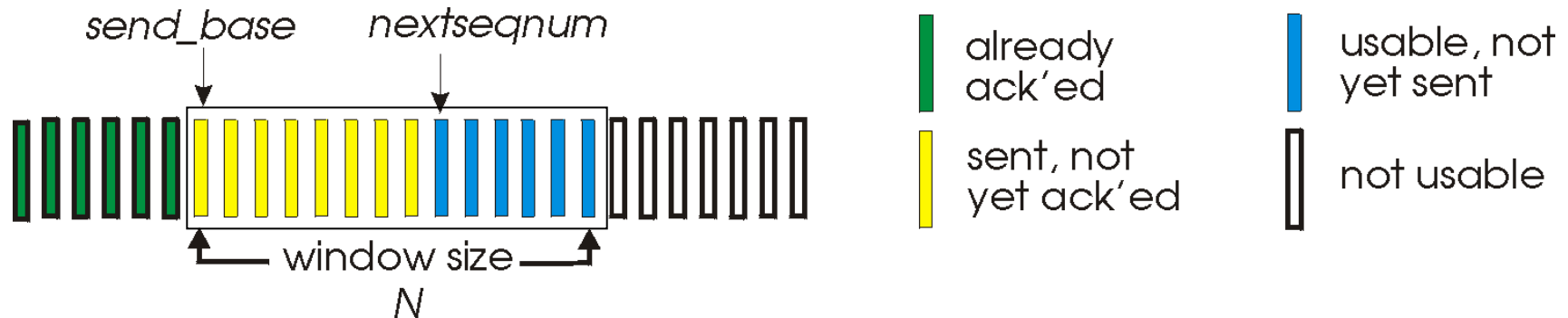


Overview of last class

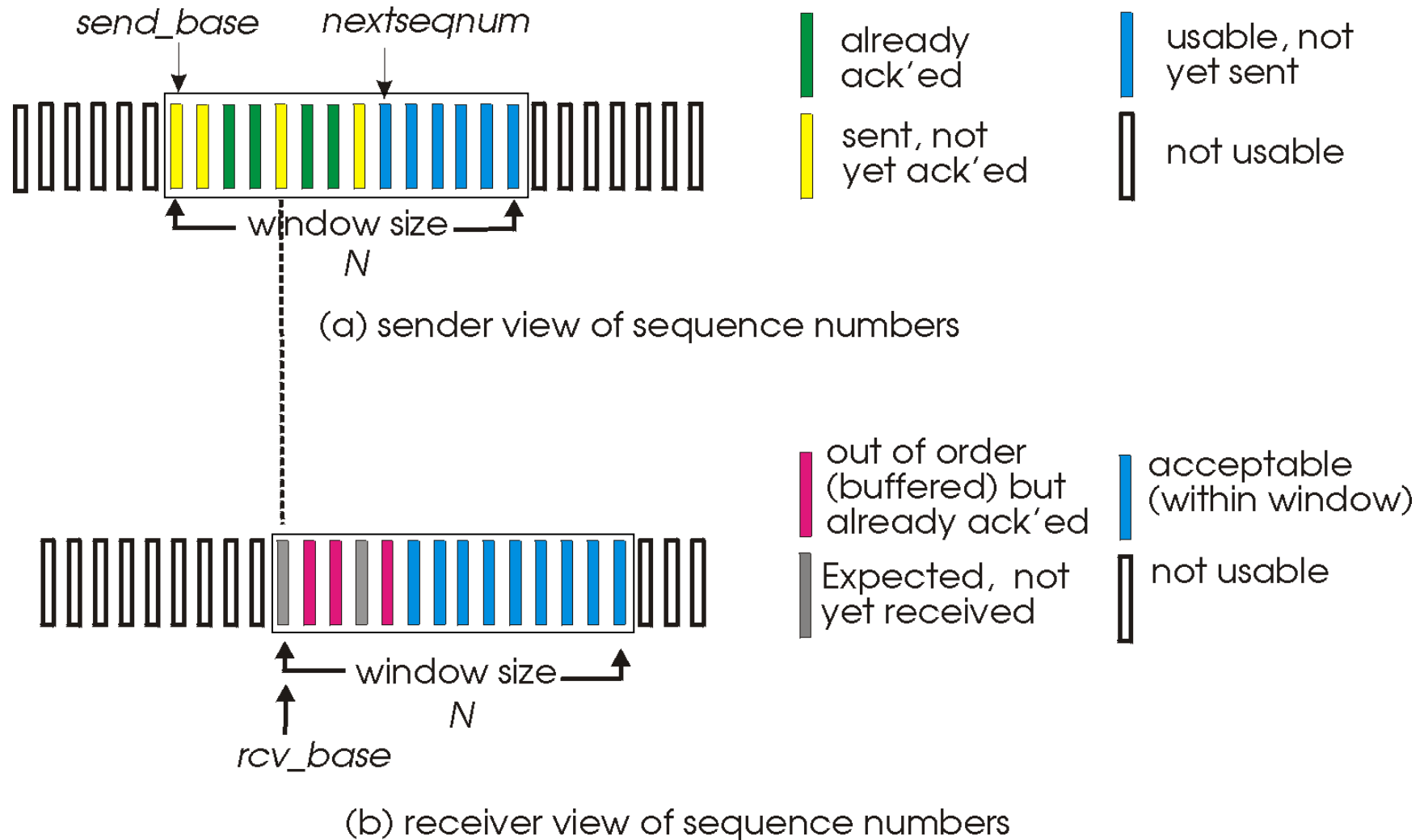
Go-Back-N: sender

- k-bit seq # in pkt header
- “window” of up to N, consecutive unack’ed pkts allowed



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

Selective repeat: sender, receiver windows



All packets on the left of the window has been successfully delivered in order !

Up to N packets in processing;
in-order, successfully delivered pkts
slide out of the window

sliding

Sequence
number

Selective repeat: dilemma

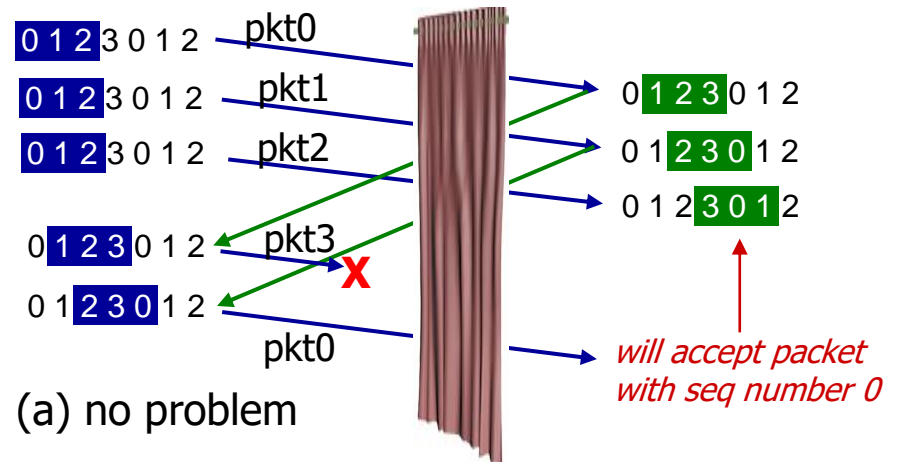
example:

- seq #'s: 0, 1, 2, 3
- window size=3
- receiver sees no difference in two scenarios!
- duplicate data accepted as new in (b)

Q: what relationship between seq # size and window size to avoid problem in (b)?

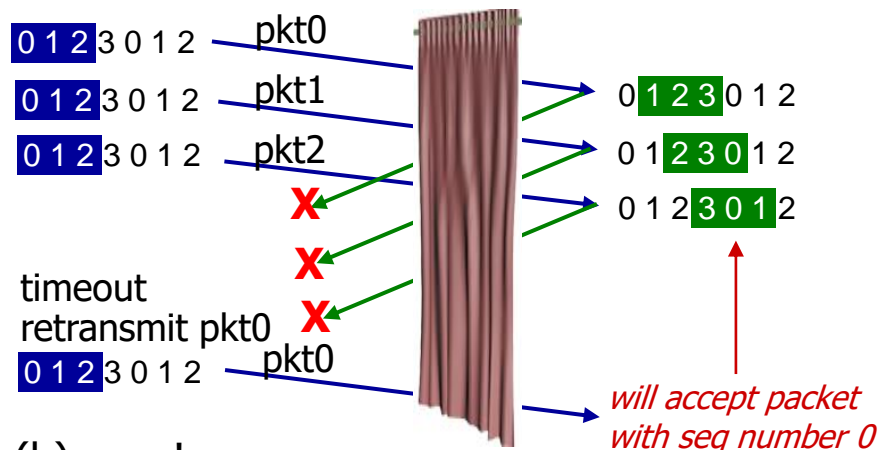
sender window
(after receipt)

receiver window
(after receipt)



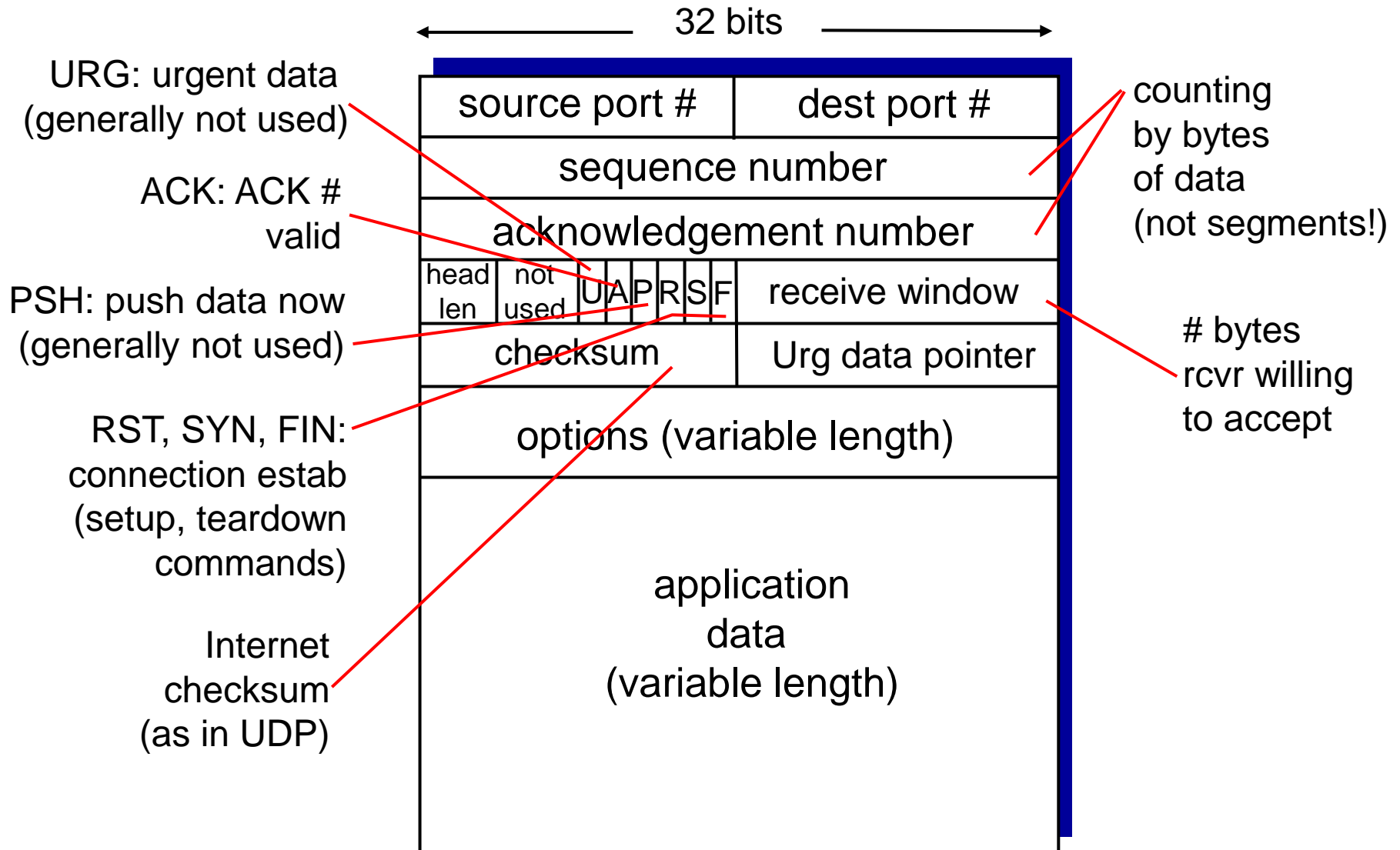
(a) no problem

*receiver can't see sender side.
receiver behavior identical in both cases!
something's (very) wrong!*

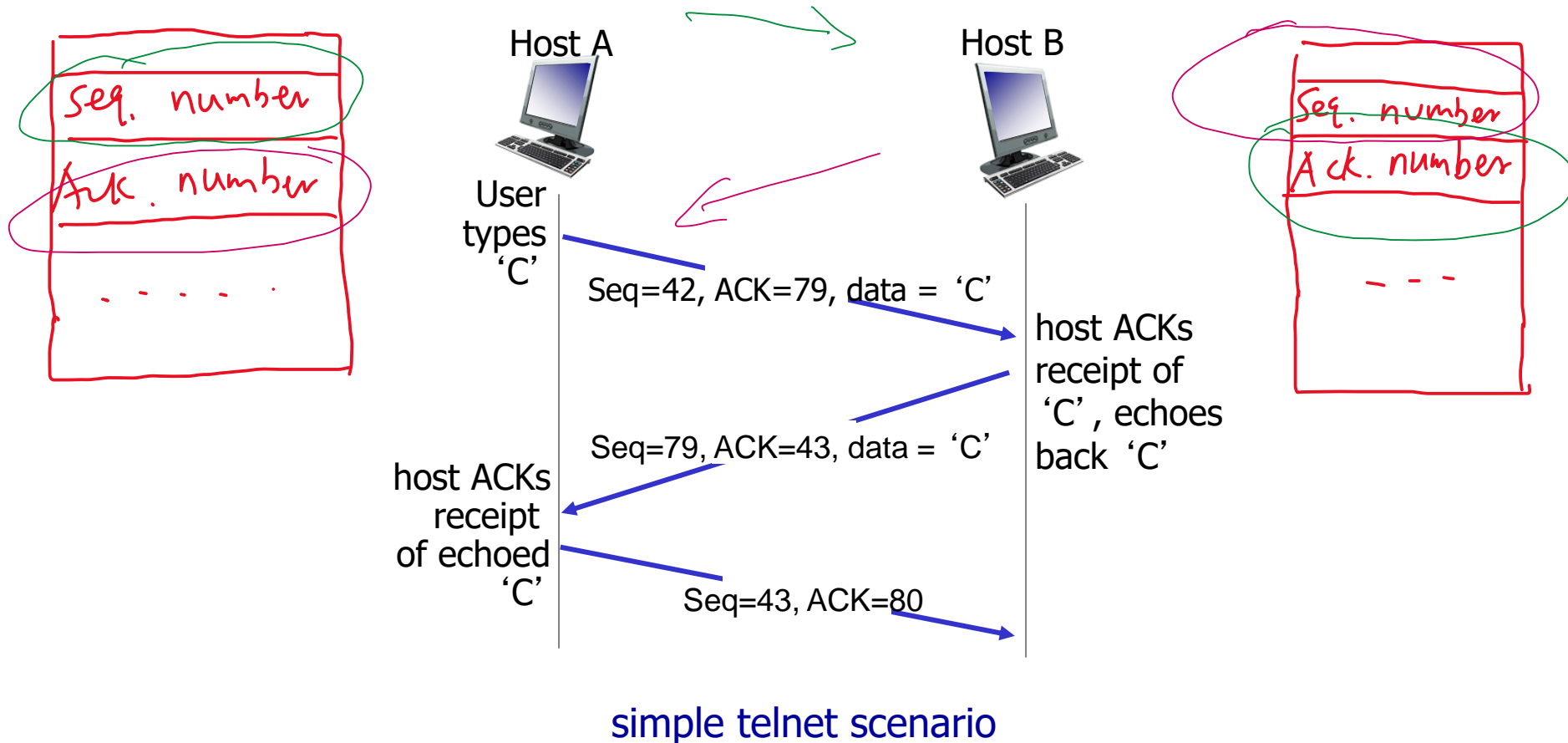


(b) oops!

TCP segment structure



TCP seq. numbers, ACKs

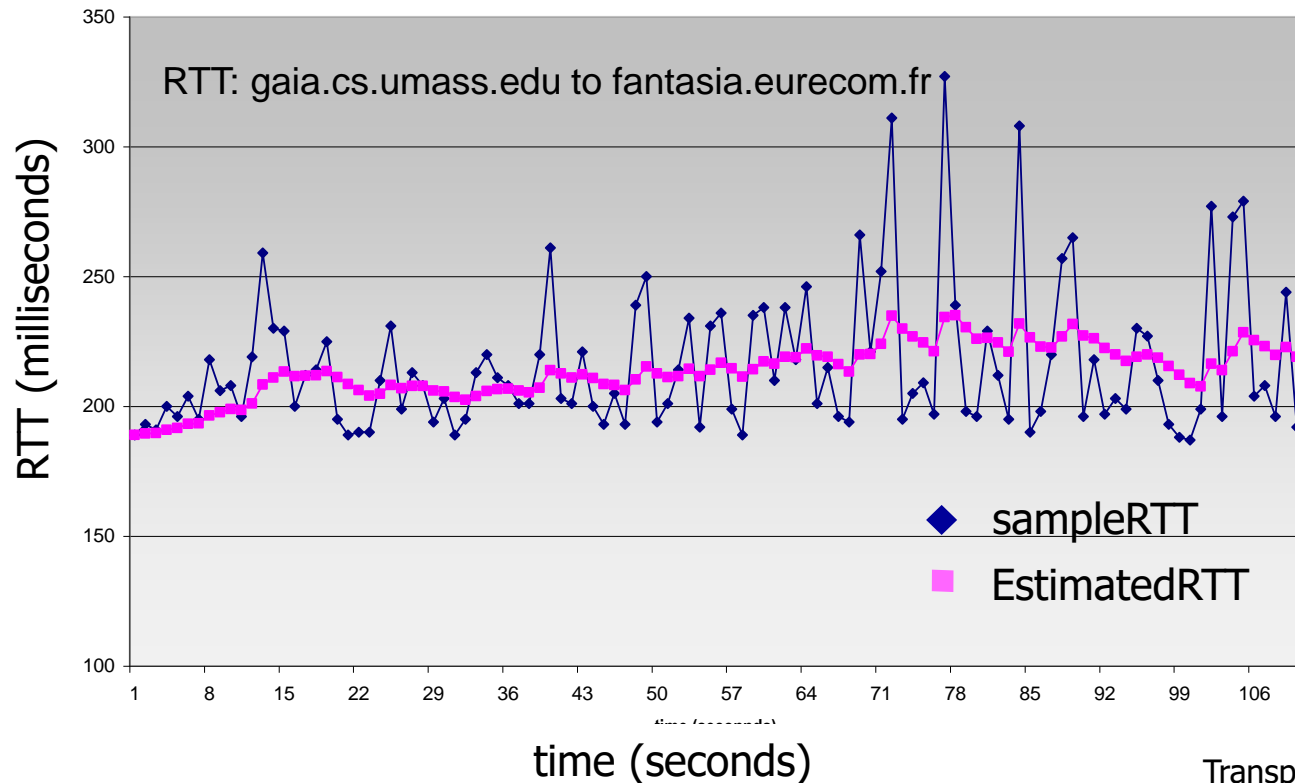


TCP round trip time, timeout

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

$\text{TimeoutInterval} = \text{EstimatedRTT} + 4 * \text{DevRTT}$

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



Class Today

Chapter 3 outline

3.1 transport-layer services

3.2 multiplexing and demultiplexing

3.3 connectionless transport: UDP

3.4 principles of reliable data transfer

3.5 connection-oriented transport: TCP

- segment structure
- reliable data transfer
- flow control
- connection management

3.6 principles of congestion control

3.7 TCP congestion control

TCP reliable data transfer

- TCP creates rdt service on top of IP' s unreliable service
 - pipelined segments
 - **cumulative acks**
 - **single retransmission timer**
- retransmissions triggered by:
 - timeout events
 - duplicate acks

let' s initially consider simplified TCP sender:

- ignore duplicate acks
- ignore flow control, congestion control

TCP sender events:

data rcvd from app:

- create segment with seq #
- seq # is byte-stream number of first data byte in segment
- start timer if not already running
 - think of timer as for oldest unacked segment
 - expiration interval: `TimeoutInterval`

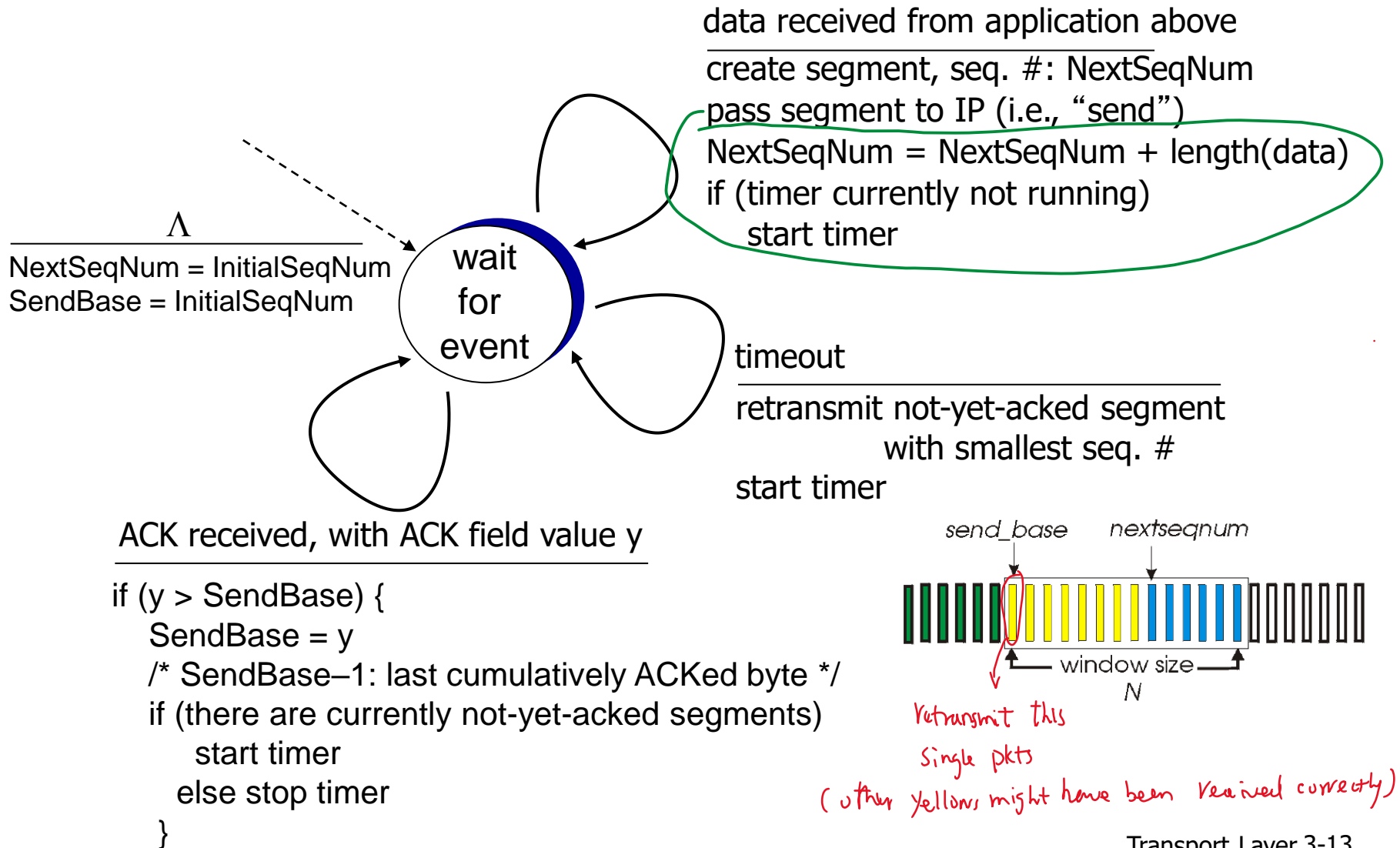
timeout:

- retransmit segment that caused timeout
- restart timer

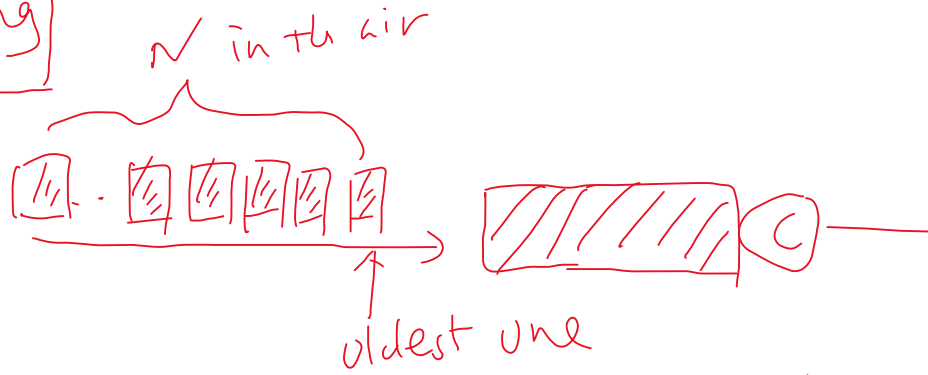
ack rcvd:

- if ack acknowledges previously unacked segments
 - update what is known to be ACKed
 - start timer if there are still unacked segments

TCP sender (simplified)



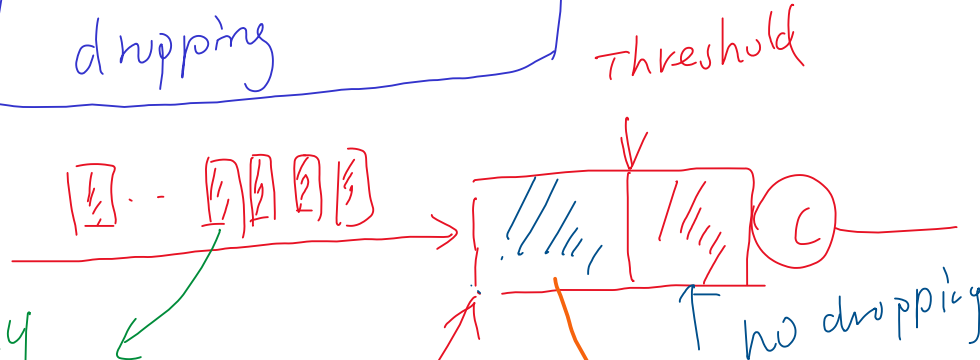
Tail-dropping



All dropped when buffer is full

Random early detection (RED)

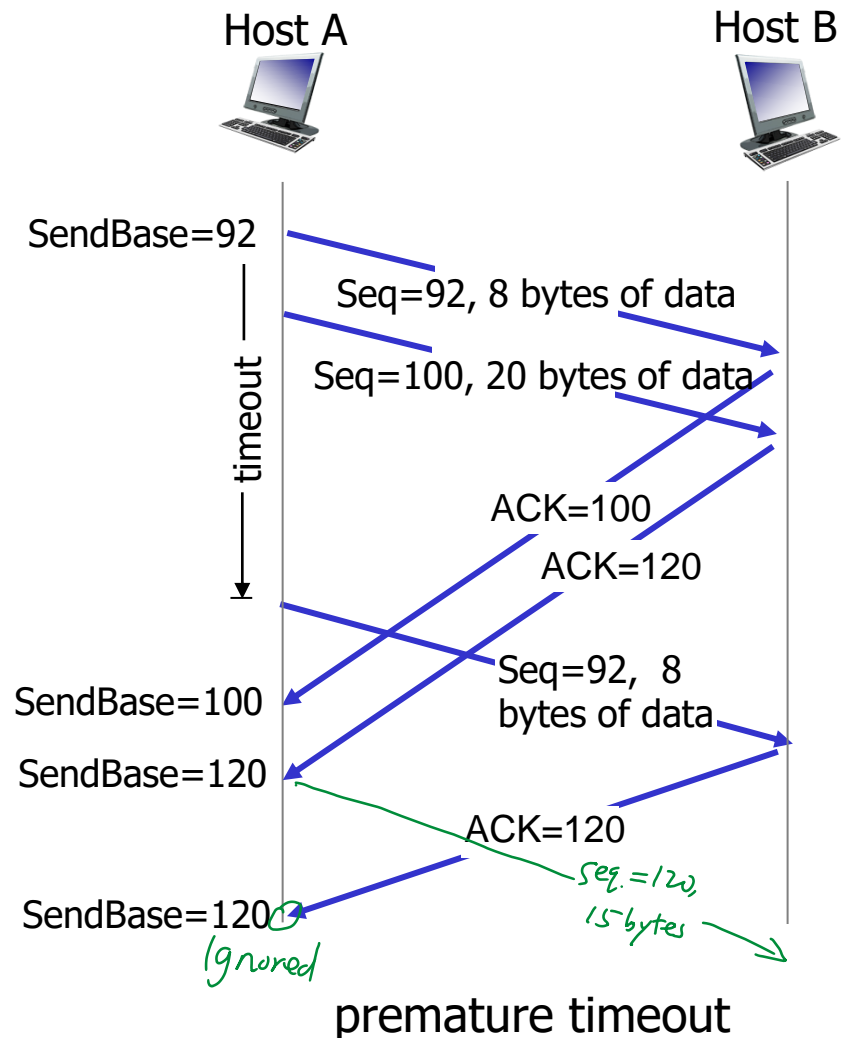
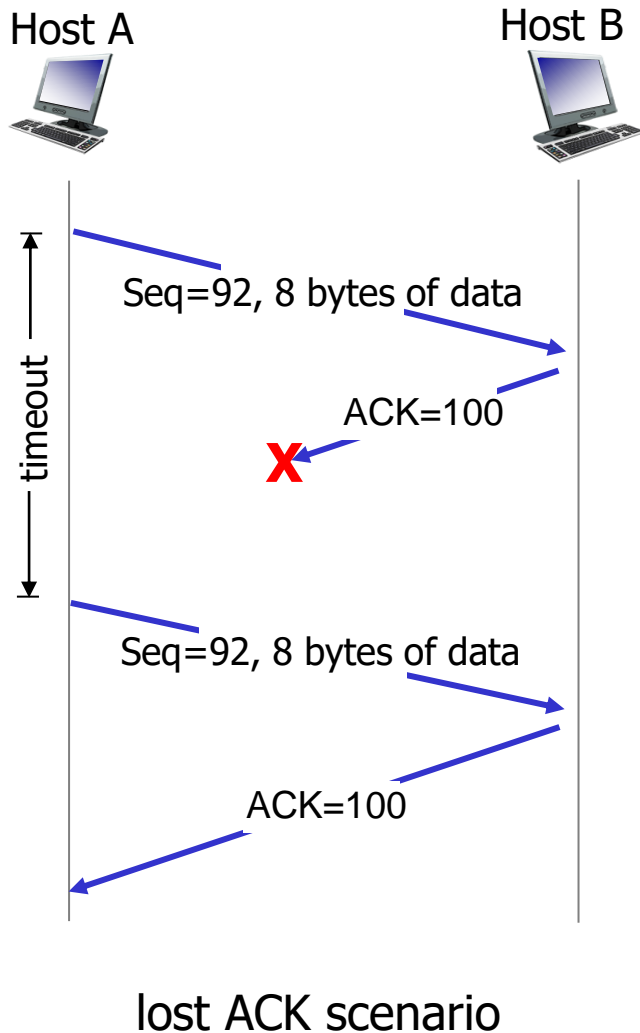
https://en.wikipedia.org/wiki/Random_early_detection



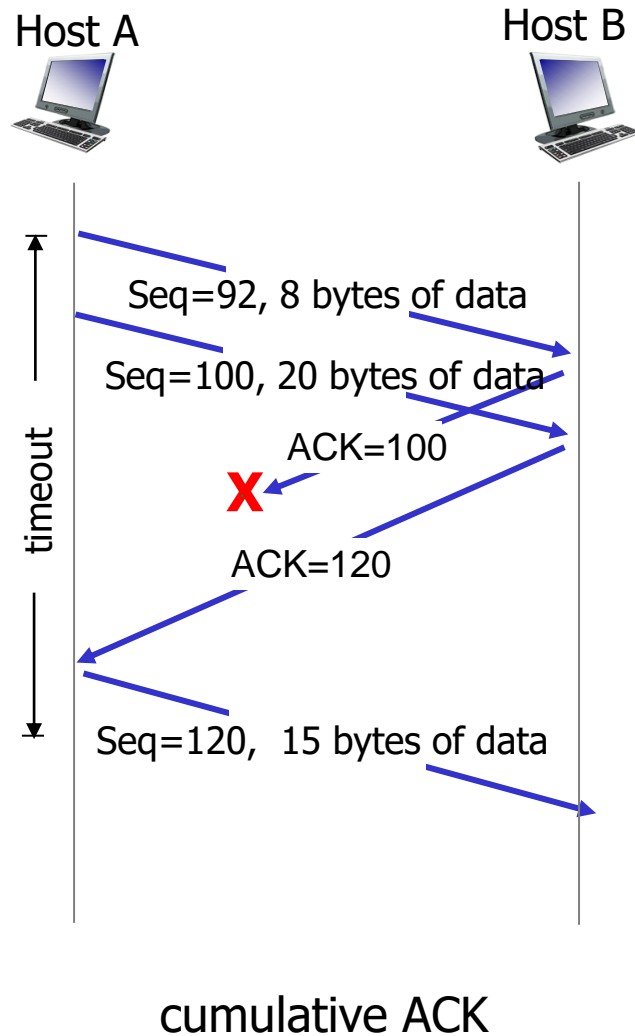
each may have a drop prob. independently

Dropping with a prob. depending on queue length

TCP: retransmission scenarios



TCP: retransmission scenarios



TCP ACK generation [RFC 1122, RFC 2581]

event at receiver

arrival of in-order segment with expected seq #. All data up to expected seq # already ACKed

arrival of in-order segment with expected seq #. One other segment has ACK pending

arrival of out-of-order segment higher-than-expect seq. # . Gap detected

arrival of segment that partially or completely fills gap
(Out of order buffered)

TCP receiver action

delayed ACK. Wait up to 500ms for next segment. If no next segment, send ACK

immediately send single cumulative ACK, ACKing both in-order segments

immediately send *duplicate ACK*, indicating seq. # of next expected byte

immediate send ACK, provided that segment starts at lower end of gap

TCP fast retransmit

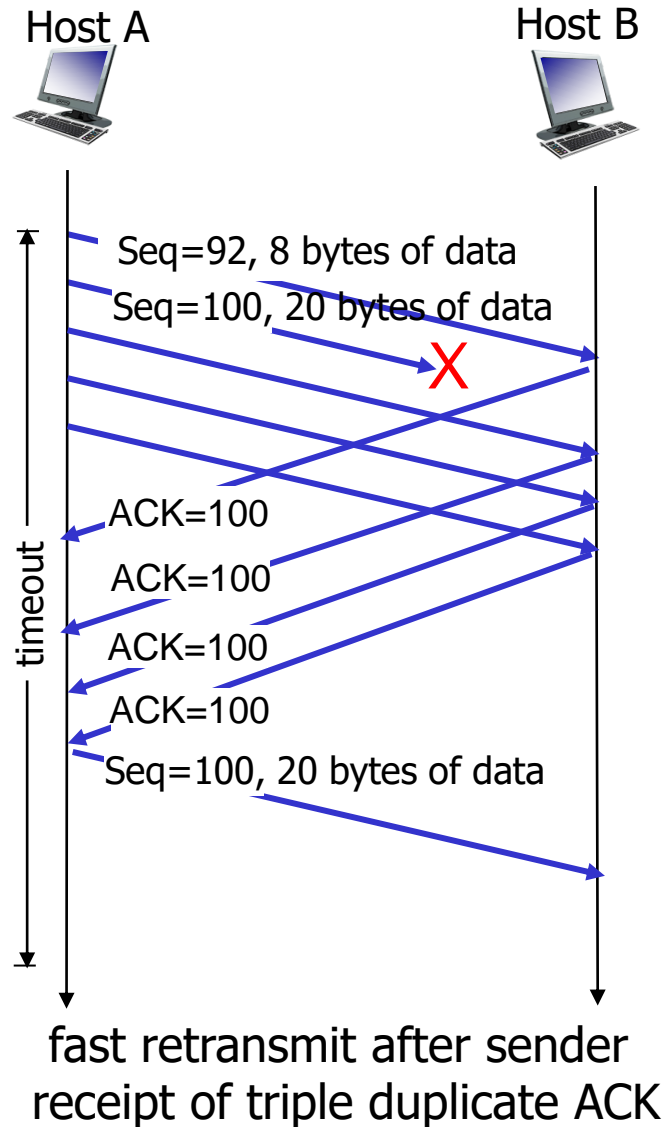
- time-out period often relatively long:
 - long delay before resending lost packet
- detect lost segments via duplicate ACKs.
 - sender often sends many segments back-to-back
 - if segment is lost, there will likely be many duplicate ACKs.

TCP fast retransmit

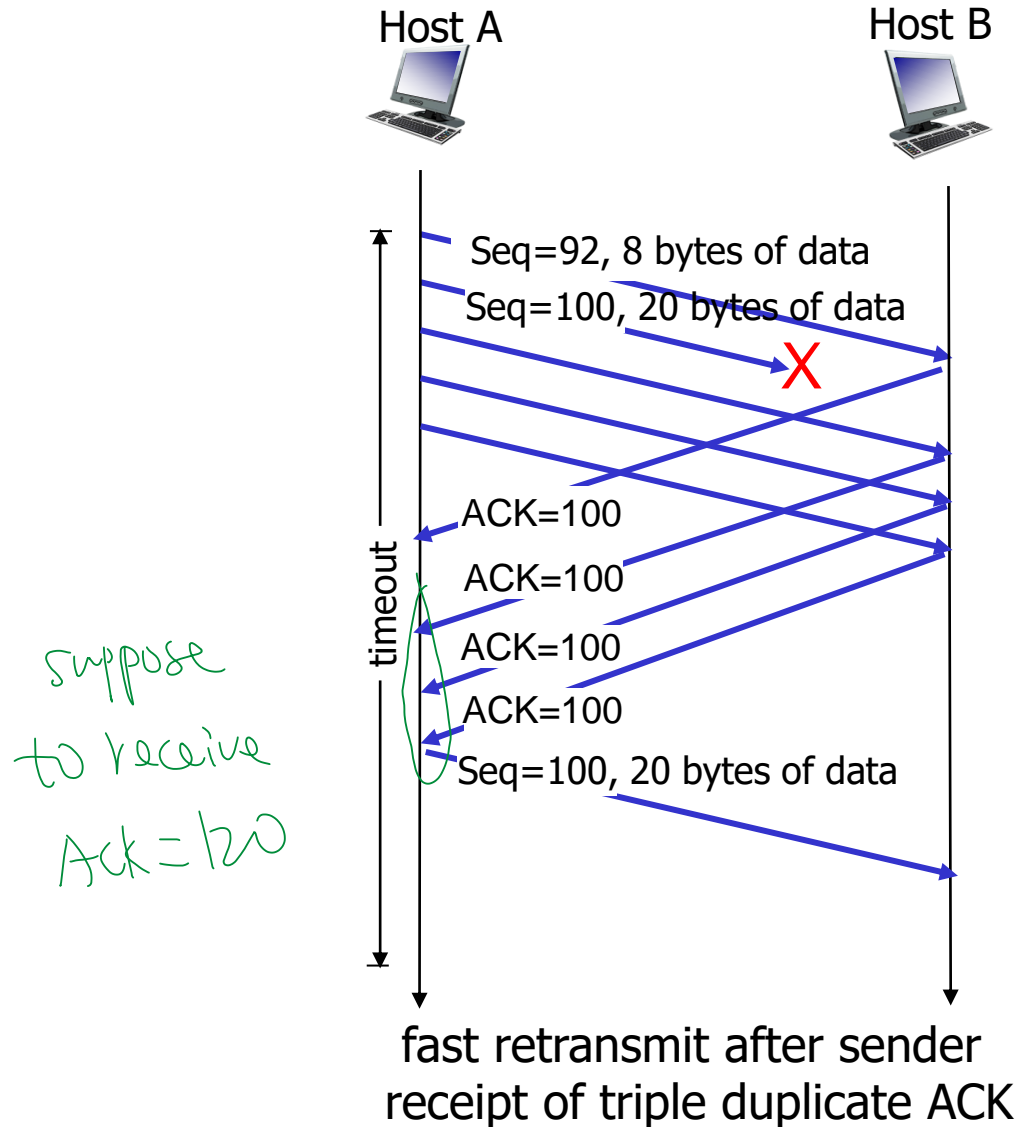
if sender receives 3 ACKs for same data (“triple duplicate ACKs”), resend unacked segment with smallest seq #

- likely that unacked segment lost, so don't wait for timeout

TCP fast retransmit



TCP fast retransmit



Chapter 3 outline

3.1 transport-layer services

3.2 multiplexing and demultiplexing

3.3 connectionless transport: UDP

3.4 principles of reliable data transfer

3.5 connection-oriented transport: TCP

- segment structure
- reliable data transfer
- flow control
- connection management

3.6 principles of congestion control

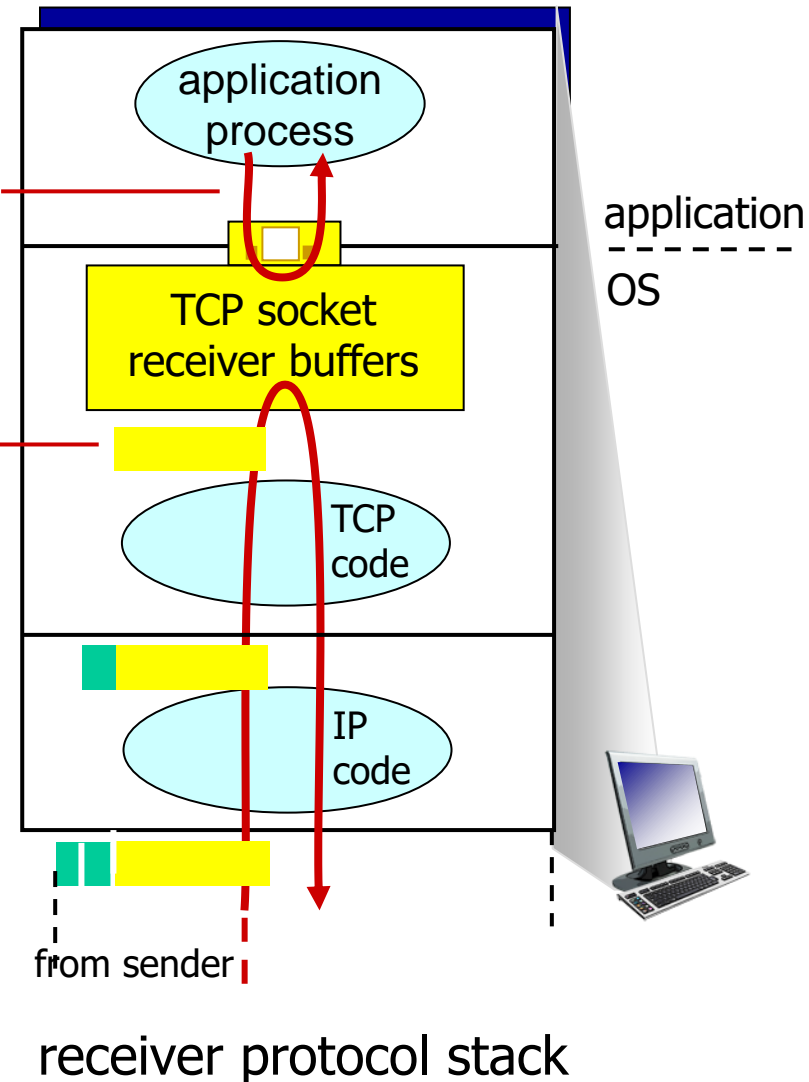
3.7 TCP congestion control

TCP flow control

application may
remove data from
TCP socket buffers

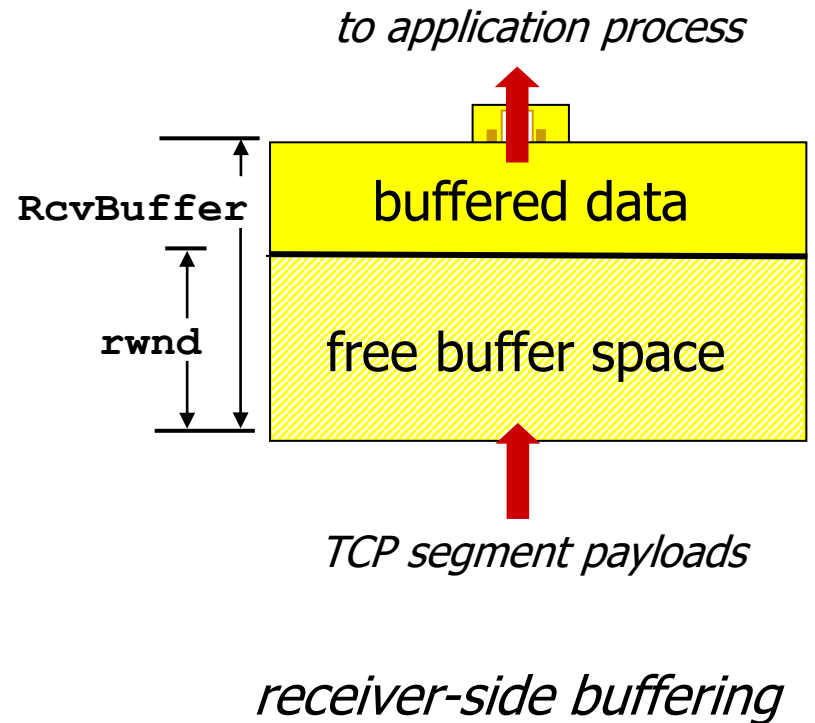
... slower than TCP
receiver is delivering
(sender is sending)

flow control
receiver controls sender, so
sender won't overflow
receiver's buffer by transmitting
too much, too fast



TCP flow control

- receiver “advertises” free buffer space by including **rwnd** value in TCP header of receiver-to-sender segments
 - **RcvBuffer** size set via socket options (typical default is 4096 bytes)
 - many operating systems autoadjust **RcvBuffer**
- sender limits amount of unacked (“in-flight”) data to receiver’s **rwnd** value
- guarantees receive buffer will not overflow



Chapter 3 outline

3.1 transport-layer services

3.2 multiplexing and demultiplexing

3.3 connectionless transport: UDP

3.4 principles of reliable data transfer

3.5 connection-oriented transport: TCP

- segment structure
- reliable data transfer
- flow control
- connection management

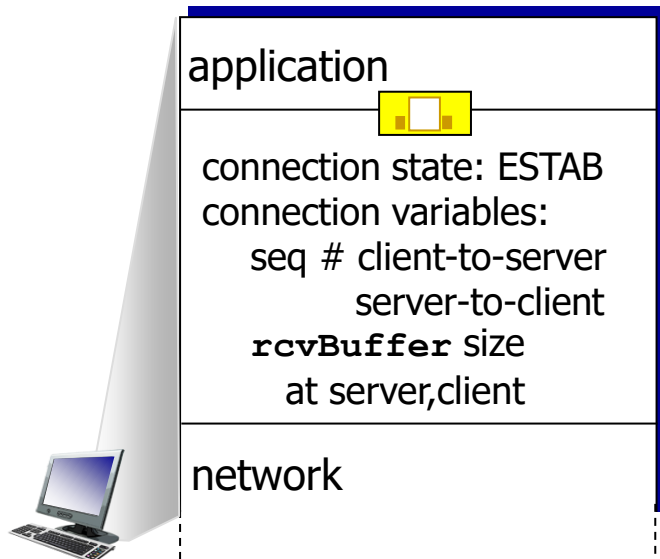
3.6 principles of congestion control

3.7 TCP congestion control

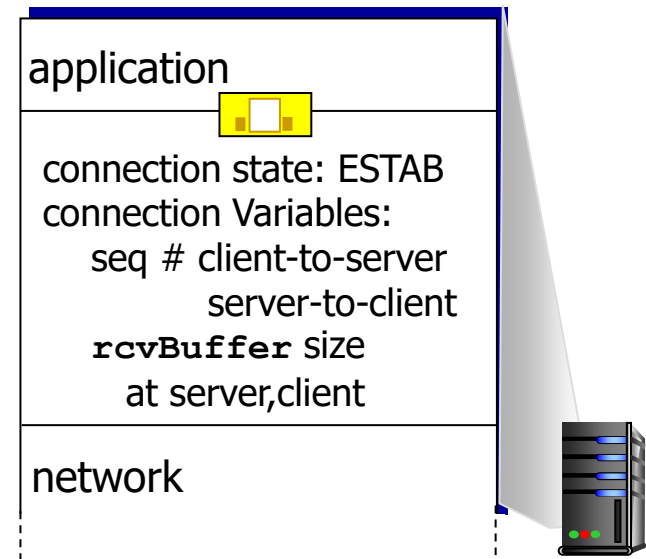
Connection Management

before exchanging data, sender/receiver “handshake”:

- agree to establish connection (each knowing the other willing to establish connection)
- agree on connection parameters



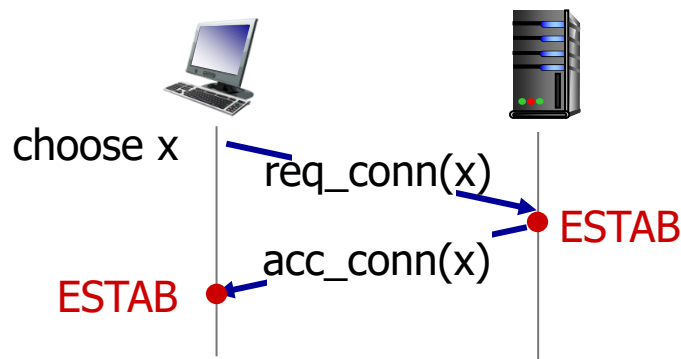
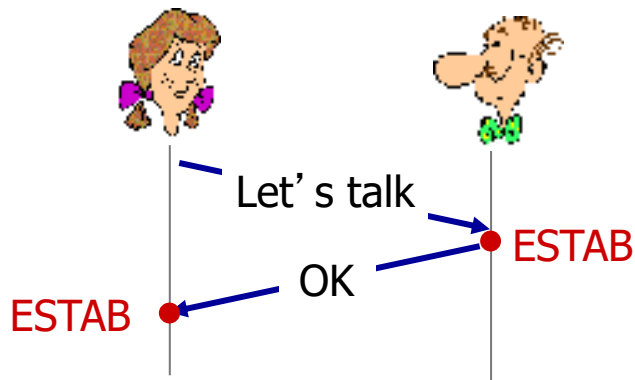
```
Socket clientSocket =  
    newSocket("hostname", "port number");  
clientSocket.connect((serverName, serverPort))
```



```
Socket connectionSocket =  
    welcomeSocket.accept();
```

Agreeing to establish a connection

2-way handshake:

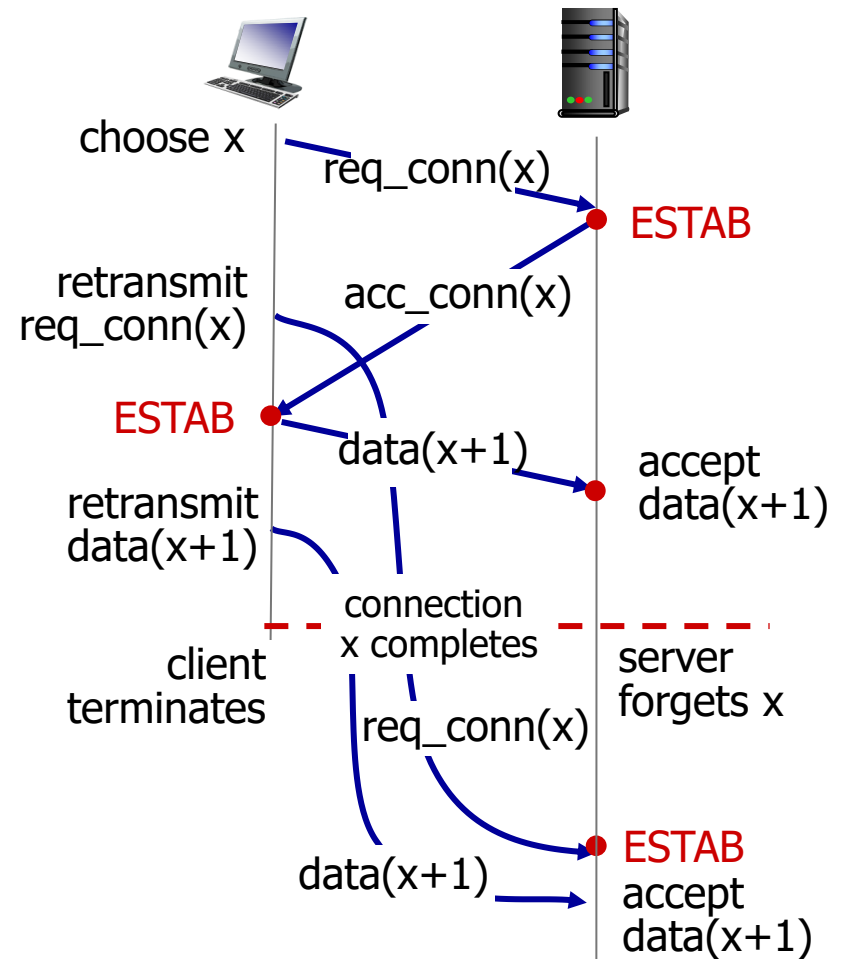
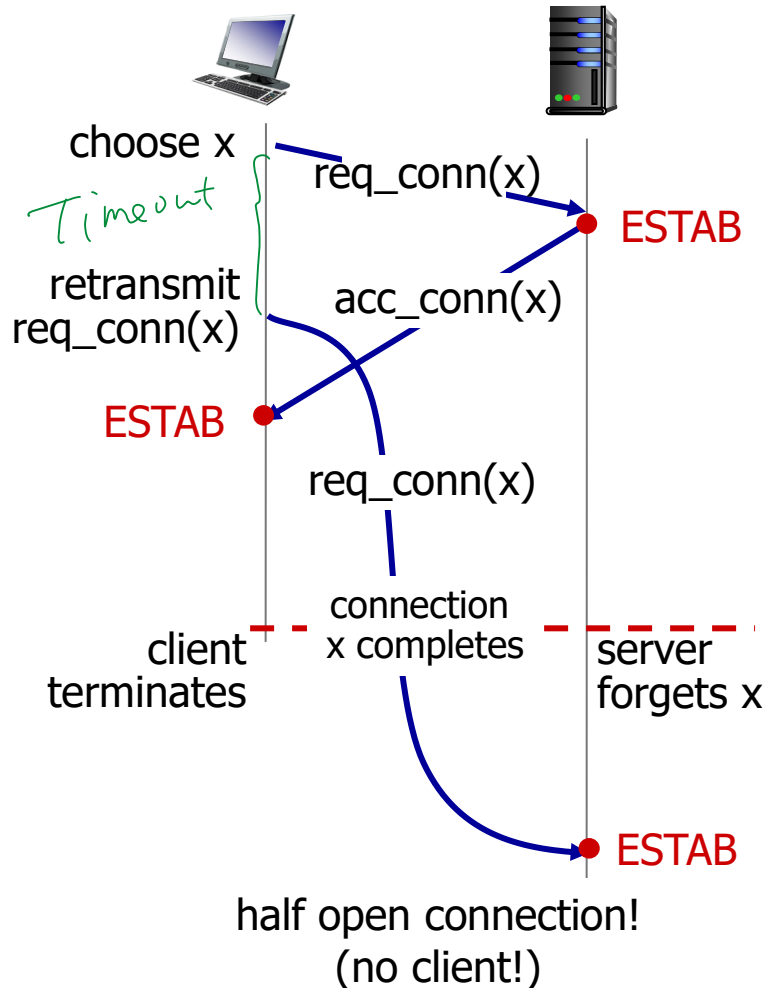


Q: will 2-way handshake always work in network?

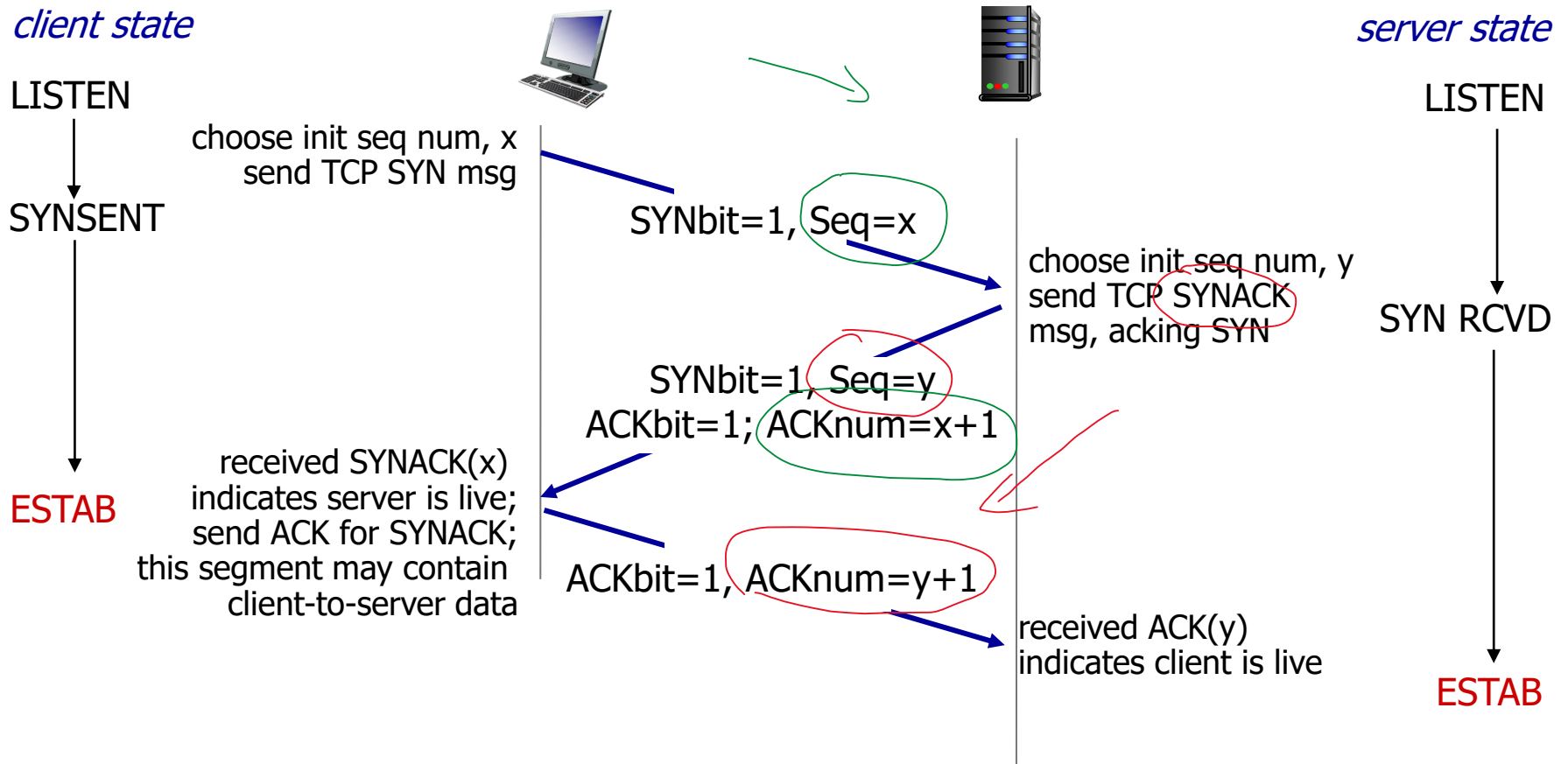
- variable delays
- retransmitted messages (e.g. req_conn(x)) due to message loss
- message reordering
- can't "see" other side

Agreeing to establish a connection

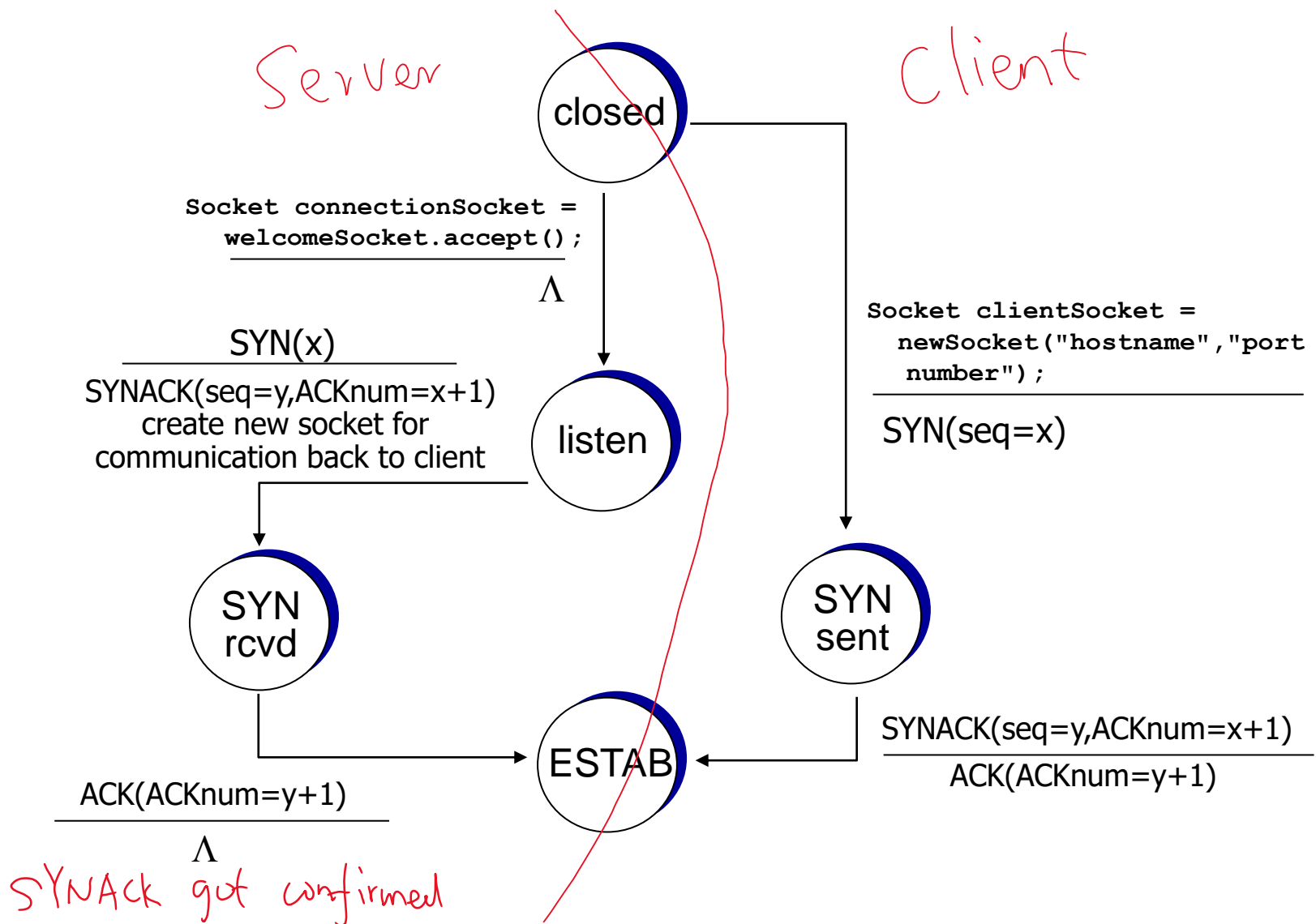
2-way handshake failure scenarios:



TCP 3-way handshake



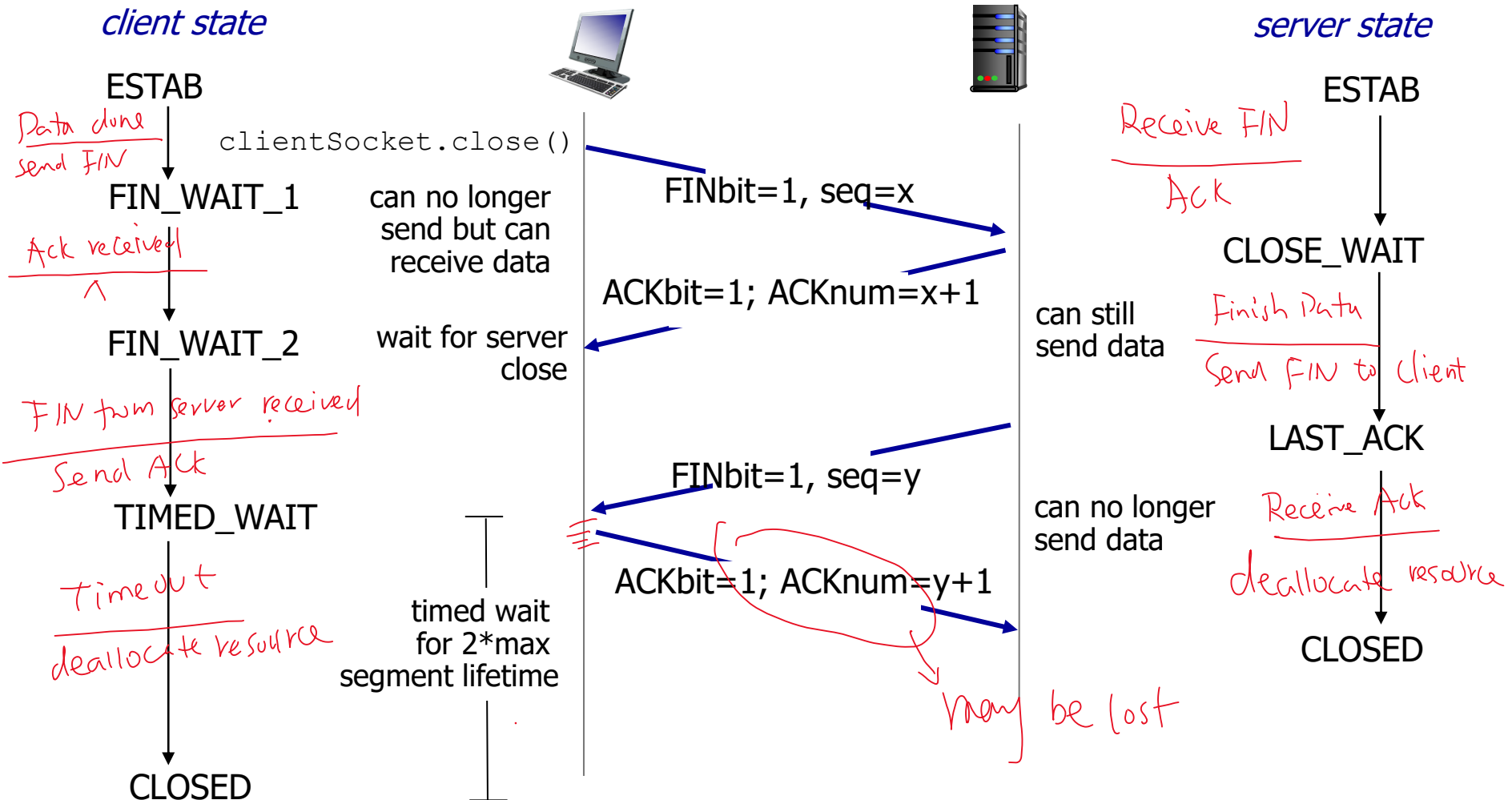
TCP 3-way handshake: FSM



TCP: closing a connection

- client, server each close their side of connection
 - send TCP segment with FIN bit = 1
- respond to received FIN with ACK
 - on receiving FIN, ACK can be combined with own FIN
- simultaneous FIN exchanges can be handled

TCP: closing a connection



Chapter 3 outline

3.1 transport-layer services

3.2 multiplexing and demultiplexing

3.3 connectionless transport: UDP

3.4 principles of reliable data transfer

3.5 connection-oriented transport: TCP

- segment structure
- reliable data transfer
- flow control
- connection management

3.6 principles of congestion control

3.7 TCP congestion control

Principles of congestion control

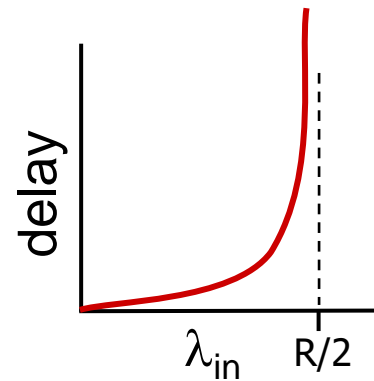
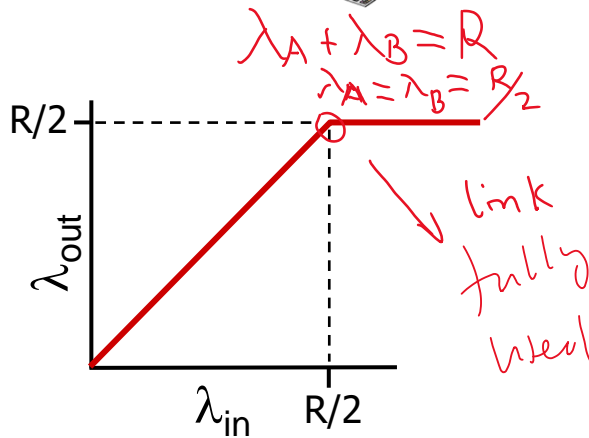
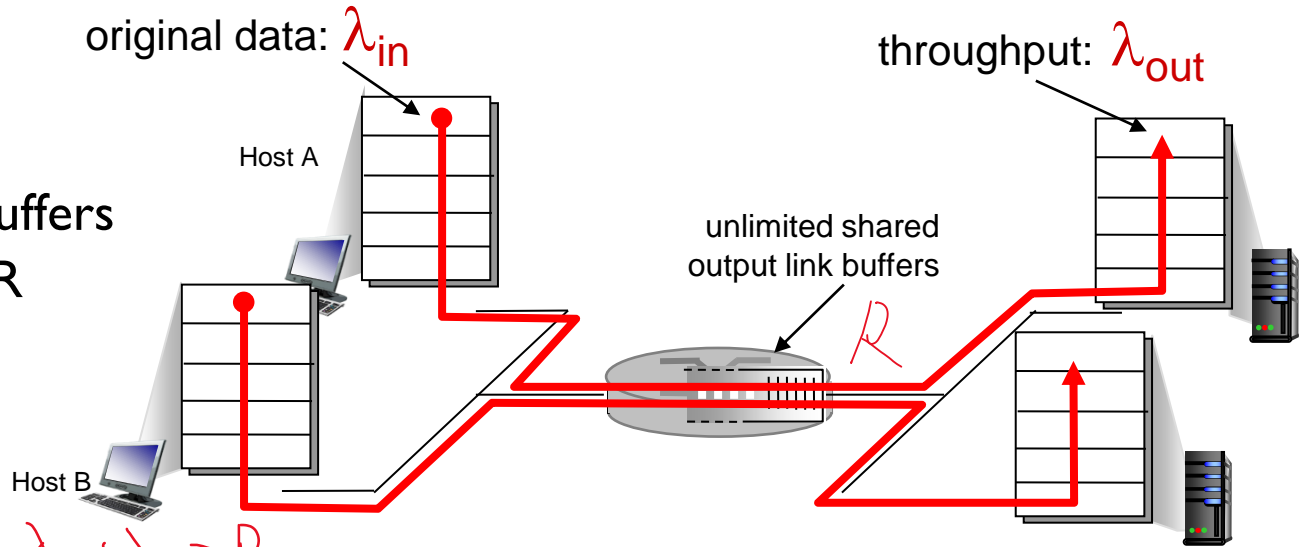
congestion:

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

Causes/costs of congestion: scenario I

- two senders, two receivers
- one router, infinite buffers
- output link capacity: R
- no retransmission

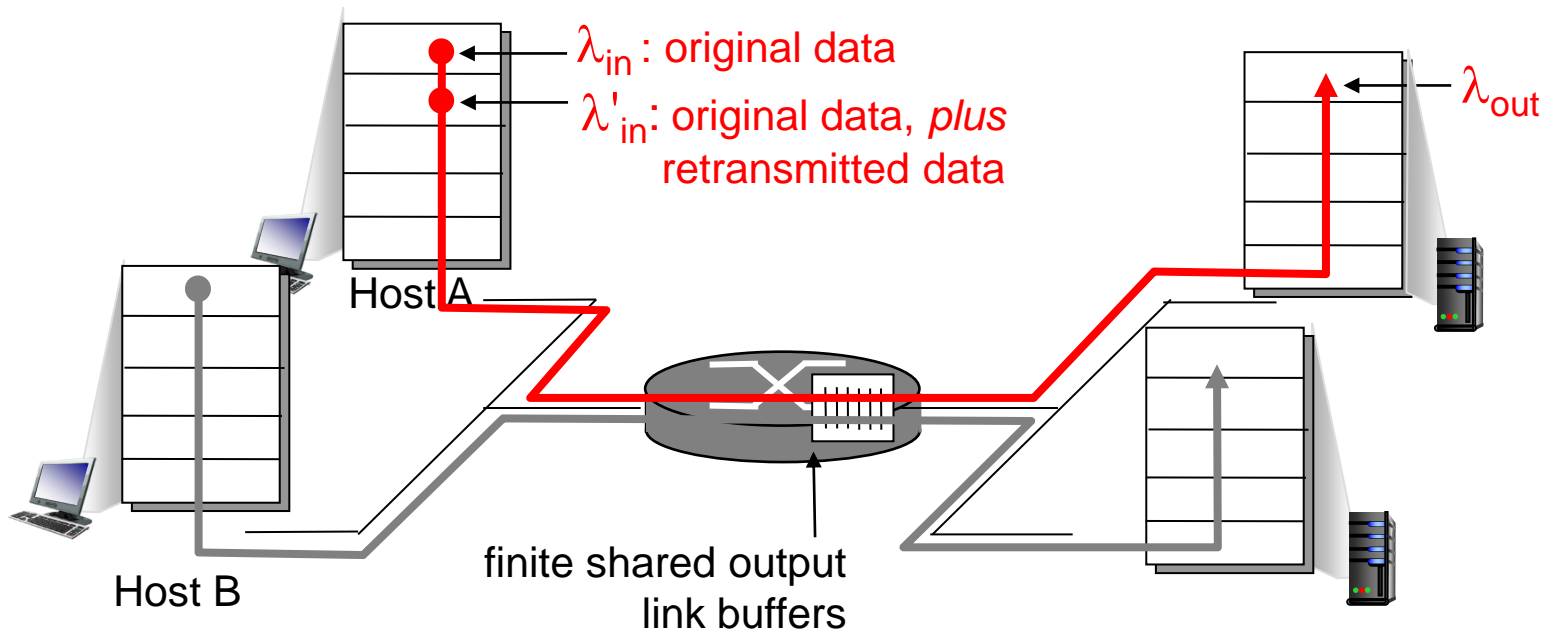
*Too much
traffic
cause
delay*



- maximum per-connection throughput: $R/2$
- ❖ large delays as arrival rate, λ_{in} , approaches capacity

Causes/costs of congestion: scenario 2

- one router, *finite* buffers
- sender retransmission of timed-out packet
 - application-layer input = application-layer output: $\lambda_{in} = \lambda_{out}$
 - transport-layer input includes *retransmissions* : $\lambda'_{in} \geq \lambda_{in}$

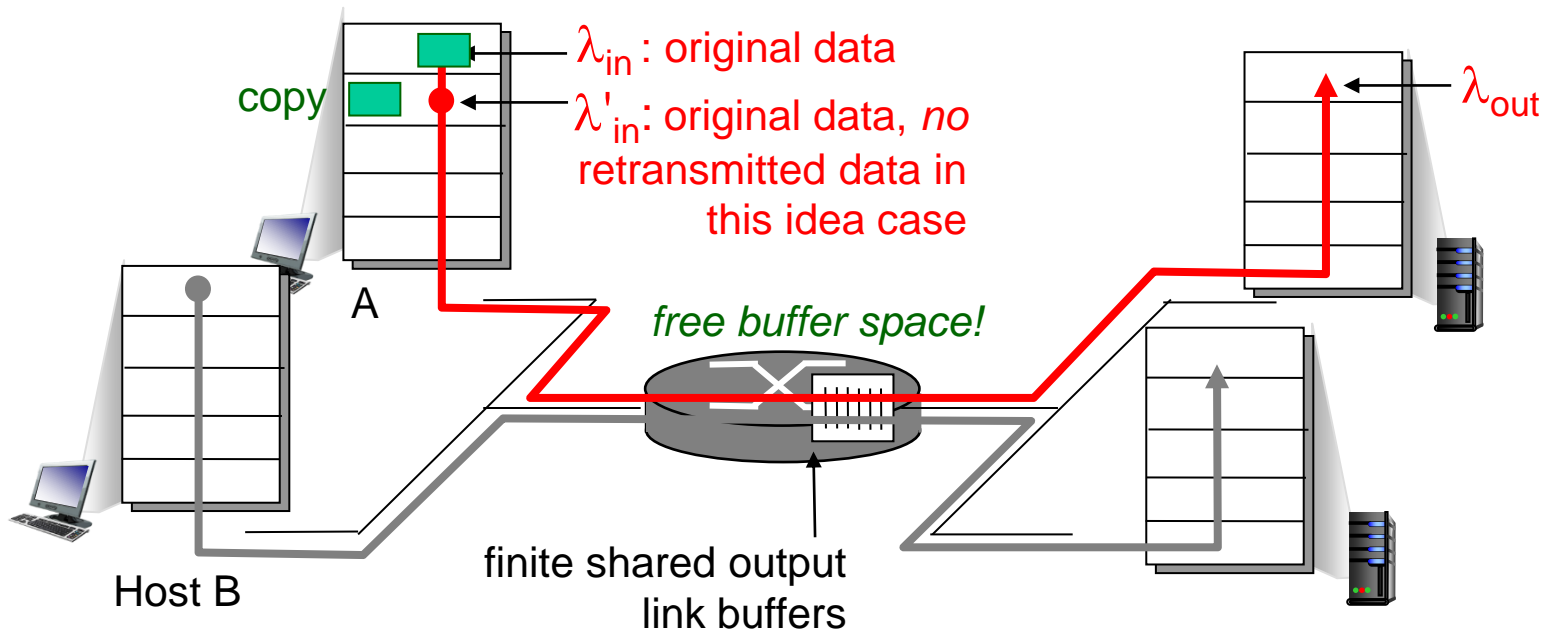
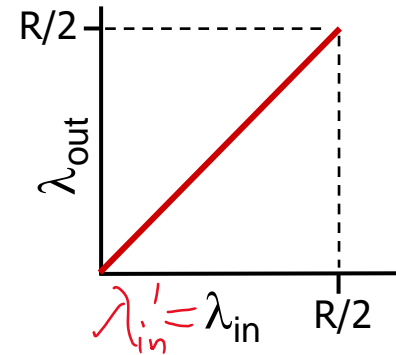


Causes/costs of congestion: scenario 2

idealization: perfect knowledge

- sender sends only when router buffers available

no loss, ignore timeout



Causes/costs of congestion: scenario 2

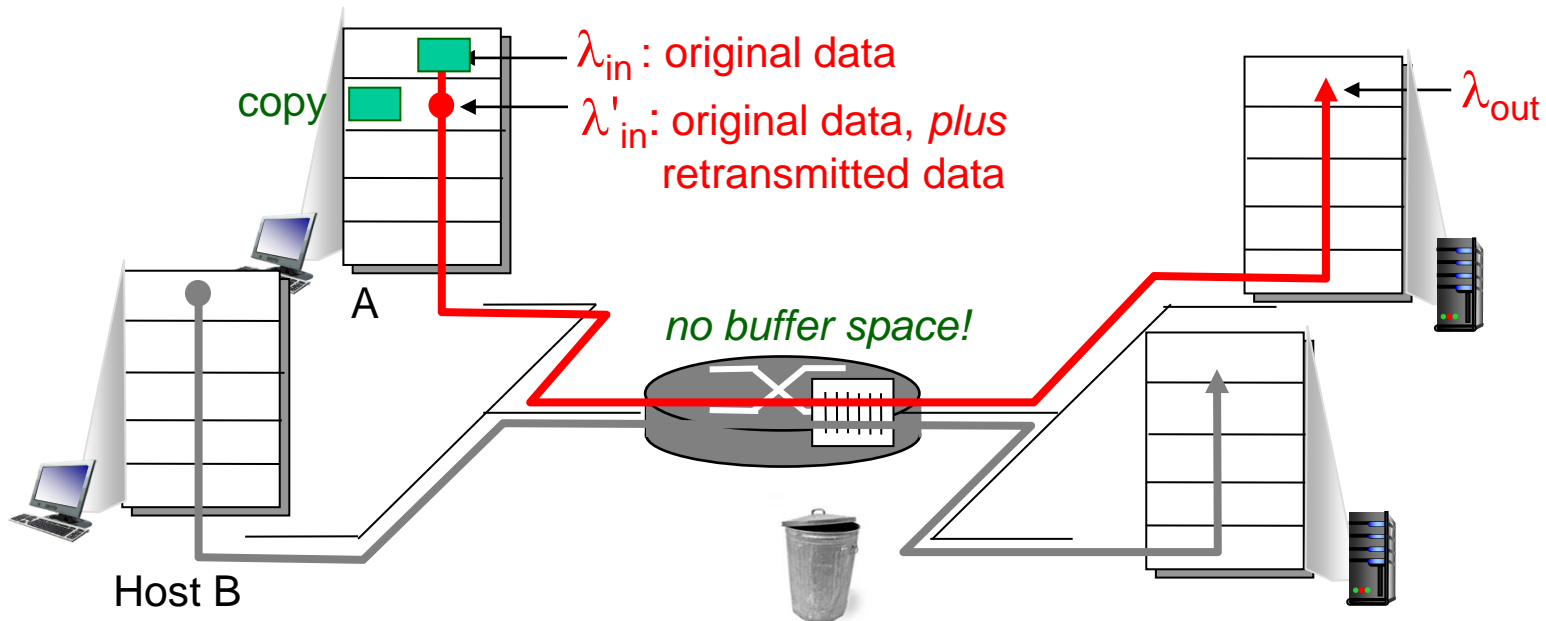
Idealization: known loss

packets can be lost,
dropped at router due
to full buffers

- sender only resends if
packet *known* to be lost

$$\lambda_{in}' = \lambda_{in} + \lambda_{retransmission} > \lambda_{in}$$

ignore timeout

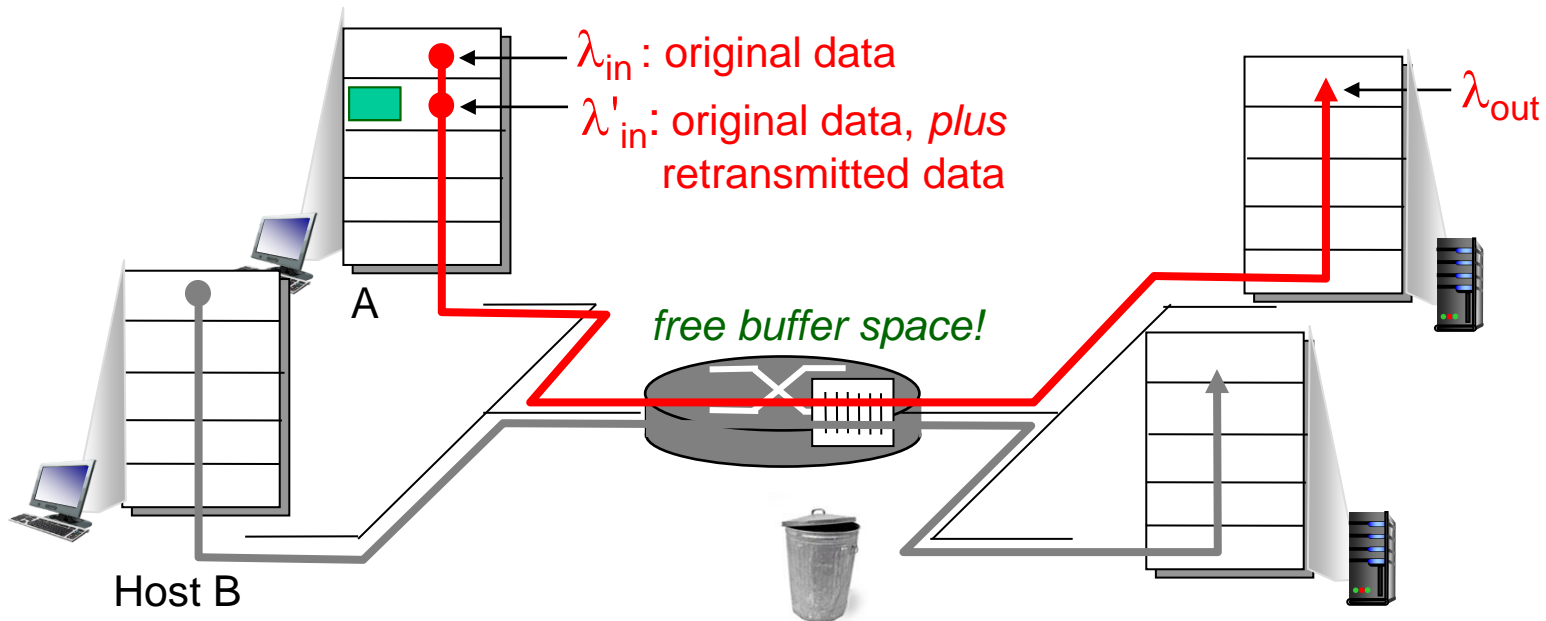
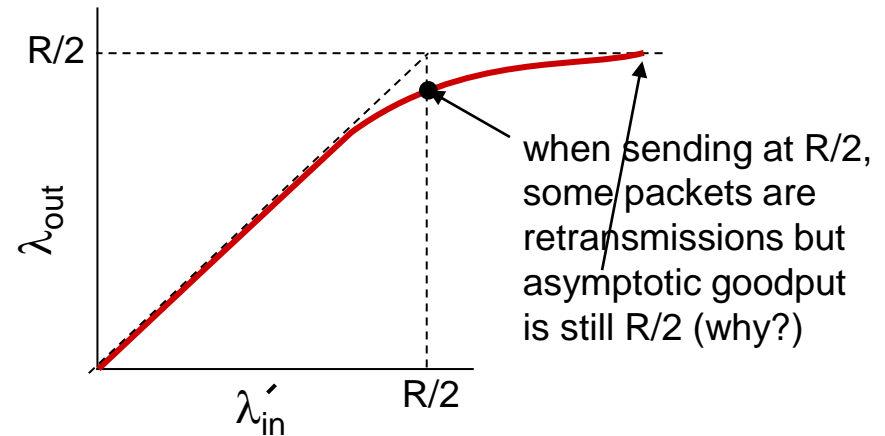


Causes/costs of congestion: scenario 2

Idealization: known loss

packets can be lost,
dropped at router due
to full buffers

- sender only resends if
packet *known* to be lost



Realistic: *duplicates*

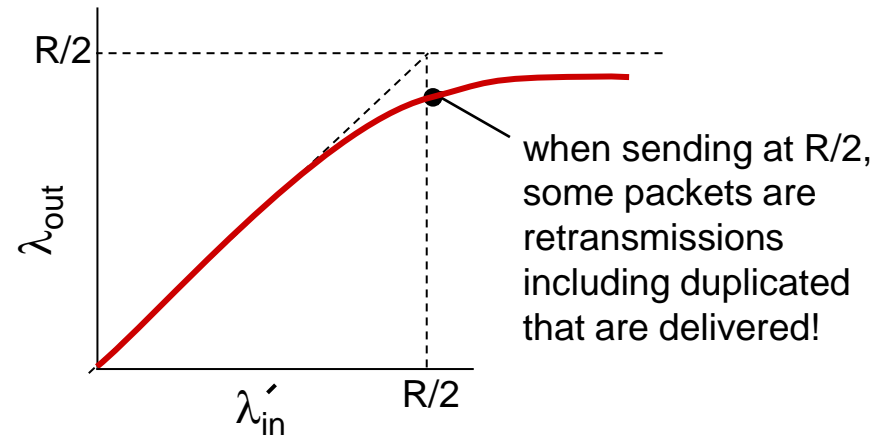
-
- A graph showing the relationship between the input rate λ'_{in} (x-axis) and the output rate λ_{out} (y-axis). The curve starts at the origin and increases, eventually saturating at a value of $R/2$ on the y-axis. A dashed line indicates the point where the input rate is $R/2$ and the output rate is $R/2$. A text box points to this point, stating: "when sending at $R/2$, some packets are retransmissions including duplicated that are delivered!".



Causes/costs of congestion: scenario 2

Realistic: duplicates

- packets can be lost, dropped at router due to full buffers
- sender times out prematurely, sending *two* copies, both of which are delivered



“costs” of congestion:

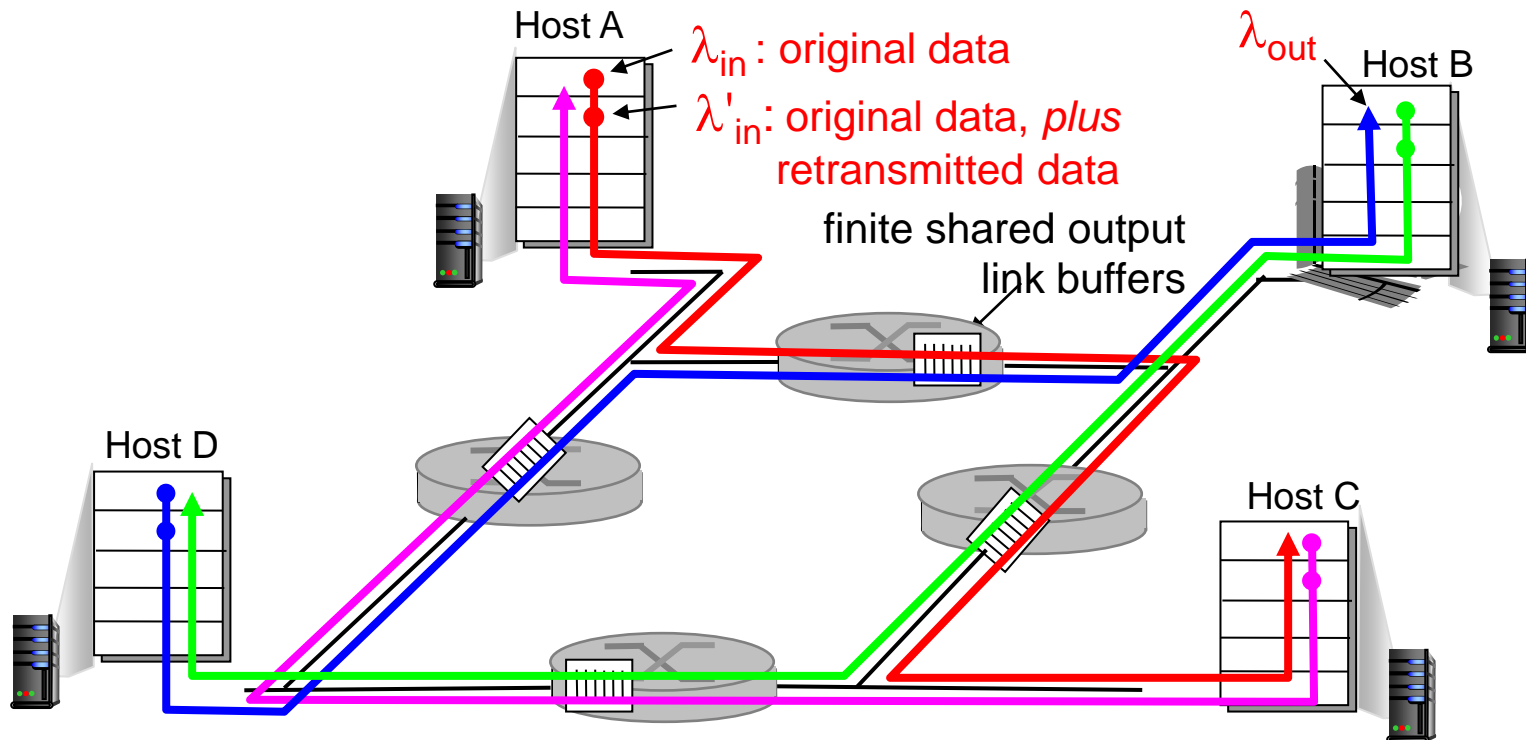
- more work (retrans) for given “goodput”
- unneeded retransmissions: link carries multiple copies of pkt
 - decreasing goodput

Causes/costs of congestion: scenario 3

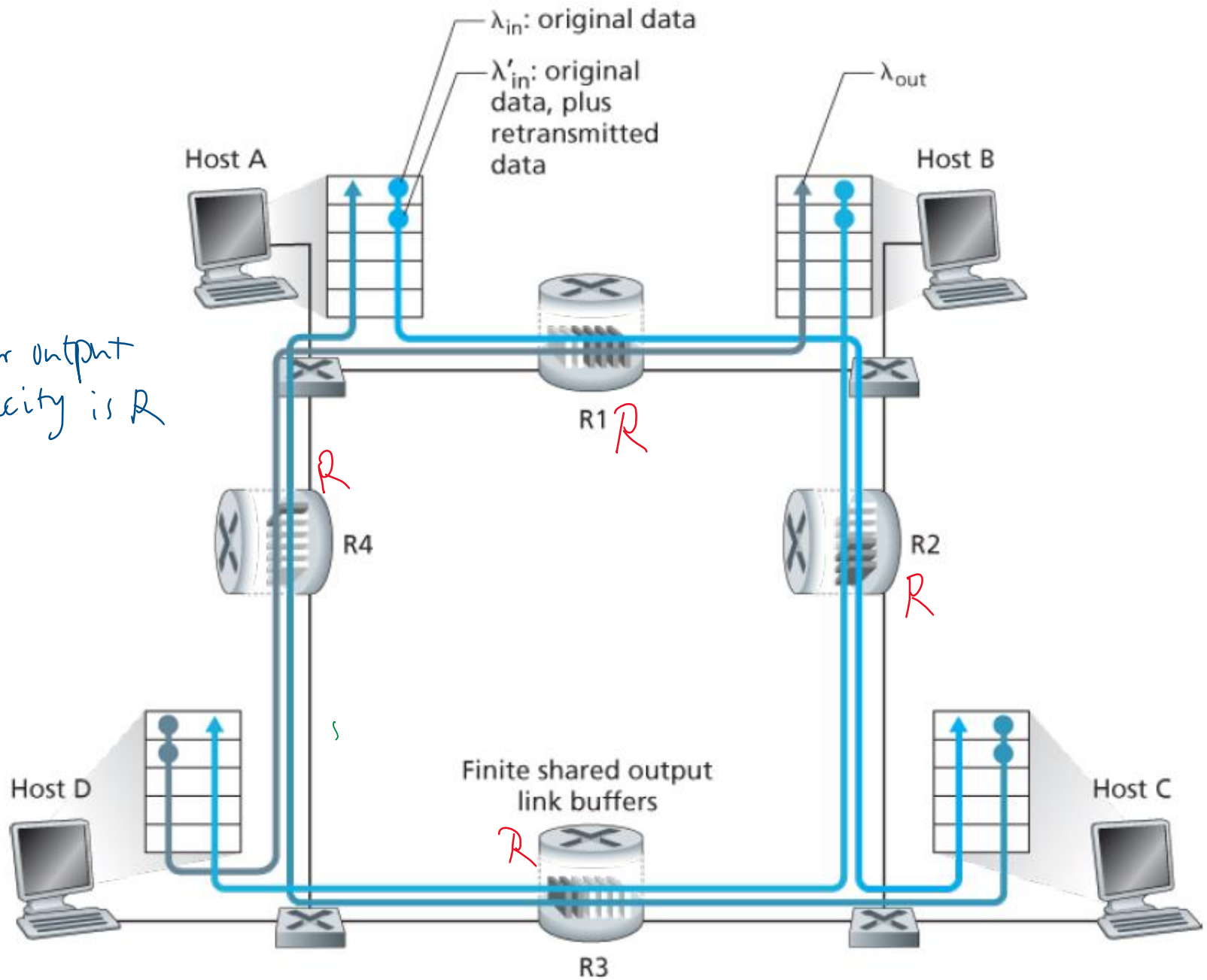
- four senders
- multihop paths
- timeout/retransmit

Q: what happens as λ_{in} and λ'_{in} increase ?

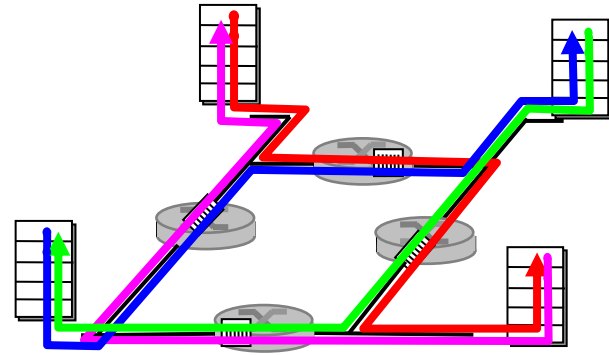
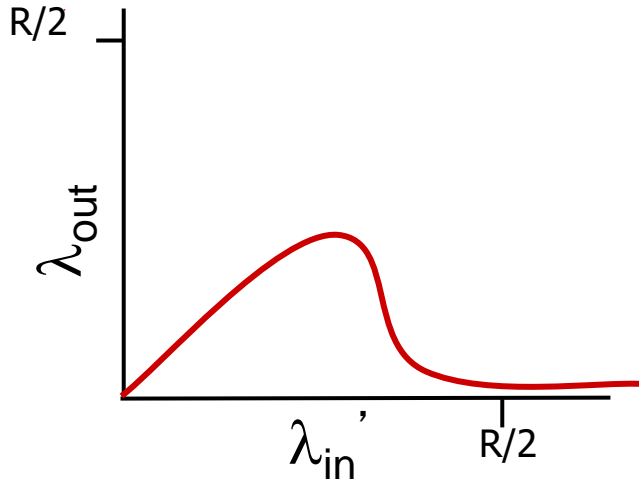
A: as red λ'_{in} increases, all arriving blue pkts at upper queue are dropped, blue throughput $\rightarrow 0$



Router output capacity is R



Causes/costs of congestion: scenario 3



another “cost” of congestion:

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