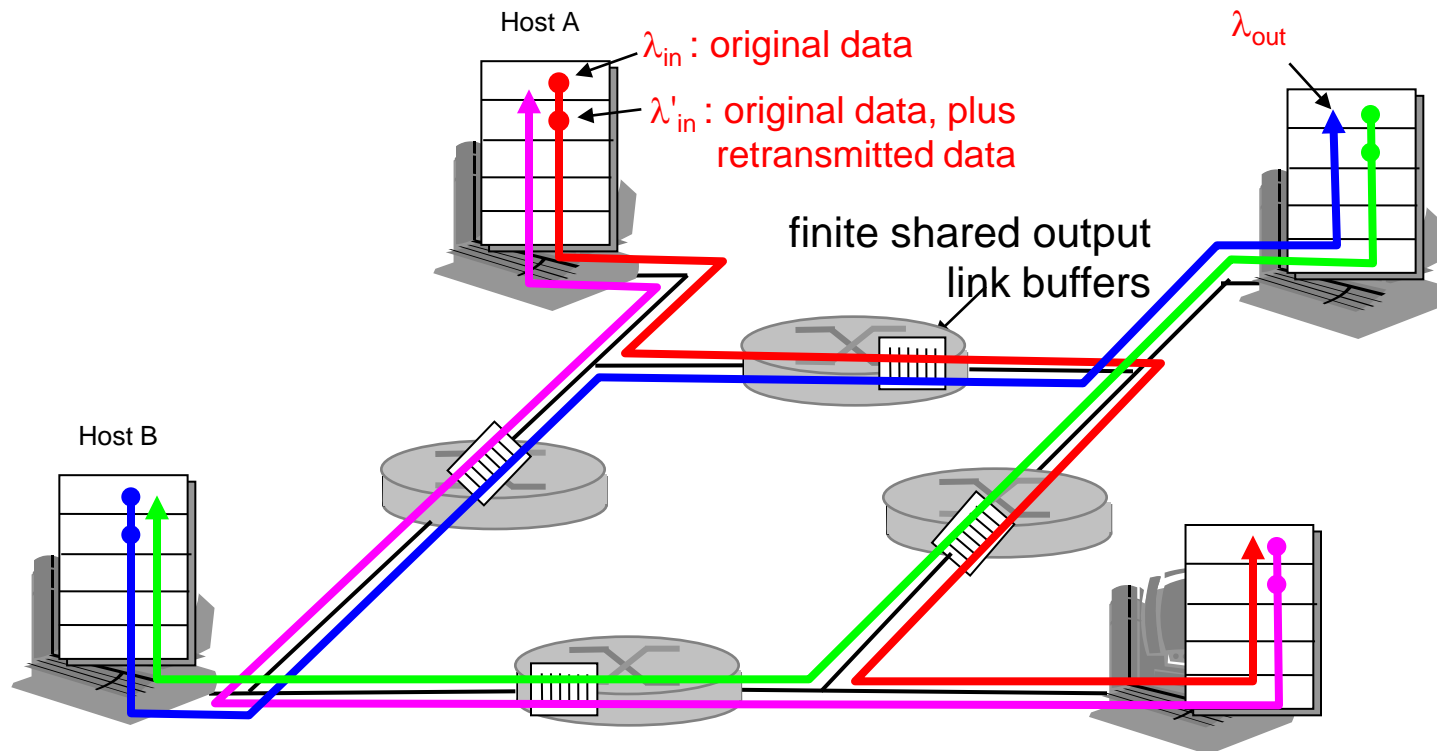
## “costs” of congestion:

- r more work (retrans) for given “goodput”
- r unneeded retransmissions: link carries multiple copies of pkt

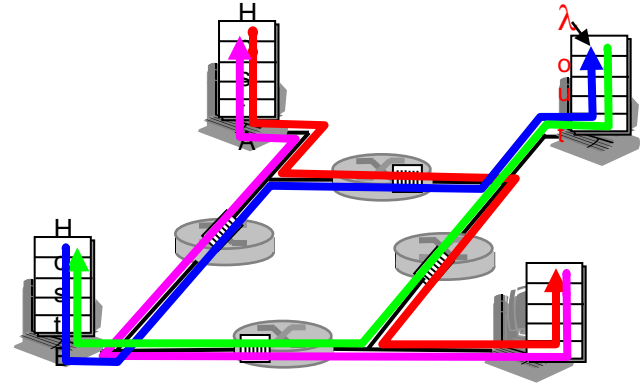
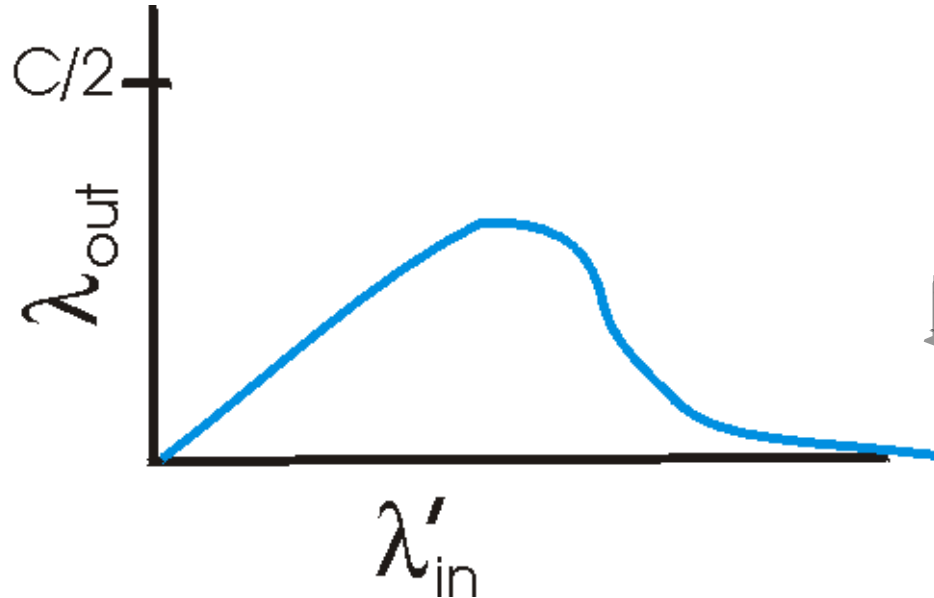
# Causes/costs of congestion: scenario 3

- r four senders
- r multihop paths
- r timeout/retransmit

**Q:** what happens as  $\lambda_{in}$  and  $\lambda'_{in}$  increase ?



## Causes/costs of congestion: scenario 3



another “cost” of congestion:

- r when packet dropped, any “upstream transmission capacity used for that packet was wasted!

# Chapter 3 outline

- r 3.1 Transport-layer services
- r 3.2 Multiplexing and demultiplexing
- r 3.3 Connectionless transport: UDP
- r 3.4 Principles of reliable data transfer
- r 3.5 Connection-oriented transport: TCP
  - m segment structure
  - m reliable data transfer
  - m flow control
  - m connection management
- r 3.6 Principles of congestion control
- r 3.7 TCP congestion control

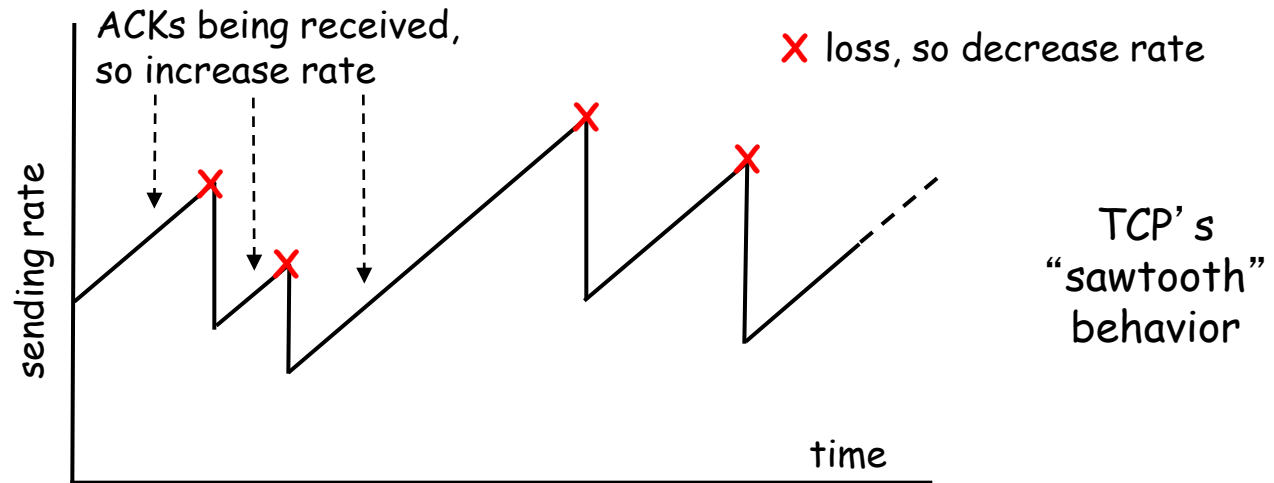
# TCP congestion control:

- r *goal:* TCP sender should transmit as fast as possible, but without congesting network
  - m Q: how to find rate *just* below congestion level
- r decentralized: each TCP sender sets its own rate, based on *implicit* feedback:
  - m *ACK:* segment received (a good thing!), network not congested, so increase sending rate
  - m *lost segment:* assume loss due to congested network, so decrease sending rate

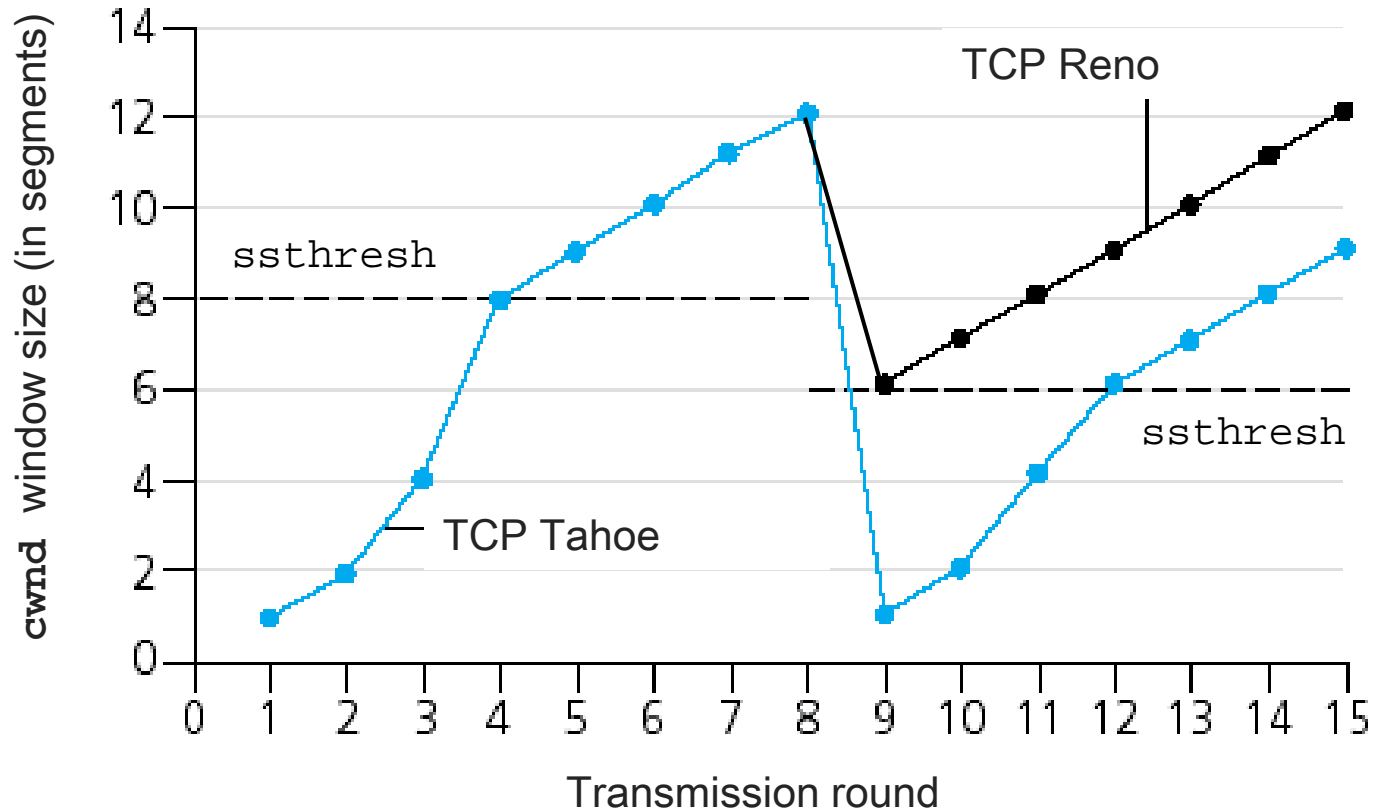


# TCP congestion control: bandwidth probing

- r “probing for bandwidth”: increase transmission rate on receipt of ACK, until eventually loss occurs, then decrease transmission rate
- m continue to increase on ACK, decrease on loss (since available bandwidth is changing, depending on other connections in network)



# Popular “flavors” of TCP



## Summary: TCP Congestion Control

- r when  $\text{cwnd} < \text{ssthresh}$ , sender in **slow-start** phase, window grows exponentially.
- r when  $\text{cwnd} \geq \text{ssthresh}$ , sender is in **congestion-avoidance** phase, window grows linearly.
- r when **triple duplicate ACK** occurs,  $\text{ssthresh}$  set to  $\text{cwnd}/2$ ,  $\text{cwnd}$  set to  $\sim \text{ssthresh}$
- r when **timeout** occurs,  $\text{ssthresh}$  set to  $\text{cwnd}/2$ ,  $\text{cwnd}$  set to 1 MSS.

# Chapter 3: Summary

- r principles behind transport layer services:
  - m multiplexing, demultiplexing
  - m reliable data transfer
  - m flow control
  - m congestion control
- r instantiation and implementation in the Internet
  - m UDP
  - m TCP

## Next:

- r leaving the network “edge” (application, transport layers)
- r into the network “core”