

Chapter 3

Transport Layer

A note on the use of these ppt slides:

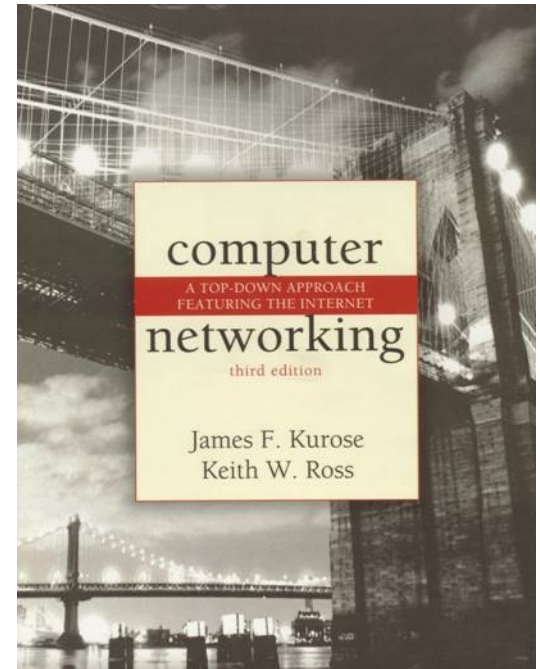
We're making these slides freely available to all (faculty, students, readers). They're in powerpoint form so you can add, modify, and delete slides (including this one) and slide content to suit your needs. They obviously represent a *lot* of work on our part. In return for use, we only ask the following:

- ❑ If you use these slides (e.g., in a class) in substantially unaltered form, that you mention their source (after all, we'd like people to use our book!)
- ❑ If you post any slides in substantially unaltered form on a www site, that you note that they are adapted from (or perhaps identical to) our slides, and note our copyright of this material.

Thanks and enjoy! JFK/KWR

All material copyright 1996-2004

J.F Kurose and K.W. Ross, All Rights Reserved



*Computer Networking:
A Top Down Approach
Featuring the Internet,*

3rd edition.

*Jim Kurose, Keith Ross
Addison-Wesley, July
2004.*

Chapter 3: Transport Layer

Our goals:

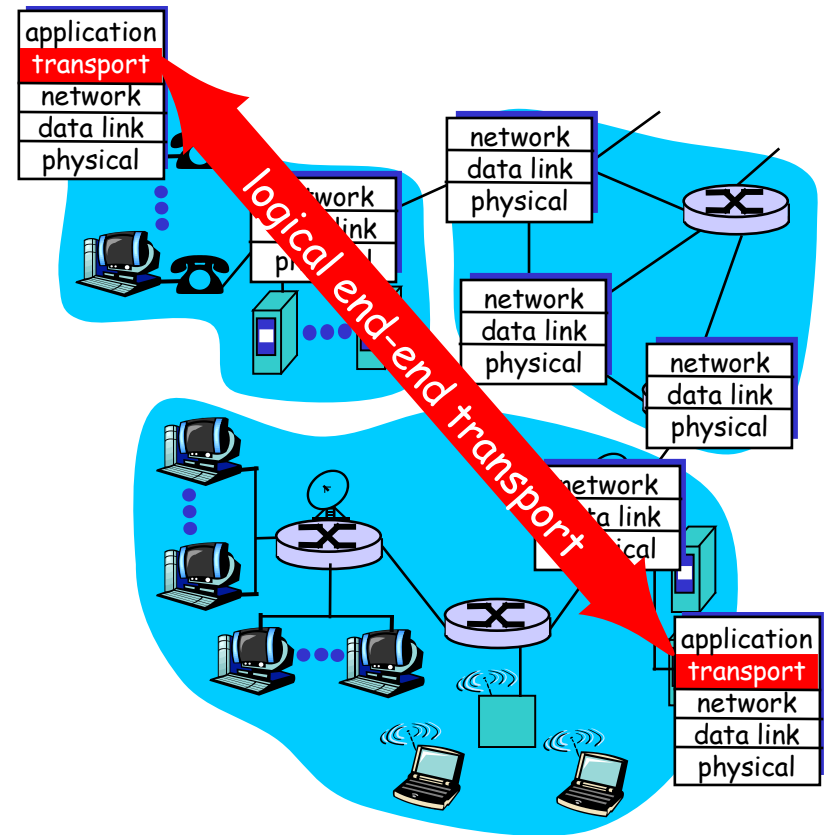
- understand principles behind transport layer services:
 - multiplexing/demultiplexing
 - reliable data transfer
 - flow control
 - congestion control
- learn about transport layer protocols in the Internet:
 - UDP: connectionless transport
 - TCP: connection-oriented transport
 - TCP congestion control

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

Transport services and protocols

- provide *logical communication* between app processes running on different hosts
- transport protocols run in end systems
 - send side: breaks app messages into *segments*, passes to network layer
 - rcv side: reassembles segments into messages, passes to app layer
- more than one transport protocol available to apps
 - Internet: TCP and UDP



Transport vs. network layer

- *network layer*: logical communication between **hosts**
- *transport layer*: logical communication between **processes**
 - relies on, enhances, network layer services

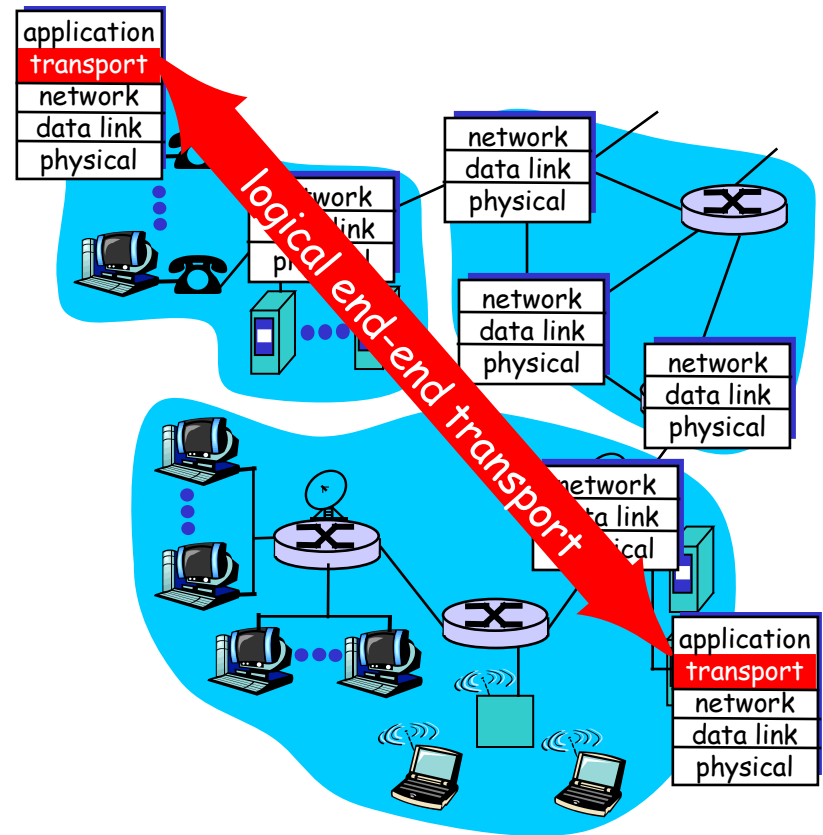
Household analogy:

12 kids sending letters to 12 kids

- processes = kids
- app messages = letters in envelopes
- hosts = houses
- transport protocol = Ann and Bill
- network-layer protocol = postal service

Internet transport-layer protocols

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



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

Multiplexing/demultiplexing

Demultiplexing at rcv host:

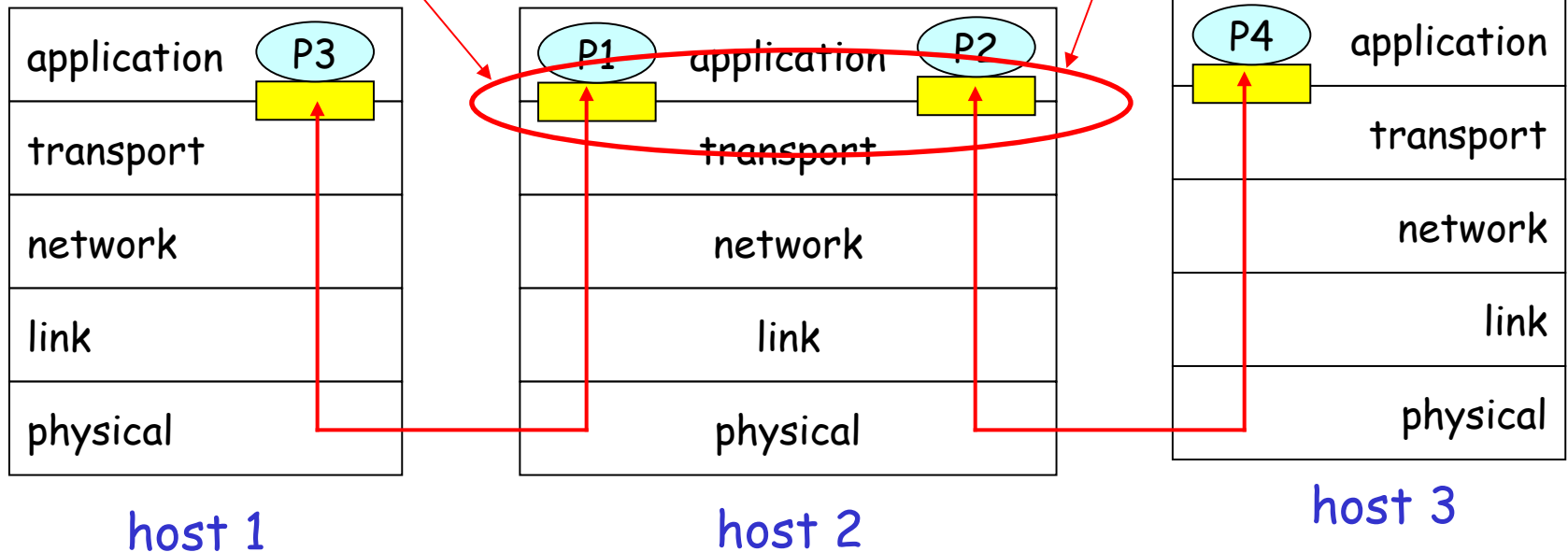
delivering received segments
to correct socket

Multiplexing at send host:

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

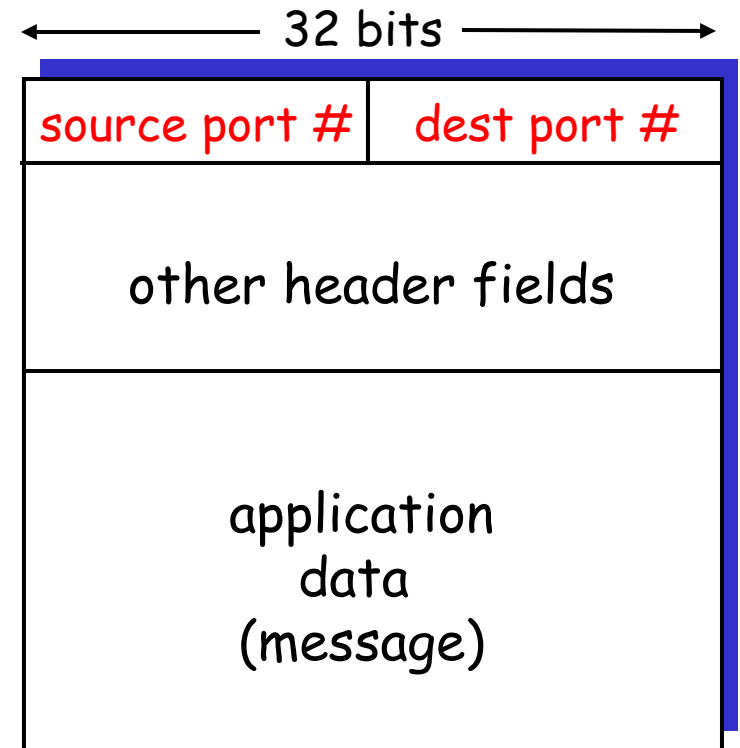
■ = socket

○ = process



How demultiplexing works

- host receives IP datagrams
 - each datagram has **source IP address, destination IP address**
 - each datagram carries 1 transport-layer segment
 - each segment has **source, destination port number** (recall: well-known port numbers for specific applications)
- host uses **IP addresses & port numbers** to direct segment to appropriate **socket**



TCP/UDP segment format

Connectionless (UDP) demultiplexing

- Create **UDP** sockets with port numbers:

```
DatagramSocket mySocket1 = new  
    DatagramSocket(99111);
```

```
DatagramSocket mySocket2 = new  
    DatagramSocket(99222);
```

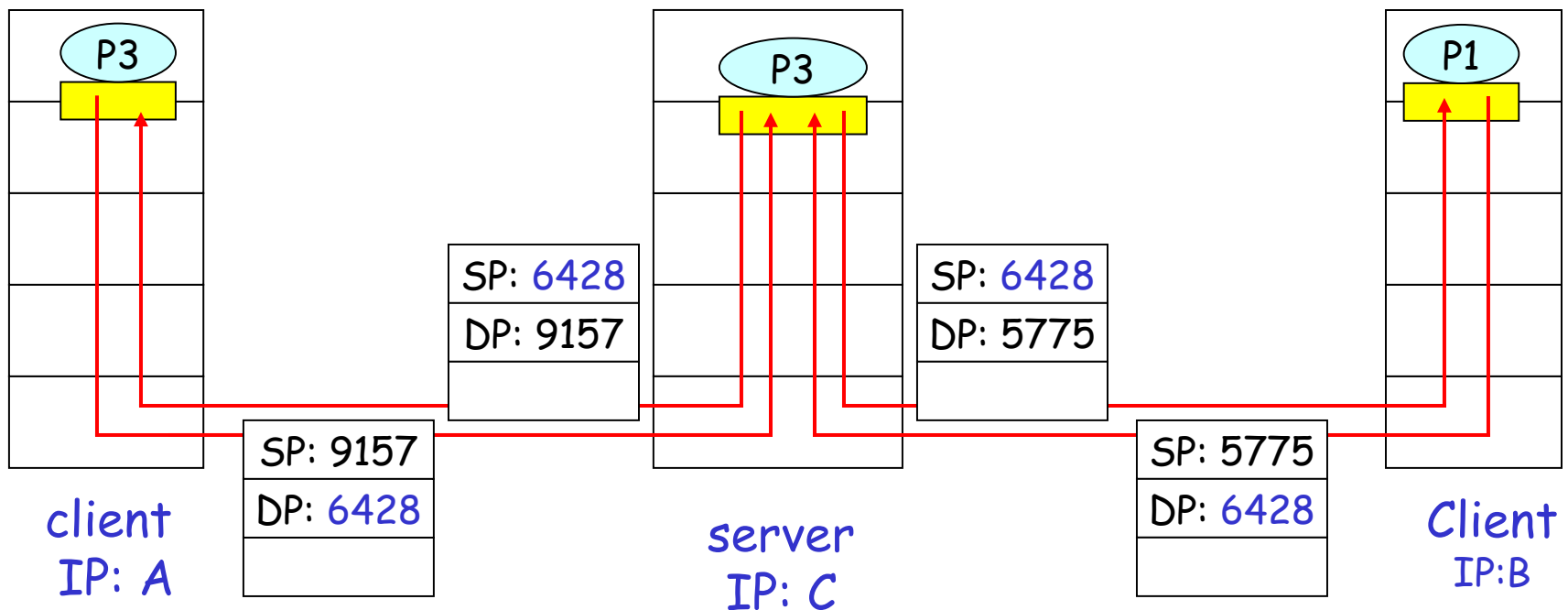
- UDP socket identified by two-tuple:

(dest IP address, dest port number)

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

Connectionless demux (cont)

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

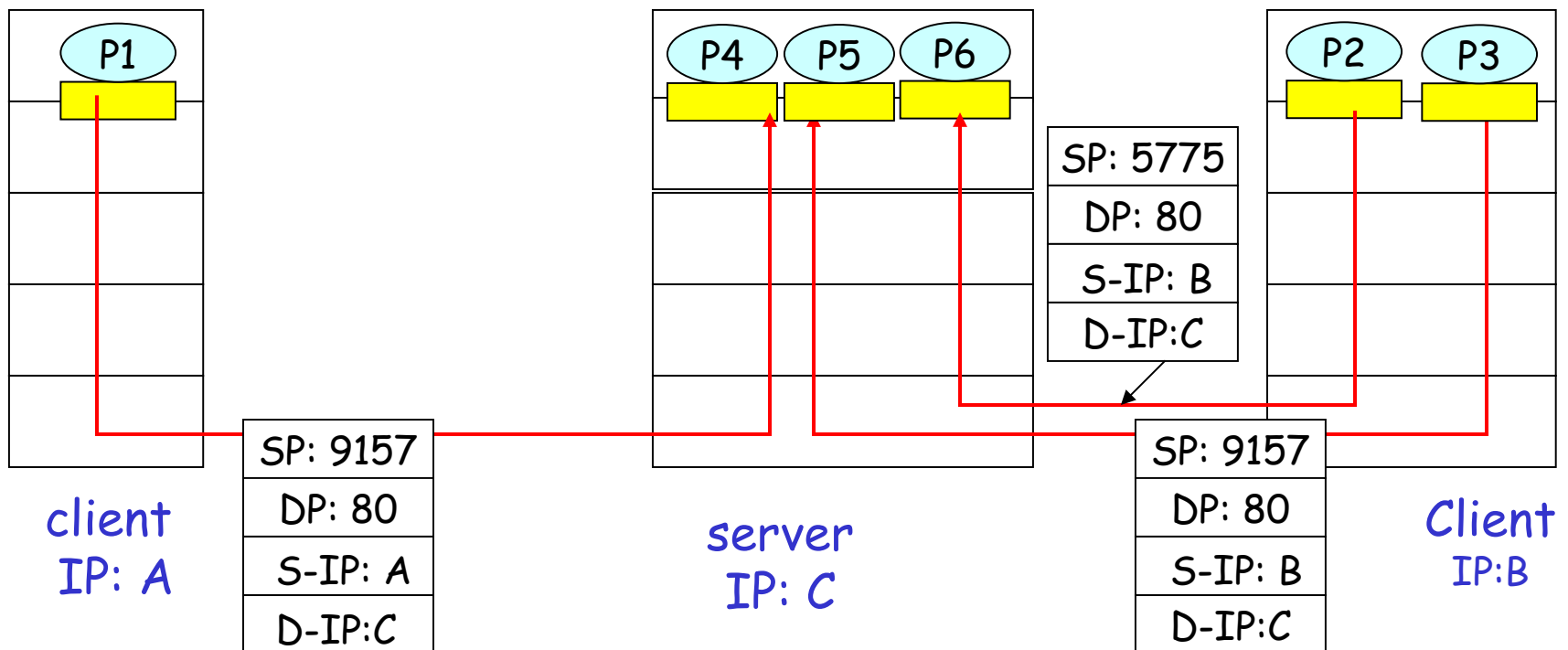


SP (source port number) provides "return address"
Complete return address: (source IP address, source port number)

Connection-oriented (TCP) demux

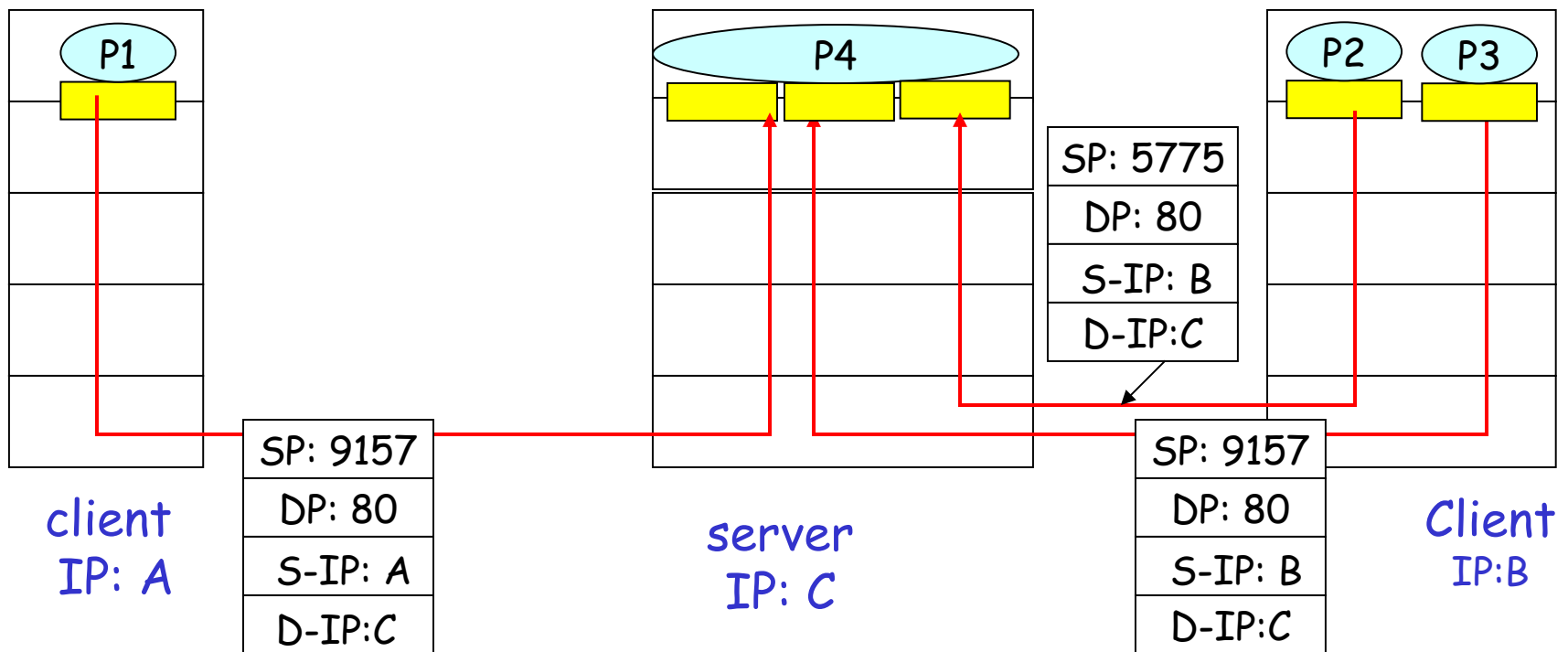
- TCP socket identified by 4-tuple:
 - source IP address
 - source port number
 - dest IP address
 - dest port number
- recv host uses all four values to direct segment to appropriate socket
- Server host may support many simultaneous TCP sockets:
 - each socket identified by its own 4-tuple
- Web servers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

Connection-oriented demux (cont)



The server spawn one process for each new client connection.

Connection-oriented demux: Threaded Web Server



The server create a new thread for each new client 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

UDP: User Datagram Protocol [RFC 768]

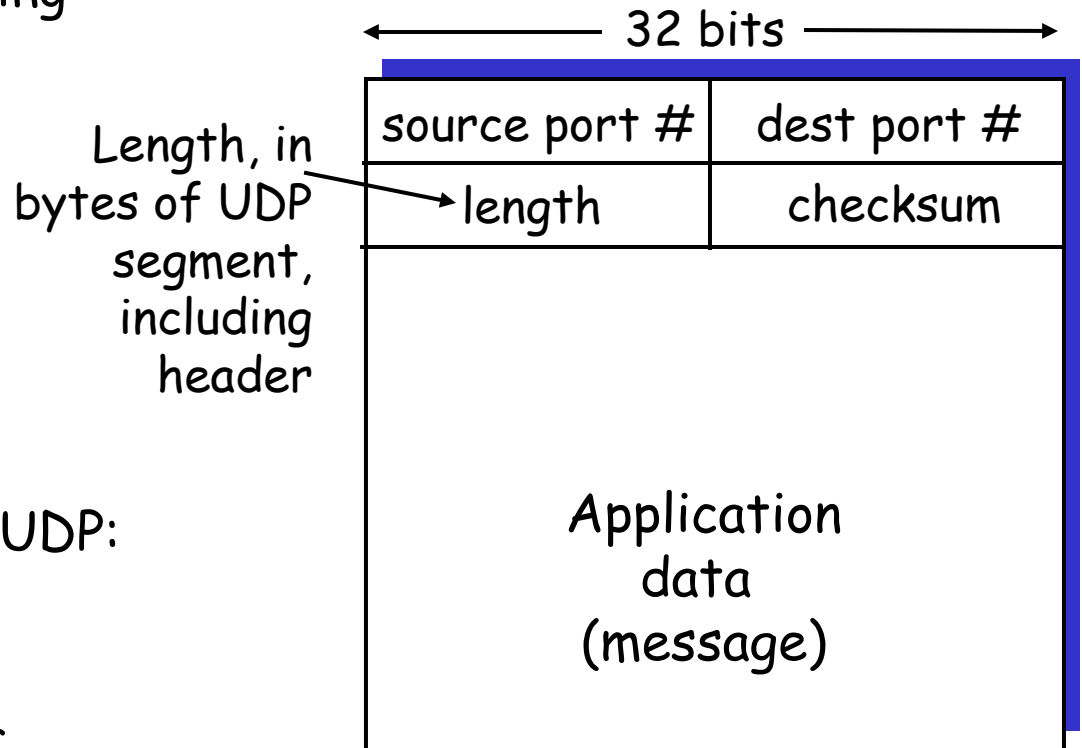
- “no frills,” “bare bones” Internet transport protocol
- “best effort” service, UDP segments may be:
 - lost
 - delivered out of order to app
- *connectionless*:
 - no handshaking between UDP sender, receiver
 - each UDP segment handled independently of others

Why is there a UDP?

- no connection establishment (which can add delay)
- simple: no connection state at sender, receiver
- small segment header
- no congestion control:
 - UDP can blast away as fast as desired
 - Applications have Better control over what and when to sent

UDP: more

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



UDP segment format

UDP checksum

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

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? More later

Internet Checksum Example

- Note
 - When adding numbers, a carryout from the most significant bit needs to be added to the result
- Example: add two 16-bit integers

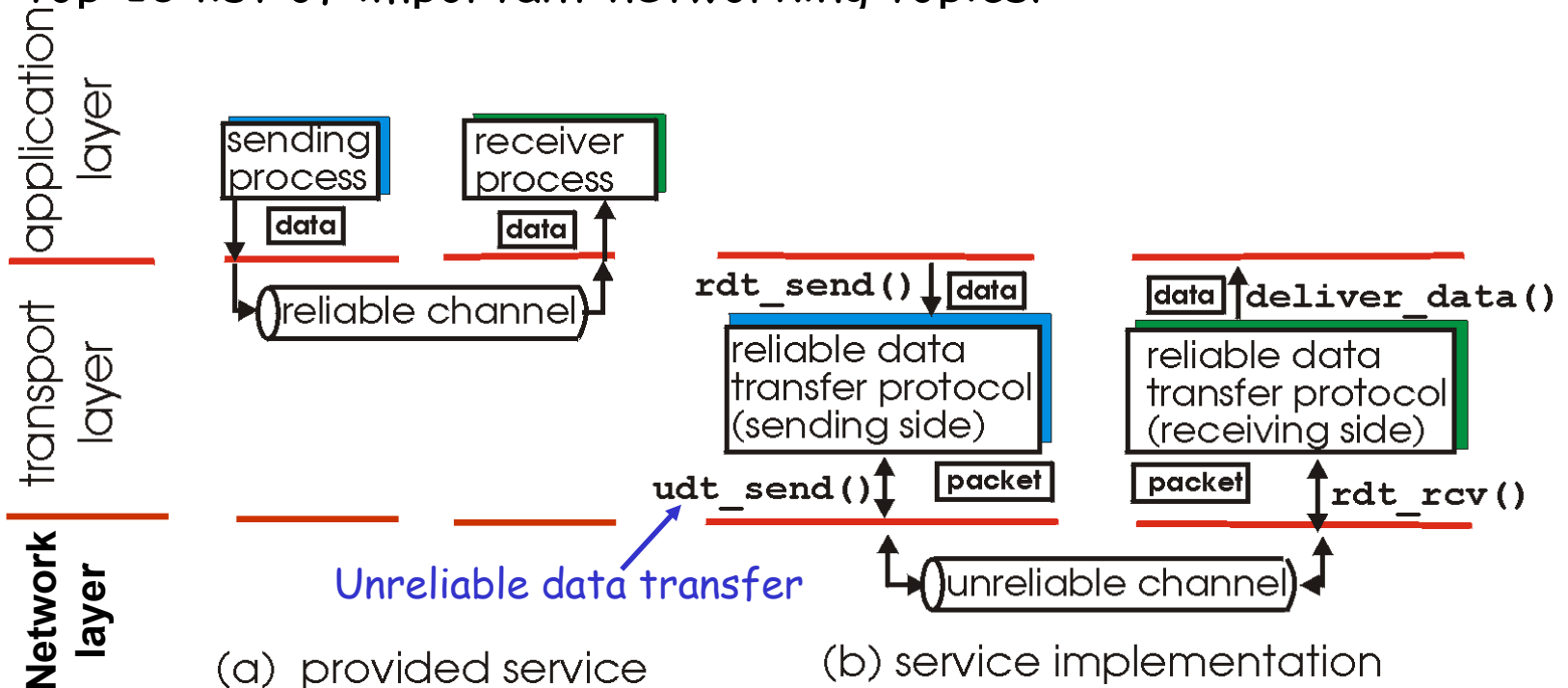
	1	1	1	0	0	1	1	0	0	1	1	0	0	1	1	0
	1	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1
<hr/>																
wraparound	1	1	0	1	1	1	0	1	1	1	0	1	1	1	0	1
<hr/>																
sum	1	0	1	1	1	0	1	1	1	0	1	1	1	1	0	0
checksum	0	1	0	0	0	1	0	0	0	1	0	0	0	0	1	1

Chapter 3 outline

- 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 Reliable data transfer

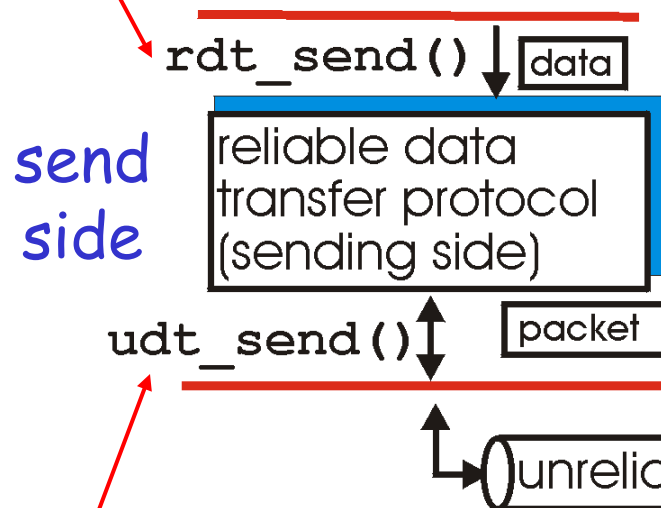
- important in implementing reliable data transfer in application, transport, link layers
- top-10 list of important networking topics!



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

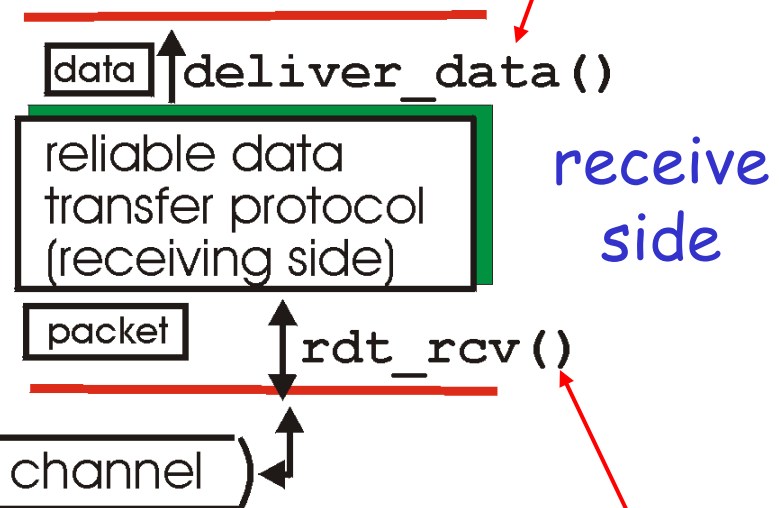
Reliable data transfer: getting started

rdt_send() : called from above, (e.g., by app.). Pass data to be delivered to receiver upper layer



udt_send() : called by rdt, to transfer packet over unreliable channel to receiver

deliver_data() : called by rdt to deliver data to upper layer



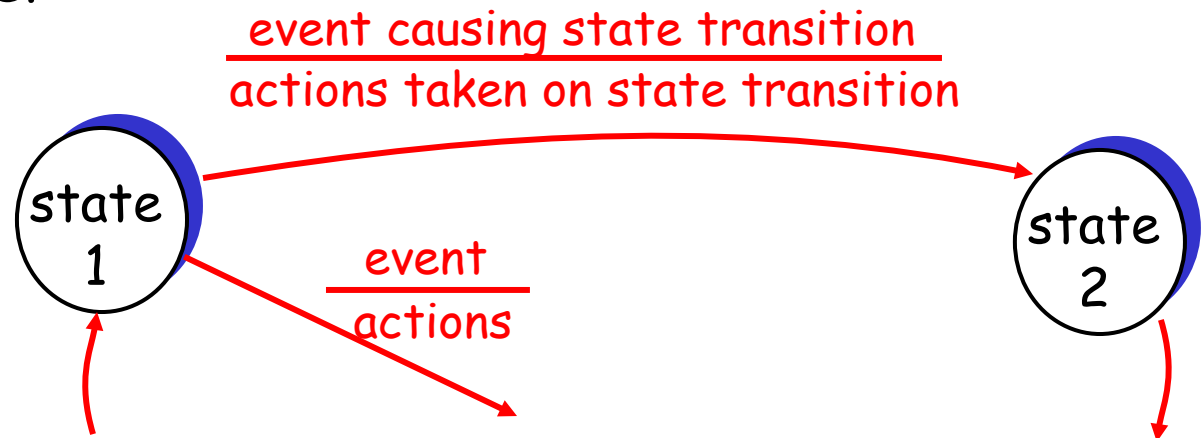
rdt_rcv() : called when packet arrives on rcv-side of channel

Reliable data transfer: getting started

We'll:

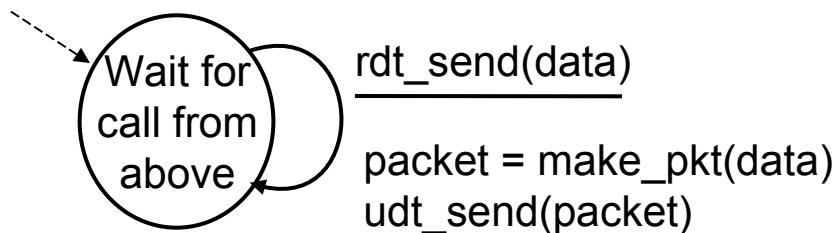
- incrementally develop sender, receiver sides of reliable data transfer protocol (rdt)
- consider only unidirectional data transfer
 - but control information will flow on both directions!
- use finite state machines (FSM) to specify sender, receiver

state: when in this "state" next state uniquely determined by next event

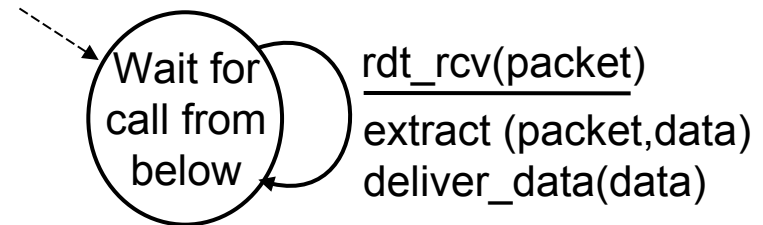


Rdt1.0: reliable transfer over a reliable channel

- Assume that the underlying channel is perfectly reliable
 - no bit errors
 - no loss of packets
- separate FSMs for sender, receiver:
 - sender sends data into underlying channel
 - receiver read data from underlying channel
 - No need for the receiver side to provide any feedback to the sender since the channel is perfectly reliable



sender

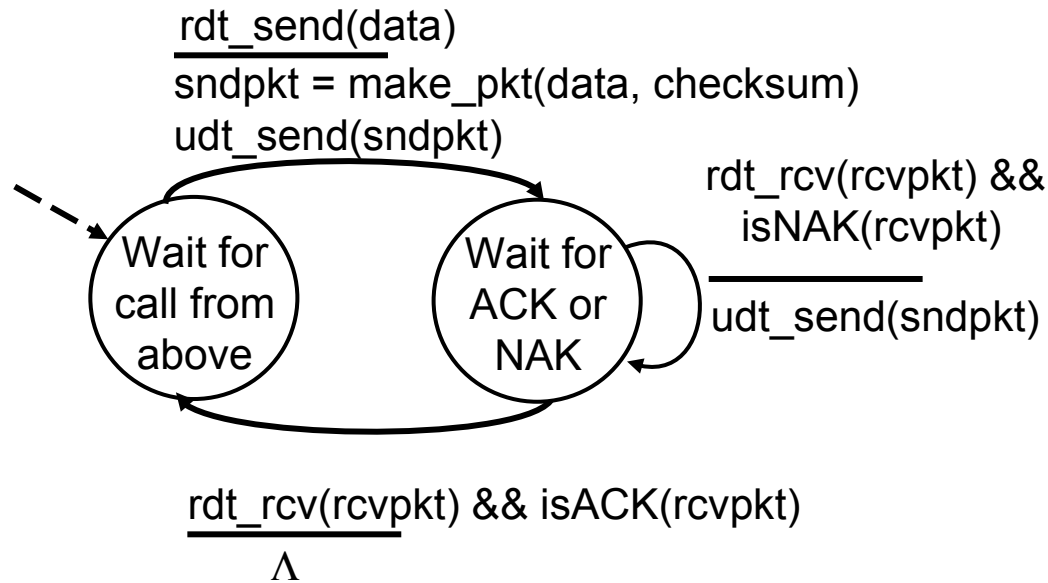


receiver

Rdt2.0: assume channel with bit errors

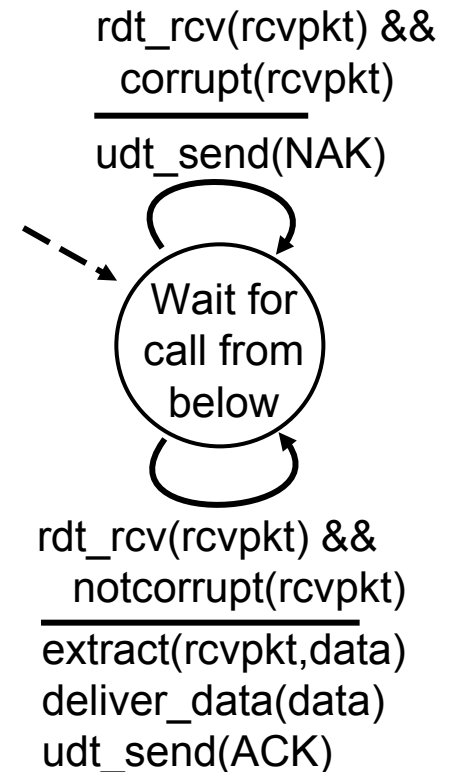
- underlying channel may flip bits in packet
 - recall: UDP checksum to detect bit errors
- *the question: how to recover from errors:*
 - *acknowledgements (ACKs):* receiver explicitly tells sender that packet received OK
 - *negative acknowledgements (NAKs):* receiver explicitly tells sender that packet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

rdt2.0: FSM specification

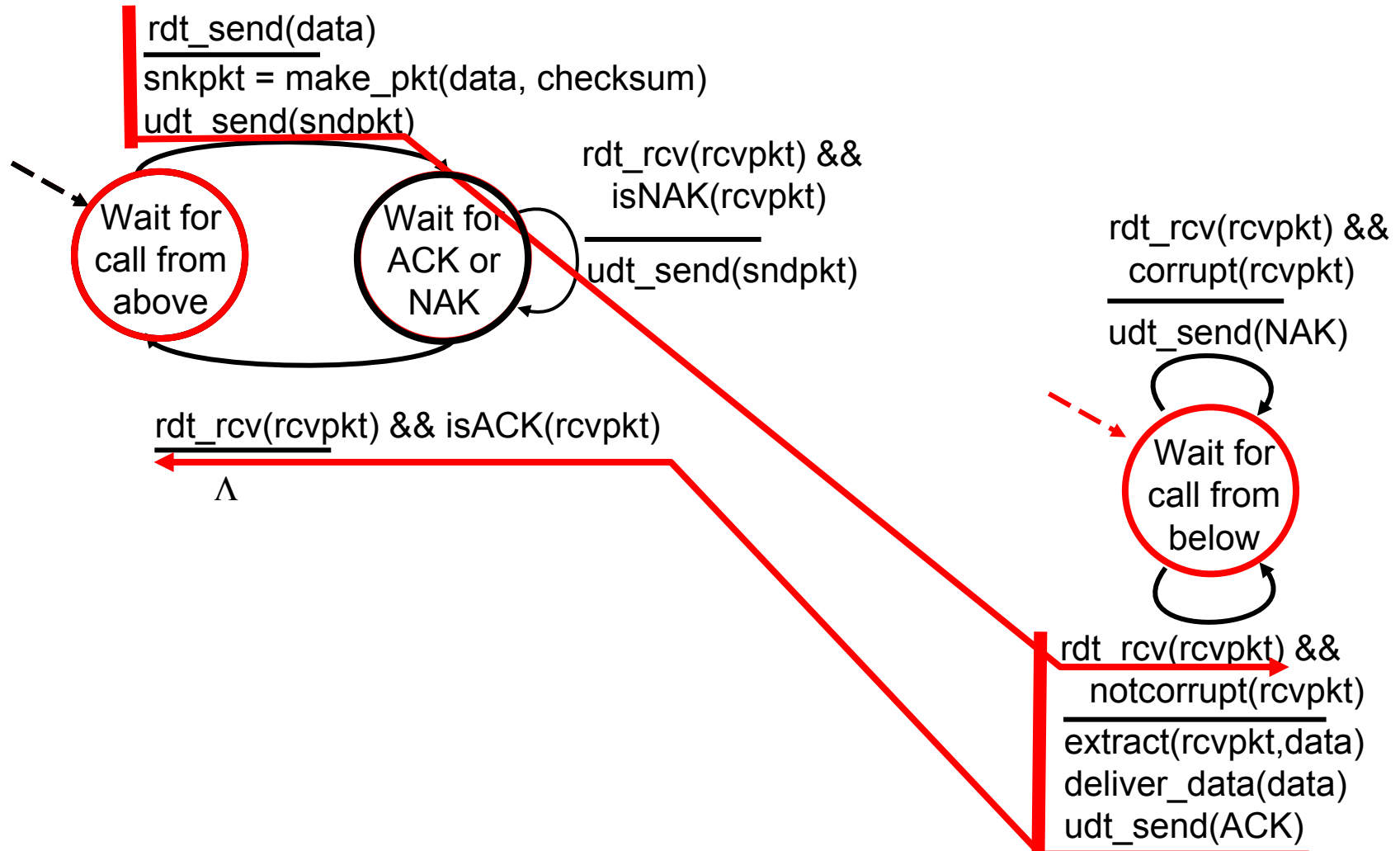


sender

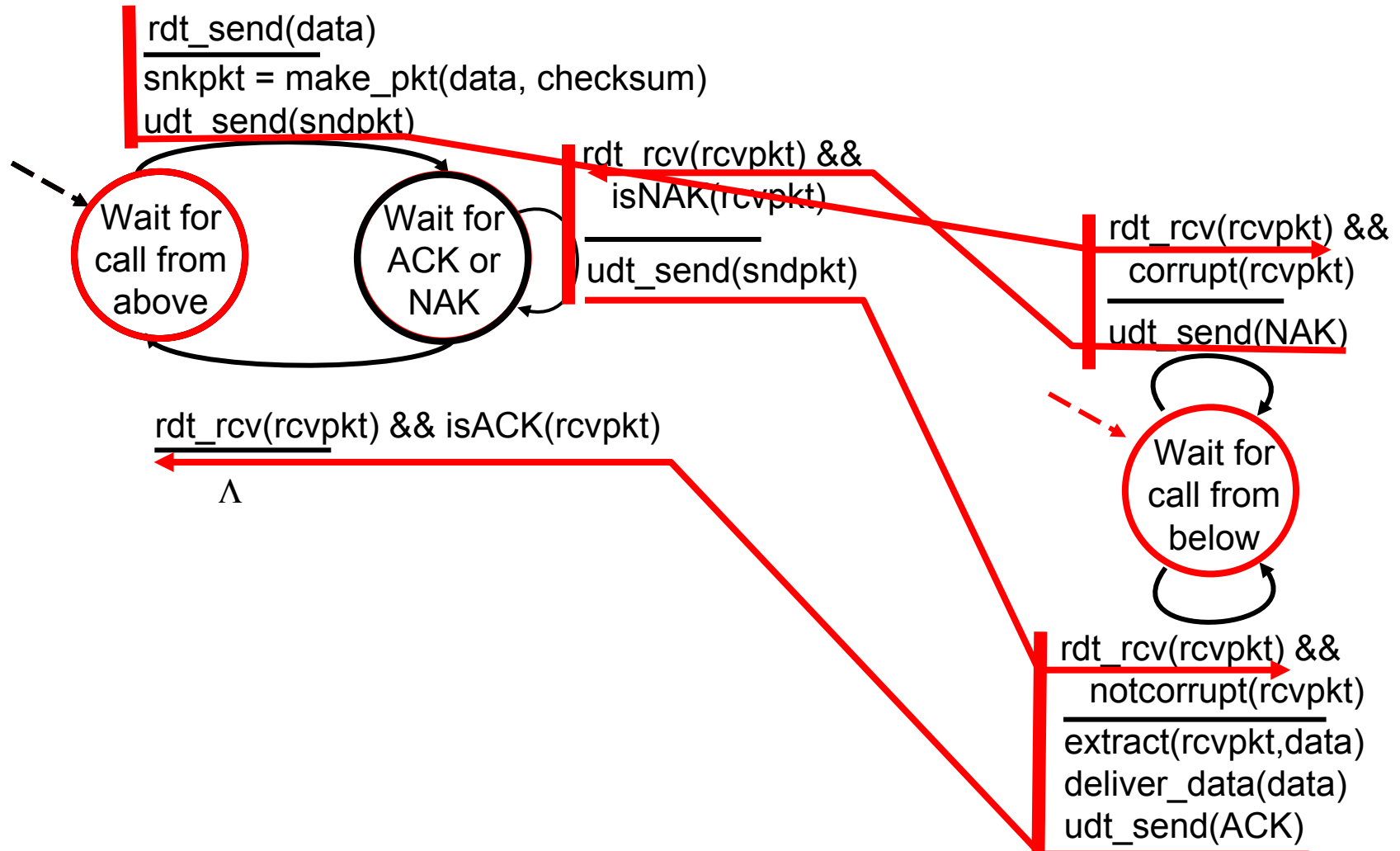
receiver



rdt2.0: operation with no errors



rdt2.0: error scenario



rdt2.0 has a fatal flaw!

What happens if ACK/NAK corrupted?

- sender doesn't know what happened at receiver!
- can't just retransmit: possible duplicate
 - Corrupted ACK - retransmitted packet is a duplicate
 - Corrupted NAK - retransmitted packet is not a duplicate
- Receiver needs to know which packet is a duplicate

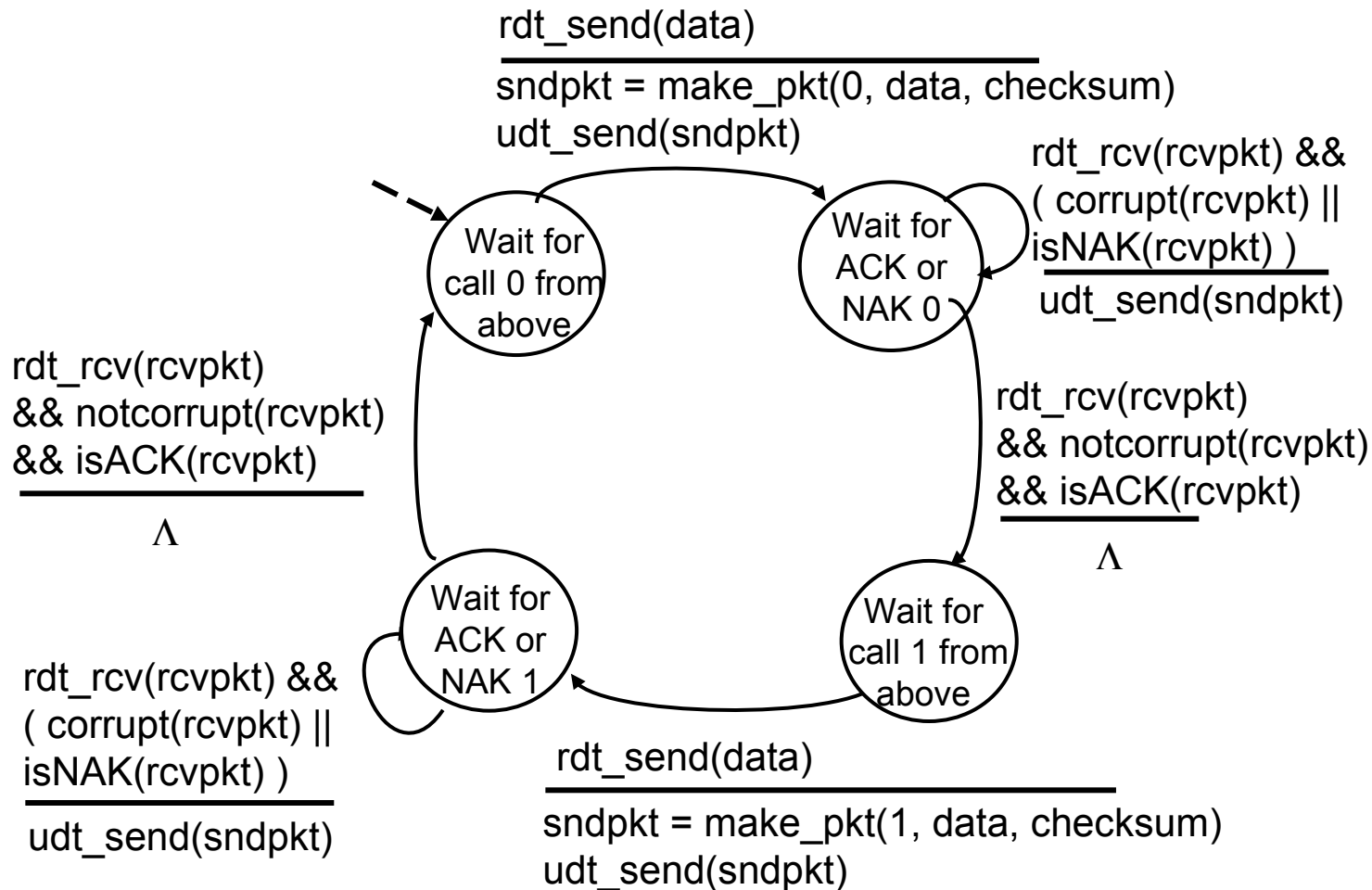
Handling duplicates:

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

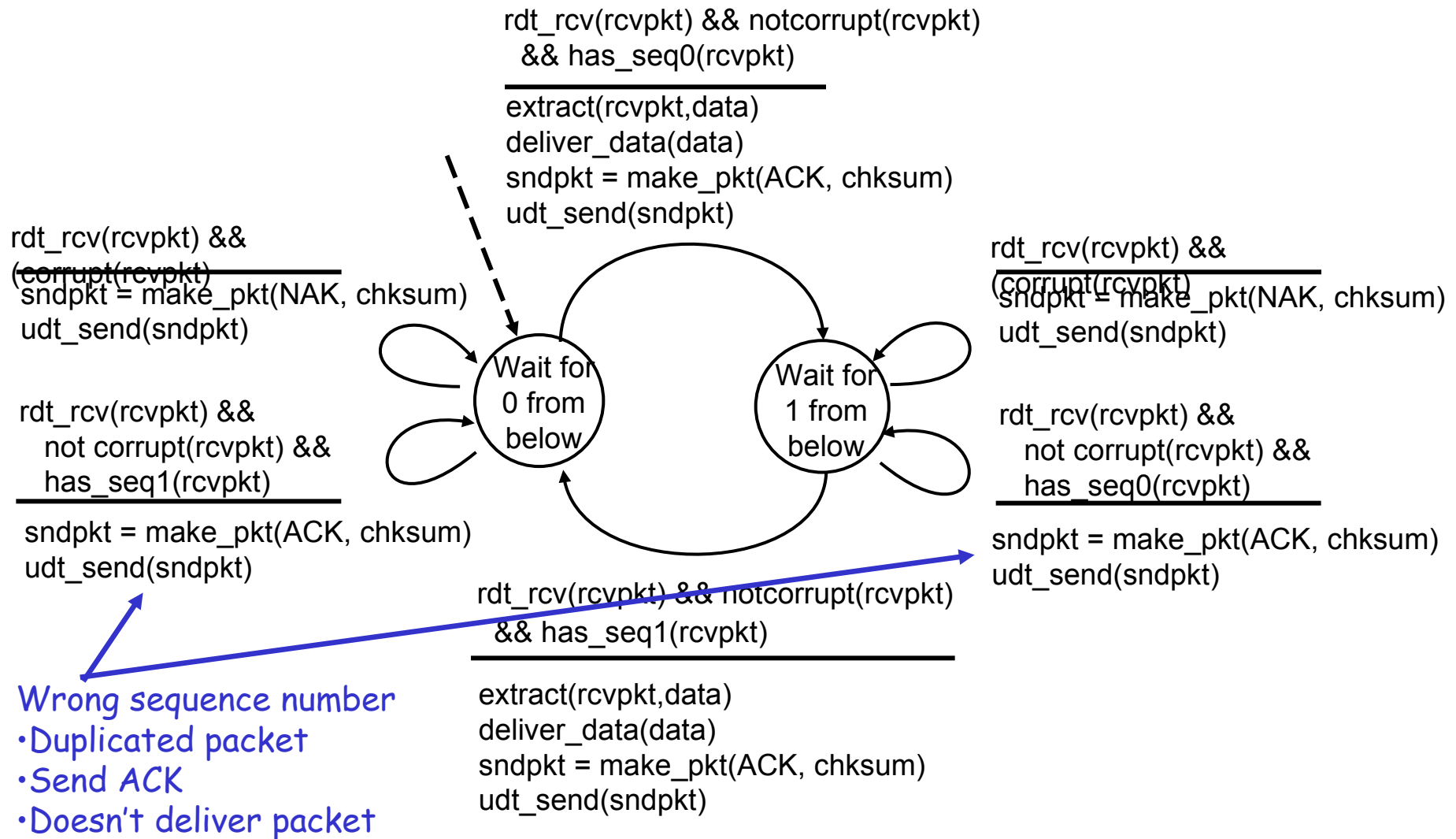
stop and wait

Sender sends one packet, then waits for receiver response

rdt2.1: sender, handles garbled ACK/NAKs



rdt2.1: receiver, handles garbled ACK/NAKs



rdt2.1: discussion

Sender:

- seq # added to packet
- two seq. #'s (0,1) will suffice. Why?
- must check if received ACK/NAK corrupted
- twice as many states
 - state must "remember" whether "current" packet has 0 or 1 seq. #

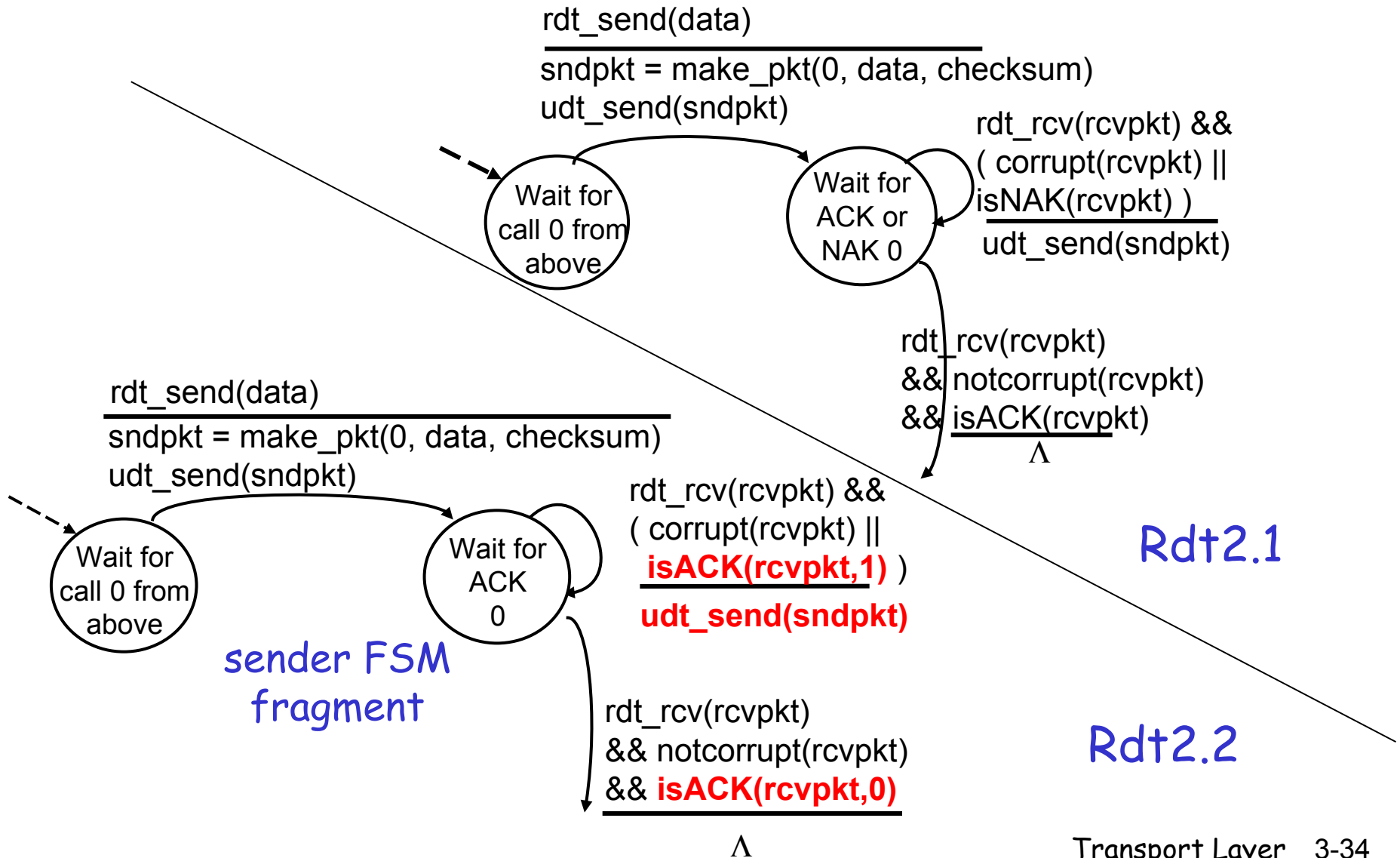
Receiver:

- must check if received packet is duplicate
 - state indicates whether 0 or 1 is expected pkt seq #
- note: receiver can *not* know if its last ACK/NAK received OK at sender

rdt2.2: a NAK-free protocol

- same functionality as rdt2.1, using ACKs only
- instead of NAK, receiver sends ACK for last pkt received OK
 - receiver must *explicitly* include seq # of packet being ACKed
- duplicate ACK at sender means that the packet following the packet being ACKed twice is corrupted
- duplicate ACK at sender results in the same action as NAK: *retransmit current packet*

rdt2.2: sender fragments



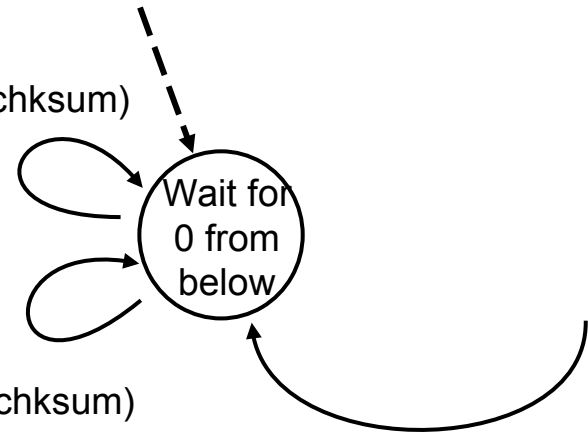
rdt2.2: receiver fragments

Rdt2.1

```
rdt_rcv(rcvpkt) &&
(corrupt(rcvpkt))
sndpkt = make_pkt(NAK, chksum)
udt_send(sndpkt)
```

```
rdt_rcv(rcvpkt) &&
not corrupt(rcvpkt) &&
has_seq1(rcvpkt)
```

```
sndpkt = make_pkt(ACK, chksum)
udt_send(sndpkt)
```



Rdt2.2

```
rdt_rcv(rcvpkt) &&
(corrupt(rcvpkt) ||
has_seq1(rcvpkt))
udt_send(sndpkt)
```



Same packet

```
rdt_rcv(rcvpkt) && notcorrupt(rcvpkt)
&& has_seq1(rcvpkt)
```

```
extract(rcvpkt,data)
deliver_data(data)
sndpkt = make_pkt(ACK1, chksum)
udt_send(sndpkt)
```

```
rdt_rcv(rcvpkt) && notcorrupt(rcvpkt)
&& has_seq1(rcvpkt)
```

```
extract(rcvpkt,data)
deliver_data(data)
sndpkt = make_pkt(ACK, chksum)
udt_send(sndpkt)
```

rdt3.0: channels with errors and loss

New assumption:

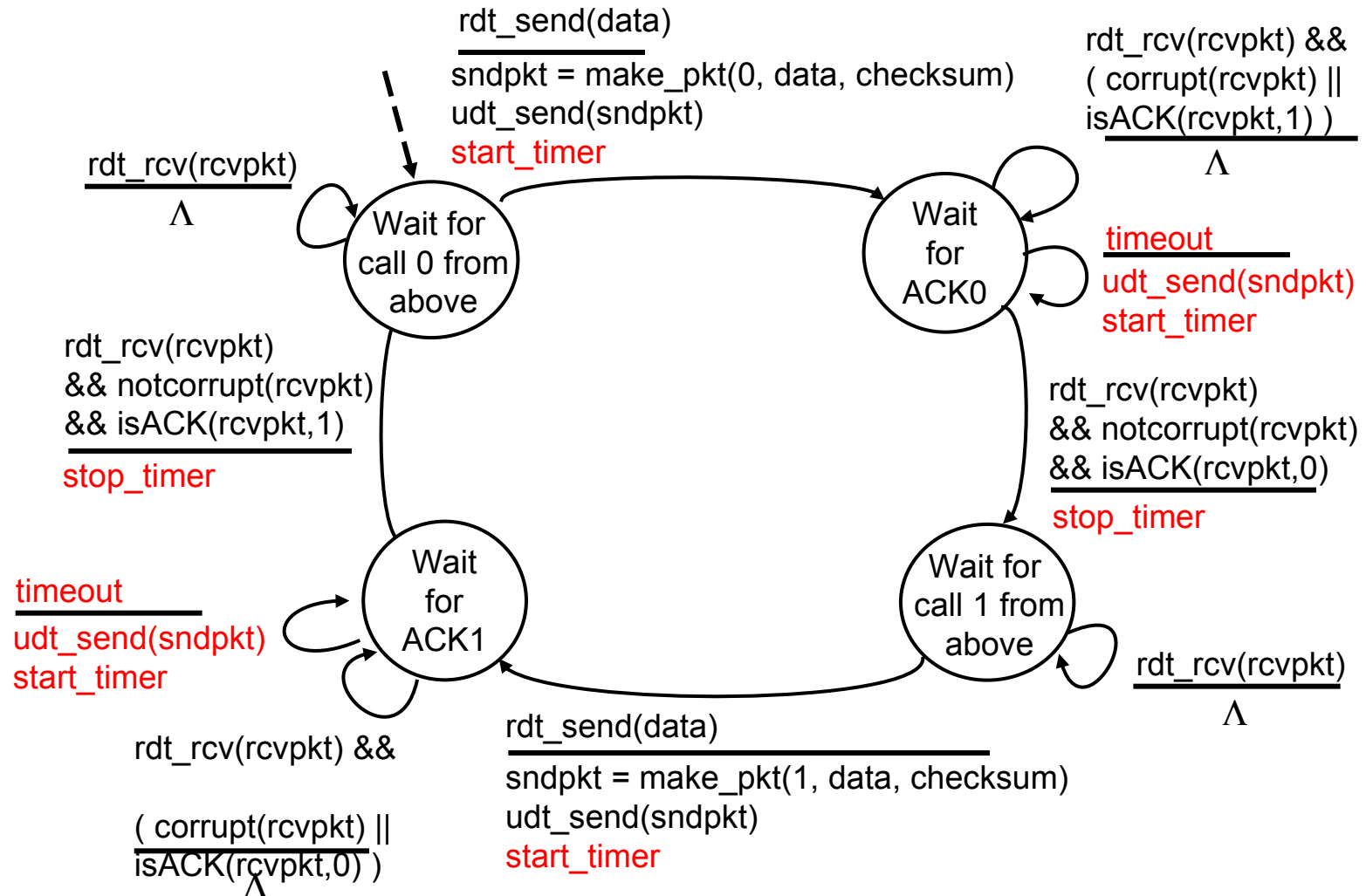
underlying channel can also lose packets (data or ACKs)

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

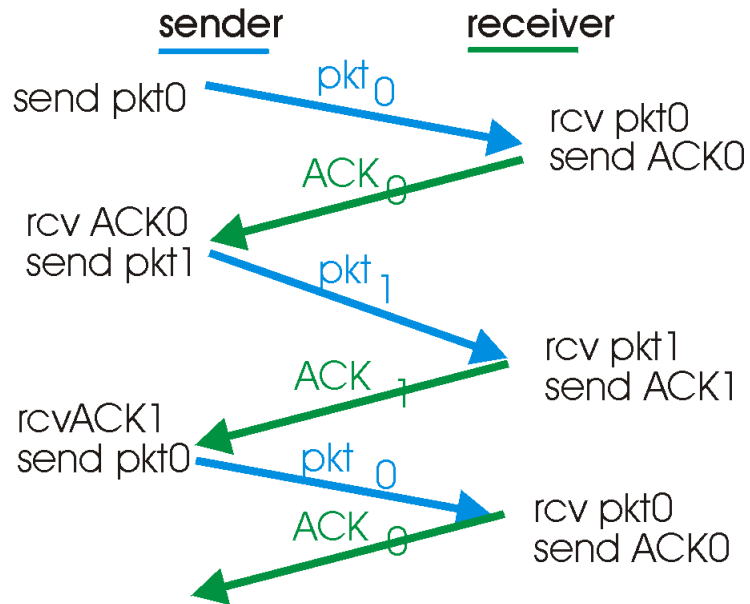
Approach: sender waits “reasonable” amount of time for ACK

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

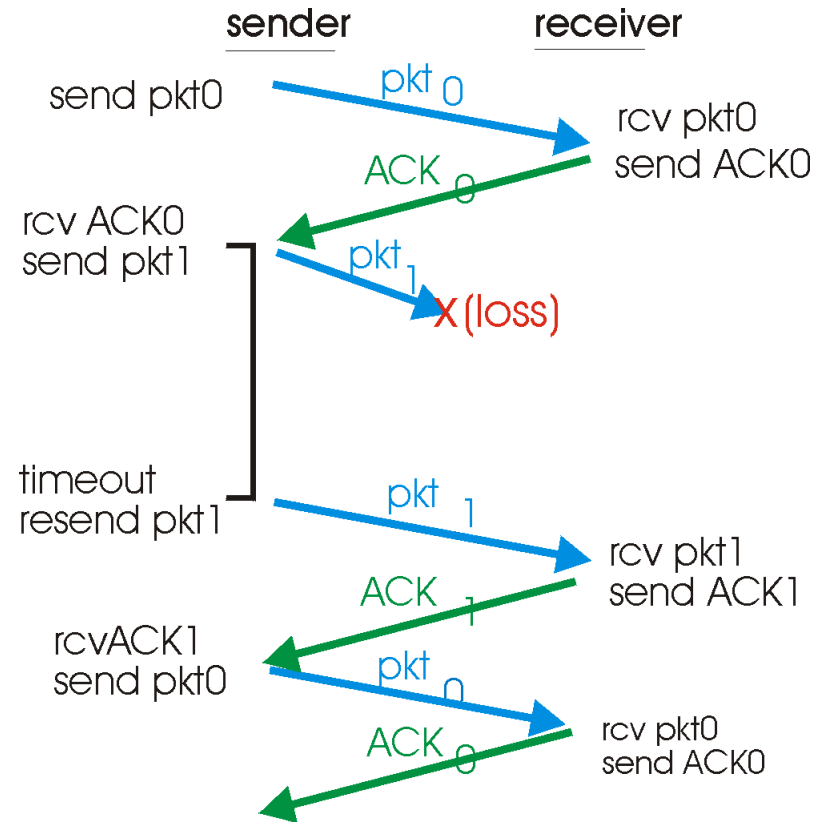
rdt3.0 sender



rdt3.0 in action

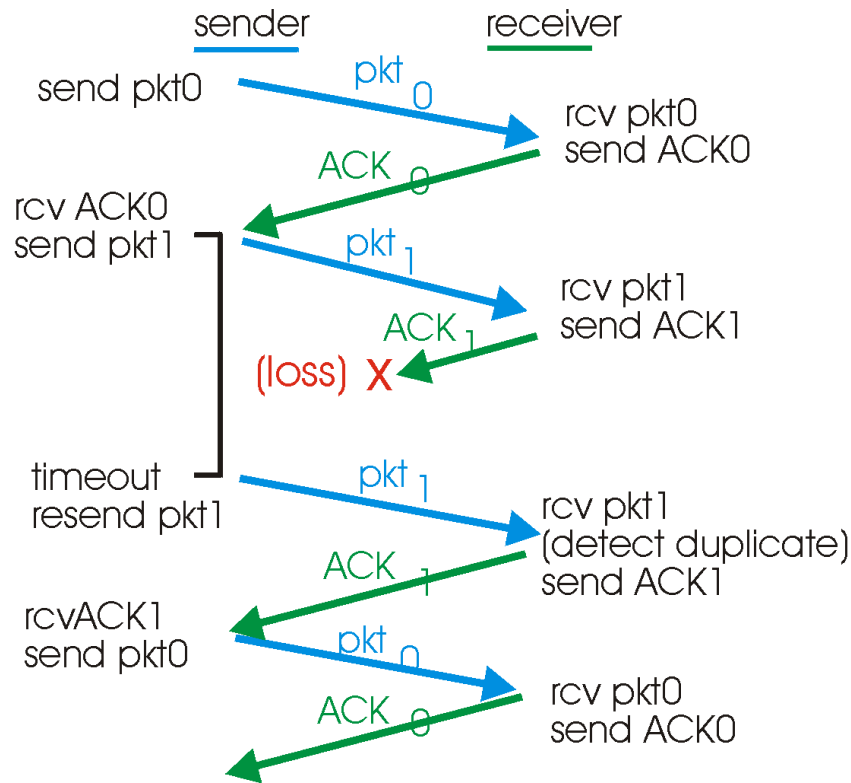


(a) operation with no loss

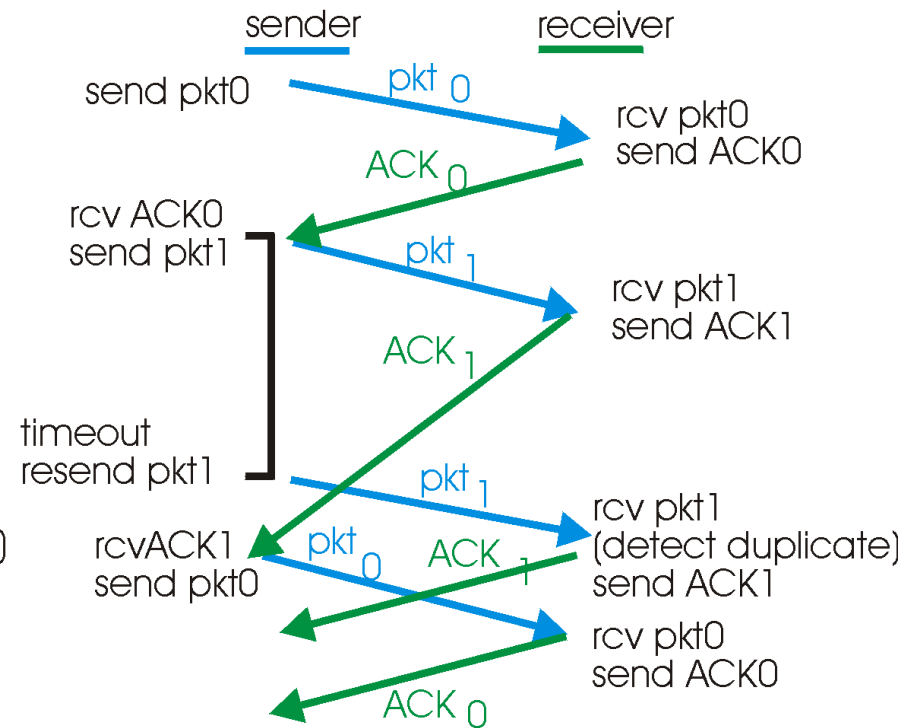


(b) lost packet

rdt3.0 in action



(c) lost ACK



(d) premature timeout

Performance of rdt3.0

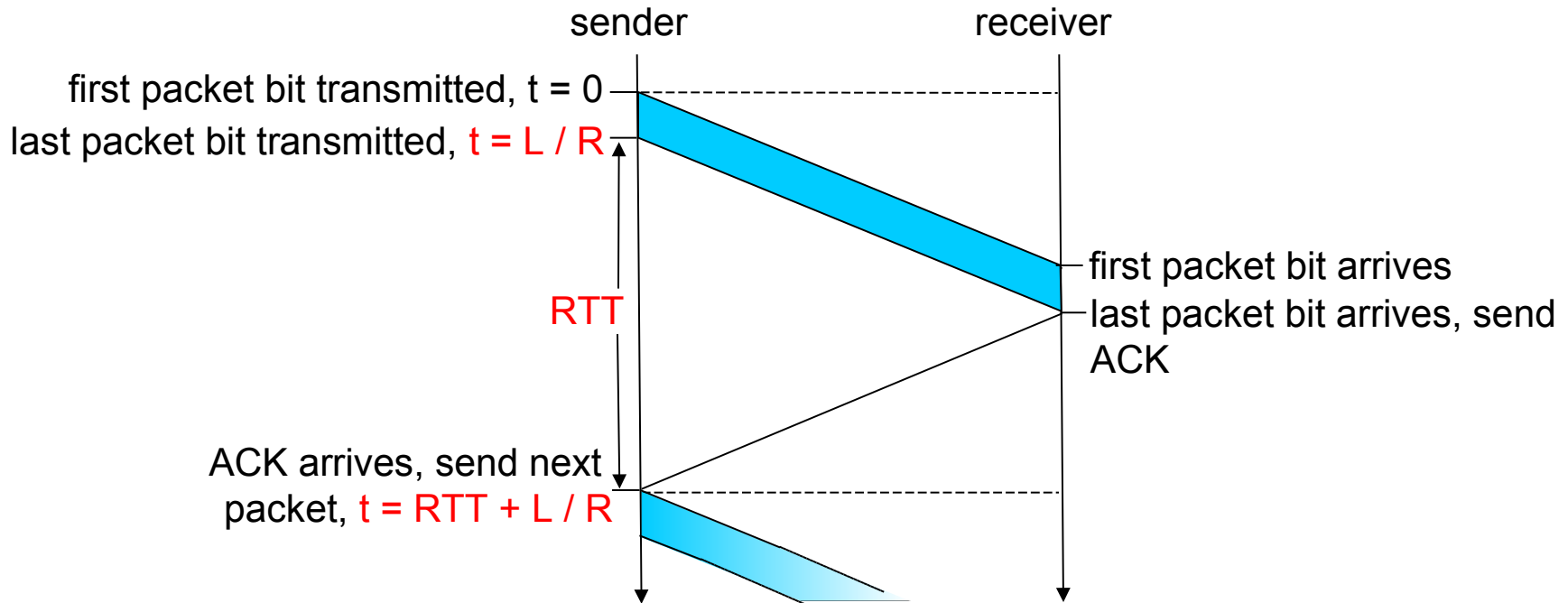
- rdt3.0 works, but performance stinks
- example: 1 Gbps link, 1KB packet
15 ms end-to-end propagation delay,

$$T_{\text{transmit}} = \frac{L \text{ (packet length in bits)}}{R \text{ (transmission rate, bps)}} = \frac{8\text{kb/pkt}}{10^9 \text{ b/sec}} = 8 \text{ microsec}$$

$$U_{\text{sender}} = \frac{L / R}{RTT + L / R} = \frac{.008}{30.008} = 0.00027$$

- U_{sender} : **utilization** - fraction of time sender busy sending
- 1KB pkt every 30 msec → 33kB/sec throughput over 1 Gbps link
- network protocol limits use of physical resources!

rdt3.0: stop-and-wait operation



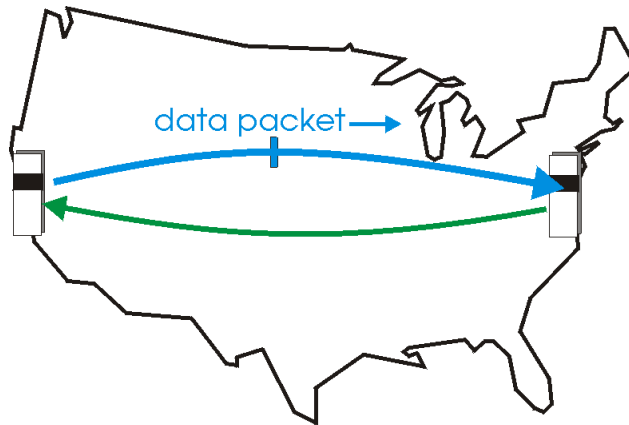
Utilization:

$$U_{\text{sender}} = \frac{L / R}{RTT + L / R} = \frac{.008}{30.008} = 0.00027$$

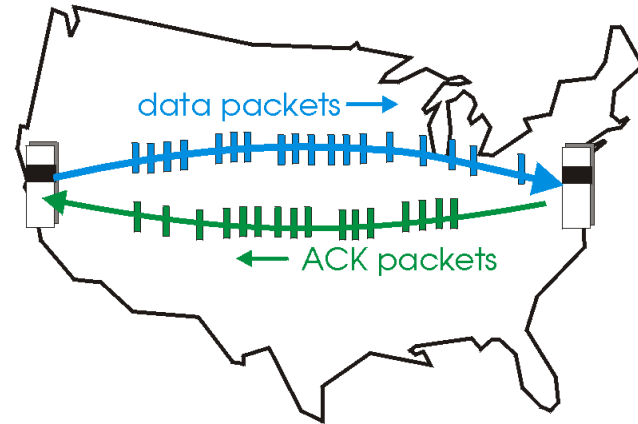
Pipelined protocols

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

- range of sequence numbers must be increased
- buffering at sender and/or receiver



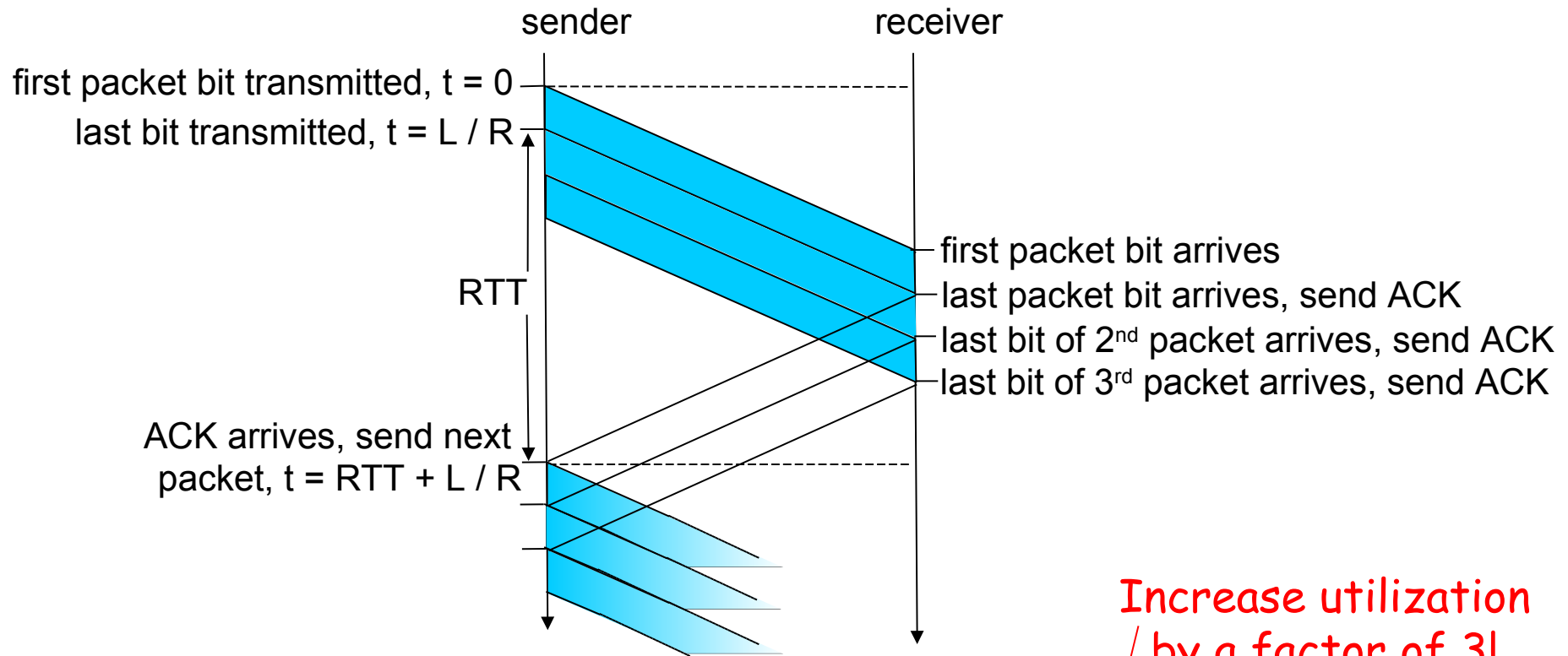
(a) a stop-and-wait protocol in operation



(b) a pipelined protocol in operation

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

Pipelining: increased utilization



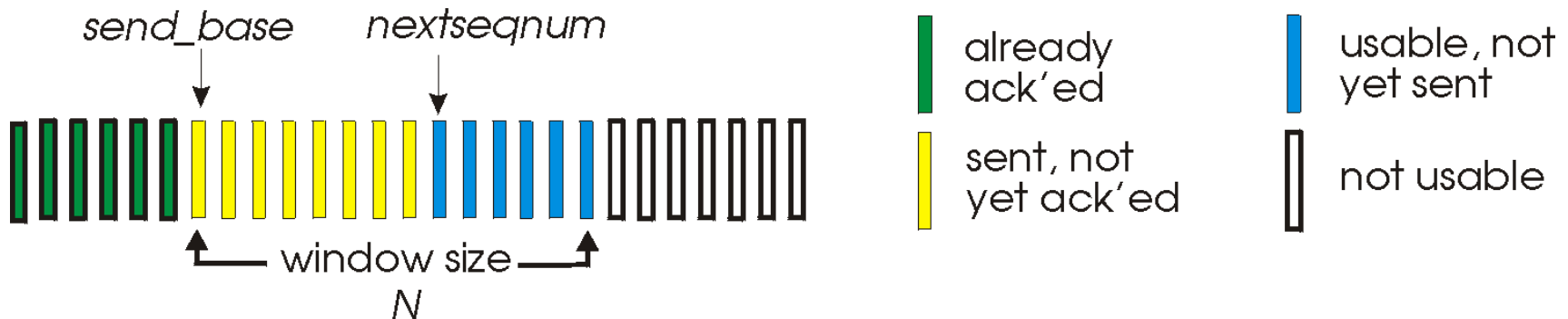
$$U_{\text{sender}} = \frac{3 * L / R}{RTT + L / R} = \frac{.024}{30.008} = 0.0008$$

Increase utilization
by a factor of 3!

Go-Back-N

Sender:

- k-bit seq # in pkt header
- "window" of up to N, consecutive unack'ed pkts allowed



- send_base: seq # of the oldest unack'ed packet
- nextseqnum: smallest unused sequence number
- ACK(n): ACKs all pkts up to, including seq # n - "cumulative ACK"
 - may receive duplicate ACKs (see receiver)
- timer for each in-flight pkt
- timeout(n): retransmit pkt n and all higher seq # pkts in window

GBN: sender extended FSM

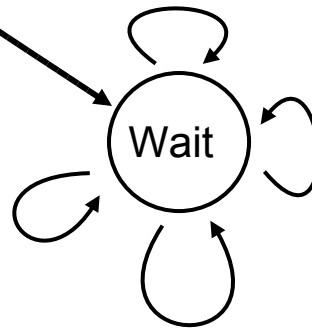
Λ

base=1
nextseqnum=1

rdt_send(data)

```

if (nextseqnum < base+N) {
    sndpkt[nextseqnum] = make_pkt(nextseqnum,data,chksum)
    udt_send(sndpkt[nextseqnum])
    if (base == nextseqnum) ← the first unack'ed packet?
        start_timer
    nextseqnum++
}
else
    refuse_data(data)
    
```



rdt_rcv(rcvpkt)
&& corrupt(rcvpkt)

timeout

```

start_timer
udt_send(sndpkt[base])
udt_send(sndpkt[base+1])
...
udt_send(sndpkt[nextseqnum-
1])
    
```

rdt_rcv(rcvpkt) &&
notcorrupt(rcvpkt)

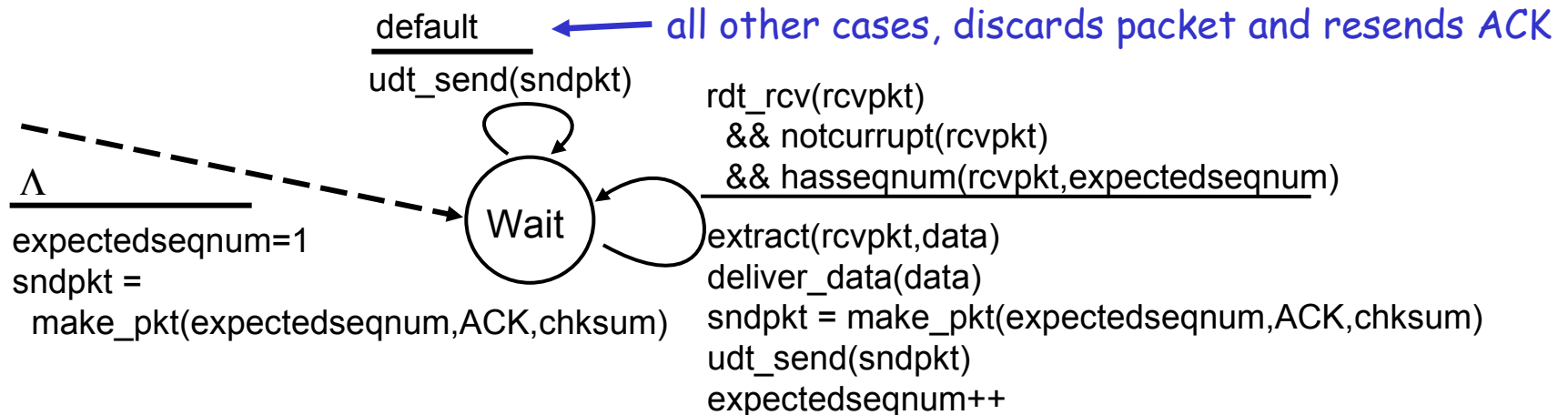
base = getacknum(rcvpkt)+1

If (base == nextseqnum) ← no more unack'ed packet?

```

stop_timer
else
    start_timer
    
```

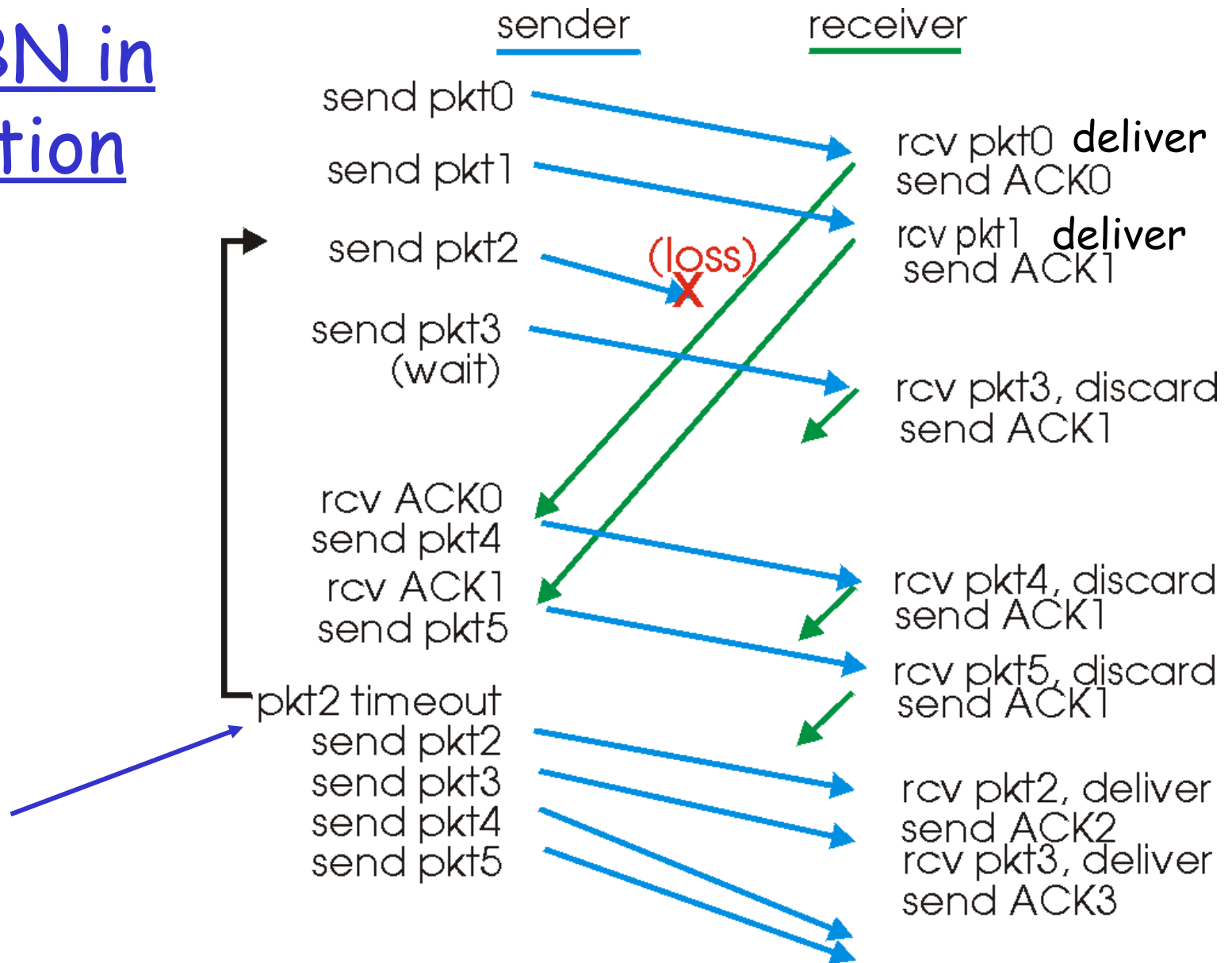
GBN: receiver extended FSM



ACK-only: always send ACK for correctly-received pkt with highest *in-order* seq #

- may generate duplicate ACKs
- need only remember **expectedseqnum**
- out-of-order pkt:
 - discard (don't buffer) -> **no receiver buffering!**
 - Re-ACK pkt with highest in-order seq #

GBN in action



Selective Repeat

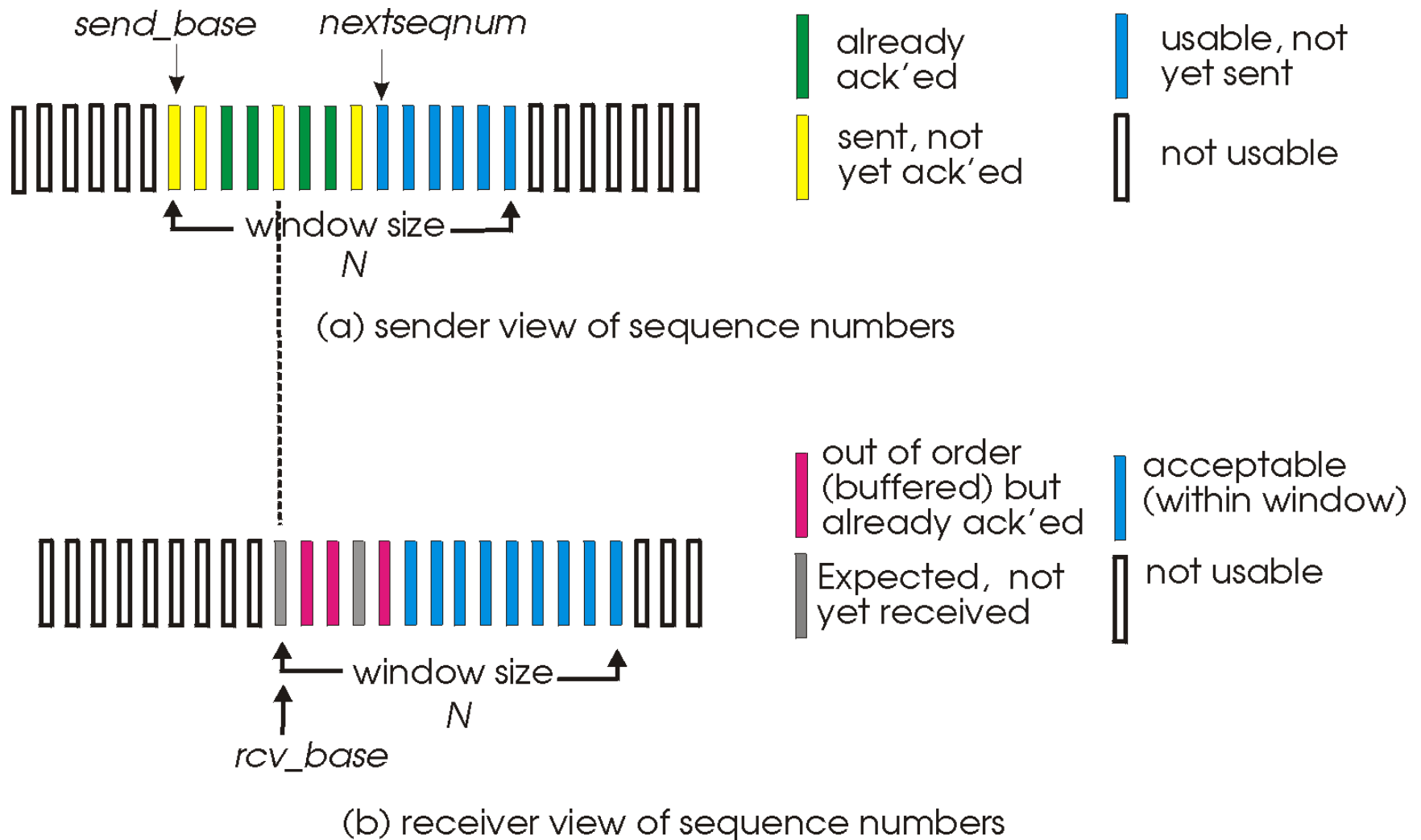
Drawback of Go-Back-N protocol

- A single packet error can cause the protocol to retransmit a large number of packets

Selected repeat

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

Selective repeat: sender, receiver windows



Selective repeat

sender

data from above :

- if next available seq # in window, send packet

timeout(n):

- resend pkt n, restart timer

receives ACK(n) in
[sendbase, sendbase+N-1]:

- mark pkt n as received
- if n is the smallest unACKed packet, advance window base to next unACKed seq #

receiver

receives packet n in
[rcvbase, rcvbase+N-1]

- send ACK(n)
- out-of-order: buffer the packet
- in-order: deliver (also deliver buffered, in-order packets), advance window to next not-yet-received packet

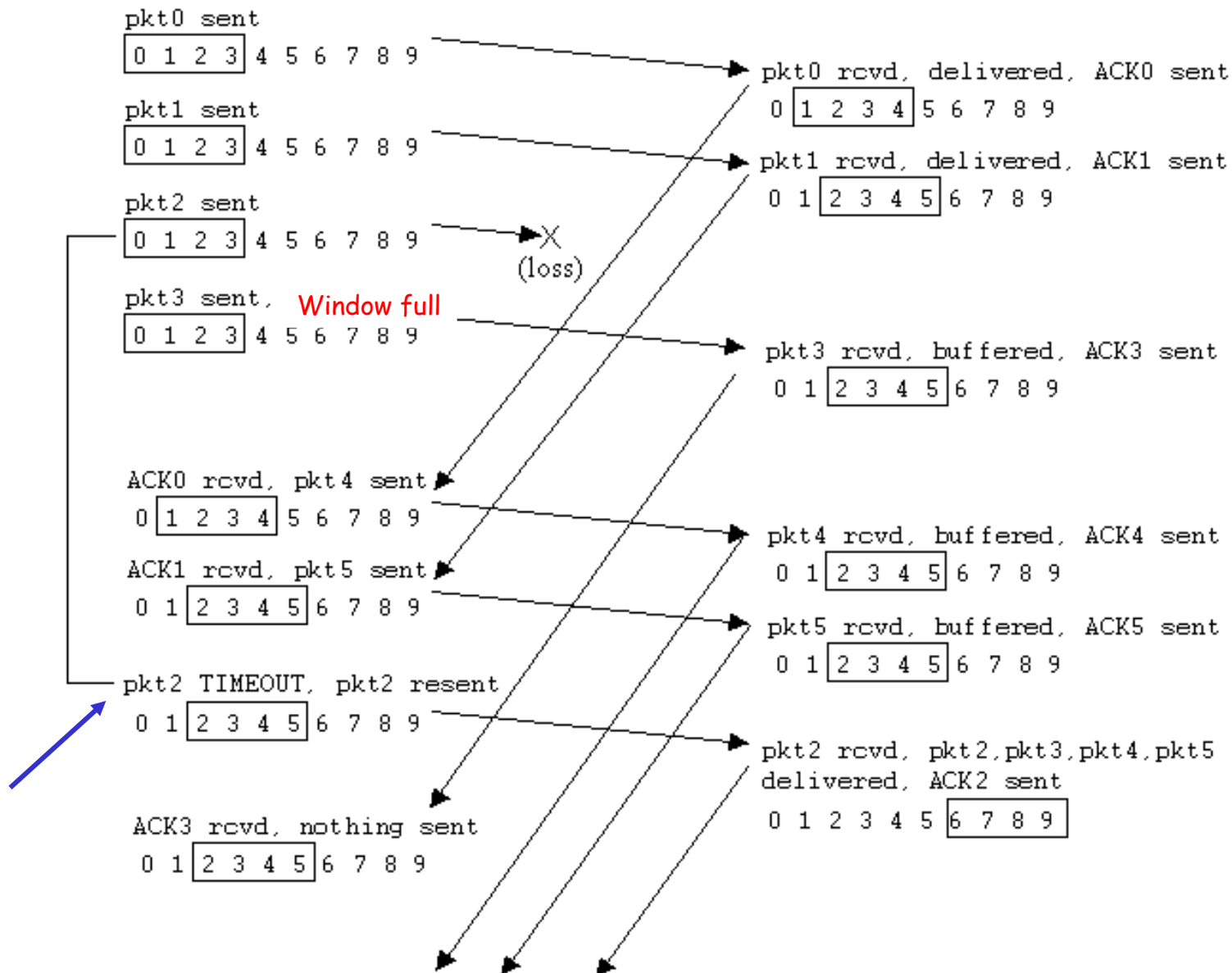
receive packet n in
[rcvbase-N, rcvbase-1]

- Already ACKed before, send ACK(n) again

otherwise:

- ignore the packet

Selective repeat in action



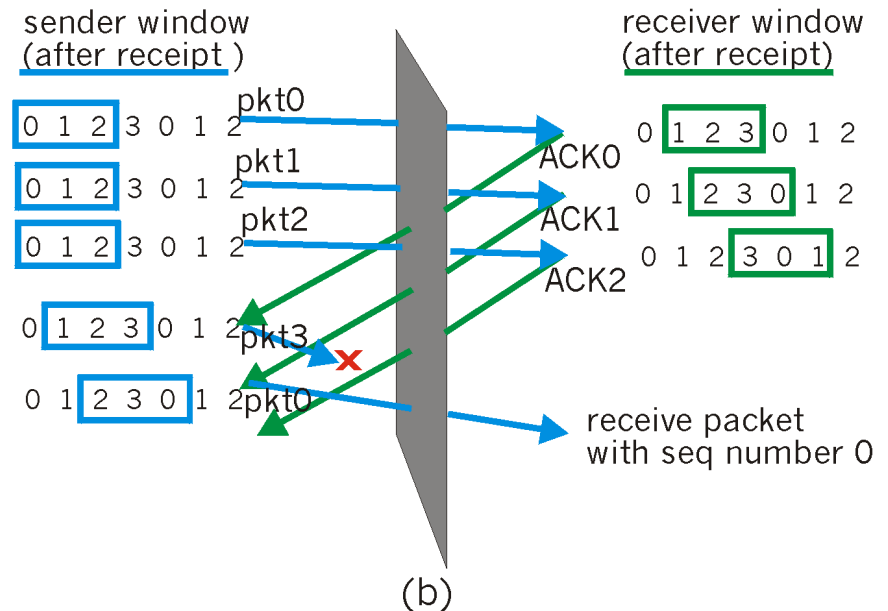
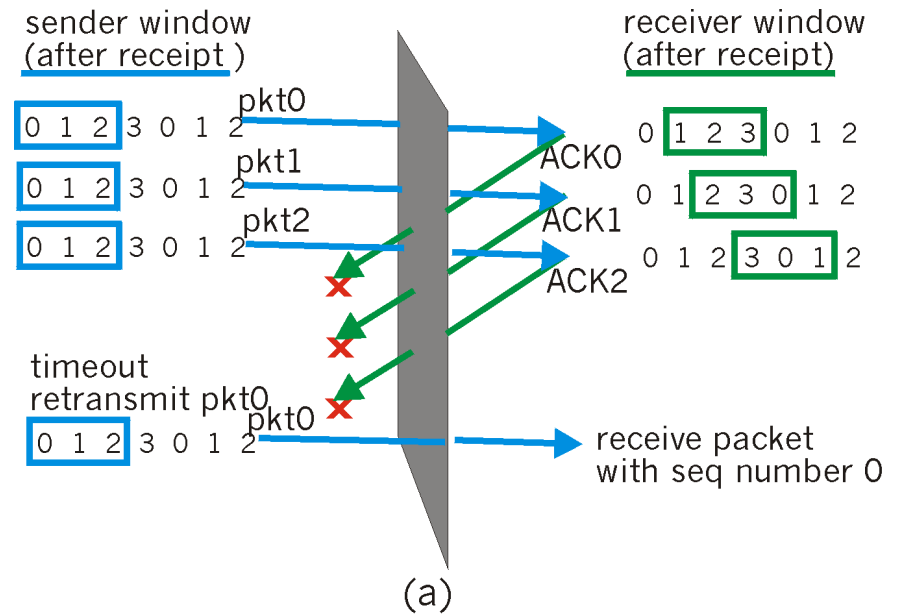
Selective repeat: dilemma

Example:

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

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

window size \leq (size of seq #)/2



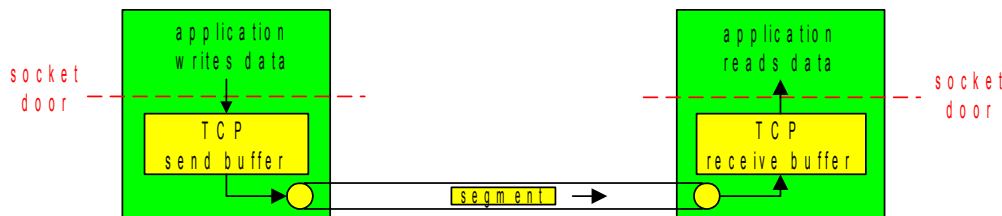
Chapter 3 outline

- 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

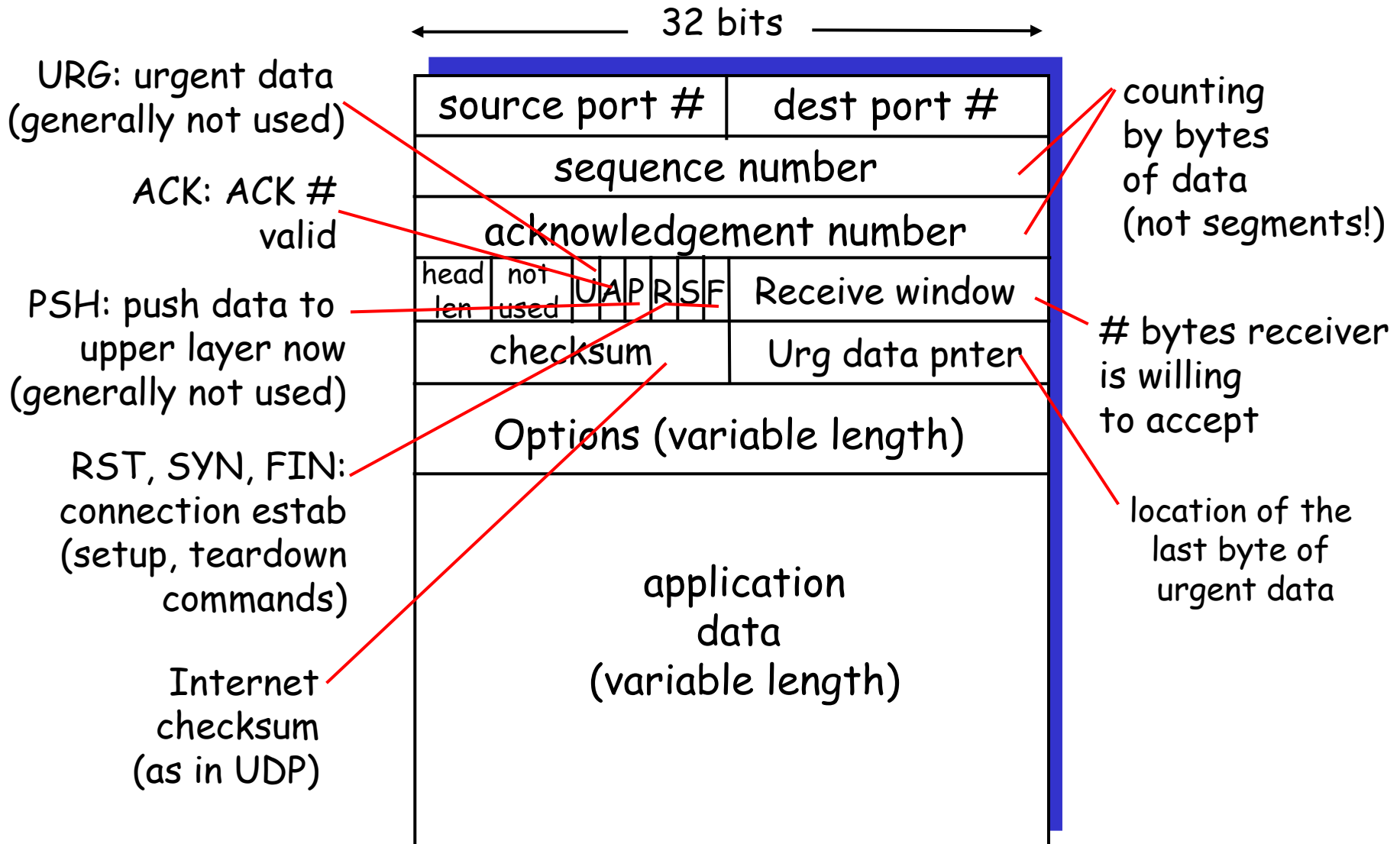
TCP: Overview

RFCs: 793, 1122, 1323, 2018, 2581

- **point-to-point:**
 - one sender, one receiver
- **reliable, in-order *byte stream*:**
 - no “message boundaries”
- **pipelined:**
 - TCP congestion and flow control set window size
- ***send & receive buffers***
- **full duplex data:**
 - bi-directional data flow in same connection
 - MSS: maximum segment size
- **connection-oriented:**
 - handshaking (exchange of control msgs) init's sender, receiver state before data exchange
- **flow controlled:**
 - sender will not overwhelm receiver



TCP segment structure



TCP seq. #'s and ACKs

Seq. #'s:

- byte-stream
"number" of first byte in segment's data

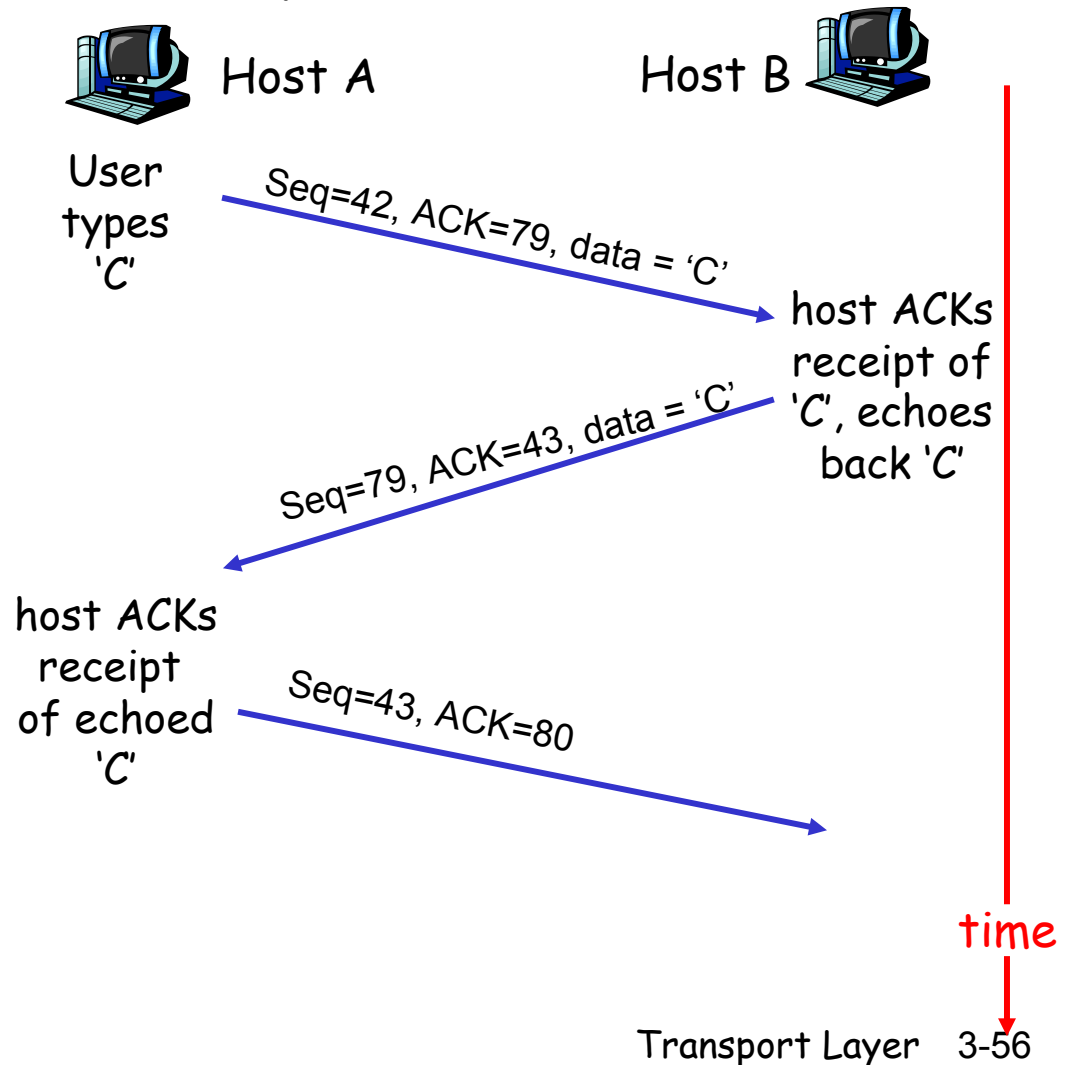
ACKs:

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

Q: how receiver handles out-of-order segments

- A: TCP spec doesn't say, - up to implementor
- two options:
discard or buffer

simple telnet scenario



TCP Round Trip Time and Timeout

Q: how to set TCP timeout value?

- longer than RTT
 - but RTT varies
- too short: premature timeout
 - unnecessary retransmissions
- too long: slow reaction to segment loss

Q: how to estimate RTT?

- **SampleRTT**: measured time from segment transmission until ACK receipt
 - Takes one **SampleRTT** measurement at a time
 - obtain a new value approximately once every RTT
 - Never computes a **SampleRTT** for a retransmitted segment
- **SampleRTT** will vary, want **estimated RTT** "smoother"
 - average several recent measurements, not just current **SampleRTT**
 - the average is called **EstimatedRTT**

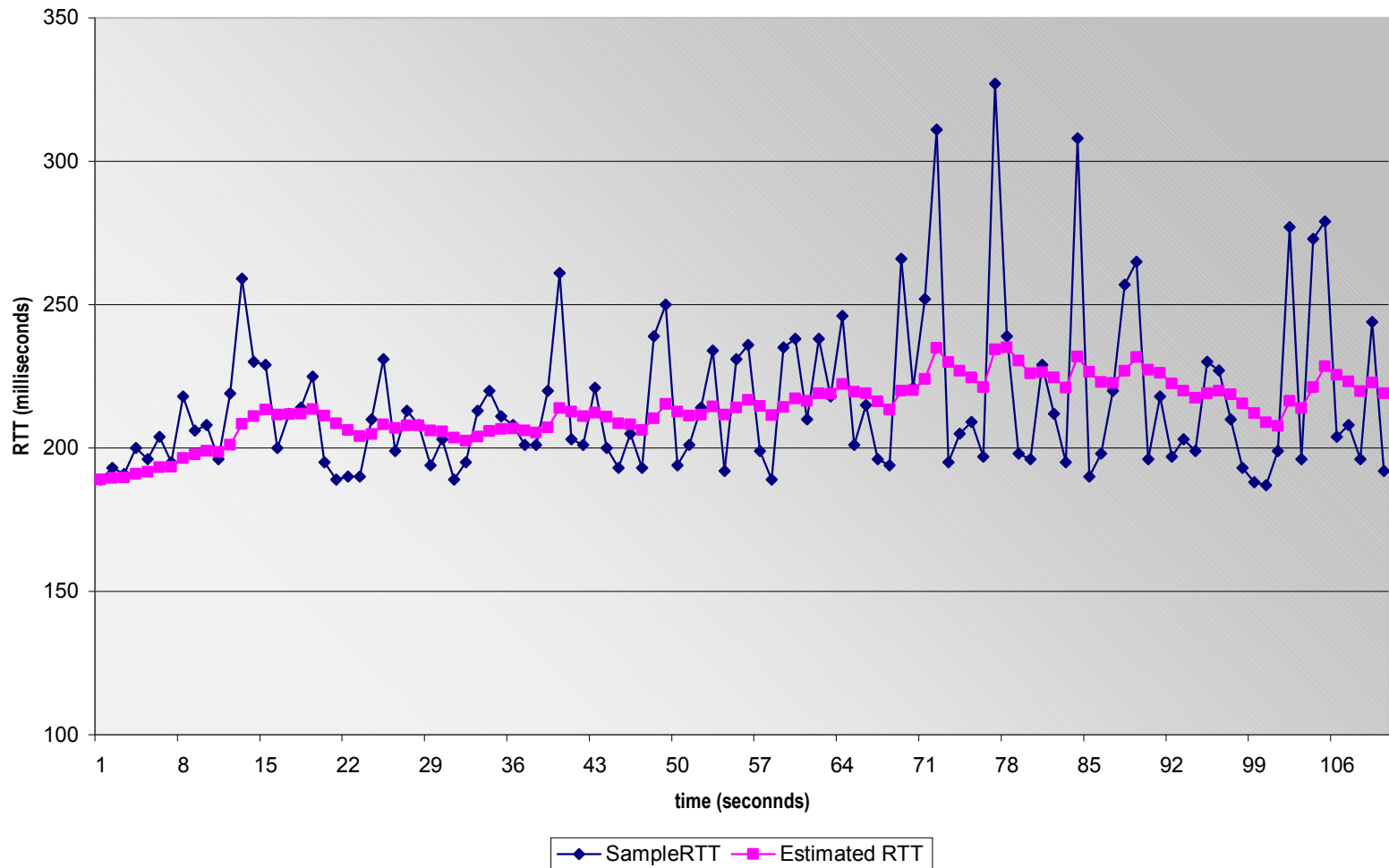
TCP Round Trip Time and Timeout

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

- Exponential weighted moving average (EWMA)
- influence of past sample decreases exponentially fast
- typical value: $\alpha = 0.125$ i.e. $1/8$

Example RTT estimation:

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



TCP Round Trip Time and Timeout

Setting the timeout

- EstimatedRTT plus "safety margin"
 - large variation in EstimatedRTT -> larger safety margin
- first estimate of how much SampleRTT deviates from EstimatedRTT:

$$\text{DevRTT} = (1-\beta) * \text{DevRTT} + \beta * |\text{SampleRTT} - \text{EstimatedRTT}|$$

(typically, $\beta = 0.25$)

Then set timeout interval:

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

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 reliable data transfer (rdt) service on top of IP's unreliable service
- Pipelined segments
- Cumulative acks
- TCP uses single retransmission timer
- Retransmissions are triggered by:
 - timeout events
 - duplicate acks
- Initially consider simplified TCP sender:
 - ignore duplicate acks
 - ignore flow control, congestion control

TCP sender events:

Three major events related to data transmission and retransmission:

1) 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`

2) timeout:

- retransmit segment that caused timeout
- restart timer

3) Ack rcvd:

- If acknowledges previously unacked segments
 - update what is known to be acked
 - start timer if there are outstanding segments

NextSeqNum = InitialSeqNum

SendBase = InitialSeqNum

loop (forever) {

 switch(event)

event: data received from application above

 create TCP segment with sequence number NextSeqNum

 if (timer currently not running)

 start timer

 pass segment to IP

 NextSeqNum = NextSeqNum + length(data)

event: timer timeout

 retransmit not-yet-acknowledged segment with

 smallest sequence number

 start timer

event: ACK received, with ACK field value of y

 if (y > SendBase) {

 SendBase = y

 if (there are currently not-yet-acknowledged segments)

 start timer

 }

 } /* end of loop forever */

TCP sender (simplified)

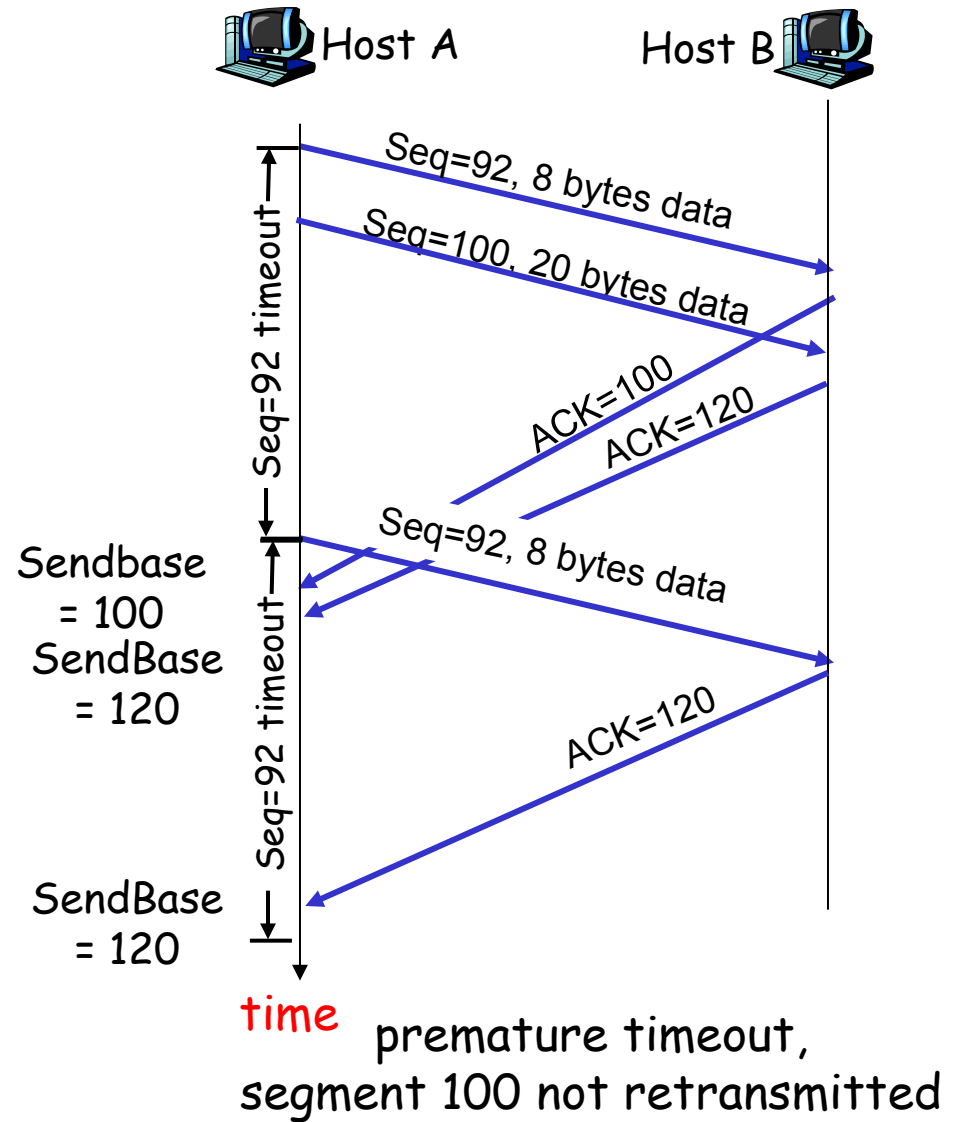
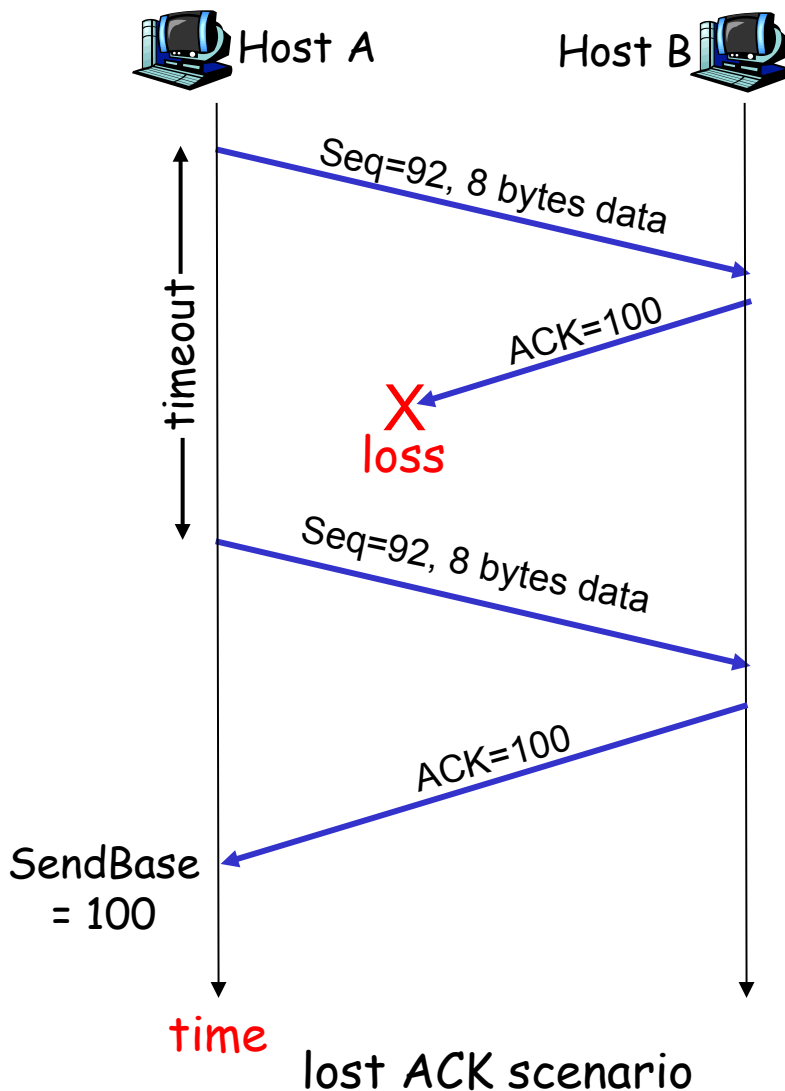
Comment:

- SendBase-1: last cumulatively ack'ed byte

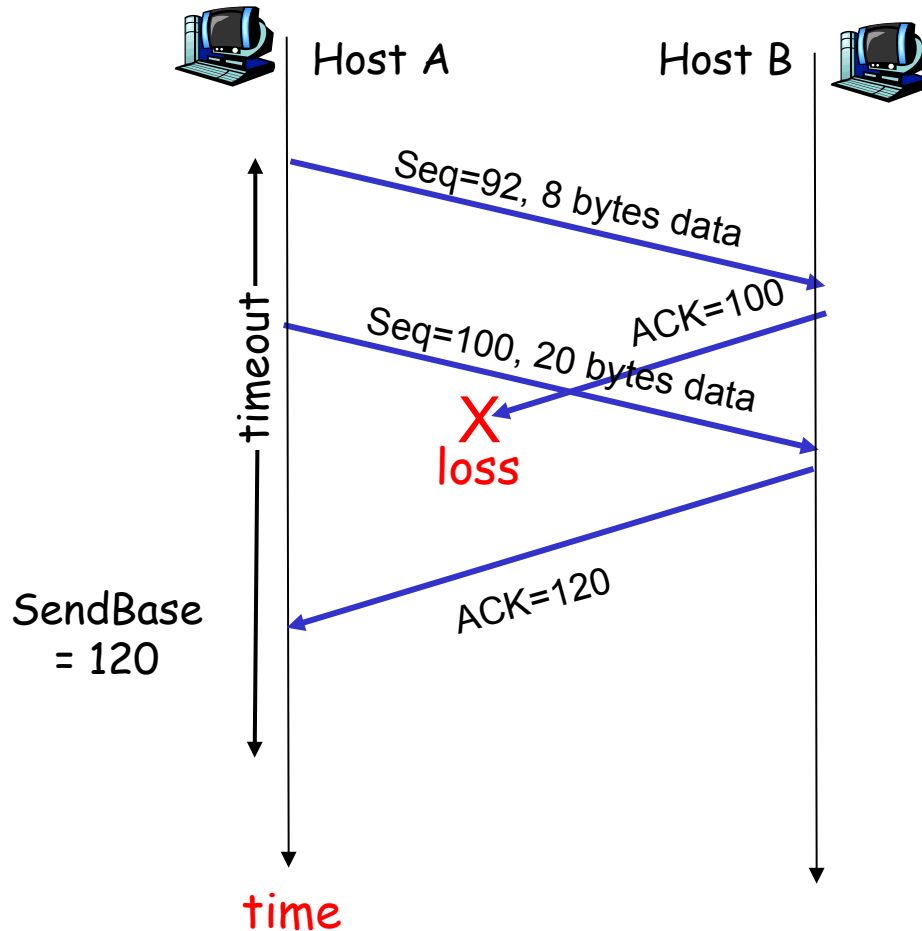
Example:

- SendBase-1 = 71;
y = 73, so the rcvr wants 73+ ;
y > SendBase, so that new data is acked

TCP: retransmission scenarios



TCP retransmission scenarios (more)



Cumulative ACK scenario,
avoid retransmission of segment 92

TCP ACK generation [RFC 1122, RFC 2581]

Event at Receiver	TCP Receiver action
Arrival of in-order segment with expected seq #. All data up to expected seq # already ACKed	Delayed ACK. Wait up to 500ms for next segment. If no next segment, send ACK
Arrival of in-order segment with expected seq #. One other in-order segment waiting for ACK transmission	Immediately send single cumulative ACK, ACKing both in-order segments
Arrival of out-of-order segment higher-than-expected seq. # . Gap detected	Immediately send duplicate ACK , indicating seq. # of next expected byte
Arrival of segment that partially or completely fills gap	Immediate send ACK, provided that segment starts at lower end of gap

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.
- If sender receives 3 ACKs for the same data, it supposes that segment after ACKed data was lost:
 - fast retransmit: resend segment before timer expires

Fast retransmit algorithm:

```
event: ACK received, with ACK field value of y
    if (y > SendBase) {
        SendBase = y
        if (there are currently not-yet-acknowledged segments)
            start timer
    }
    else {
        increment count of dup ACKs received for y
        if (count of dup ACKs received for y = 3) {
            resend segment with sequence number y
        }
    }
```

a duplicate ACK for
already ACKed segment

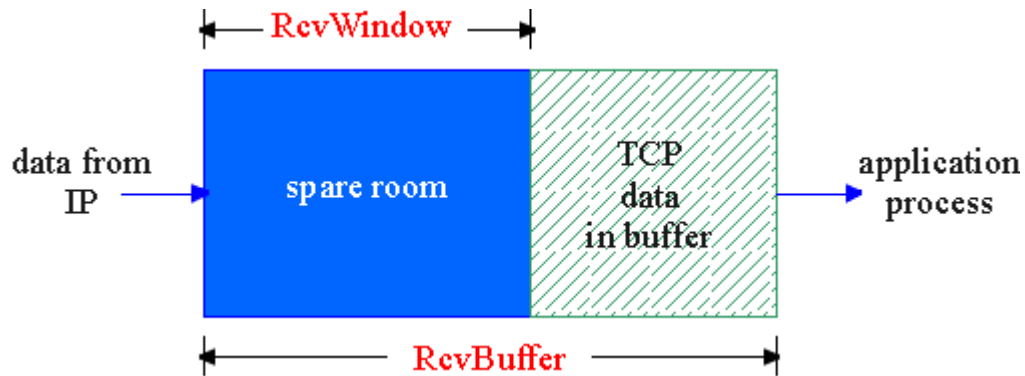
fast retransmit

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

- receive side of TCP connection has a receive buffer:

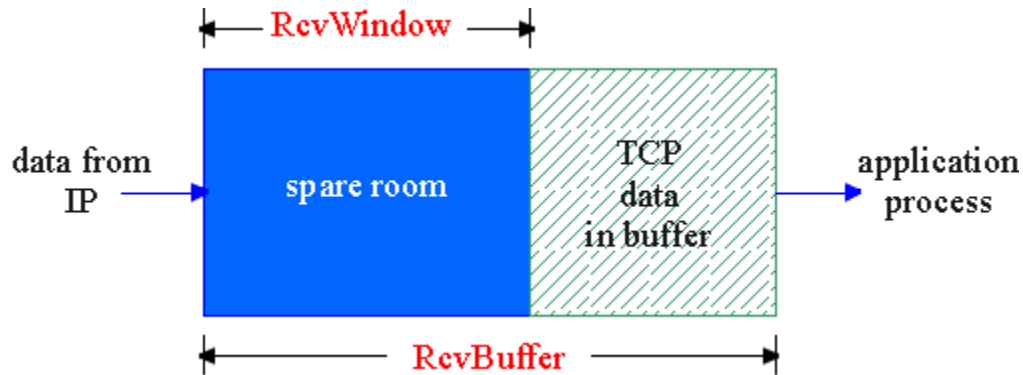


- app process may be slow at reading from buffer

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

- speed-matching service: matching the send rate to the receiving app's drain rate

TCP Flow control: how it works



(Suppose TCP receiver discards out-of-order segments)

□ spare room in buffer

= RcvWindow

= $\text{RcvBuffer} - [\text{LastByteRcvd} - \text{LastByteRead}]$

- Rcvr advertises spare room by including value of RcvWindow in segments
- Sender limits unACKed data to RcvWindow
 - guarantees receive buffer doesn't overflow

TCP Flow control: how it works

- Is sender blocked when $RcvWindow = 0$?
- Sender must be prepared to accept from the application and send at least one byte of new data even if $RcvWindow = 0$
- Sender must regularly retransmit to the receiver even when $RcvWindow = 0$

(Two minutes is recommended for the retransmission interval)

- This retransmission is essential to guarantee that when either $RcvWindow = 0$ or re-opening of the window will be reported to sender

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 Connection Management

Recall: TCP sender, receiver establish "connection" before exchanging data segments

- initialize TCP variables:

 - seq. #s

 - buffers, flow control info (e.g. RcvWindow)

- *client*: connection initiator

```
Socket clientSocket =
```

```
    new Socket("hostname", "port number");
```

- *server*: contacted by client

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

TCP Connection Management (cont.)

Three way handshake:

Step 1: client host sends TCP SYN segment to server

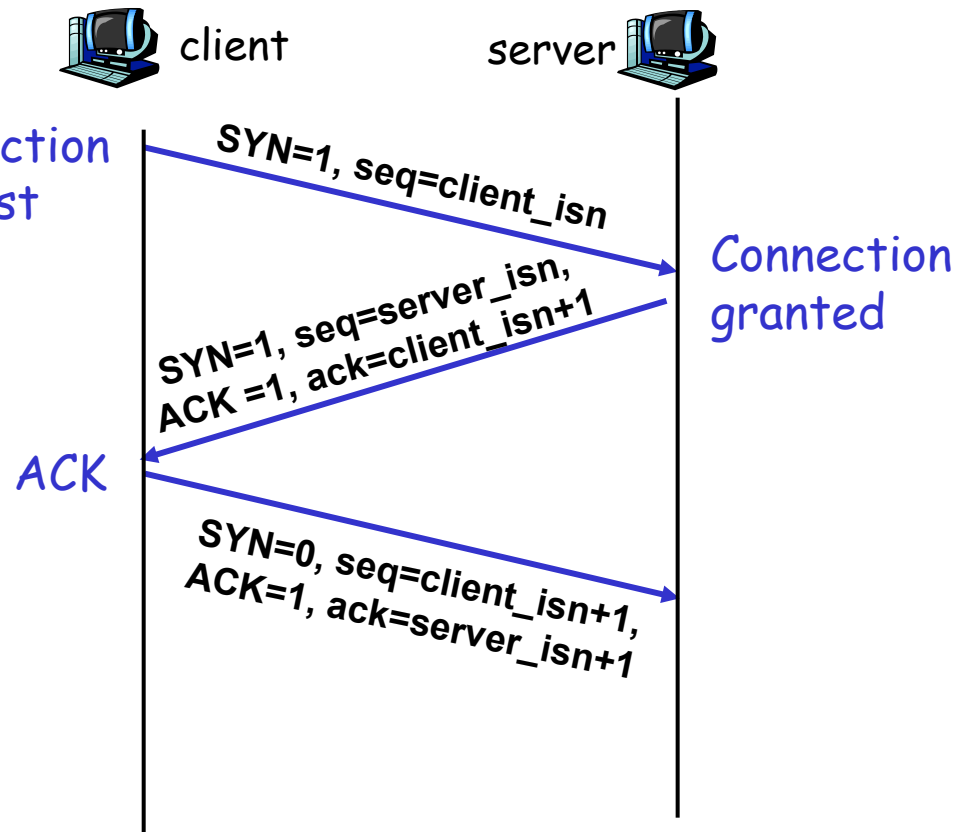
- specifies initial seq #
- no data

Connection request

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

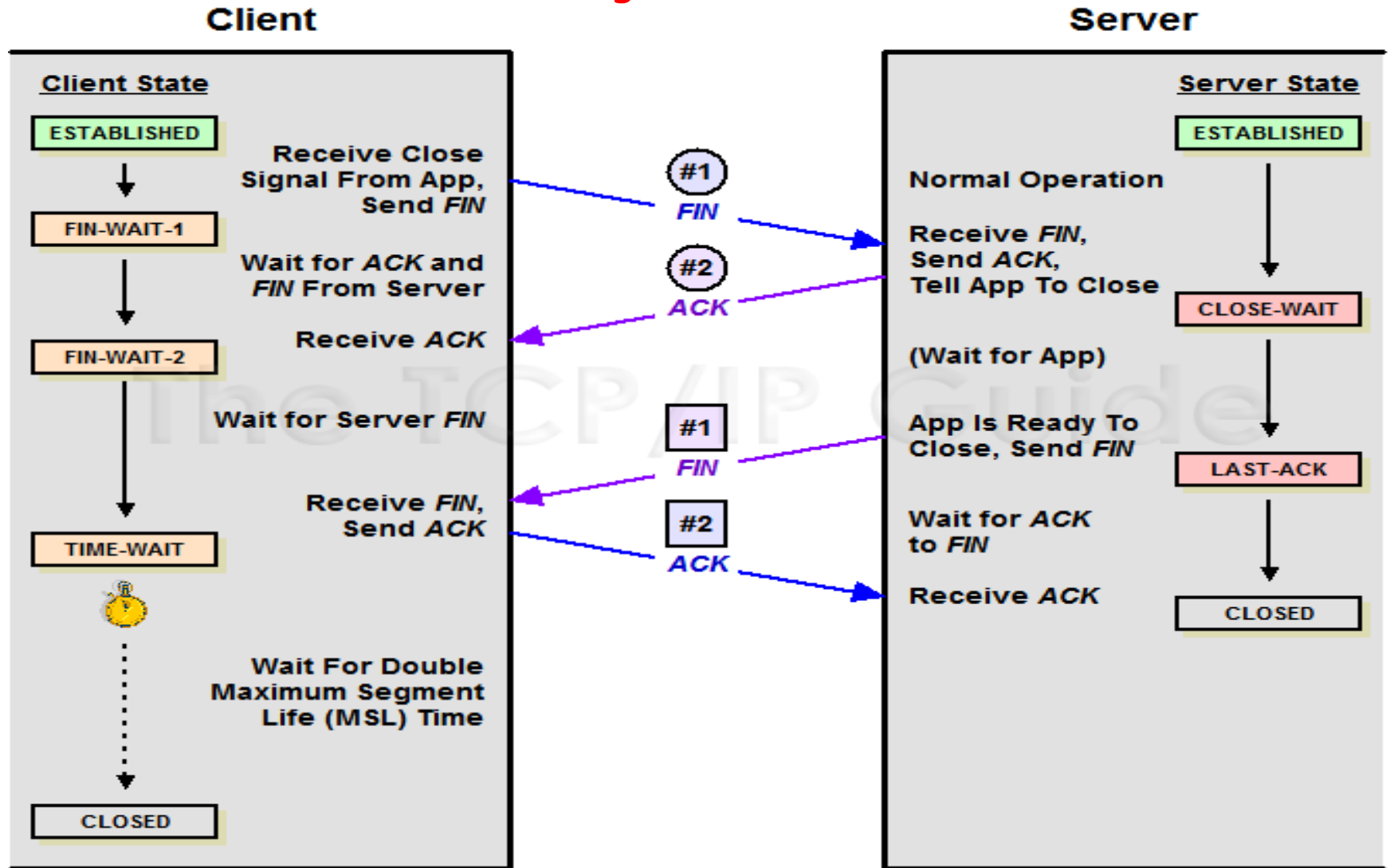
- server allocates buffers
- specifies server initial seq. #

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



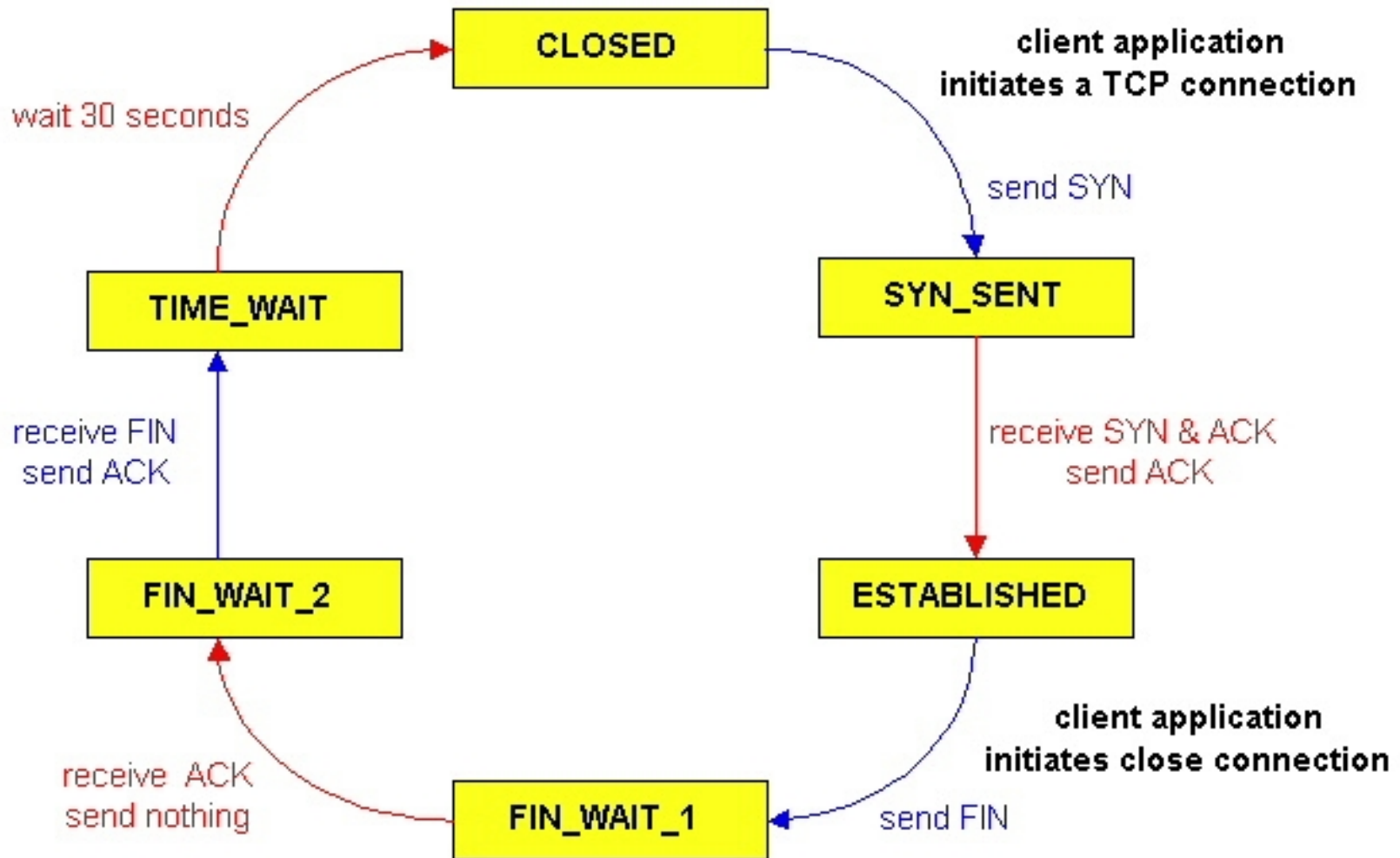
TCP Connection Management (cont.)

Closing a connection:



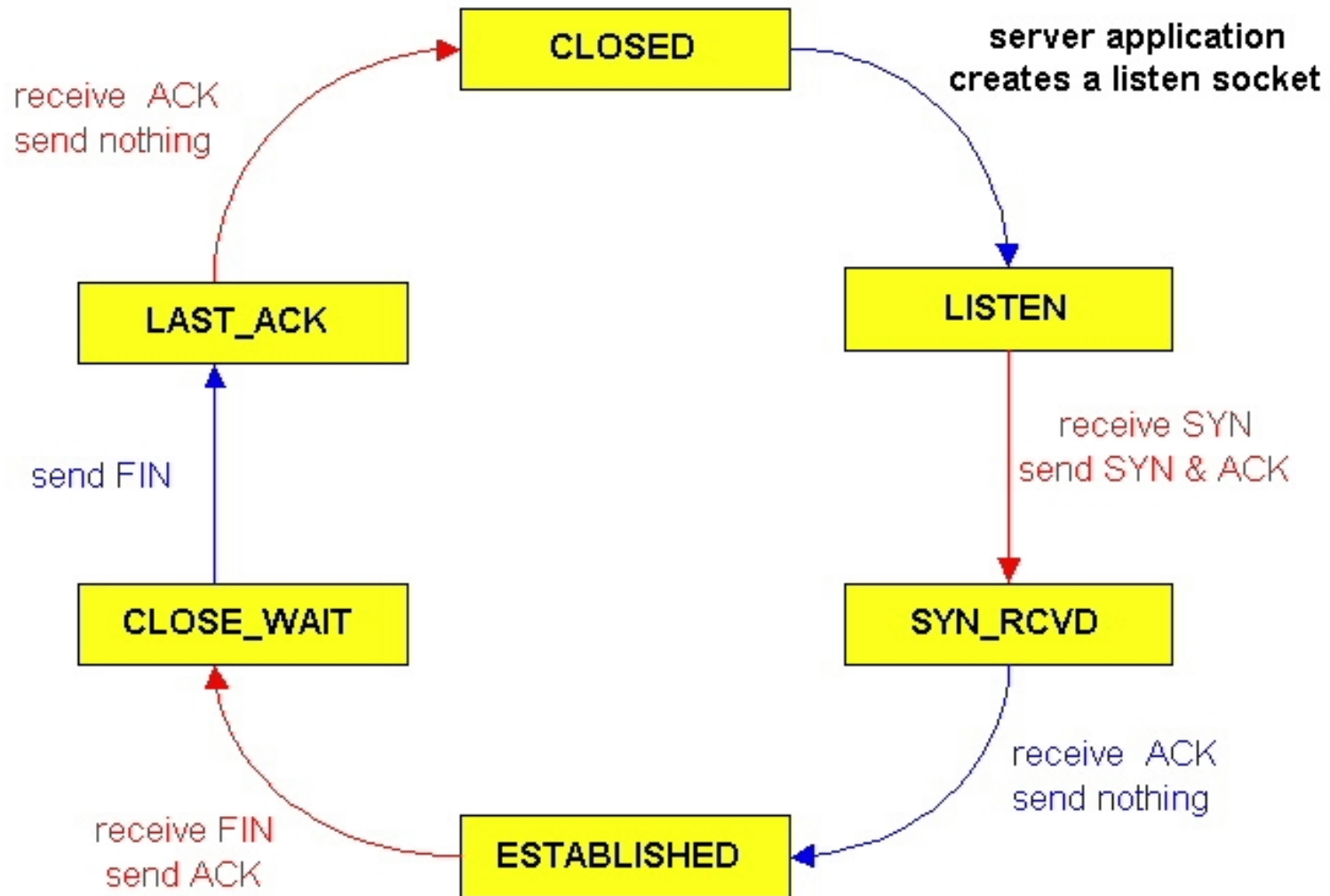
TCP Connection Management (cont)

TCP client lifecycle



TCP Connection Management (cont)

TCP server lifecycle



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

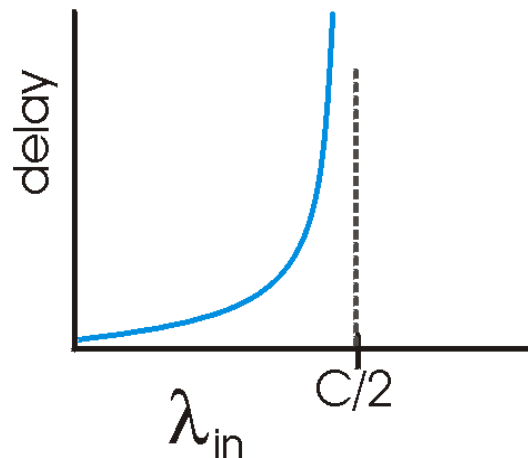
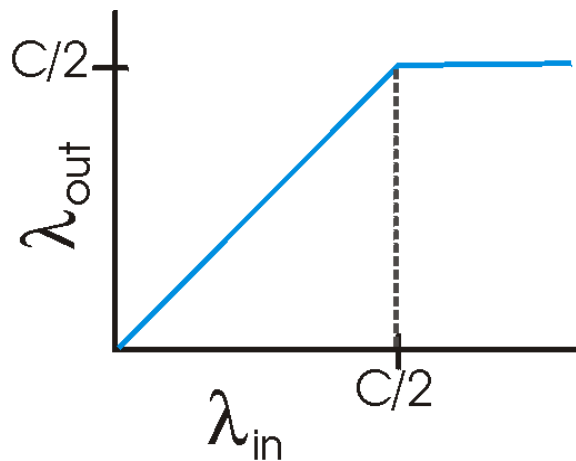
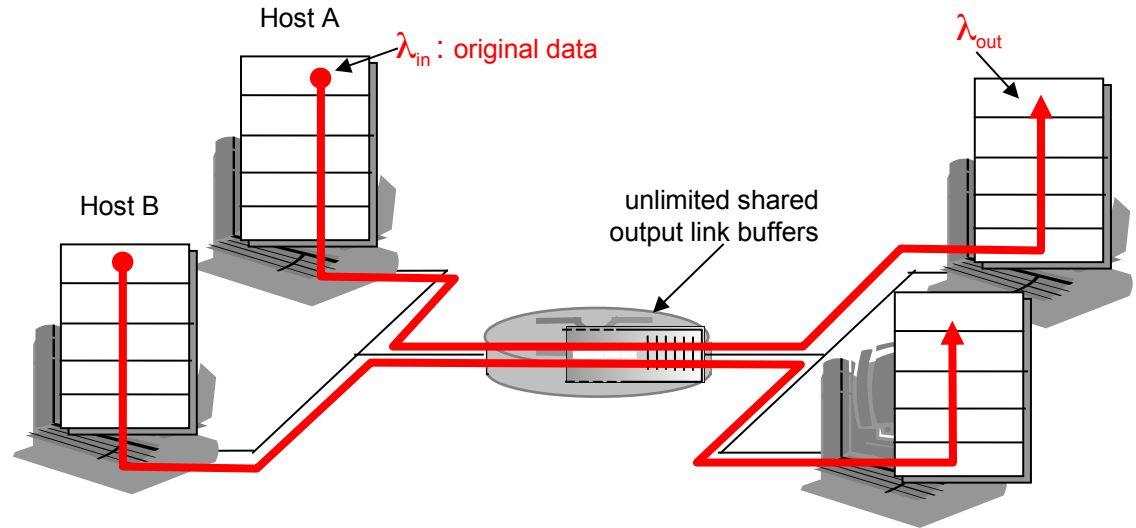
Principles of Congestion Control

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 1

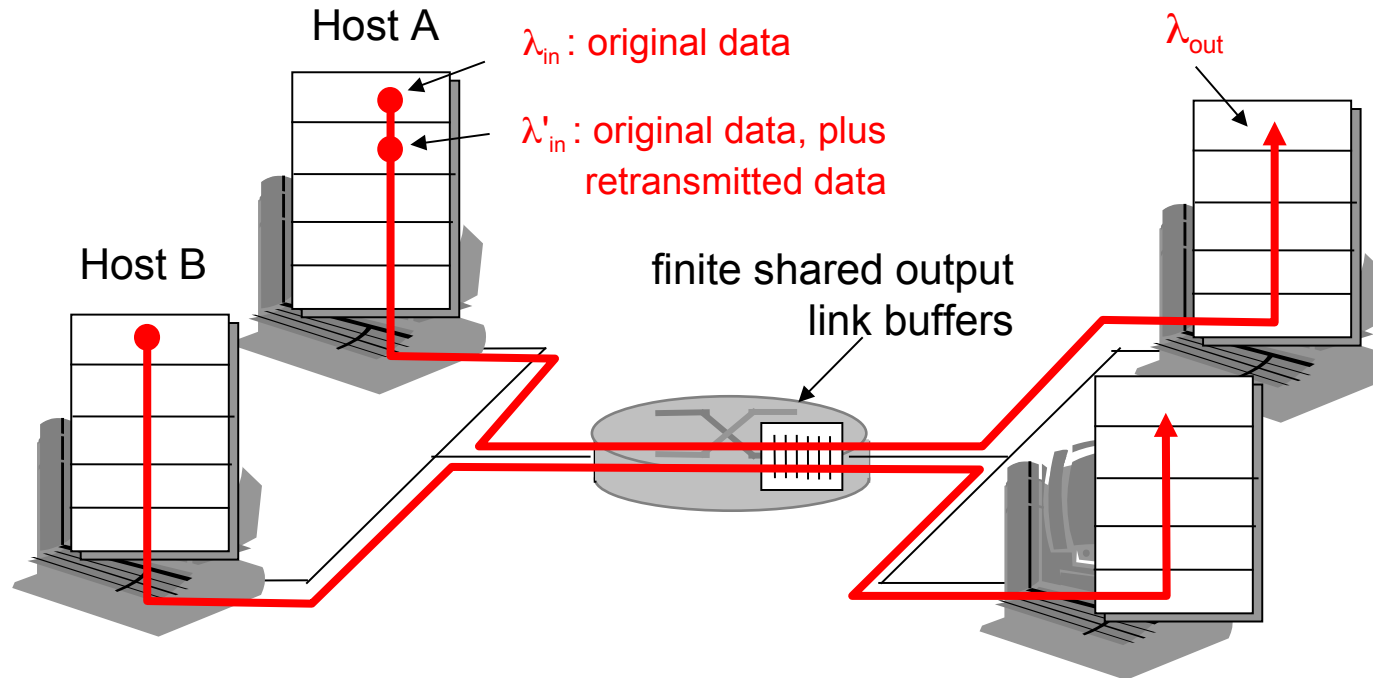
- two senders, two receivers
- one router, infinite buffers
- no retransmission



- large delays when congested
- maximum achievable throughput

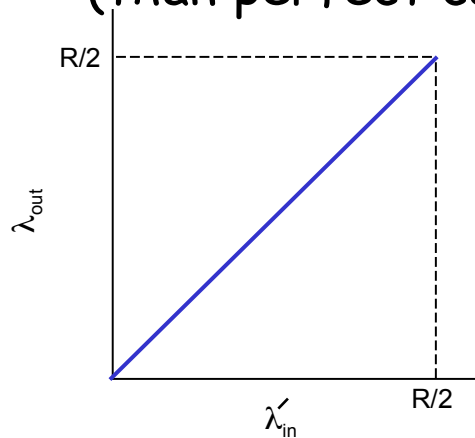
Causes/costs of congestion: scenario 2

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

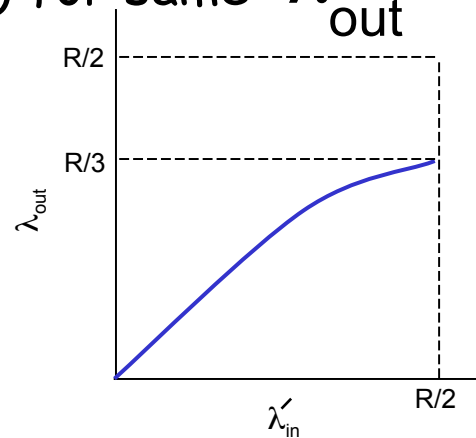


Causes/costs of congestion: scenario 2

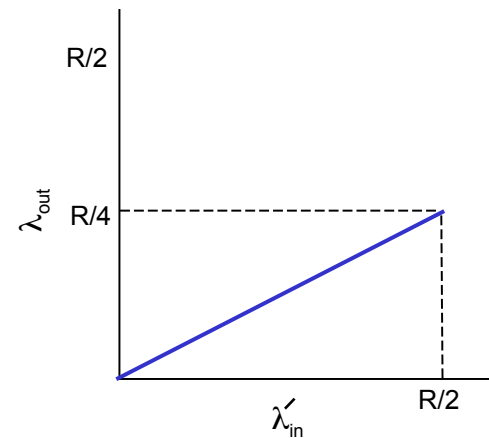
- a) always: $\lambda_{in} = \lambda_{out}$ (goodput)
- b) "perfect" retransmission only when loss: $\lambda'_{in} > \lambda_{out}$
- c) retransmission of delayed (not lost) packet makes λ'_{in} larger (than perfect case) for same λ_{out}



a.



b.



c.

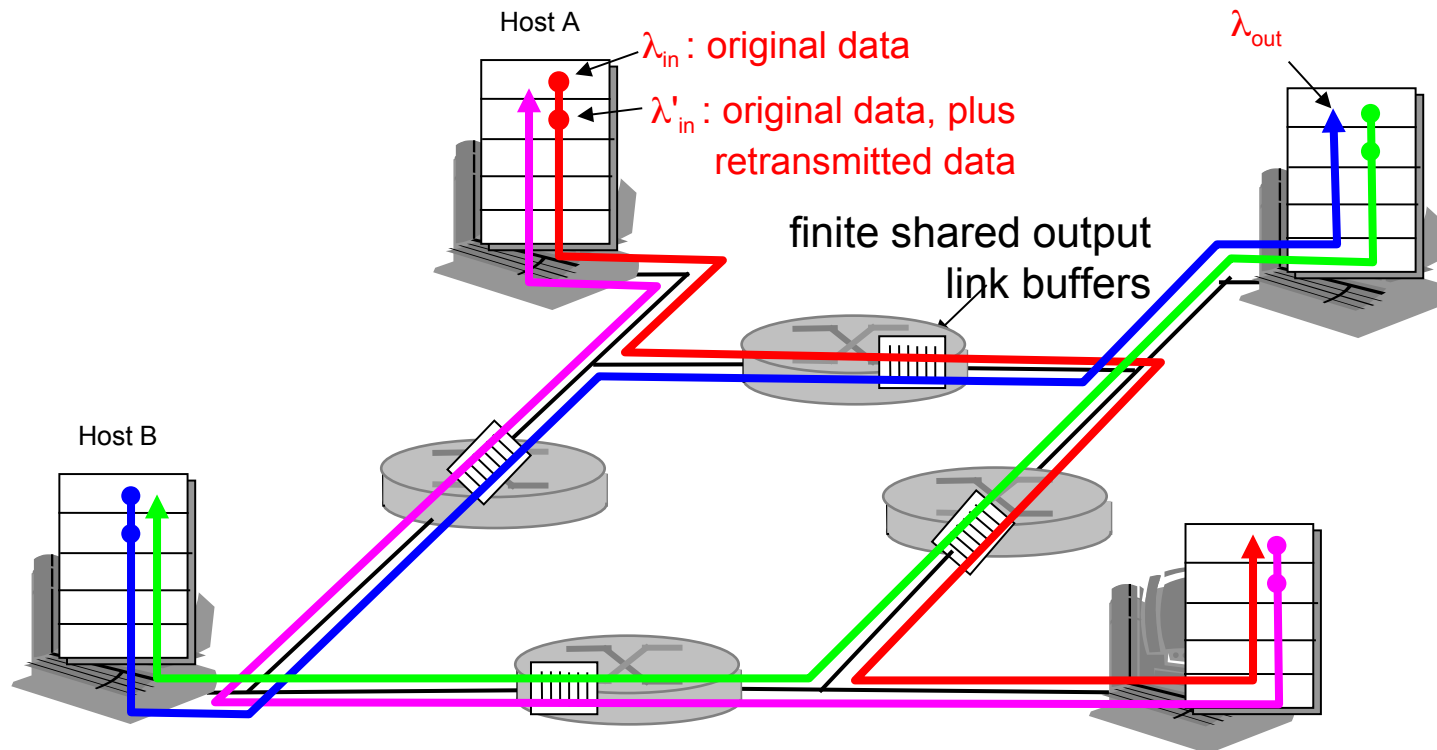
"costs" of congestion:

- more work (retransmissions) for given "goodput"
- unneeded retransmissions: link carries multiple copies of packet

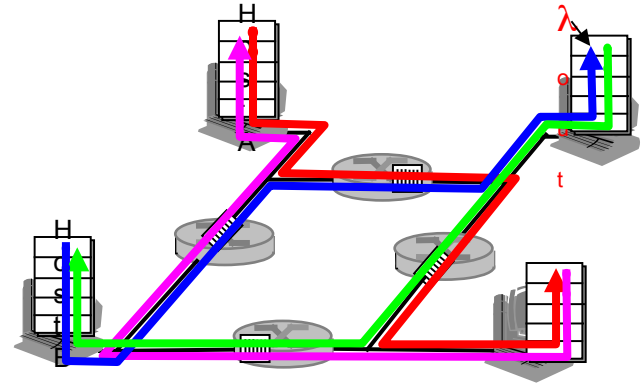
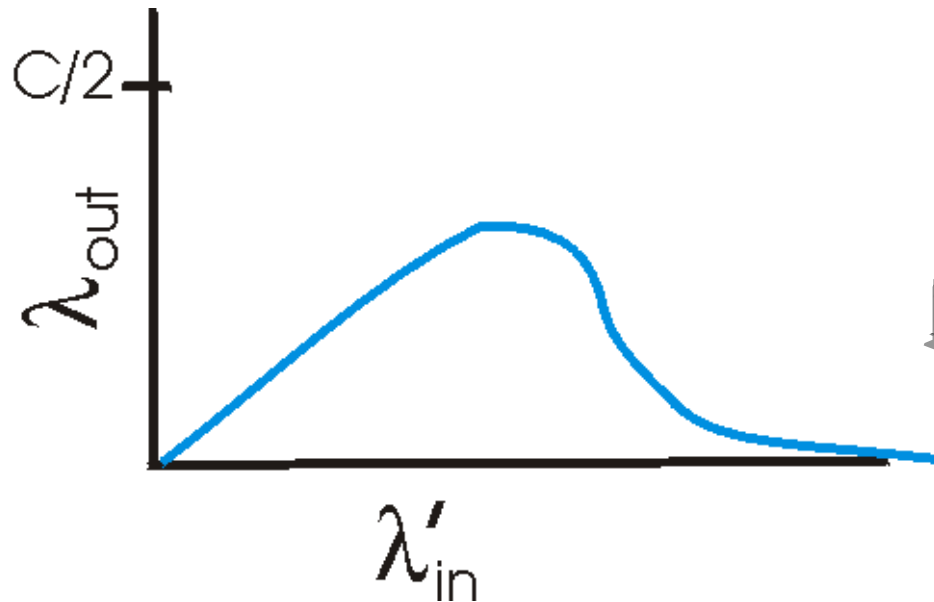
Causes/costs of congestion: scenario 3

- four senders
- multihop paths
- timeout/retransmit

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



Causes/costs of congestion: scenario 3



Another "cost" of congestion:

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

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

Network-assisted congestion control

example: ATM ABR congestion control

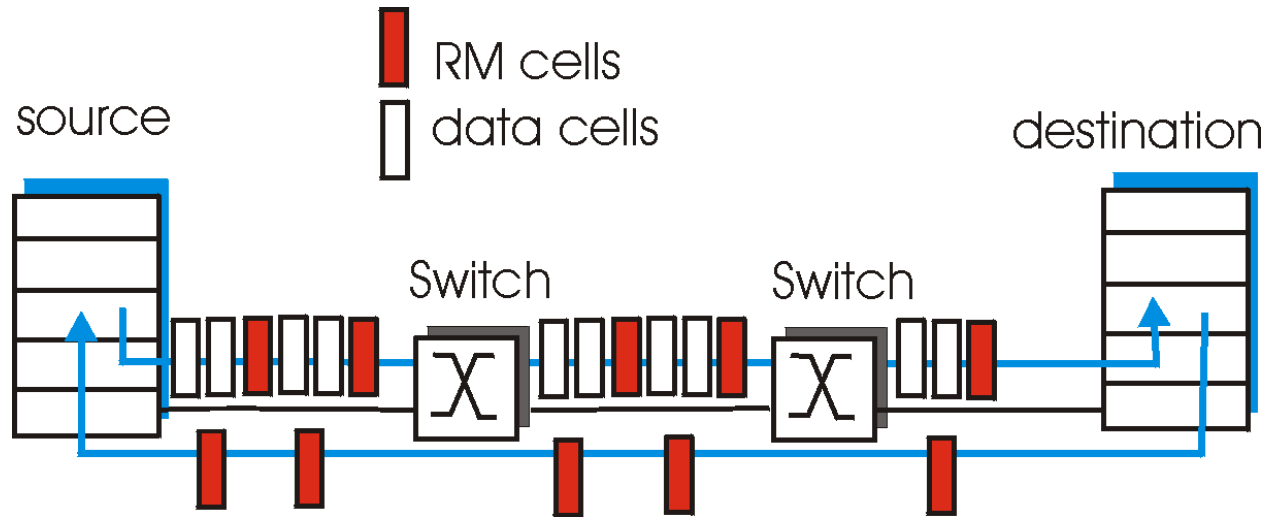
ABR: available bit rate:

- "elastic service"
- if sender's path "underloaded":
 - sender should use available bandwidth
- if sender's path congested:
 - sender throttled to minimum guaranteed rate

RM (resource management) cells:

- sent by sender, interspersed with data cells (every 32 data cells)
- bits in RM cell set by switches ("*network-assisted*")
 - NI bit: no increase in rate (mild congestion)
 - CI bit: congestion indication
- RM cells returned to sender by receiver, with NI and CI bits intact

Case study: ATM ABR congestion control



- two-byte **ER (explicit rate)** field in RM cell
 - congested switch may lower ER value in cell
 - sender's send rate at minimum supportable rate on path
- **EFCI (explicit forward congestion indication)** bit in **data cells**: set to 1 in congested switch
 - if data cell preceding RM cell has EFCI set, receiver sets CI bit in returned RM cell

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 Congestion Control

- end-end control (no network assistance)
- If the sender perceives that there is little congestion, the sender increases its send rate
- If the sender perceives that there is a congestion, the sender reduces its send rate

This approaches raises three questions:

- How does the sender limits its sent rate?
- Hoe does the sender perceives that there is a congestion?
- What algorithm should the sender use to change its send rate as a function of perceived end-to-end congestion?

TCP Congestion Control(CONT)

How does sender limits send rate?

- The amount of unACKed data at sender may not exceed $\min(\text{CongWin}, \text{RcvWindow})$; i.e.,

$$\text{LastByteSent} - \text{LastByteAcked}$$

$$\leq \min(\text{CongWin}, \text{RcvWindow})$$

- To focus on congestion control, ignore flow control.

Assume $\text{rcvWindow} > \text{CongWin}$

- Roughly,

$$\text{rate} = \frac{\text{CongWin}}{\text{RTT}} \text{ Bytes/sec}$$

- Sender adjusts send rate by adjusting CongWin

TCP Congestion Control(cont)

How does sender perceive congestion?

- loss event = timeout or 3 duplicate acks
- TCP sender reduces rate (CongWin) after loss event

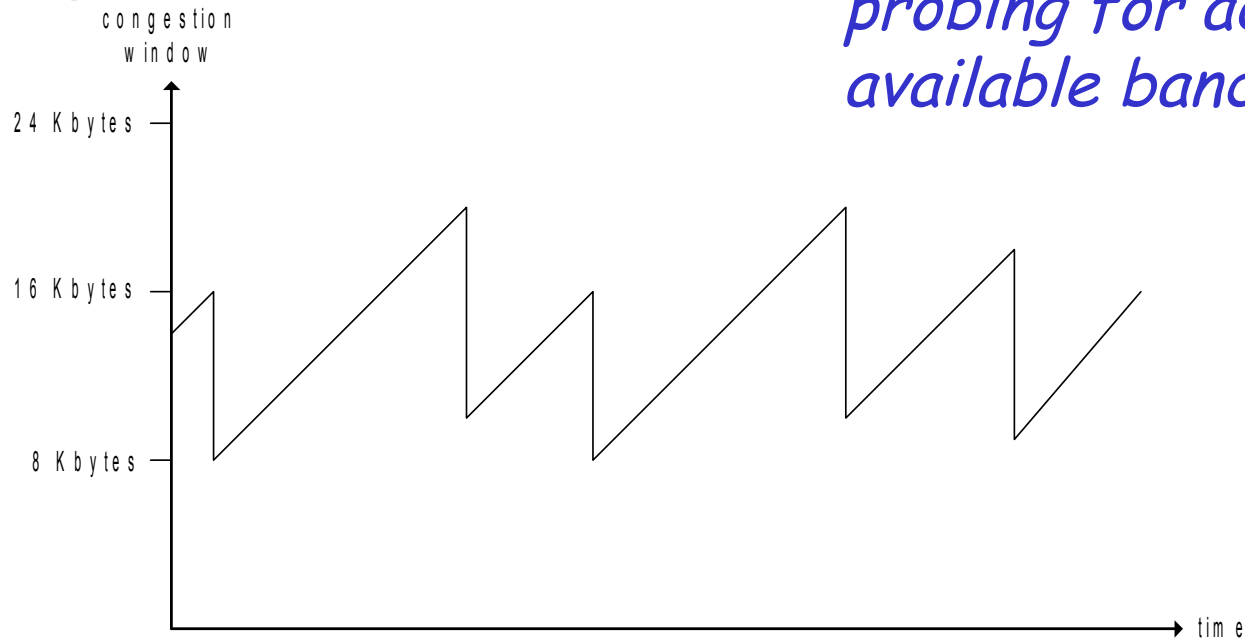
What algorithm is used to regulate send rate?

- Three components of the TCP congestion control algorithm
 - Additive-increase, multiplicative-decrease (AIMD)
 - slow start
 - Reaction to timeout events

TCP AIMD

multiplicative decrease:

- cut CongWin in half after loss event
- CongWin is not allowed to drop below 1 maximum segment size (MSS)



additive increase: increase

CongWin by 1 maximum segment size (MSS) every RTT in the absence of loss events:
probing for additional available bandwidth

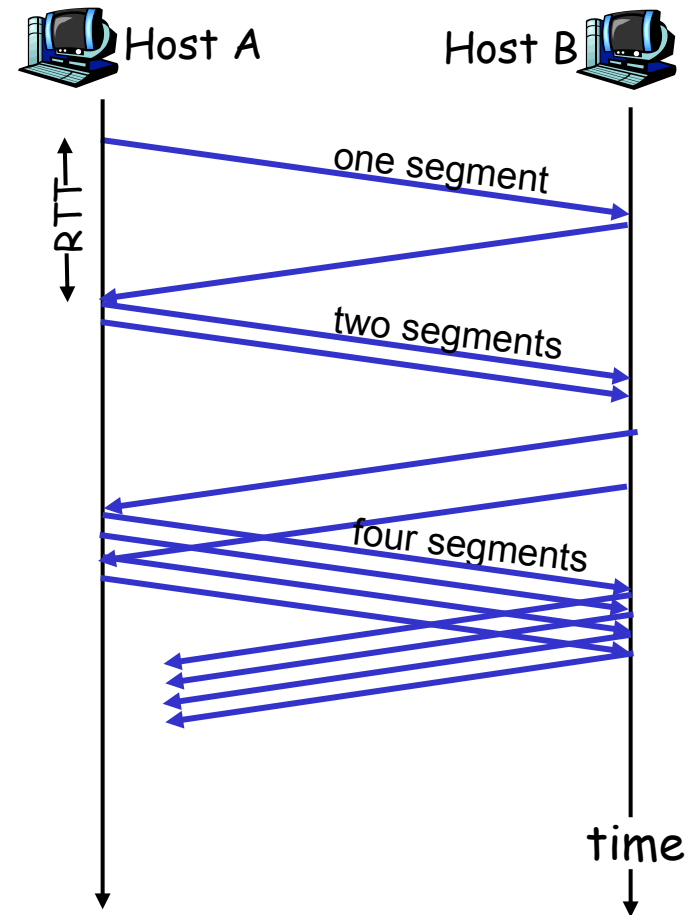
Long-lived TCP connection

TCP Slow Start

- When connection begins, $\text{CongWin} = 1 \text{ MSS}$ (slow start)
 - Example: $\text{MSS} = 500 \text{ bytes}$ & $\text{RTT} = 200 \text{ msec}$
 - initial rate = $(500 \times 8) / (200 \times 10^{-3}) = 20 \text{ kbps}$
- available bandwidth may be $\gg \text{MSS}/\text{RTT}$
 - desirable to quickly ramp up to respectable rate
- When connection begins, increase rate exponentially fast until first loss event: *Double its value of CongWin every RTT*
- When the first lost event occurs, CongWin is cut in half and then increases linearly

TCP Slow Start (more)

- When connection begins, increase rate exponentially until first loss event:
 - double CongWin every RTT
 - done by incrementing CongWin for every ACK received
- Summary: initial rate is slow but ramps up exponentially fast



Reaction to timeout events

- After 3 dup ACKs:
 - CongWin is cut in half
 - window then grows linearly
- But after timeout event:
 - CongWin set to 1 MSS;
 - window then grows exponentially
 - to a **threshold**, then grows linearly

Philosophy:

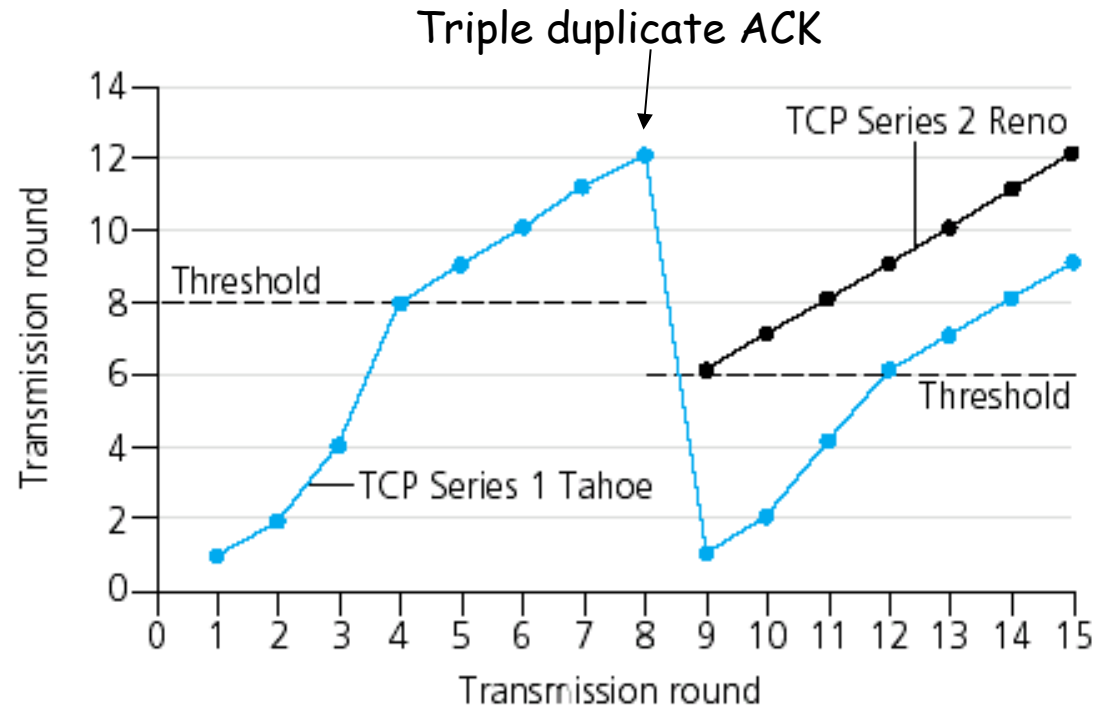
- 3 dup ACKs indicates network capable of delivering some segments
- timeout before 3 dup ACKs is “more alarming”

Reaction to timeout events (more)

- Q:** When should the exponential increase switch to linear?
- A:** When CongWin gets to 1/2 of its value before timeout.

Implementation:

- Maintains a variable.
Threshold
- At loss event, Threshold is set to 1/2 of CongWin just before loss event



- TCP Tahoe sets its congestion window to 1 MSS and enters slow start phase after either type of loss event
- TCP Reno cancels the slow start phase of TCP Tahoe after a triple duplicate ACK (**fast recovery**)

Summary: TCP Congestion Control

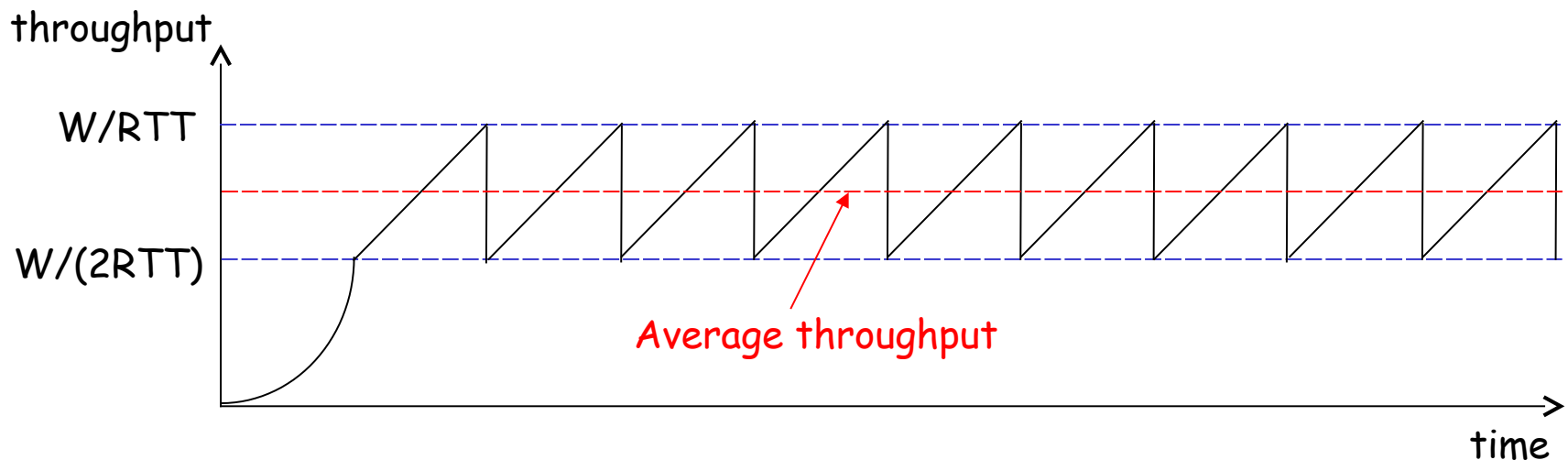
- When CongWin is below Threshold, sender in **slow-start** phase, window grows exponentially.
- When CongWin is above Threshold, sender is in **congestion-avoidance** phase, window grows linearly.
- When a **triple duplicate ACK** occurs, Threshold set to $\text{CongWin}/2$ and CongWin set to Threshold.
- When **timeout** occurs, Threshold set to $\text{CongWin}/2$ and CongWin is set to 1 MSS.

TCP sender congestion control

Event	State	TCP Sender Action	Commentary
ACK receipt for previously unacked data	Slow Start (SS)	$\text{CongWin} = \text{CongWin} + \text{MSS}$, If ($\text{CongWin} > \text{Threshold}$) set state to "Congestion Avoidance"	Resulting in a doubling of CongWin every RTT
ACK receipt for previously unacked data	Congestion Avoidance (CA)	$\text{CongWin} = \text{CongWin} + \text{MSS} * (\text{MSS} / \text{CongWin})$	Additive increase, resulting in increase of CongWin by 1 MSS every RTT
Loss event detected by triple duplicate ACK	SS or CA	$\text{Threshold} = \text{CongWin} / 2$, $\text{CongWin} = \text{Threshold}$, Set state to "Congestion Avoidance"	Fast recovery, implementing multiplicative decrease. CongWin will not drop below 1 MSS.
Timeout	SS or CA	$\text{Threshold} = \text{CongWin} / 2$, $\text{CongWin} = 1 \text{ MSS}$, Set state to "Slow Start"	Enter slow start
Duplicate ACK	SS or CA	Increment duplicate ACK count for segment being acked	CongWin and Threshold not changed

Average throughput

- W : the window size when a lost event occurs
- RTT: round-trip time
- Assume W and RTT are approximately constant
- Ignore the slow start phases that occur after timeout events



$$\text{Average throughput} = (0.75W)/RTT$$

TCP Futures

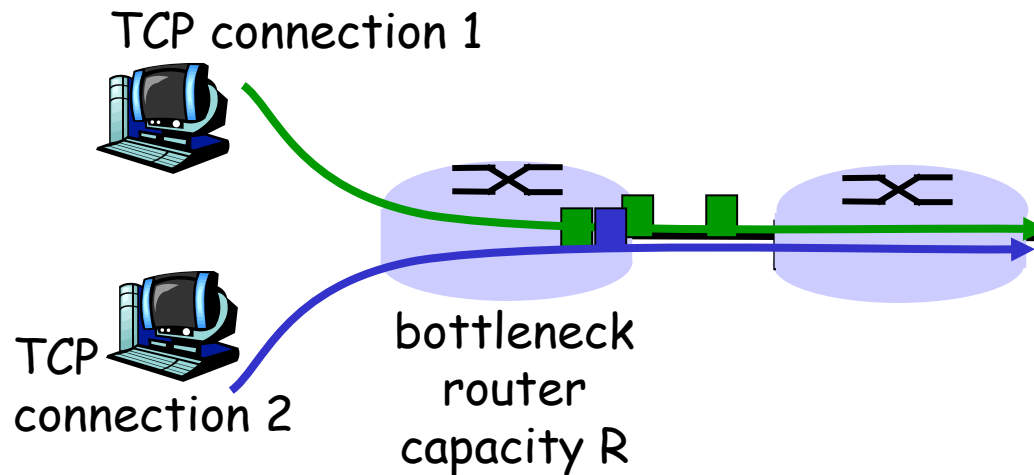
- Example: 1500 byte segments, 100ms RTT, want 10 Gbps throughput
- Requires window size $W = 83,333$ in-flight segments
- Throughput in terms of loss rate:

$$\frac{1.22 \cdot MSS}{RTT \sqrt{L}}$$

- Segment loss probability : $\rightarrow L = 2 \cdot 10^{-10}$ *Wow*
- New versions of TCP for high-speed needed!

TCP Fairness

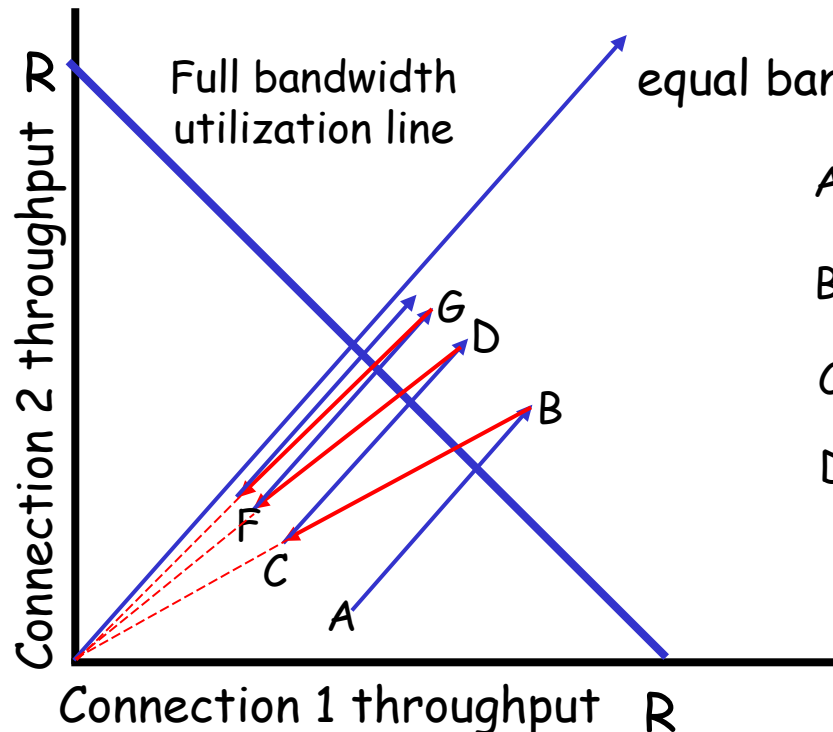
Fairness goal: if K TCP sessions share same bottleneck link of bandwidth R , each should have average rate of R/K



Why is TCP fair?

Two competing sessions:

- Additive increase gives slope of 1, as throughput increases
- multiplicative decrease decreases throughput by factor of 2



- A: congestion avoidance: both connections additive increase
- B: loss: both connections decrease window by factor of 2
- C: congestion avoidance: both connections additive increase
- D: loss: both connections decrease window by factor of 2

⋮
Finally, fluctuates along the equal bandwidth share line

Fairness (more)

Fairness and UDP

- Multimedia apps often do not use TCP
 - do not want rate throttled by congestion control
- Instead use UDP:
 - pump audio/video at constant rate, tolerate packet loss
- Research area: congestion control for UDP

Fairness and parallel TCP connections

- nothing prevents app from opening parallel connections between 2 hosts.
- Web browsers do this
- Example: link of rate R supporting 9 connections:
 - new app asks for 1 TCP, gets rate $R/10$
 - new app asks for 11 TCPs, gets $R/2$

Delay modeling

Q: How long does it take to receive an object from a Web server after sending a request?

Ignoring congestion, delay is influenced by:

- TCP connection establishment
- data transmission delay
- TCP slow start

Notation, assumptions:

- Assume one link between client and server of rate R
- S : MSS (bits)
- O : object size (bits)
- no retransmissions (no loss, no corruption)

Window size:

- First assume: fixed congestion window, W segments
- Then dynamic window, modeling slow start

Fixed congestion window (1)

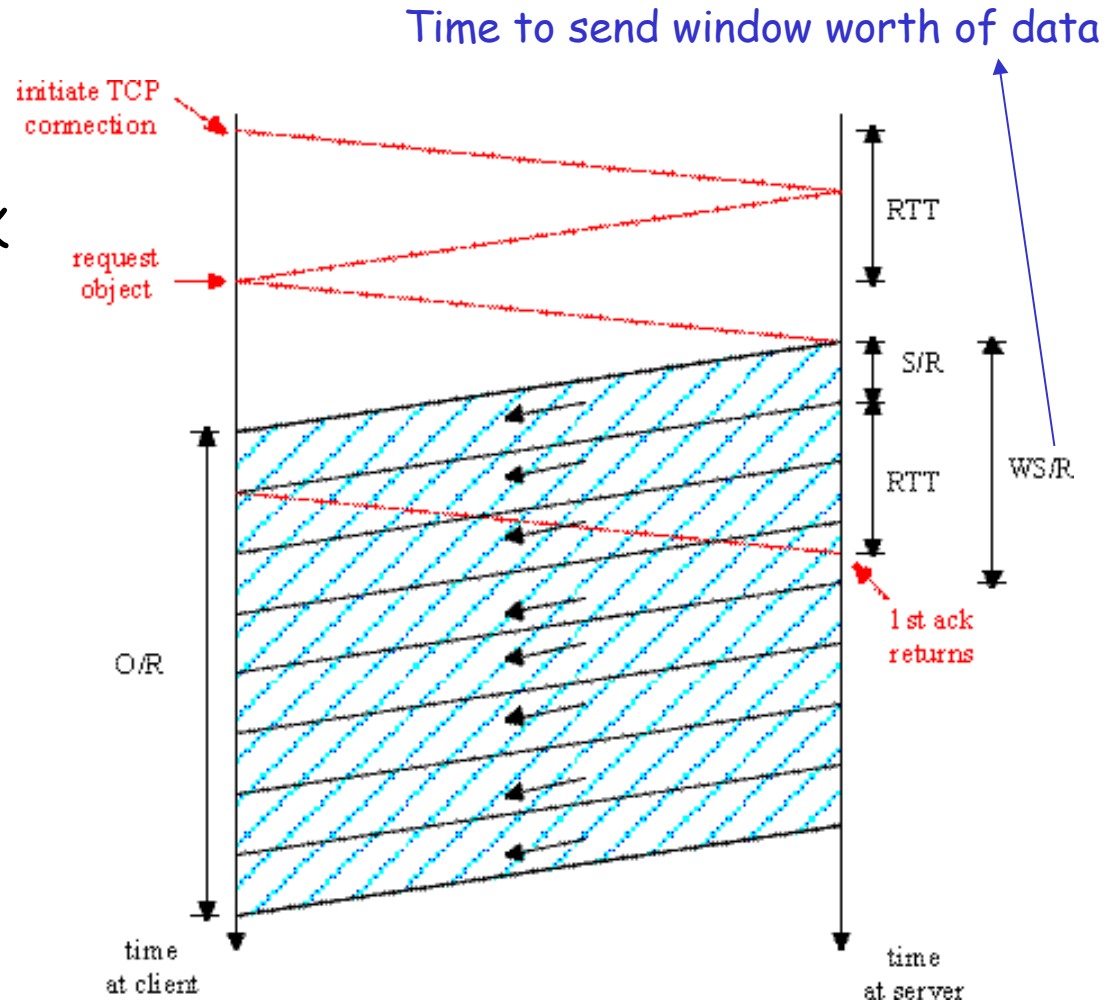
Two cases:

First case:

- $WS/R \leq RTT + S/R$: ACK for first segment in window returns before window's worth of data is sent

Second case:

- $WS/R > RTT + S/R$: wait for ACK after sending window's worth of data



Fixed congestion window (2)

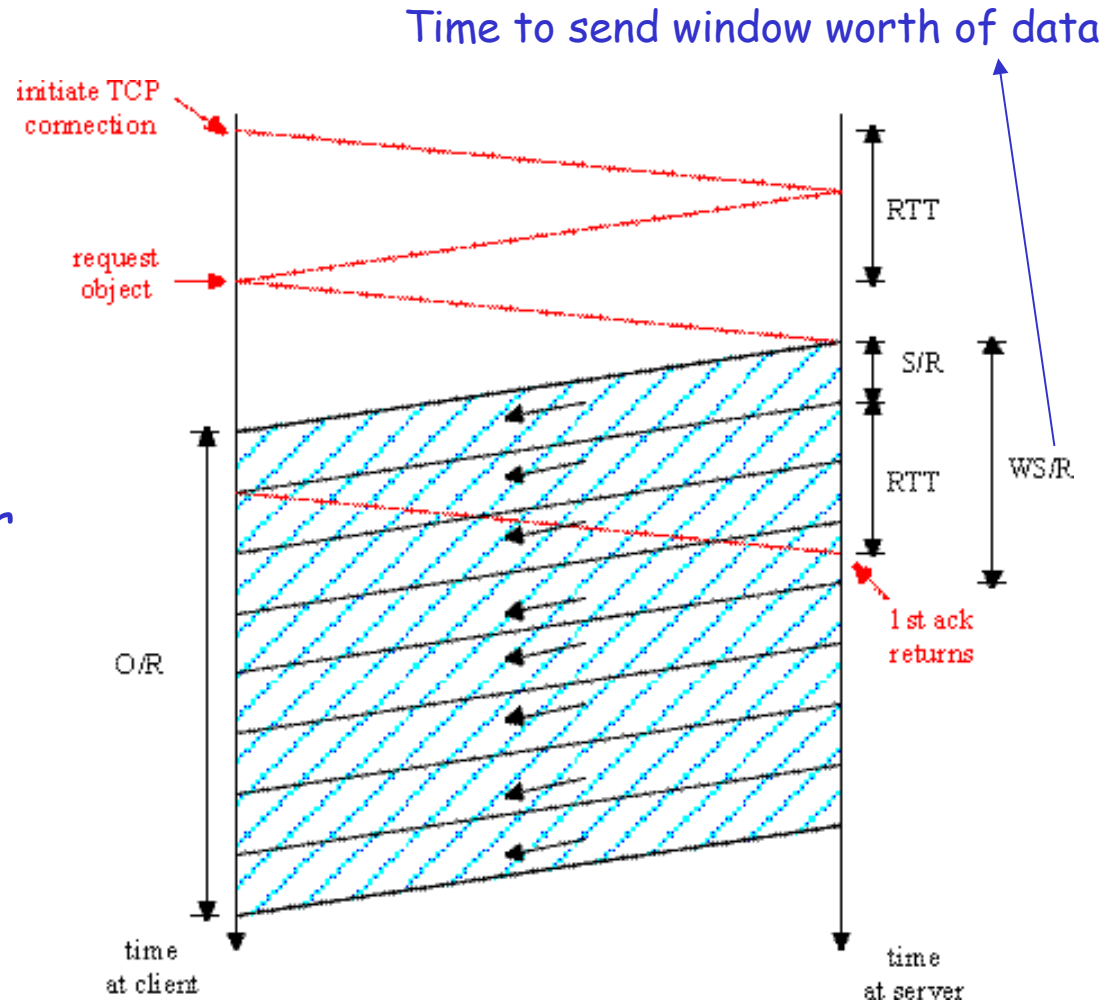
First case:

WS/R RTT + S/R: ACK
for first segment in
window returns before
window's worth of data
is sent

Time for sending one
window's data Time for
the first ACK to
return

The server can send data
without stopping

$$\text{delay} = 2\text{RTT} + \text{O/R}$$



Fixed congestion window (3)

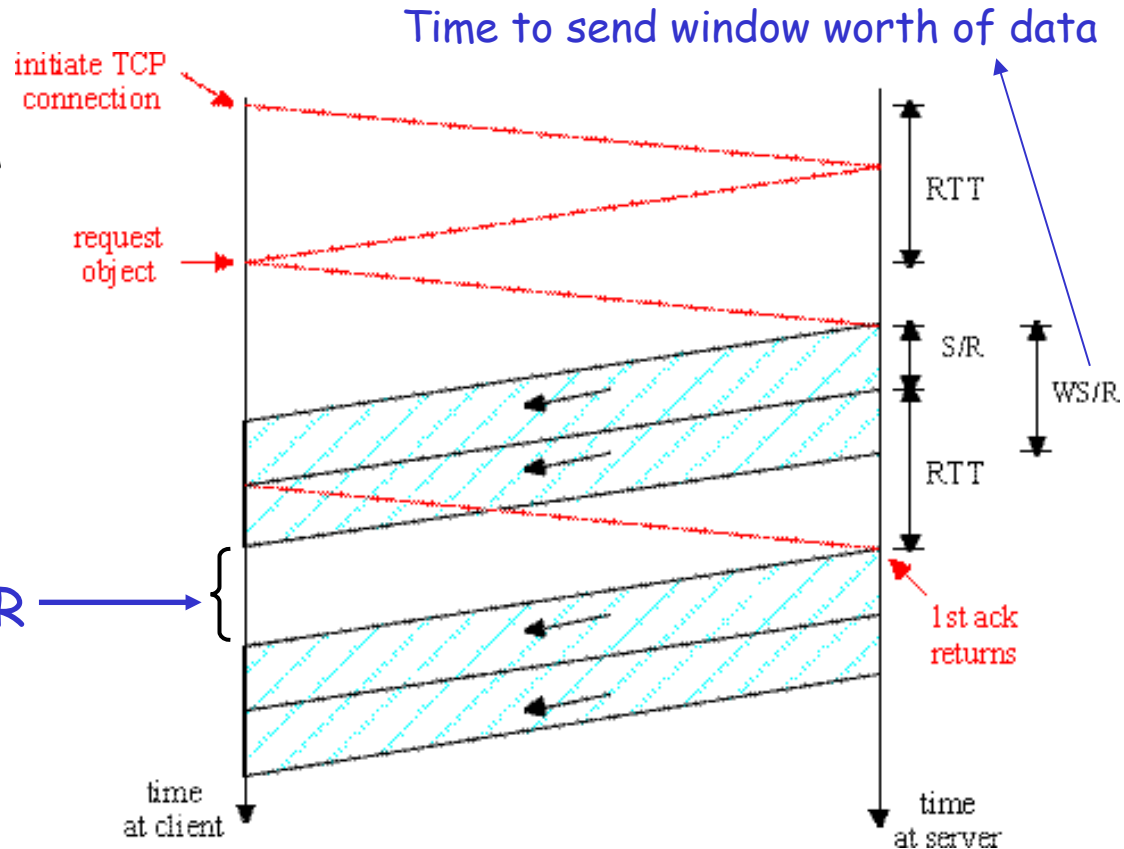
Second case:

- $WS/R < RTT + S/R$: wait for ACK after sending window's worth of data is sent

Time for sending one window's data < Time for the first ACK to return

- $K = O/WS$: the number of windows that cover the object

$$S/R + RTT - WS/R$$



$$\text{delay} = 2RTT + O/R + (K-1)[S/R + RTT - WS/R]$$

TCP Delay Modeling: Slow Start (1)

Now suppose window grows according to slow start

Will show that the delay for one object is:

$$Latency = 2RTT + \frac{O}{R} + P \left[RTT + \frac{S}{R} \right] - (2^P - 1) \frac{S}{R}$$

where P is the number of times TCP idles at server:

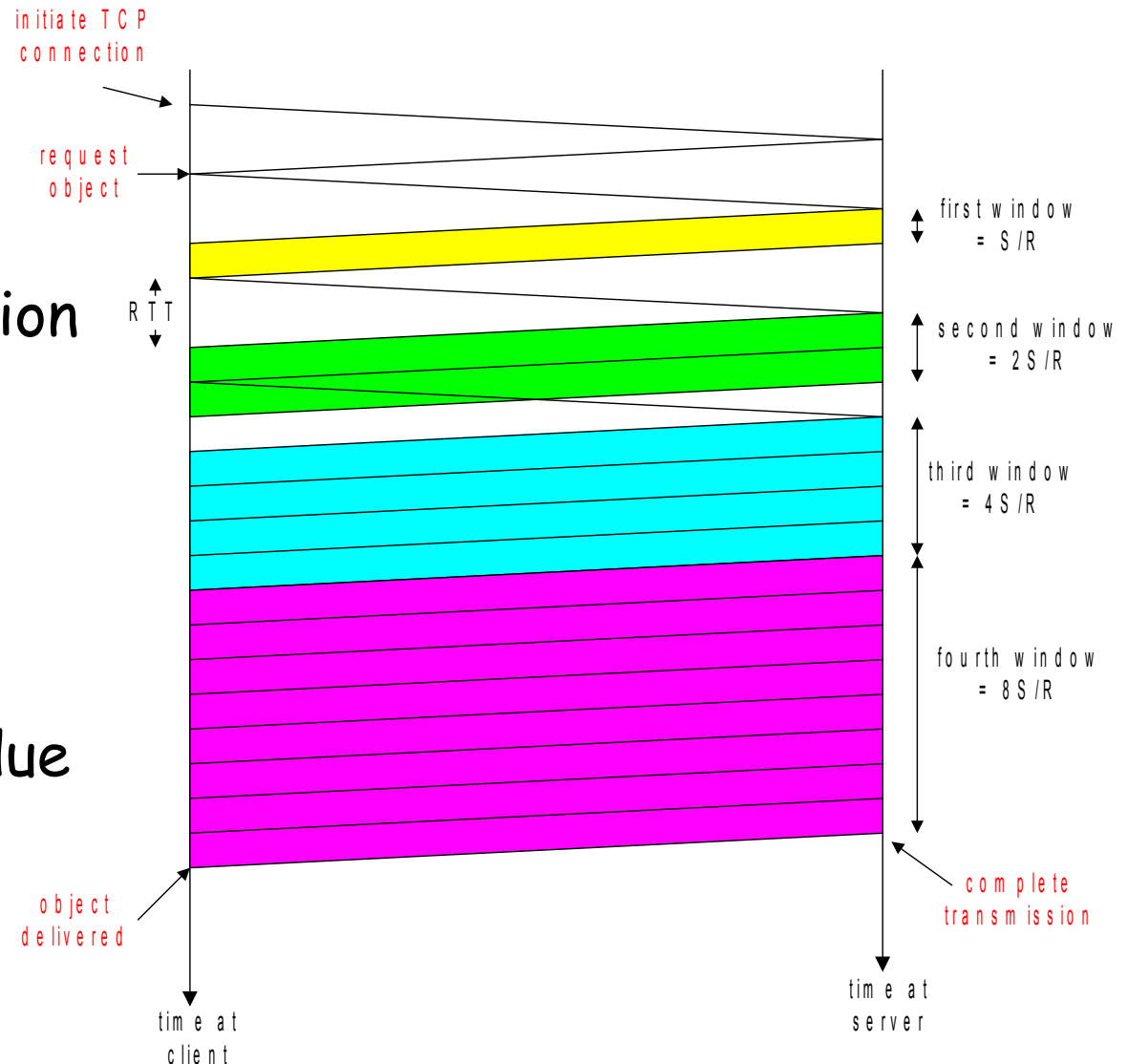
$$P = \min\{Q, K - 1\}$$

- where Q is the number of times the server idles if the object were of infinite size.
- and K is the number of windows that cover the object.

TCP Delay Modeling: Slow Start (2)

Delay components:

- ❑ 2 RTT for connection estab and request
- ❑ O/R to transmit object
- ❑ time server idles due to slow start



TCP Delay Modeling (3)

$\frac{S}{R} + RTT$ = time from when server starts to send segment
until server receives acknowledgement

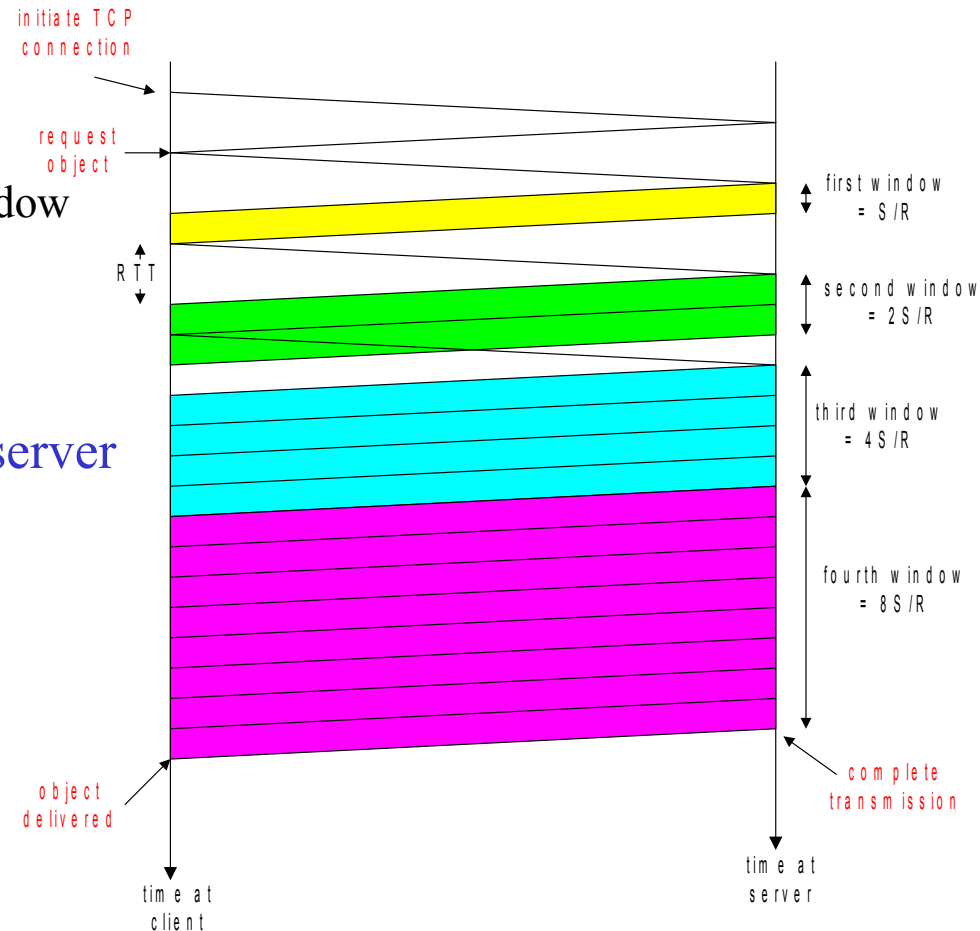
$2^{k-1} \frac{S}{R}$ = time to transmit the k th window

$\left[\frac{S}{R} + RTT - 2^{k-1} \frac{S}{R} \right]^+$ = idle time after the k th window

O/R = Time to transmit the object

P : the number of times TCP idles at server

$$\begin{aligned} \text{delay} &= \frac{O}{R} + 2RTT + \sum_{p=1}^P \text{idleTime}_p \\ &= \frac{O}{R} + 2RTT + \sum_{k=1}^P \left[\frac{S}{R} + RTT - 2^{k-1} \frac{S}{R} \right] \\ &= \frac{O}{R} + 2RTT + P \left[RTT + \frac{S}{R} \right] - (2^P - 1) \frac{S}{R} \end{aligned}$$



TCP Delay Modeling (4)

How to calculate P?

P : the number of times TCP idles at server

Q : the number of times the server idles if the object were of infinite size

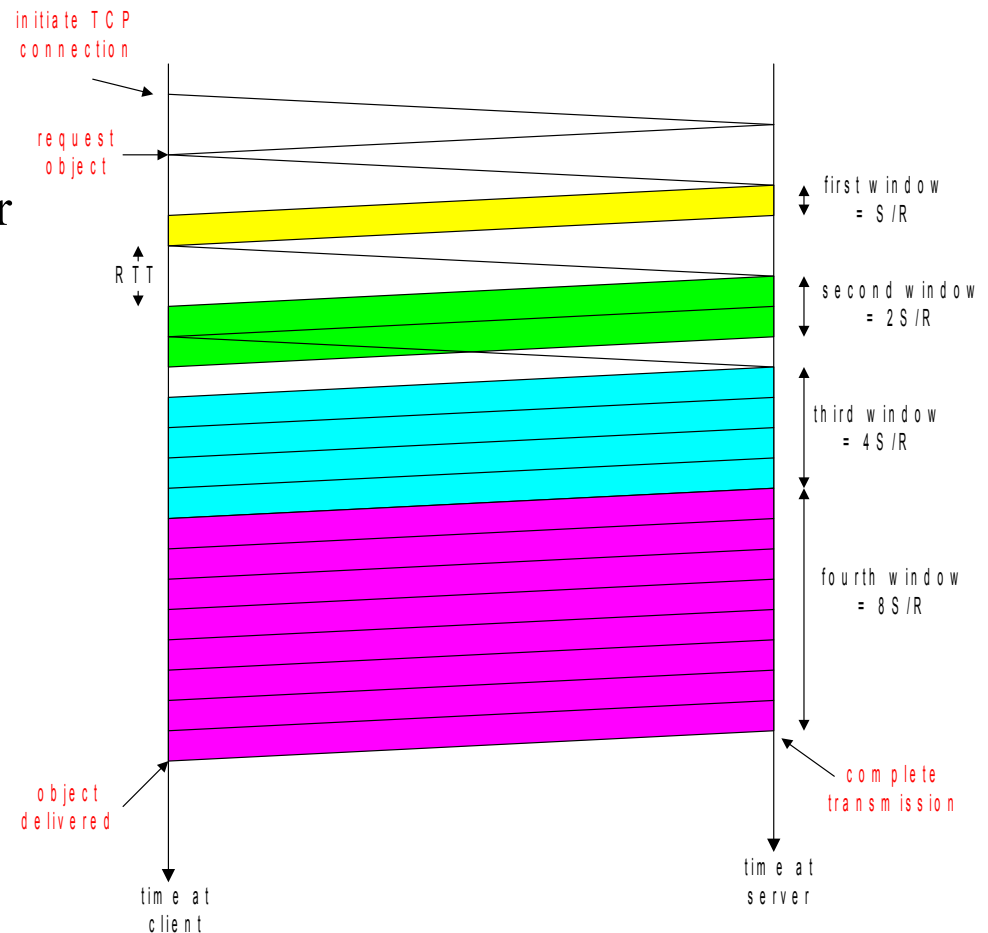
K : the number of windows that cover the object.

$$P = \min\{Q, K - 1\}$$

Example:

- $O/S = 15$ segments
- $K = 4$ windows
- $Q = 2$
- $P = \min\{K-1, Q\} = 2$

Server idles $P=2$ times



TCP Delay Modeling (5)

Recall K = number of windows that cover object

How do we calculate K ?

S : MSS (bits)

O : object size (bits)

$$\begin{aligned} K &= \min\{k : 2^0 S + 2^1 S + \cdots + 2^{k-1} S \geq O\} \\ &= \min\{k : 2^0 + 2^1 + \cdots + 2^{k-1} \geq O/S\} \\ &= \min\{k : 2^k - 1 \geq \frac{O}{S}\} \\ &= \min\{k : k \geq \log_2(\frac{O}{S} + 1)\} \\ &= \left\lceil \log_2(\frac{O}{S} + 1) \right\rceil \end{aligned}$$

TCP Delay Modeling (6)

Recall Q = the number of times the server idles if the object were of infinite size

How do we calculate Q ?

$\left[\frac{S}{R} + RTT - 2^{k-1} \frac{S}{R} \right]^+ = \text{idle time after the } k\text{th window}$

$$Q = \max \left\{ k : RTT + \frac{TS}{R} - \frac{S}{R} 2^{k-1} \geq 0 \right\}$$

$$= \max \left\{ k : 2^{k-1} \leq T + \frac{RTT}{S/R} \right\}$$

$$= \max \left\{ k : k \leq \log_2 \left(T + \frac{RTT}{S/R} \right) + 1 \right\}$$

$$= \left\lfloor \log_2 \left(T + \frac{RTT}{S/R} \right) \right\rfloor + 1$$

TCP latency - congestion control vs. no congestion control

With congestion control:

$$Latency = 2RTT + \frac{O}{R} + P \left[RTT + \frac{S}{R} \right] - (2^P - 1) \frac{S}{R}$$

Without congestion control: no congestion window constraint

$$latency = 2RTT + O/R$$

← minimum latency

$$P \left[\frac{S}{R} \right] - (2^P - 1) \frac{S}{R} < 0$$

$$\frac{latency}{minimum\ latency} \leq 1 + \frac{P}{O/(R*RTT) + 2}$$

RTT << O/R or
small R*RTT (delay•bandwidth) ⇒ TCP slow start does not significantly increase latency

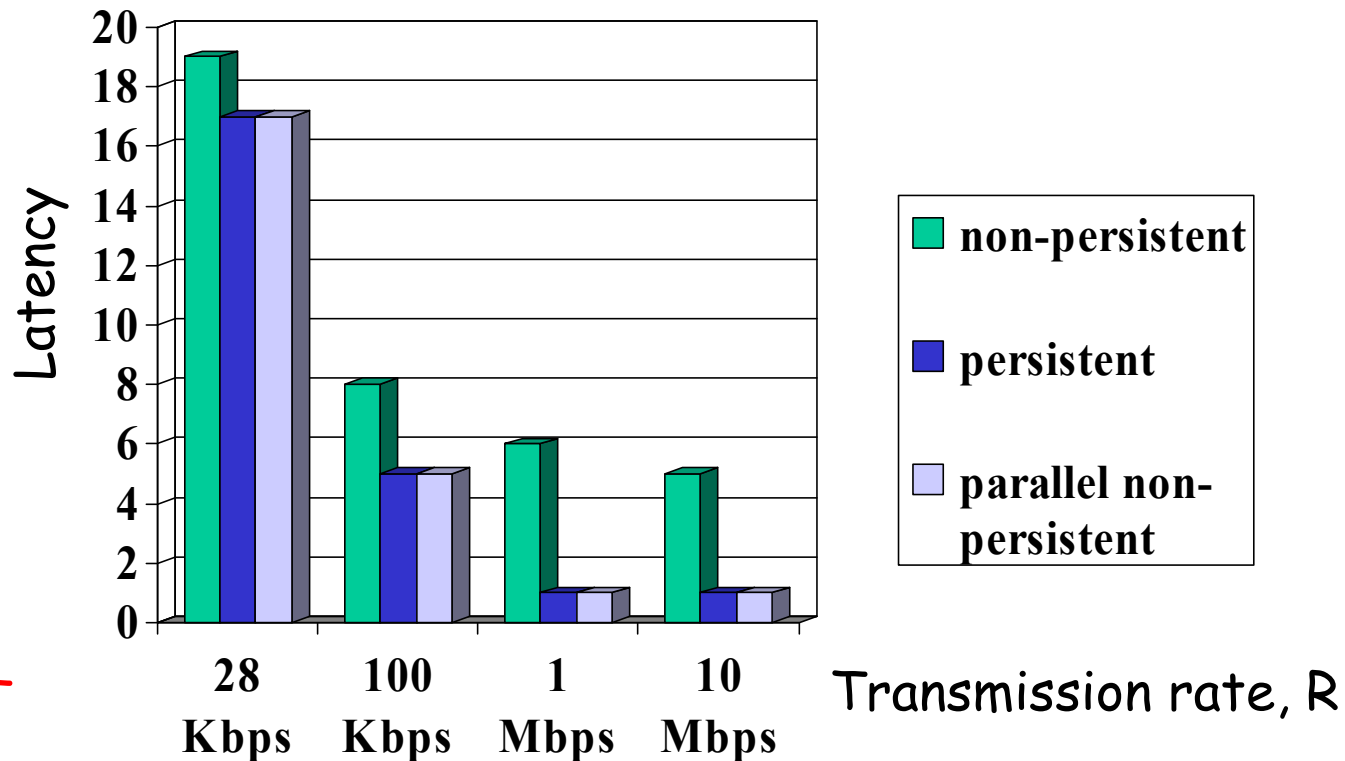
Otherwise; or
large R*RTT (delay•bandwidth) ⇒ TCP slow start significantly increases latency

HTTP Modeling

- Assume Web page consists of:
 - 1 base HTML page (of size O bits)
 - M images (each of size O bits)
- Non-persistent HTTP:
 - $M+1$ TCP connections in series
 - *Response time = $(M+1)O/R + (M+1)2RTT + \text{sum of idle times}$*
- Persistent HTTP:
 - $2 RTT$ to request and receive base HTML file
 - $1 RTT$ to request and receive M images
 - *Response time = $(M+1)O/R + 3RTT + \text{sum of idle times}$*
- Non-persistent HTTP with X parallel connections
 - Suppose M / X integer.
 - 1 TCP connection for base file
 - M / X sets of parallel connections for images.
 - *Response time = $(M+1)O/R + (M/X + 1)2RTT + \text{sum of idle times}$*

HTTP Response time (in seconds)

RTT = 100 msec. O = 5 Kbytes. M=10 and X=5



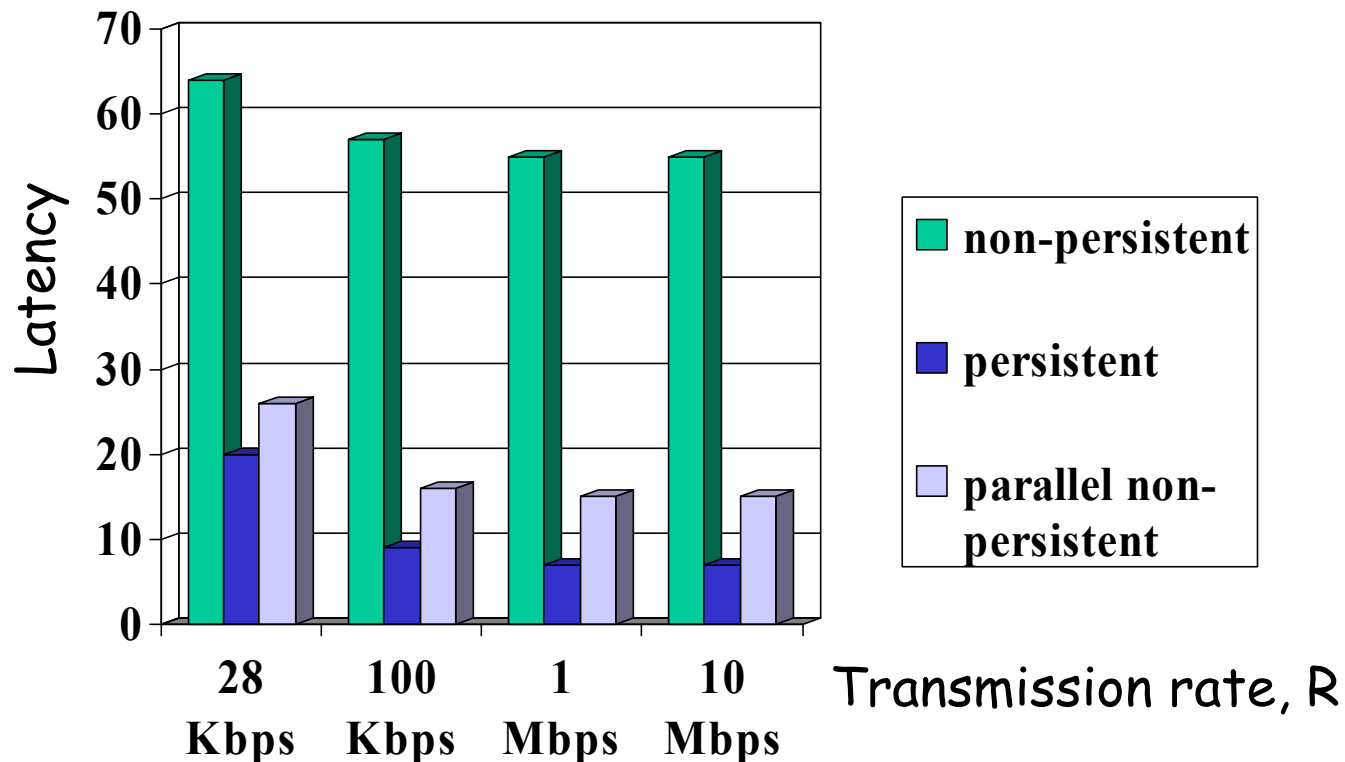
Small RTT

For low bandwidth (low $\text{delay} \cdot \text{bandwidth}$), connection & response time dominated by transmission time.

Persistent connections only give minor improvement over parallel connections.

HTTP Response time (in seconds)

RTT = 1 sec, O = 5 Kbytes, M=10 and X=5



For **larger RTT**, response time dominated by **TCP establishment & slow start delays**. Persistent connections now give important improvement: particularly in **high delay•bandwidth** networks.

Chapter 3: Summary

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

Next:

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