

# Transport Services and Protocols

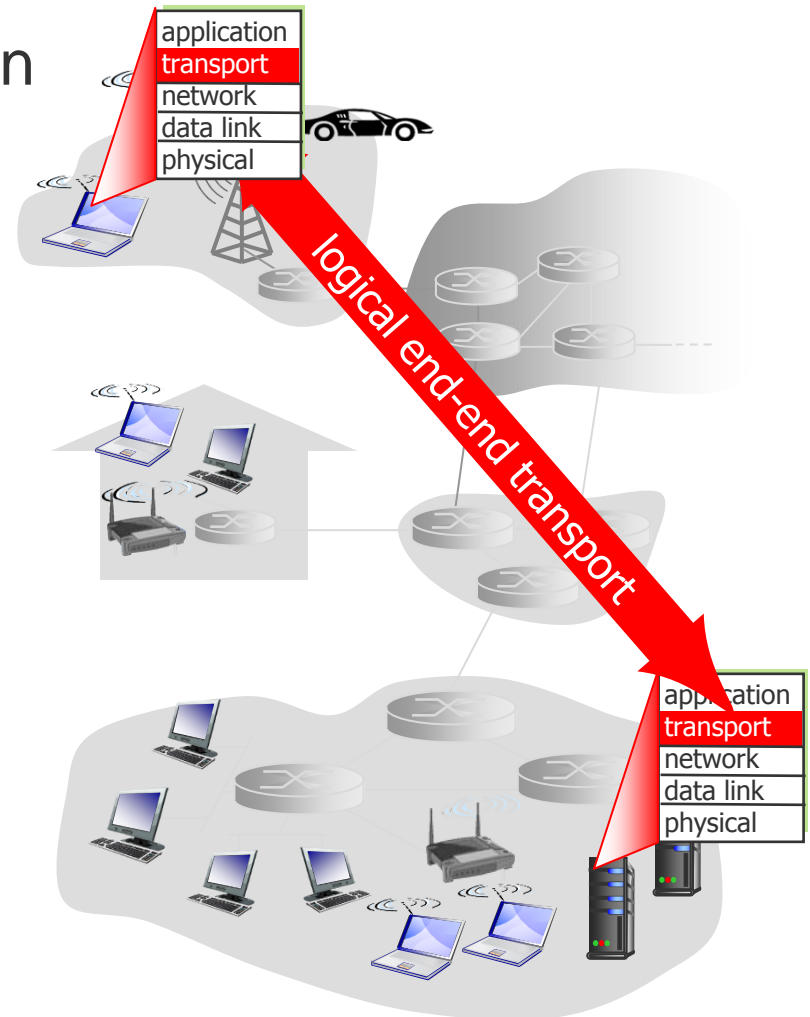
---

Andy Carpenter  
([Andy.Carpenter@manchester.ac.uk](mailto:Andy.Carpenter@manchester.ac.uk))

Elements these slides come from Kurose and Ross, authors of "Computer Networking: A Top-down Approach", and are copyright Kurose and Ross

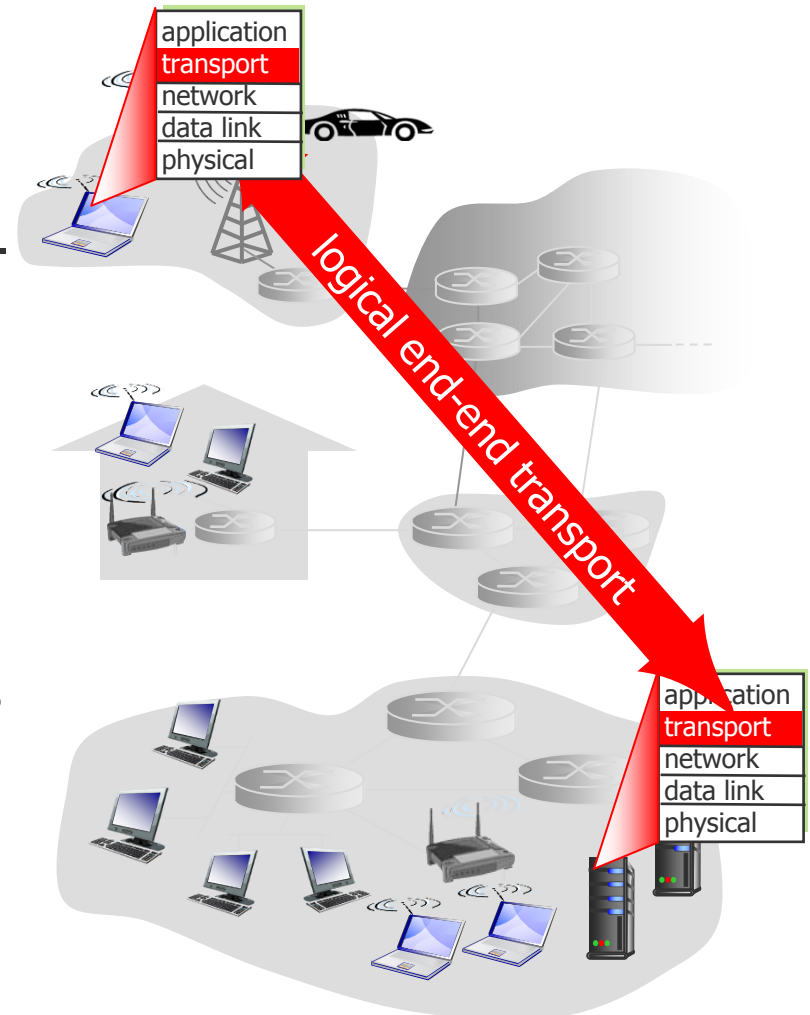
# Transport Services and Protocols

- Provides logical communication between applications
  - process-to-process
- Runs in end systems
  - send side: breaks app messages into segments, passes to network layer
  - receiver side: reassembles segments into messages, passes to app layer
- Multiplex/ id end-points
- Protocols available:
  - Internet: TCP, UDP



# Transport Service QoS Params

- Provides QoS to applications
- Options
  - Connectionless/connection-orientated
  - reliability
  - flow control
  - congestion control
- Alters underlying network QoS
- Services not available:
  - delay guarantees
  - bandwidth guarantees



# End-points (Ports/Sockets)

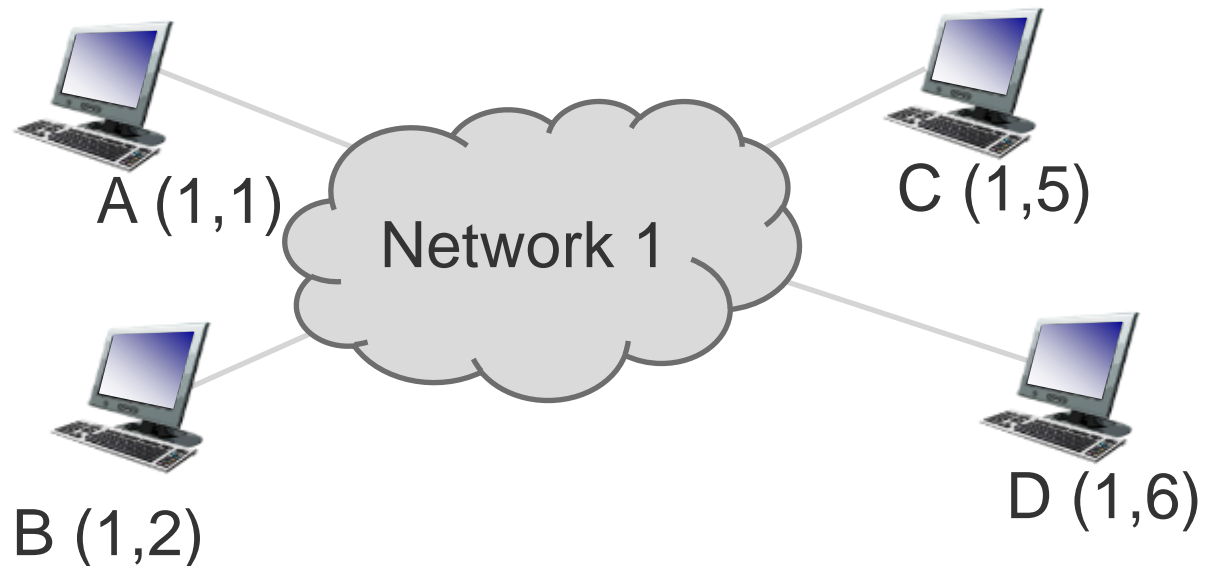
- Support multiple destinations
- Identify process via 16-bit identifier; port
- Applications use sockets which have associated port
- Well know services allocated ports 0-1024

Port	Mnemonic	Service
7	ECHO	Echo
20	FTP-DATA	FTP (default data)
21	FTP	FTP (control)
53	DOMAIN	Domain name service
80	HTTP	Hypertext Transfer protocol

Centrally  
allocated

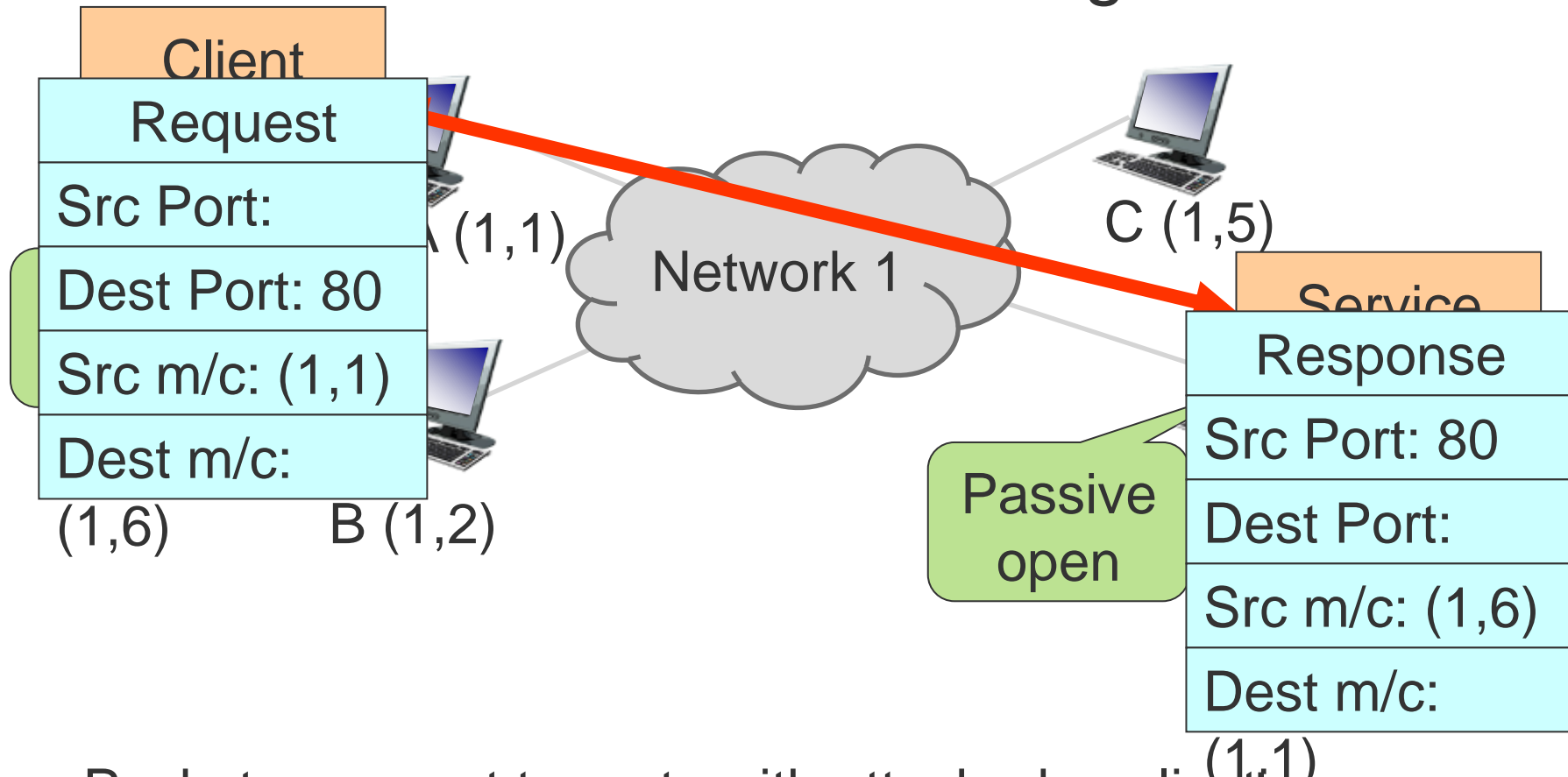
- Client ports are dynamically allocated by o/s

# End-Point: Connecting



- Packets are sent to ports with attached applications
- Data to other ports will generate an error message

# End-Point: Connecting



- Packets are sent to ports with attached applications
- Data to other ports will generate an error message

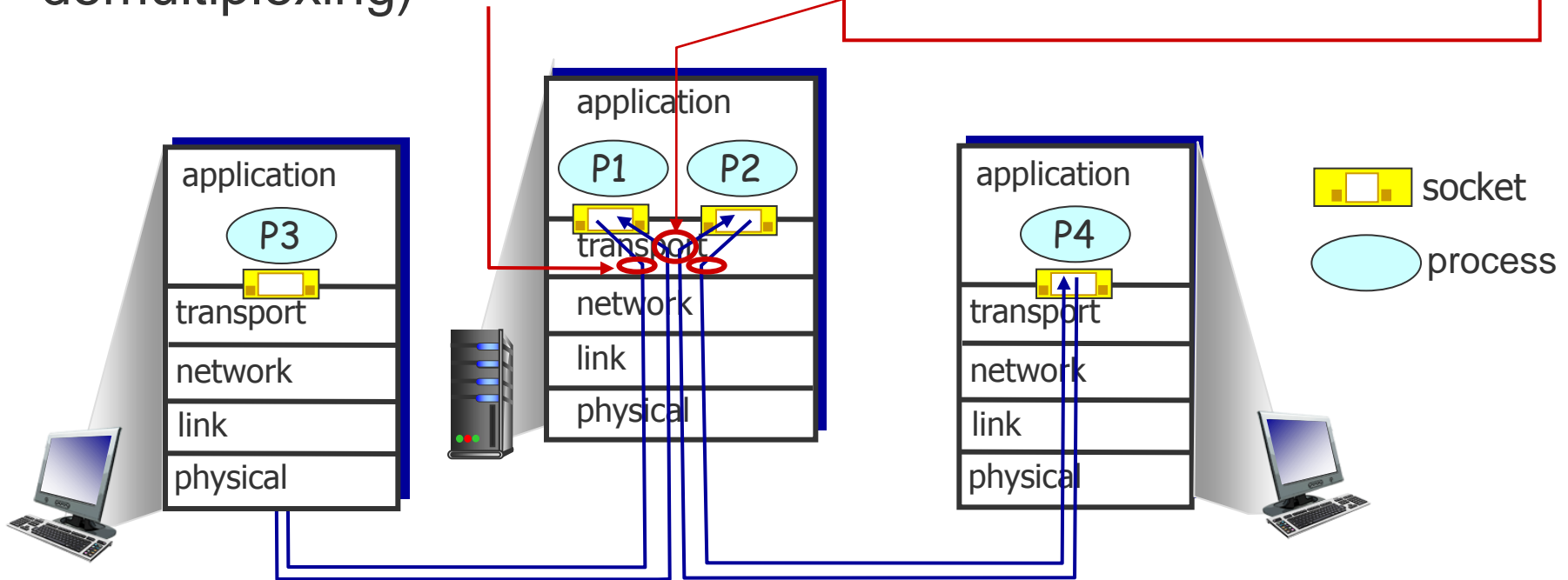
# End-Point: Multiplexing

## *multiplexing at sender:*

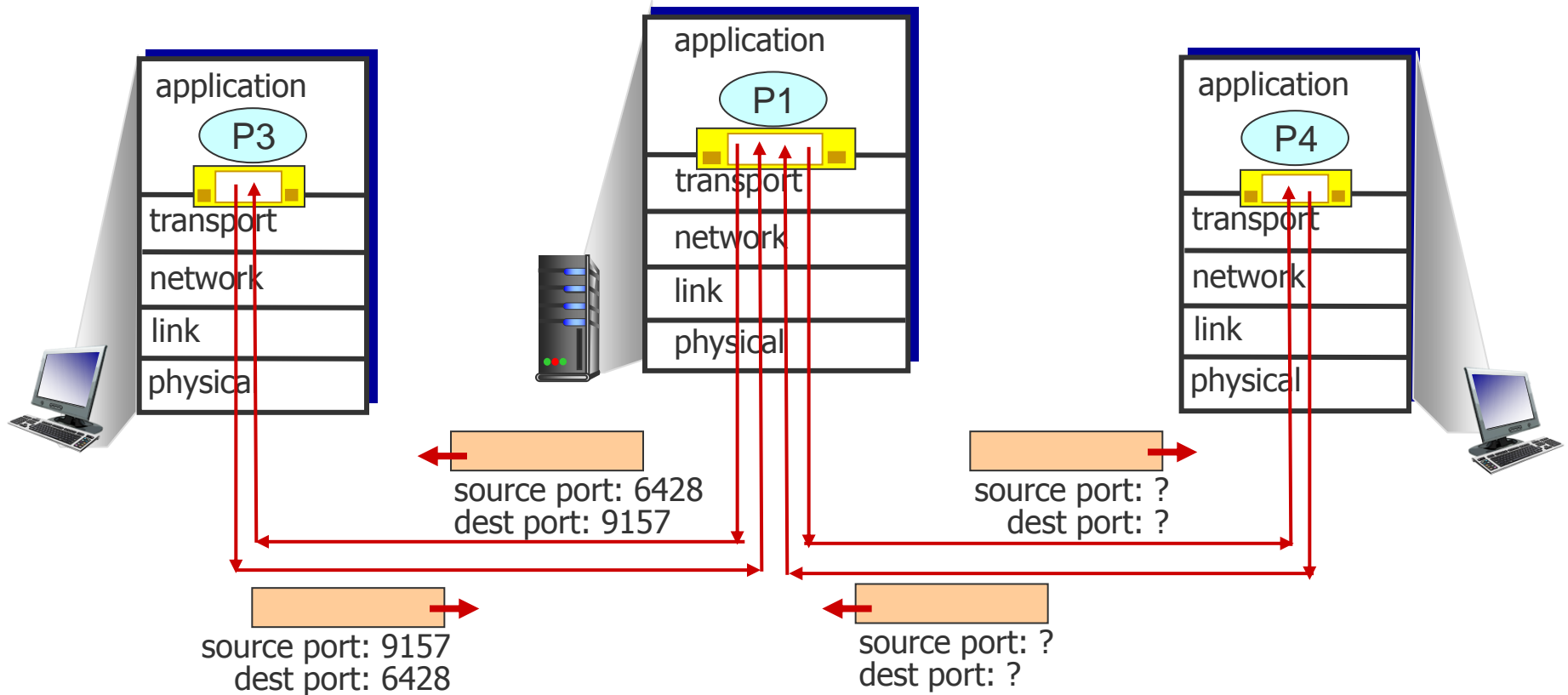
handle data from multiple sockets, add transport header (later used for demultiplexing)

## *demultiplexing at receiver:*

use header info to deliver received segments to correct socket

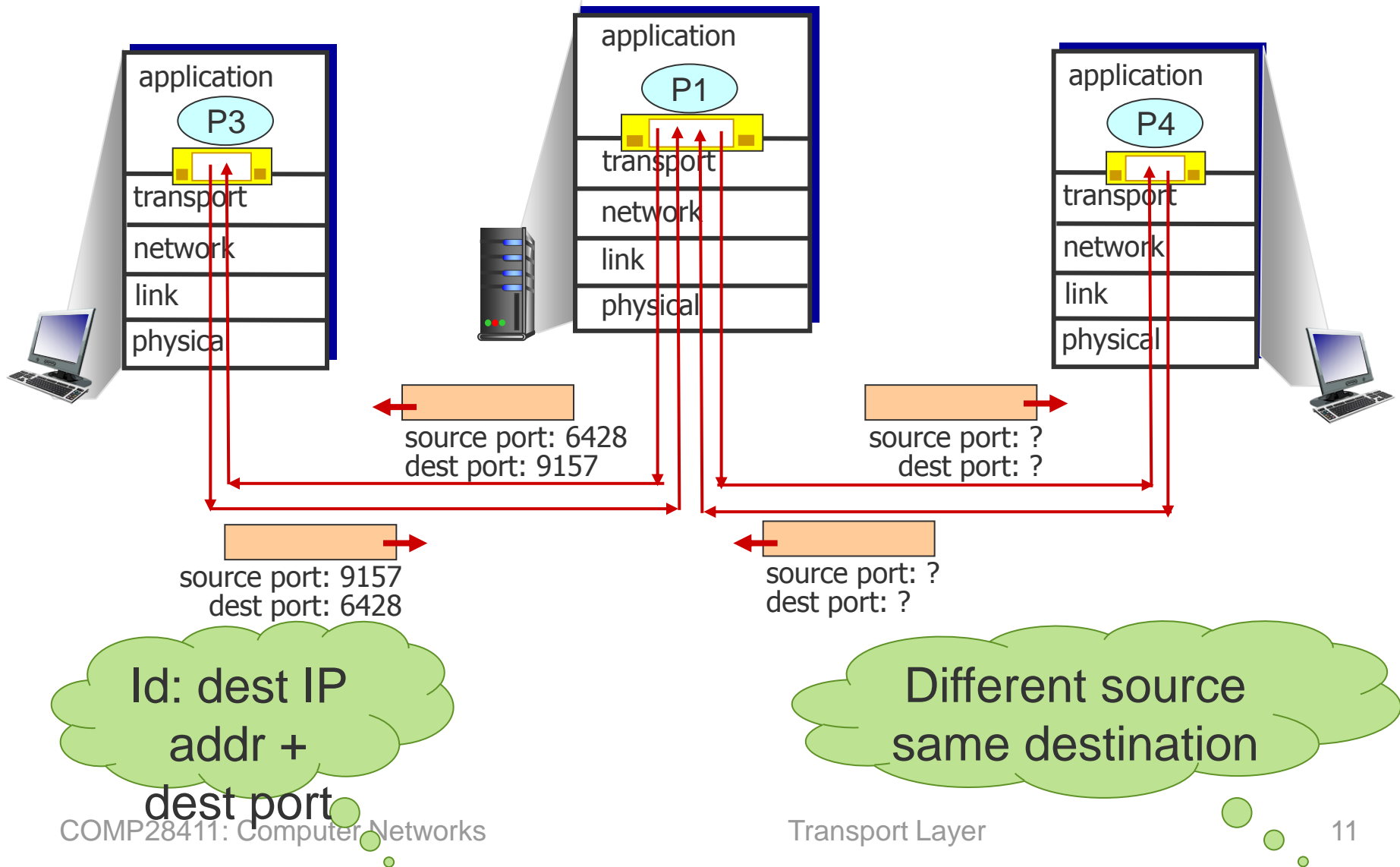


# End-Point: Id (Connectionless)

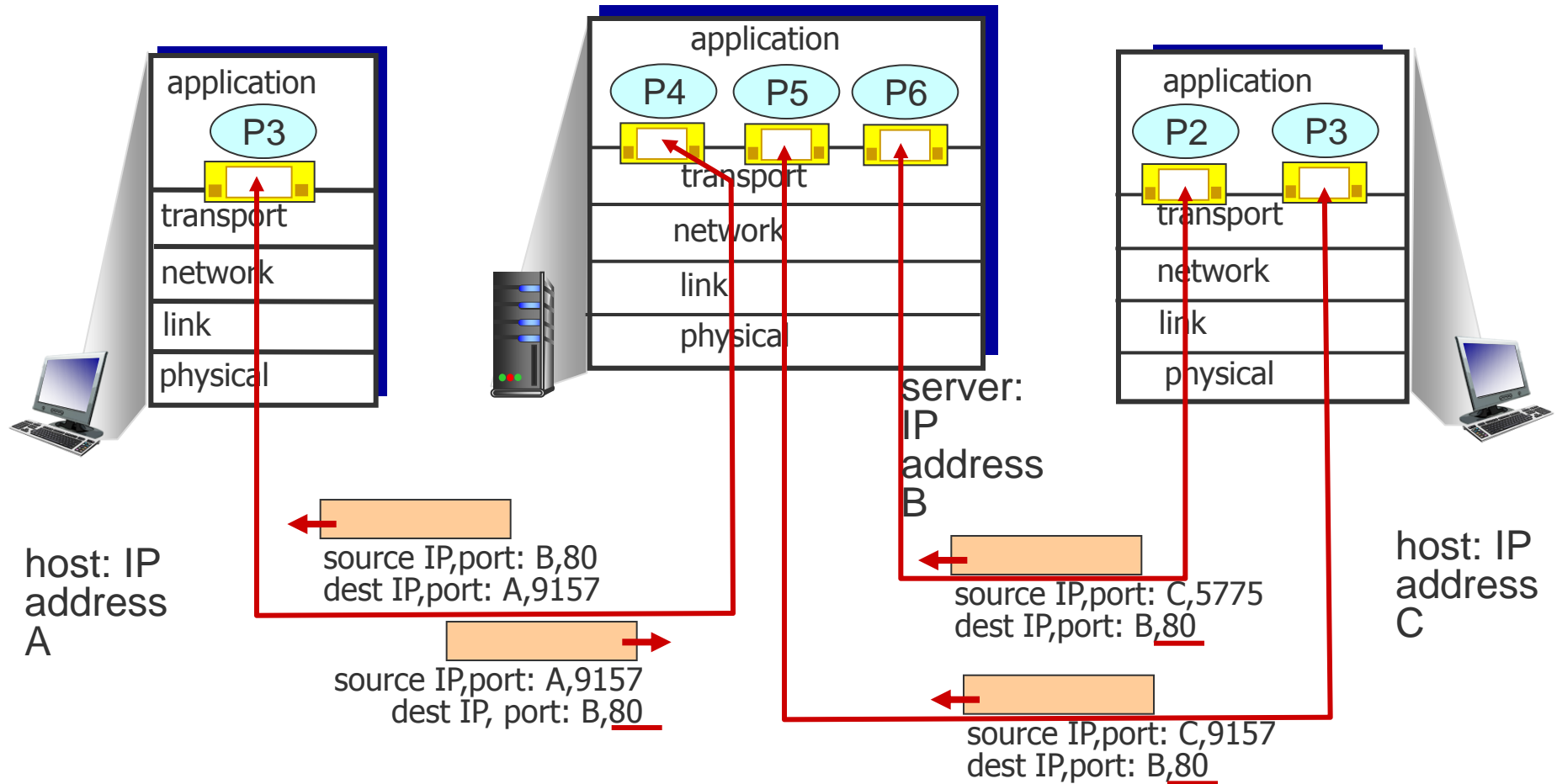




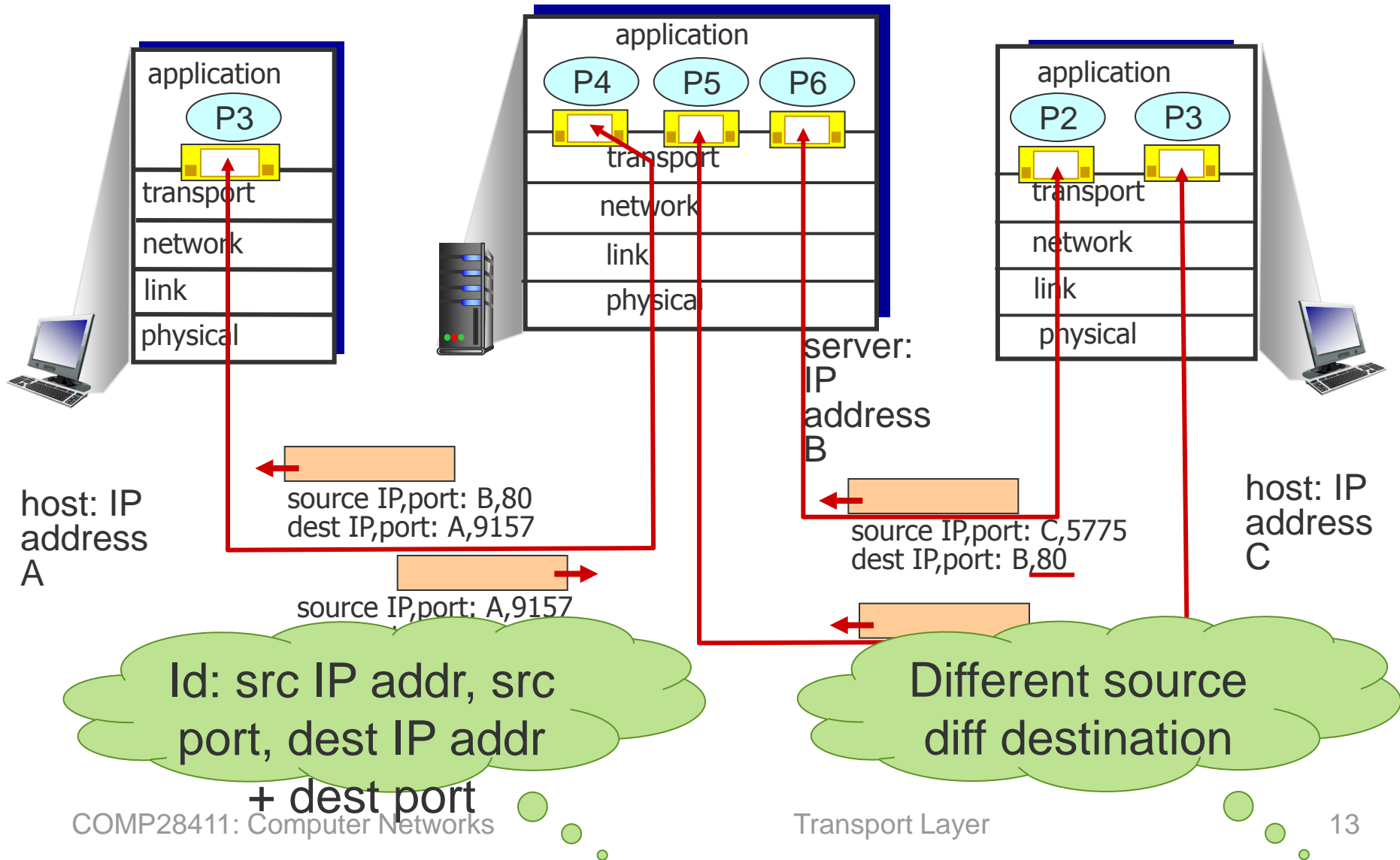
# End-Point: Id (Connectionless)



# End-Point: Id (Connection-Oriented)



# End-Point: Id (Connection-Oriented)



# User Datagram Protocol (UDP)

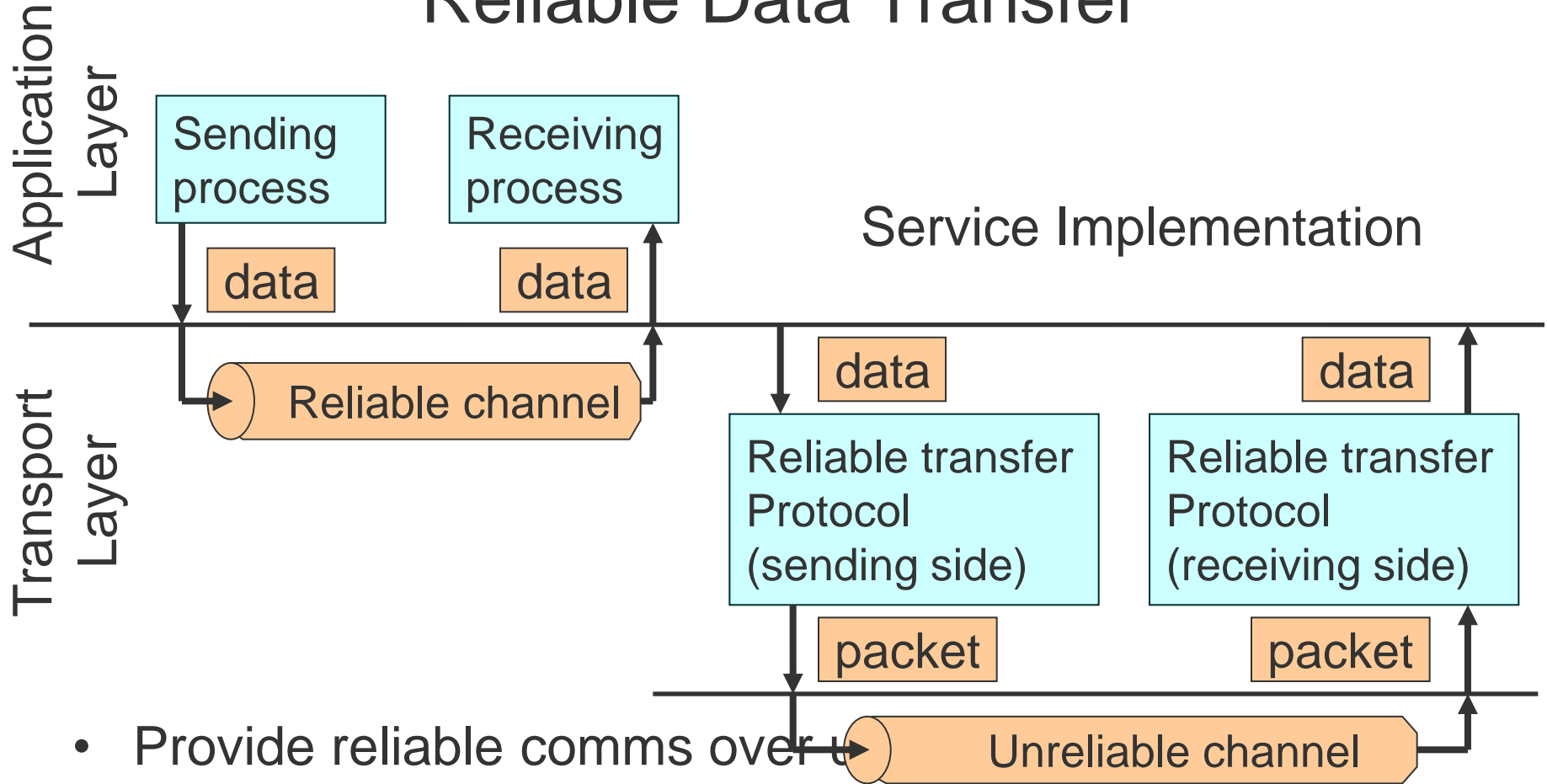
- “no frills”, “best effort” connectionless service, segments may be:
  - lost, reordered
- Segments independent
- Often used for streaming multimedia applications
  - loss tolerant
  - rate sensitive
- Other uses:
  - DNS, SNMP
  - NFS
  - o/s applications

## Why is there a UDP?

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

[RFC 768]

# Reliable Data Transfer

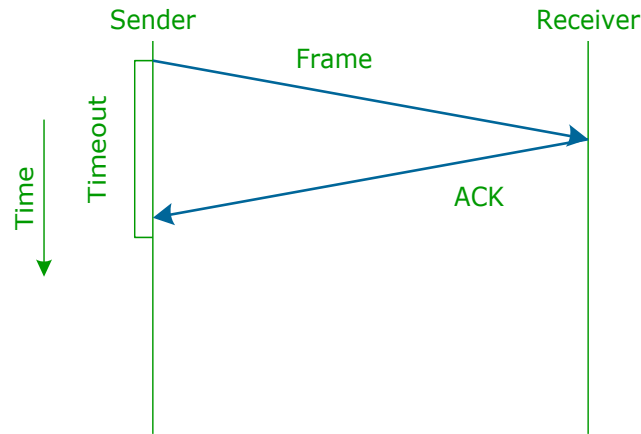


- Provide reliable comms over unreliable channel
- Complexity depends on characteristics of unreliable channel

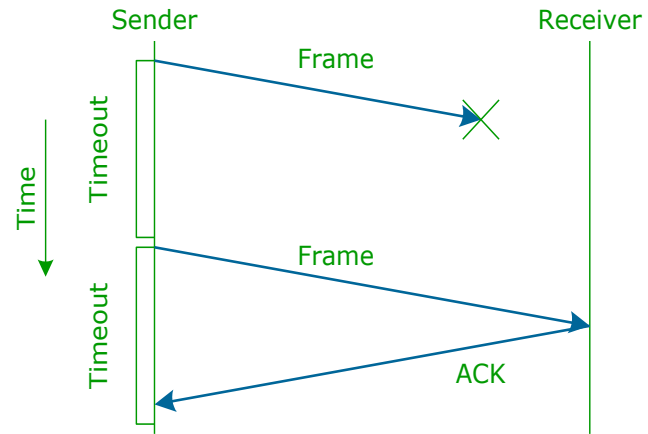
# Recovering from Errors

- Recovery uses two mechanisms:
  - acknowledgements and timeouts
- Acknowledgement (**ACK**) is control packet from
  - receiver to transmitter of data packet being ACKed
- Receipt of ACK confirms delivery of data
- If ACK not received within **timeout**:
  - transmitter of data retransmits data; needs copy
- Process called automated repeat request (ARQ)
- ARQ mechanisms: **stop-and-wait, sliding window**

# Reliability: Stop-and-Wait

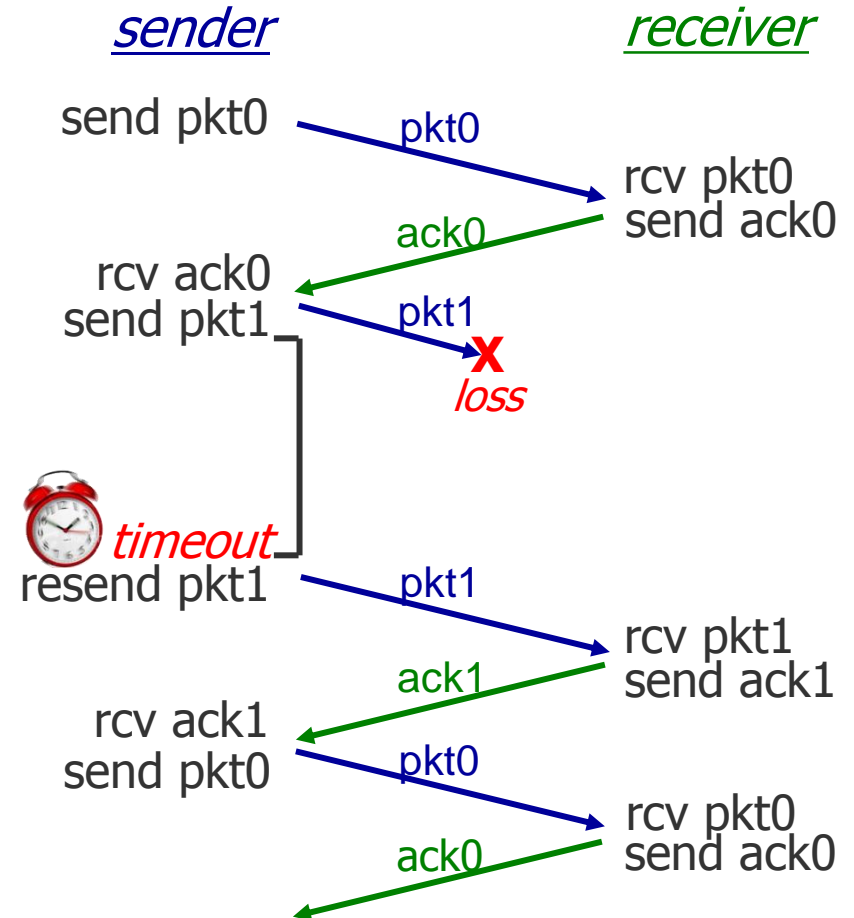
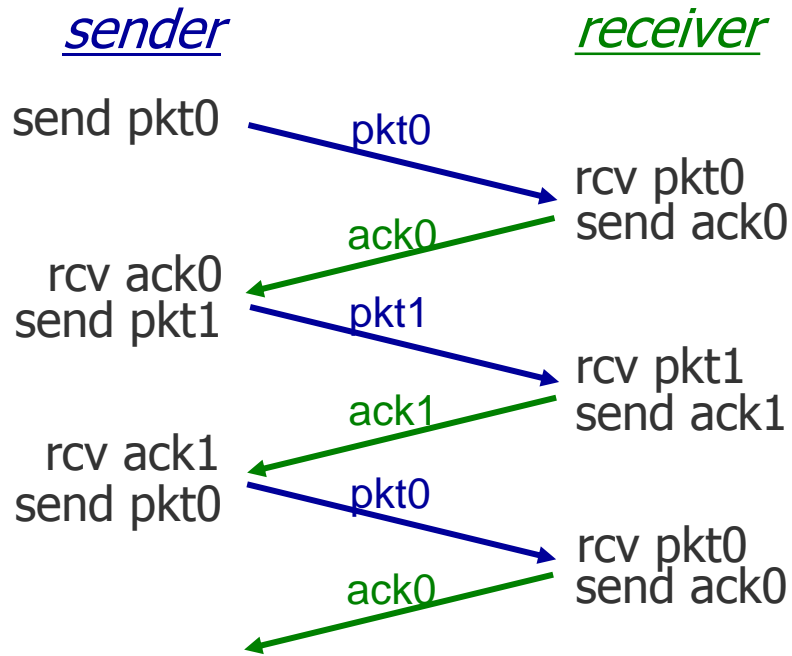


Peterson and Davie, Figure 2.19



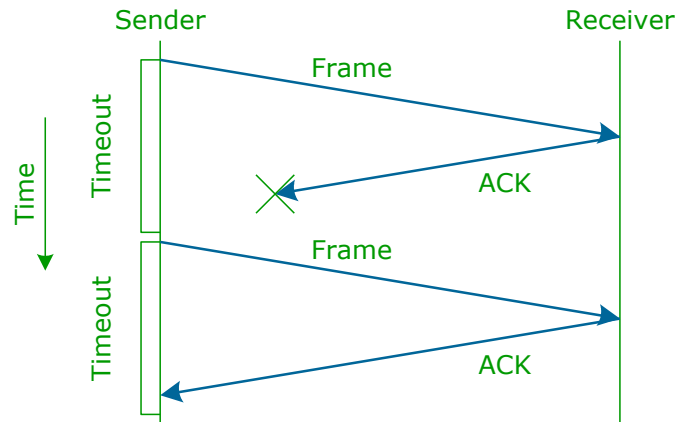
- Wait for acknowledgement before sending next packet
- Normal operation, ACK received before timeout expires
- If not, data packet is retransmitted and, hopefully, ACKed

# Reliability: Stop-and-Wait

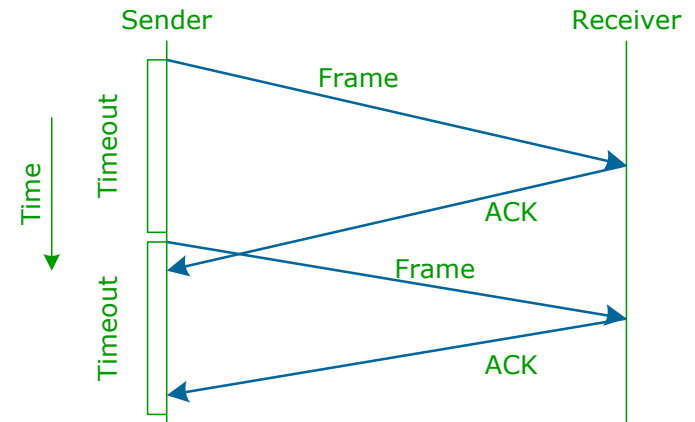




# Reliability: Stop-and-Wait

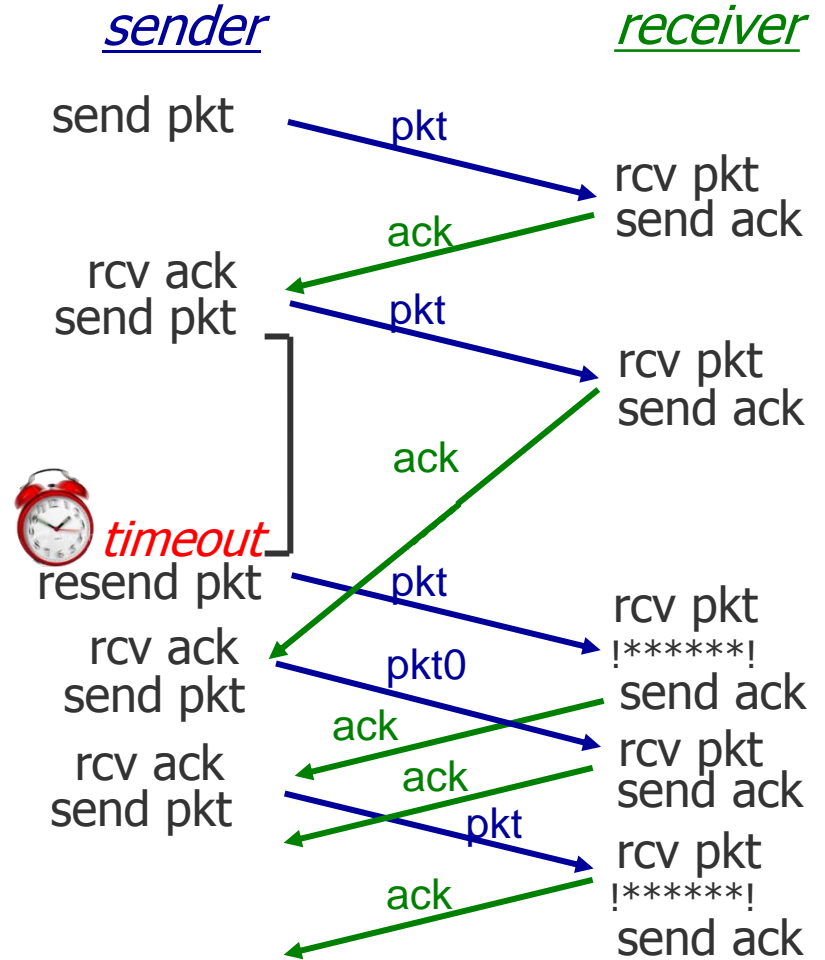
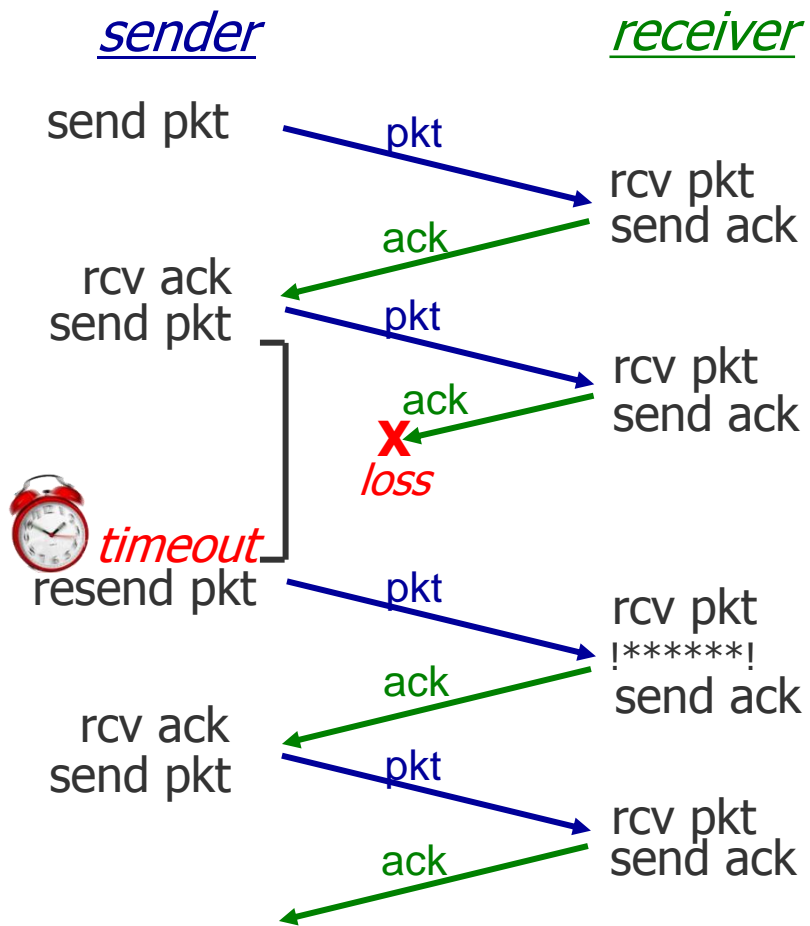


Peterson and Davie, Figure 2.19

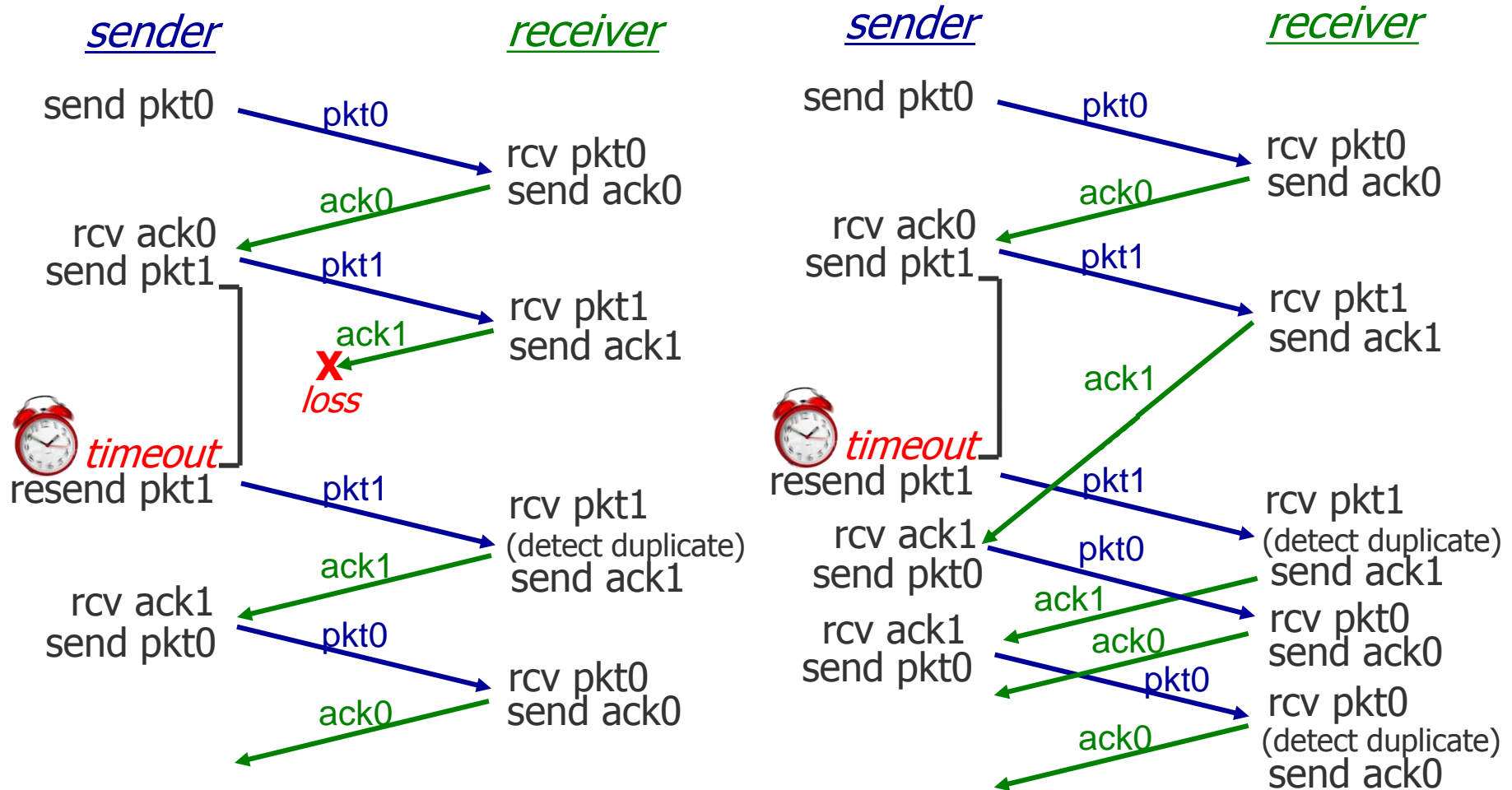


- Loss of ACK, or late arrival, causes retransmission
- Duplicate packet is received, but
- receiver believes that it is receiving a new packet
- Duplicates can be detected using sequence numbers
- One-bit sequence numbers used; 0 or 1

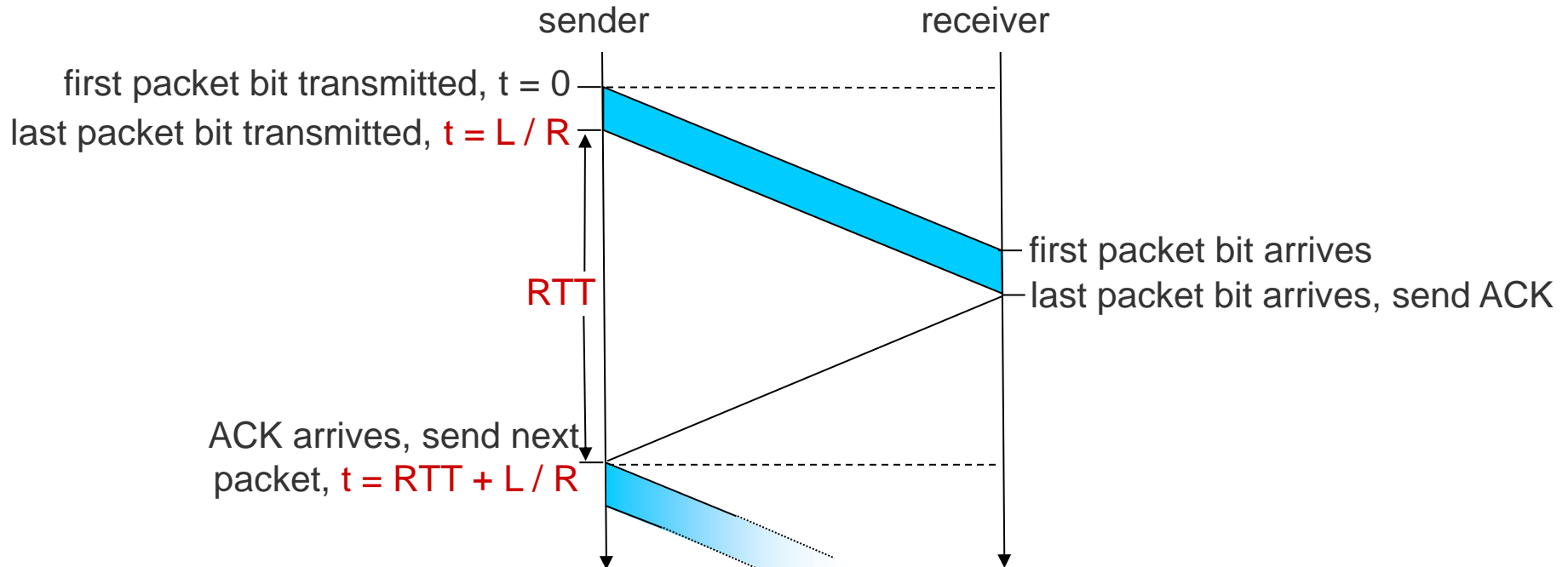
# Reliability: Stop-and-Wait



# Reliability: Stop-and-Wait



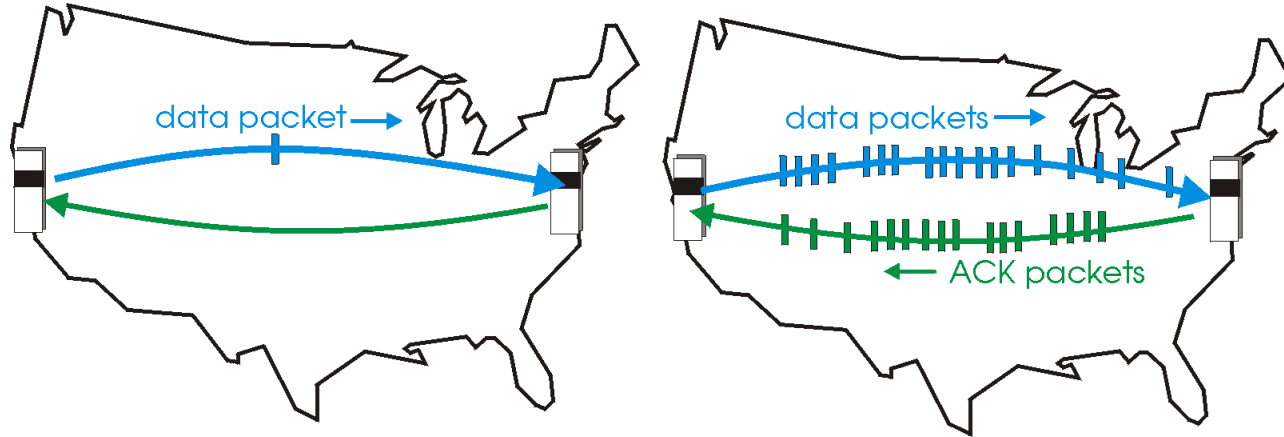
# Reliability: S-a-Wait - Utilisation



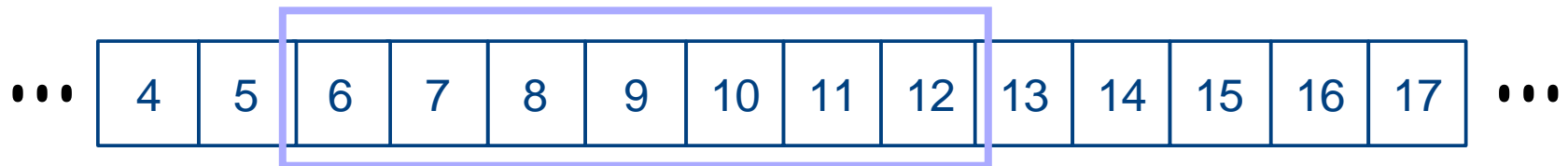
$$U_{\text{sender}} = \frac{L / R}{RTT + L / R} = \frac{.008}{30.008} = 0.00027$$

Network protocols limit use  
of physical resources

# Reliability: Pipelined Protocols

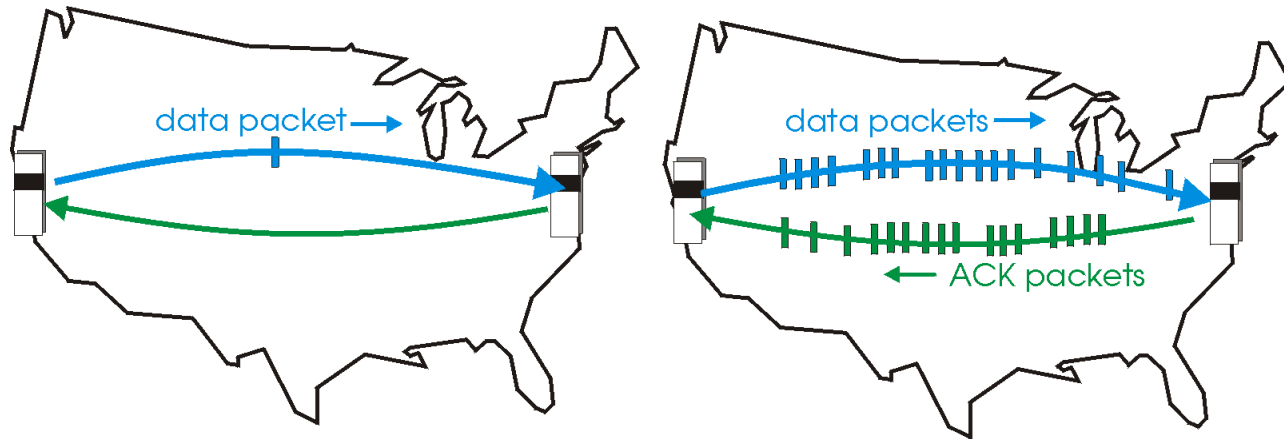


- Allow multiple, “in-flight”, yet-to-be-acknowledged pkts
- Example: received ACK 5, window 6

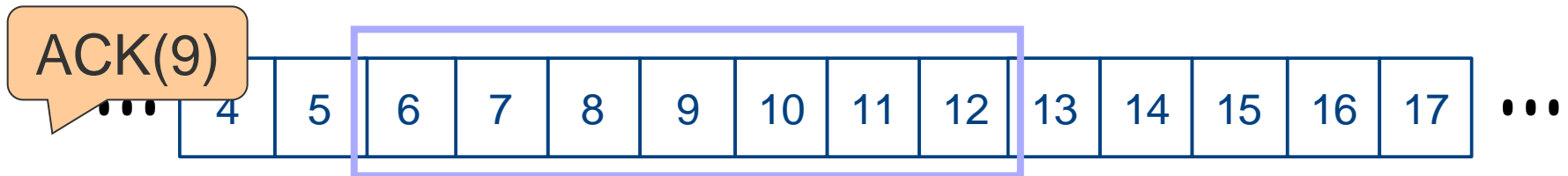


- Often referred to as sliding window

# Reliability: Pipelined Protocols

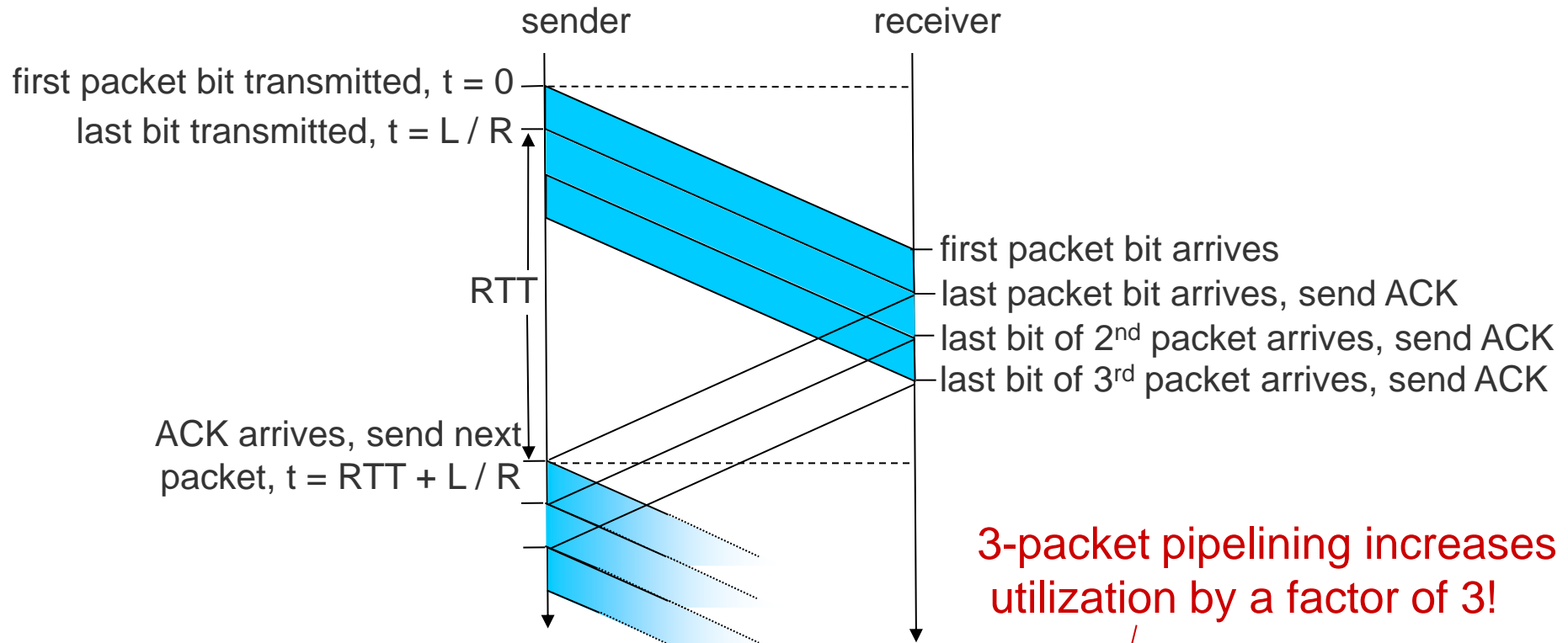


- Allow multiple, “in-flight”, yet-to-be-acknowledged pkts
- Example: received ACK 5, window 6



- Often referred to as sliding window

# Reliability: Pipelined Utilisation



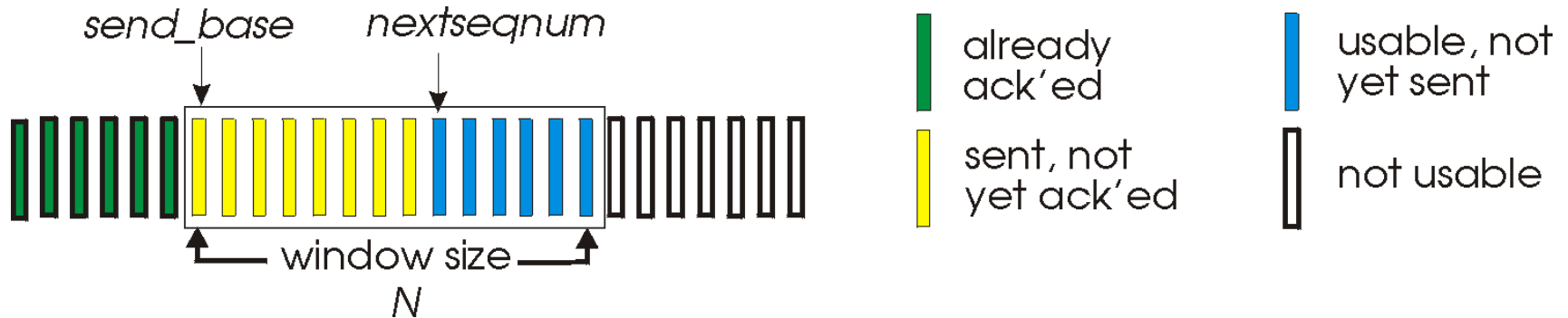
$$U_{\text{sender}} = \frac{3L / R}{RTT + L / R} = \frac{.0024}{30.008} = 0.00081$$

# Reliability: Pipelined Approaches

- What if expecting packet 8 and get packet 9?
- If packet 8 delayed, when arrives can send ACK(9)
- If lost, timeout will expire and will be resent
- **Go-Back-N**, send cumulative ACKs:
  - ACK(n) acknowledges all packets upto n
  - likely that will also get packets 9, 10, ... resent
- **Selective repeat**:
  - explicitly acknowledges all packet
  - only unsuccessfully received packets are resent
- Could **negatively acknowledge** (NACK) packet 8,
  - requests retransmission without timeout expiring



# Reliability: Go-Back-N



- k-bit seq # in pkt header
- “window” of up to  $N$ , consecutive unack'd 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

# Reliability: Go-Back-N Example (H)

sender window (N=4)

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

sender

send pkt0

send pkt1

send pkt2

send pkt3

(wait)

rcv ack0, send pkt4

rcv ack1, send pkt5

ignore duplicate ACK



*pkt 2 timeout*

send pkt2

send pkt3

send pkt4

send pkt5

receiver

receive pkt0, send ack0

receive pkt1, send ack1

receive pkt3, discard,  
(re)send ack1

receive pkt4, discard,  
(re)send ack1

receive pkt5, discard,  
(re)send ack1

rcv pkt2, deliver, send ack2

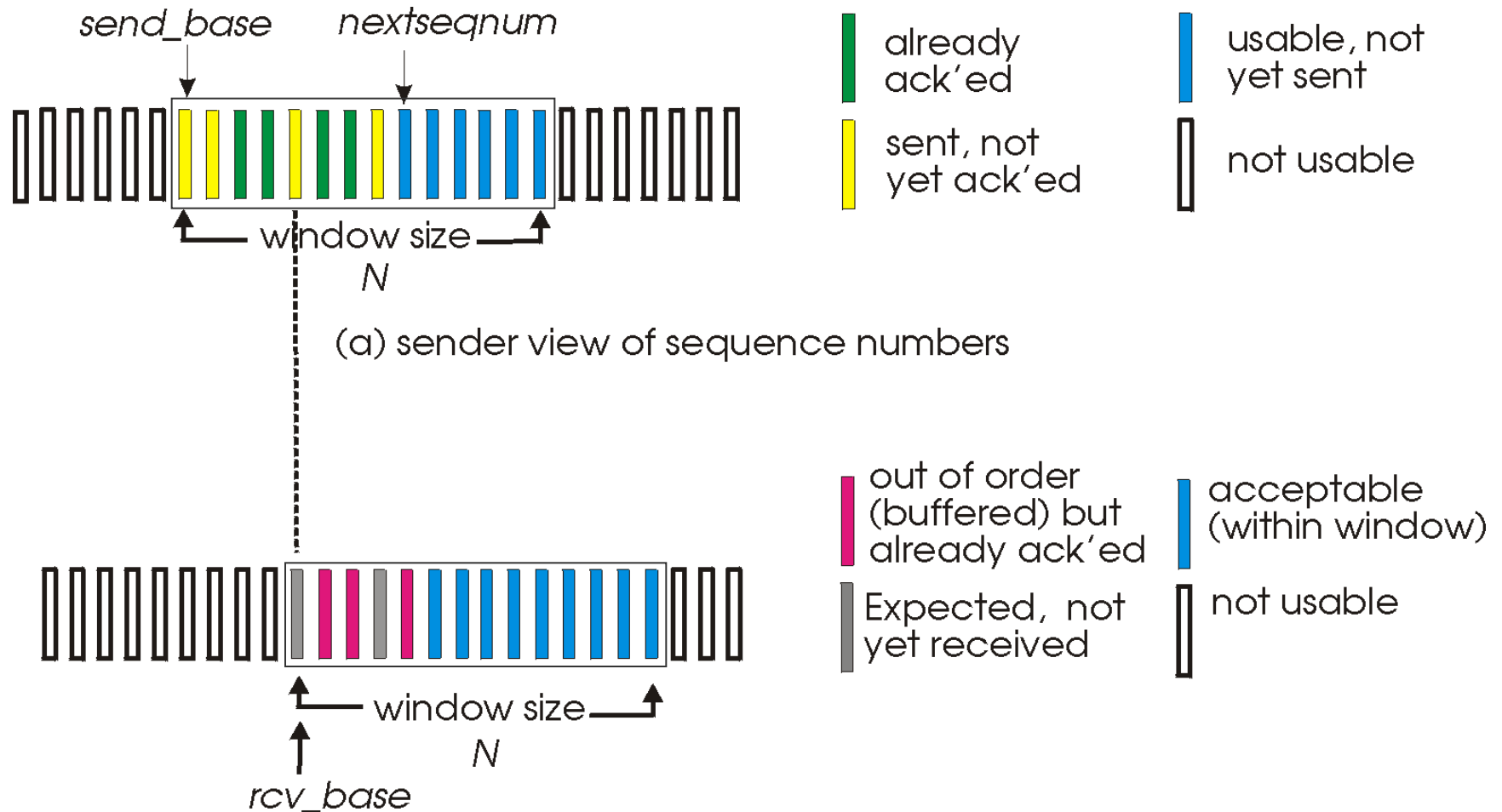
rcv pkt3, deliver, send ack3

rcv pkt4, deliver, send ack4

rcv pkt5, deliver, send ack5

*X loss*

# Reliability: Selective Repeat



# Reliability: Selective Repeat

sender window (N=4)

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

sender

send pkt0

send pkt1

send pkt2

send pkt3

(wait)

rcv ack0, send pkt4

rcv ack1, send pkt5

record ack3 arrived



*pkt 2 timeout*

send pkt2

record ack4 arrived

record ack5 arrived

receiver

receive pkt0, send ack0

receive pkt1, send ack1

receive pkt3, buffer,  
send ack3

receive pkt4, buffer,  
send ack4

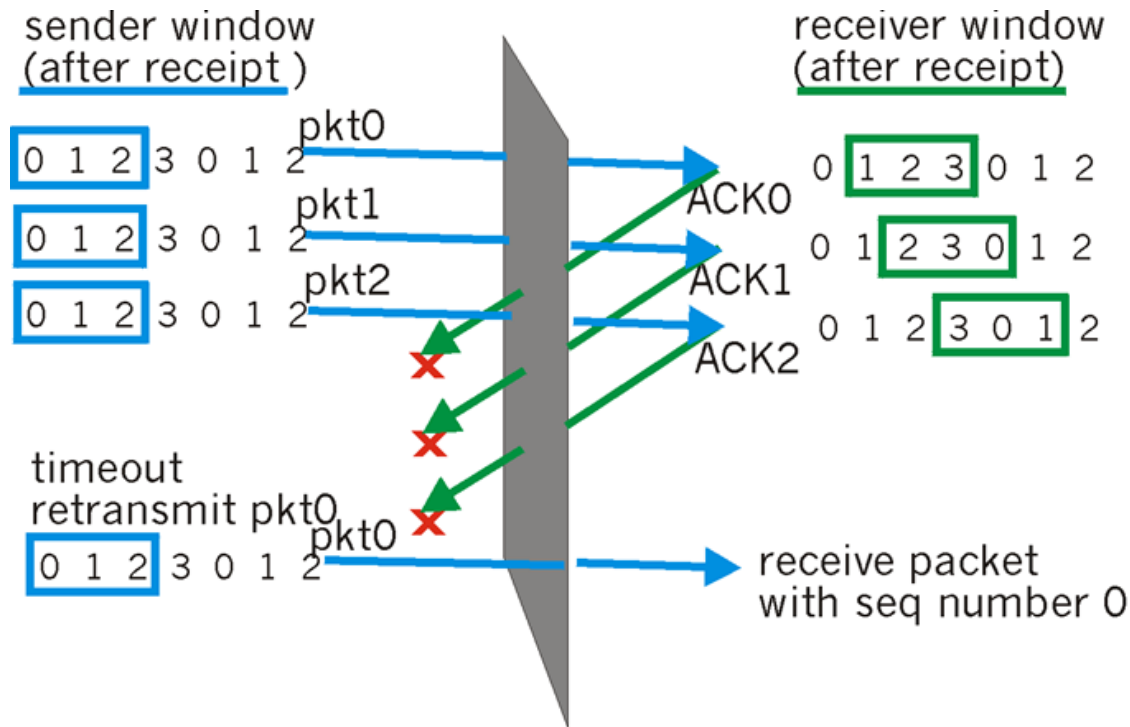
receive pkt5, buffer,  
send ack5

rcv pkt2; deliver pkt2,  
pkt3, pkt4, pkt5; send ack2

*Q: what happens when ack2 arrives?*

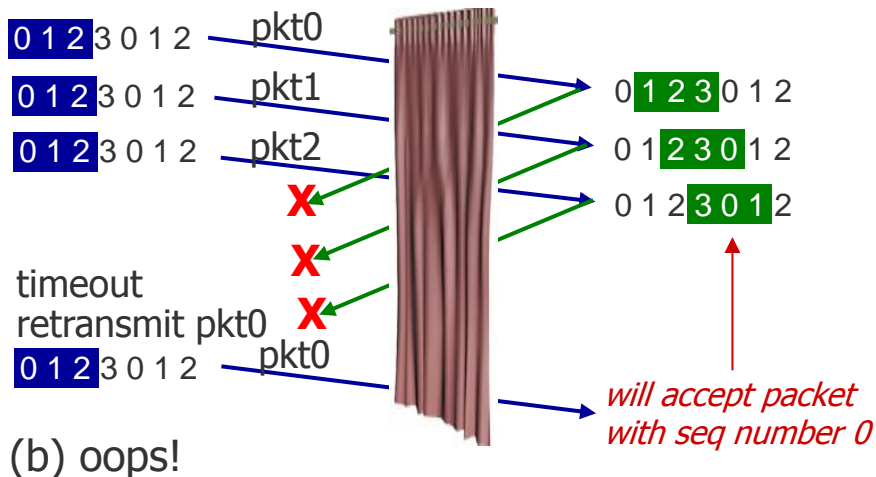
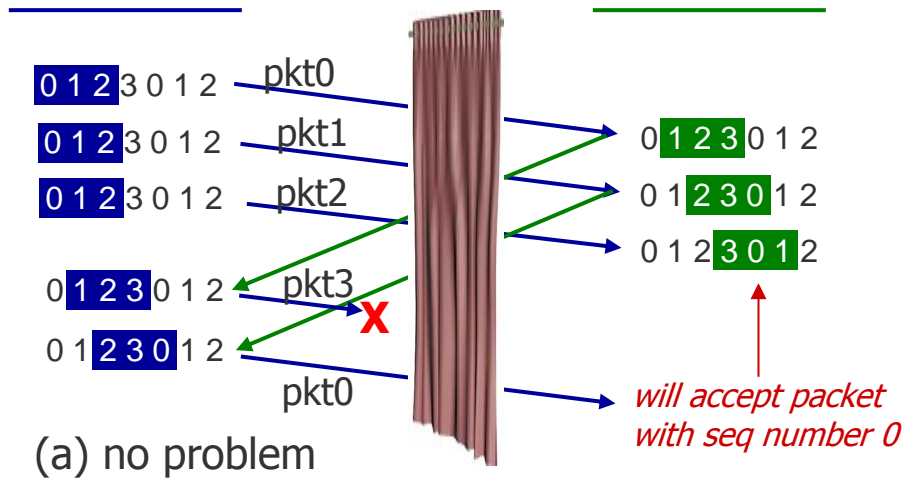
# Sliding Windows – Finite Seq. Numbers

- Sequence numbers wrap, must interpret correctly
- For example,



- Max window size =  $(\text{max sequence number} + 1)/2$

# Reliability: Window Issue



# Transport Services and Protocols (2)

---

Andy Carpenter

(Andy.Carpenter@manchester.ac.uk)

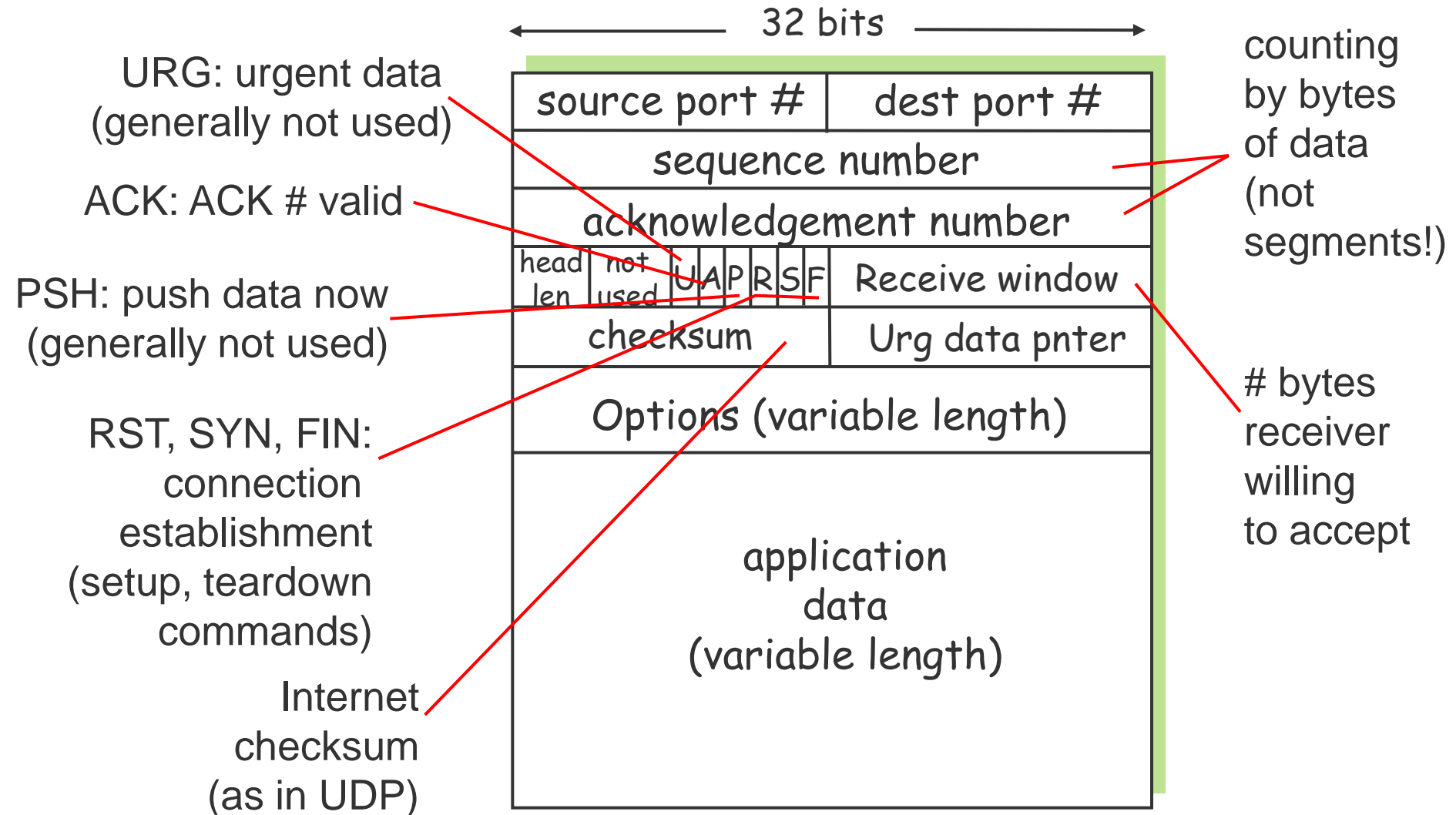
Elements these slides come from Kurose and Ross, authors of "Computer Networking: A Top-down Approach", and are copyright Kurose and Ross

# Transmission Control Protocol

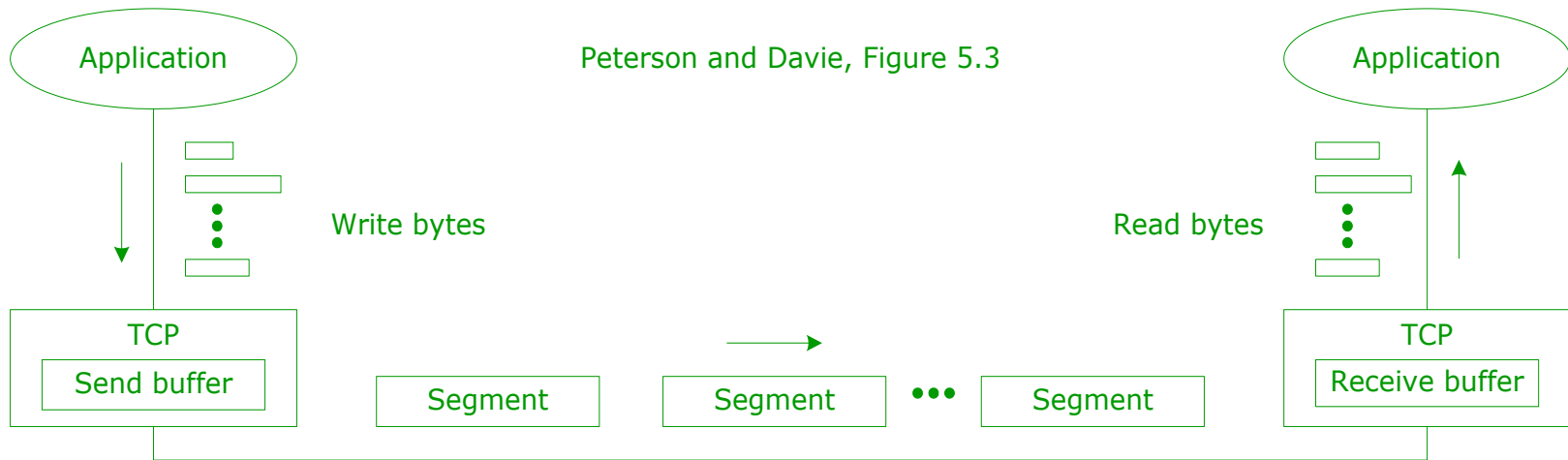
- **point-to-point:**
  - one sender, one receiver
- **reliable, in-order *byte stream*:**
  - no “message boundaries”
- **pipelined:**
  - TCP congestion and flow control set window size
- full duplex data:
- bi-directional data flow in same connection
- MSS: maximum segment size
- connection-oriented: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



# TCP: Segment Format



# TCP: Segmentation



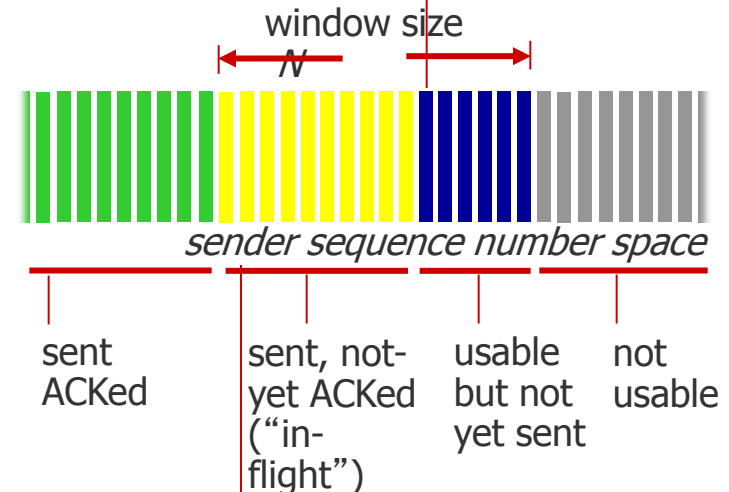
- Application writes bytes to connection as needs
- TCP collects (buffers) bytes before sending
- Sends set of bytes as a TCP segment in IP datagram
- Receiving TCP also buffers
- Application reads bytes from connection as required

# TCP seq. numbers, ACKs

- sequence numbers:
  - byte stream “number” of first byte in segment’s data
- acknowledgements:
  - seq # of next byte expected from other side
- cumulative ACK

outgoing segment from sender

source port #	dest port #
sequence number	
acknowledgement number	
	rwnd
checksum	urg pointer

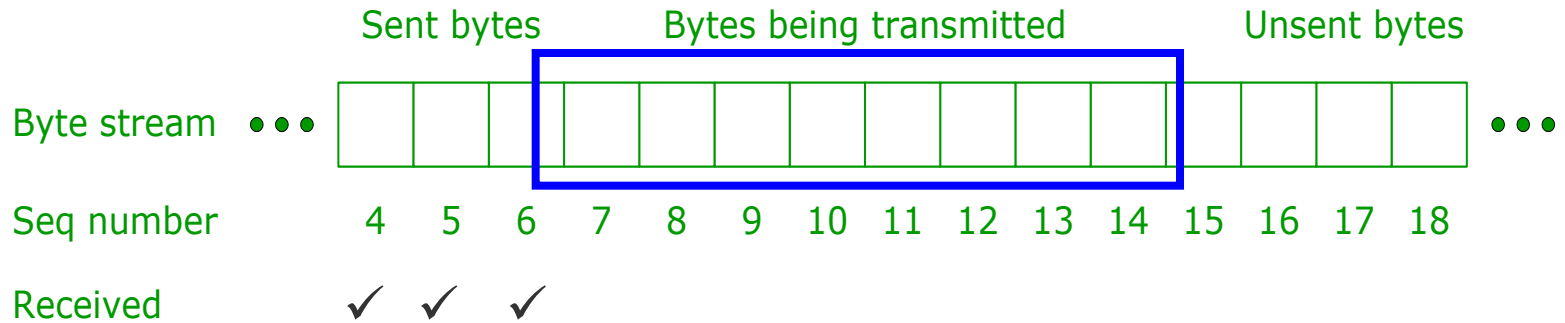


incoming segment to sender

source port #	dest port #
sequence number	
acknowledgement number	
	A
checksum	urg pointer

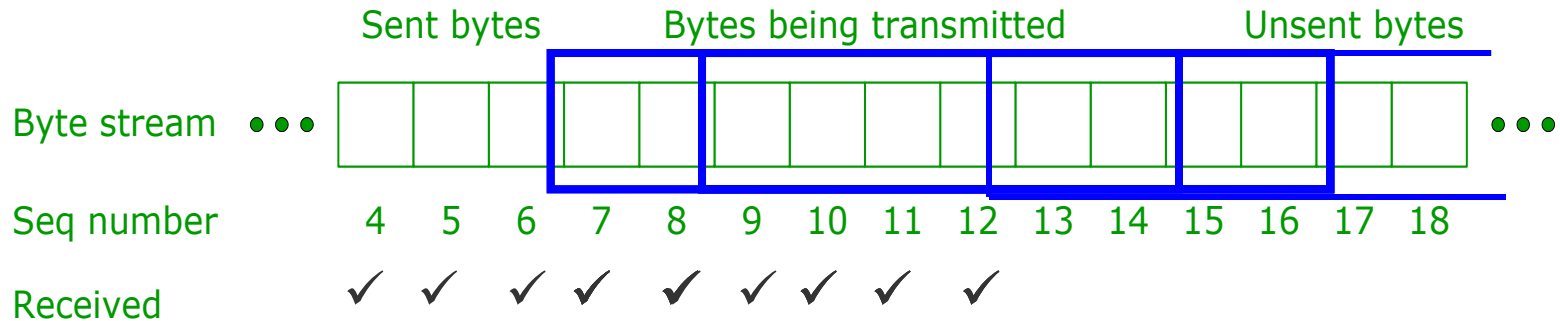
# TCP: Reliability

- Example, window 8, last ACK 7



# TCP: Reliability

- Example, window 8, last ACK 7



Receive segment containing bytes with sequence numbers 7 and 8

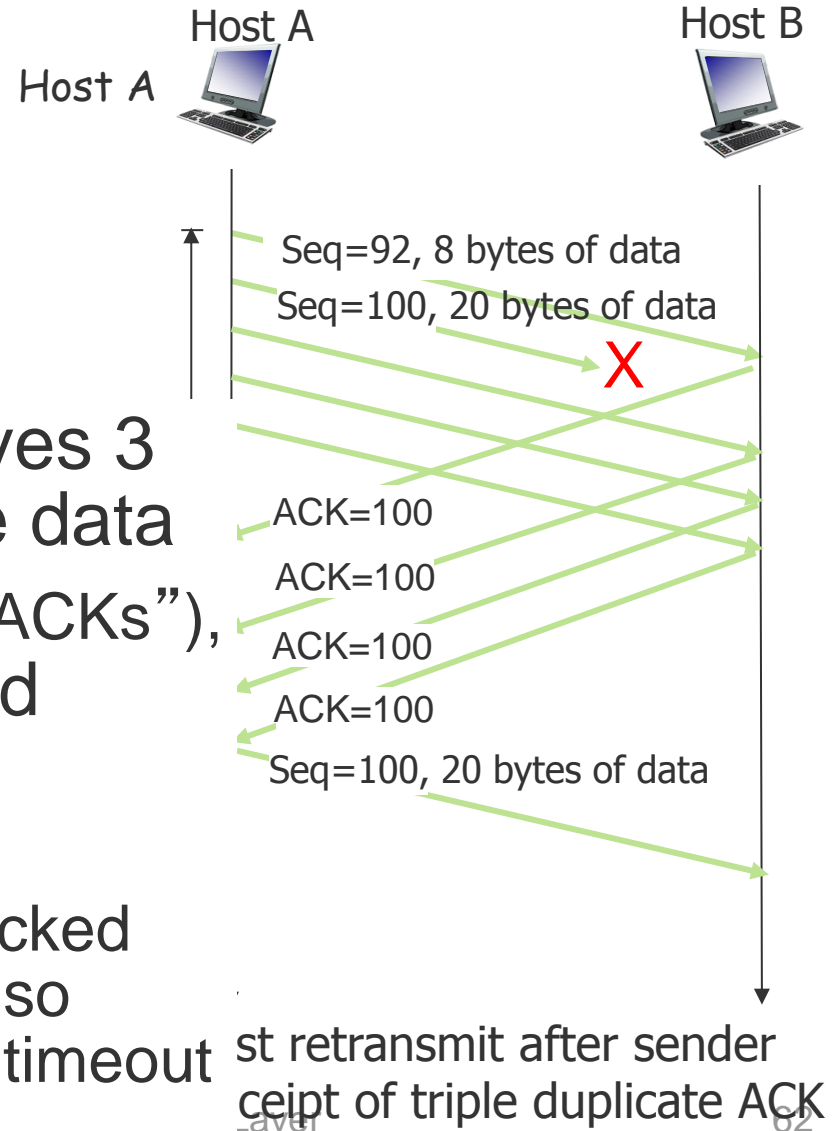
Receive segment containing bytes with sequence numbers 10 and 11

Receive segment containing byte with sequence number 9

Acknowledge byte with sequence number 9 and slide window along byte stream  
Note bytes received, but do not acknowledge or slide window  
Acknowledge byte with sequence number 13 and slide window along byte stream

# Reliability: TCP – Fast Retransmit

- ❖ time-out period often relatively long:
  - long delay before resending lost packet
- ❖ detect lost if sender receives 3 duplicate ACKs for same data
  - sender (“triple duplicate ACKs”), many senders resend unacked to-back segment with smallest seq #
  - if segment lost, so there will be many duplicate ACKs
  - likely that unacked segment lost, so don't wait for timeout



# Reliability: Retransmit Timeout

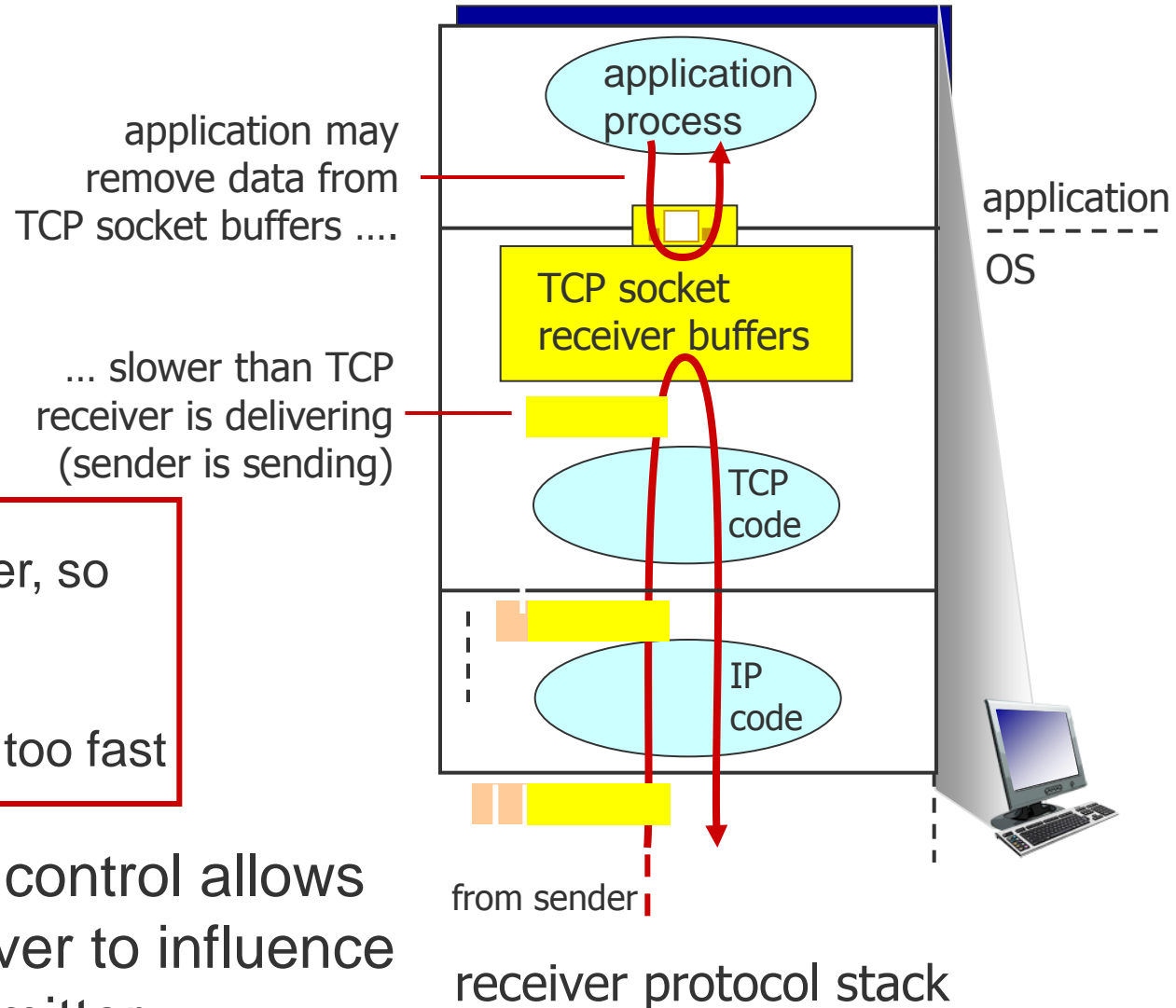
- If value too small, unnecessarily retransmit segments
- If value too large, get excessive delays before retransmit
- Appropriate value related to RTT, but that depends on:
  - pair of hosts involved, time, congestion in network
- For each connection:
  - use adaptable algorithm to determine RTT
- Basic algorithm sets timeout to twice estimated RTT
- Karn/Partridge algorithm reduces miscalculations
- Jacobson/Karels algorithm
  - copes with significant variance in real RTT

# Reliability: Estimation of RTT

- Estimate RTT as average of measured RTTs:
- When re-sending, to which transmit does ACK relate?
- Karn/Partridge algorithm:
  - only measures RTT for non-retransmitted segments
  - doubles timeout value for each retransmission
- Jacobson/Karels Algorithm:
  - if RTT variation small, average is good approx.
  - if RTT variation large, average is poor approx.
  - algorithm takes account of variation



# Flow Control



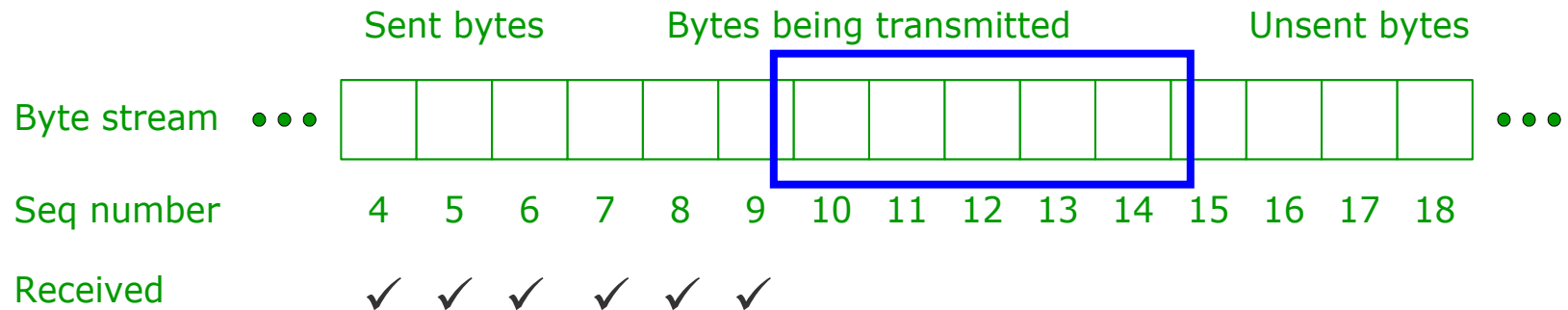
*flow control*

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

# Flow control allows receiver to influence transmitter

# Flow Control: Example

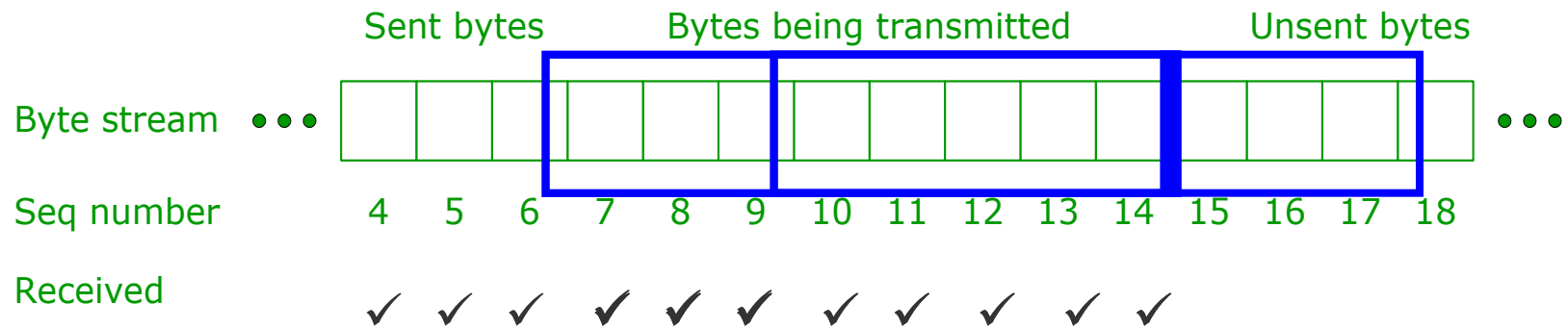
- Uses variant of sliding windows algorithm
  - window size can be changed



- Can get deadlock

# Flow Control: Example

- Uses variant of sliding windows algorithm
  - window size can be changed



Receive segment containing edge byte with sequence number 10 and reduce window size to 4

Receive segment containing edge byte with sequence number 14

Can get deadlock to Re-acknowledge byte with sequence number 15 and increase window size

# TCP 3-way handshake

*client state*

LISTEN

SYNSENT

ESTAB

choose init seq num,  $x$   
send TCP SYN msg

received SYNACK( $x$ )  
indicates server is live;  
send ACK for SYNACK;  
this segment may contain  
client-to-server data



SYNbit=1, Seq= $x$

SYNbit=1, Seq= $y$   
ACKbit=1; ACKnum= $x+1$

ACKbit=1, ACKnum= $y+1$

choose init seq num,  $y$   
send TCP SYNACK  
msg, acking SYN

received ACK( $y$ )  
indicates client is live

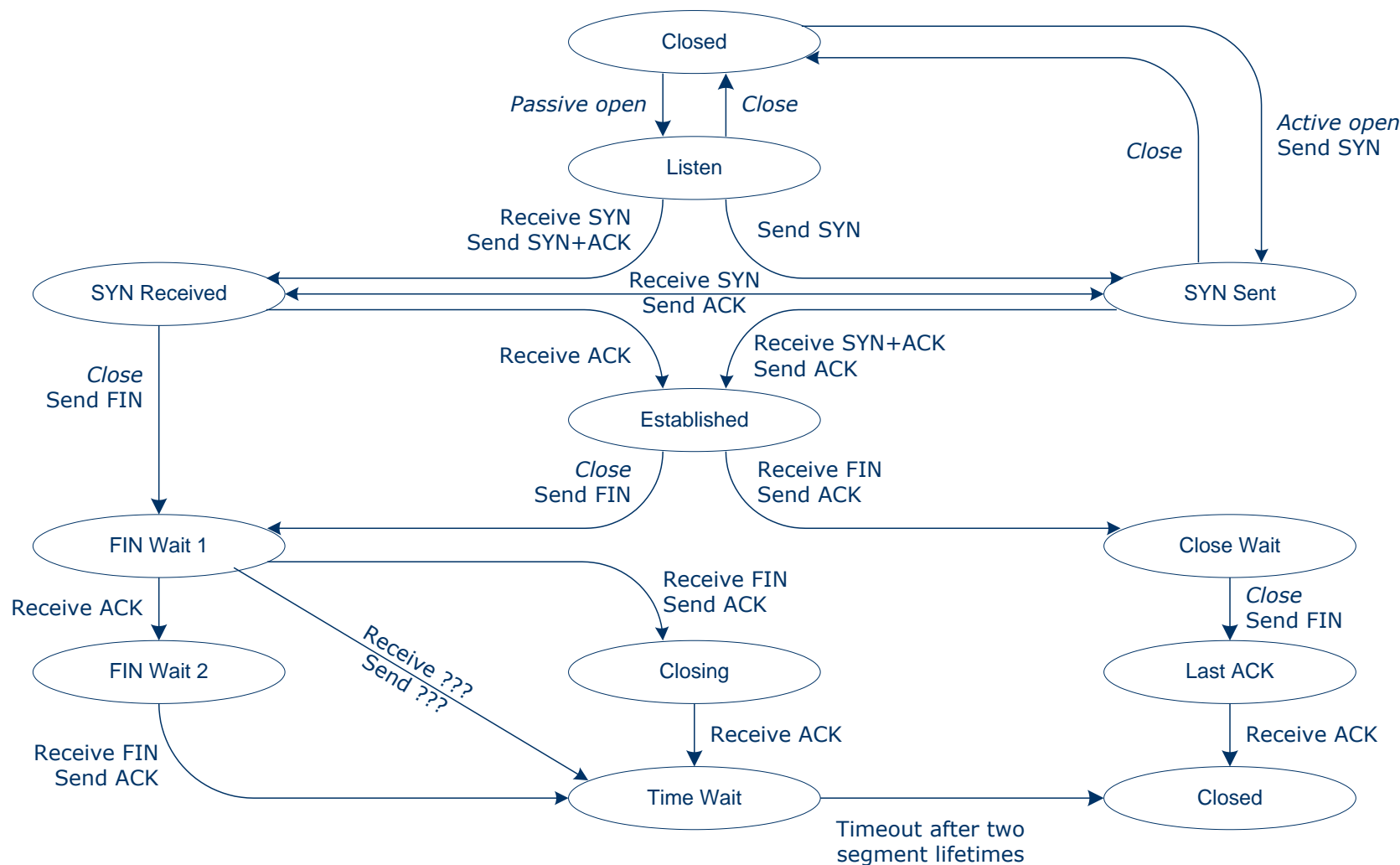
*server state*

LISTEN

SYN RCVD

ESTAB

# TCP: State Transition Diagram



# TCP: Keeping the Pipe Full

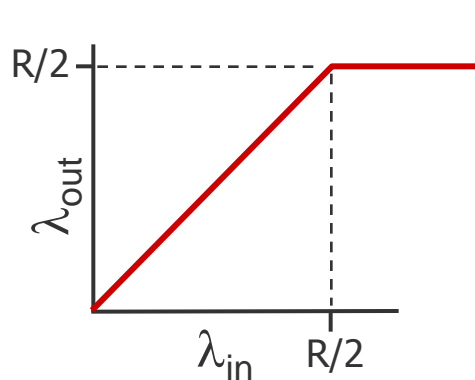
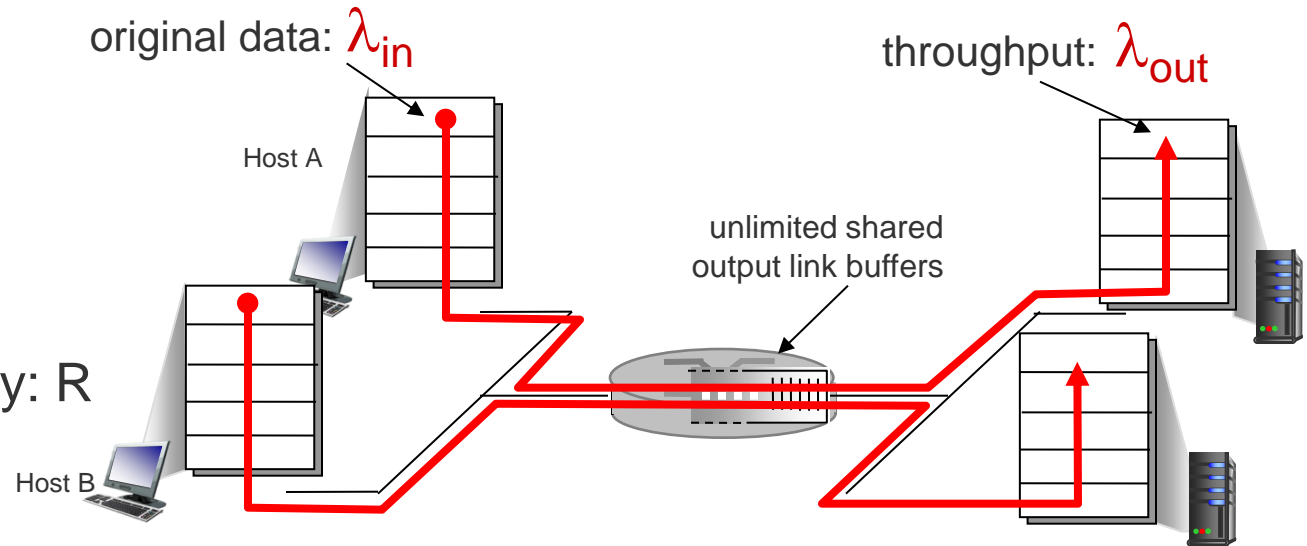
- Capacity of physical connection (100ms RTT)

T1 (1.5Mbps)	18KB	Ethernet (10Mbps)	122KB
T3 (45Mbps)	549KB	FDDI (100Mbps)	1.2MB
STS-3 (155Mbps)	1.8MB	STS-12 (622Mbps)	7.4MB
STS-24 (1.2Gbps)	14.8MB		

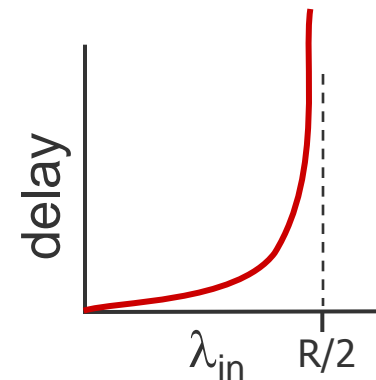
- Maximum window of single TCP connection:  $2^{16} = 64\text{KB}$
- Maximising use of cables relies on shared use
- Approaching point where:
  - capacity on cable exists but no acknowledgement
- Proposed extension allows:
  - multiplier factor for sequence number and window size

# Congestion: Causes/Costs 1

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

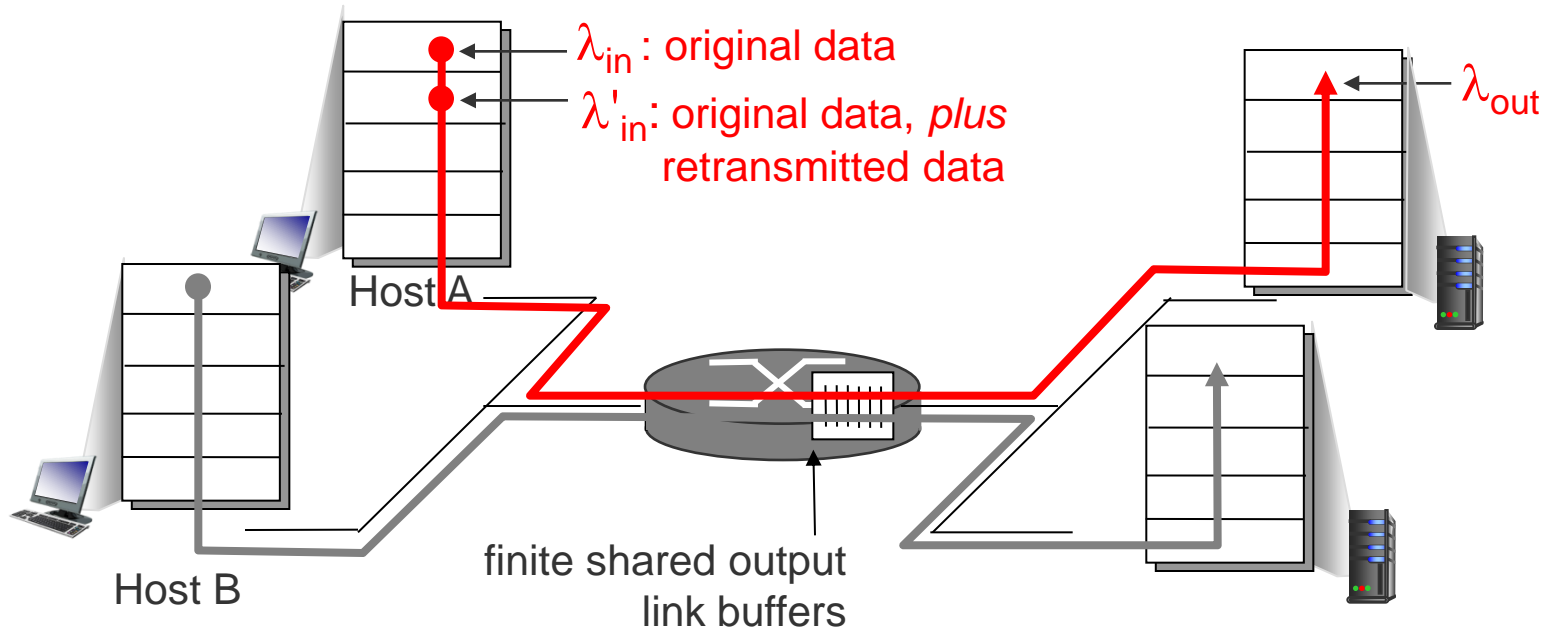


- ❖ maximum per-connection throughput:  $R/2$



- ❖ large delays as arrival rate,  $\lambda_{in}$ , approaches capacity

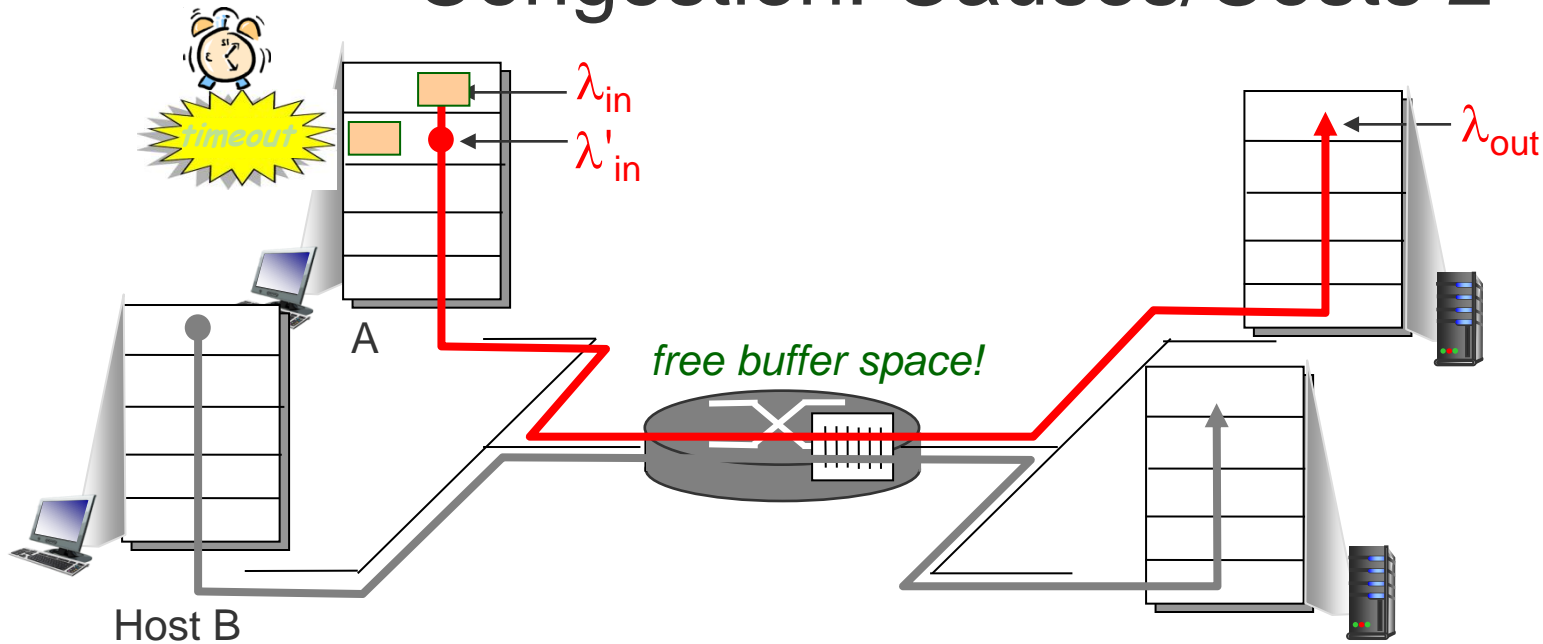
# Congestion: Causes/Costs 2



- ❖ 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* :  $\lambda_{in}$   $\lambda_{in}$

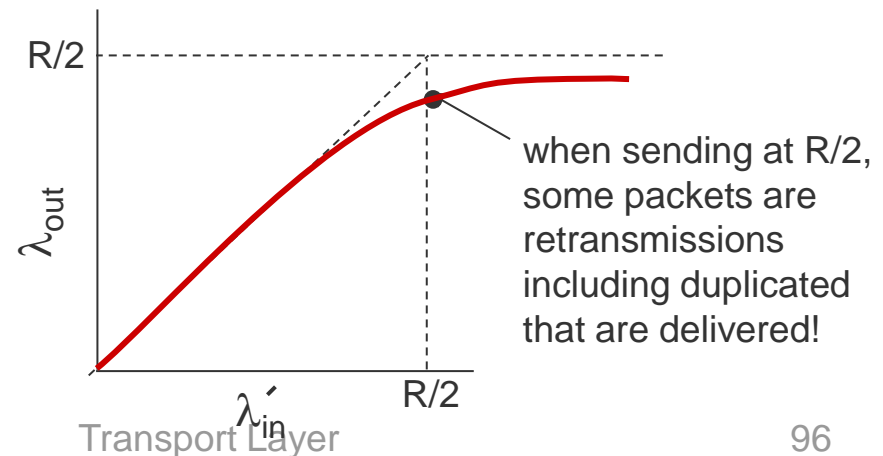


# Congestion: Causes/Costs 2

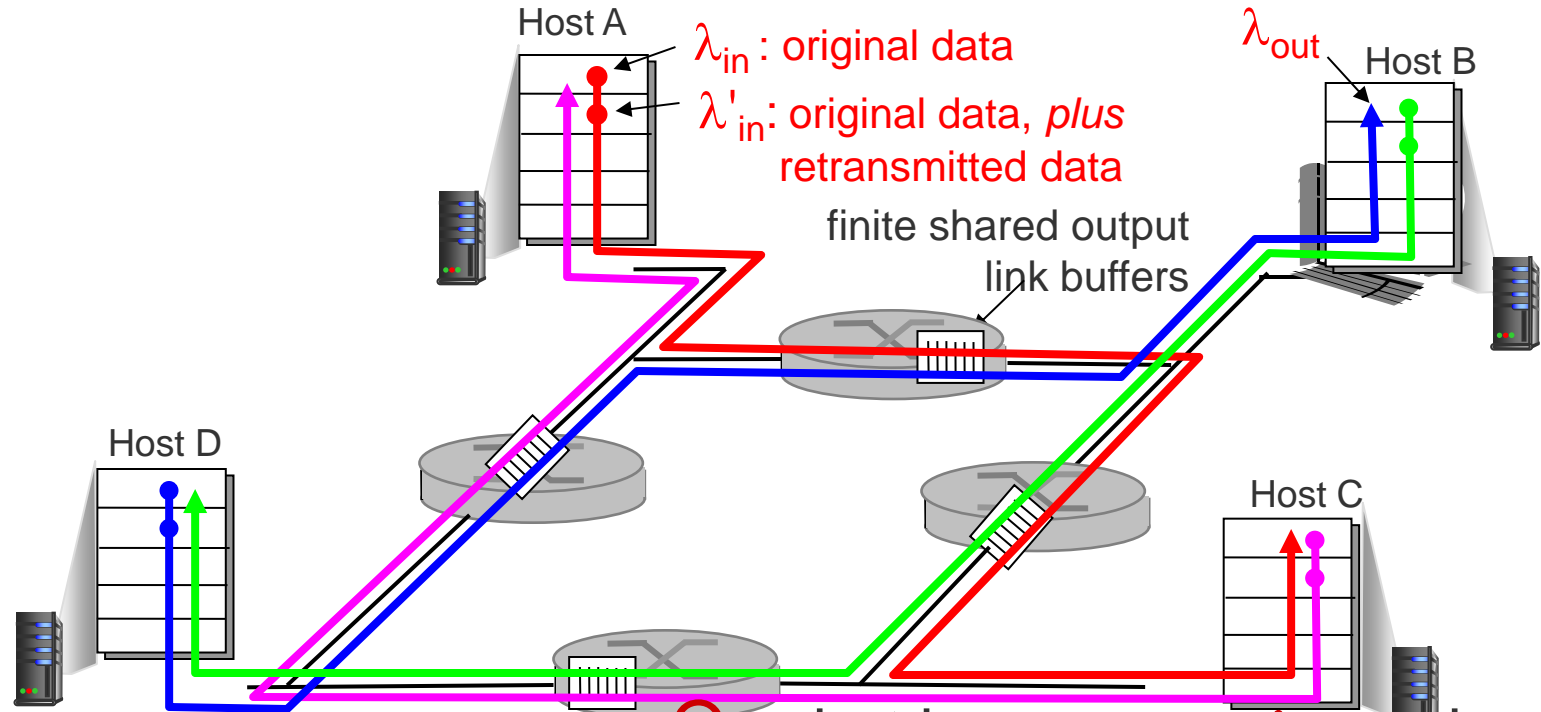


## Realistic: *duplicates*

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



# Congestion: Causes/Costs 3

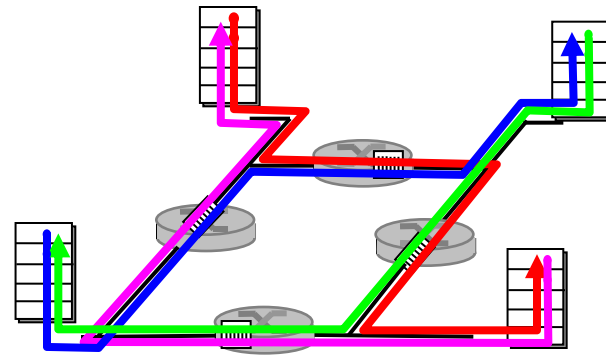
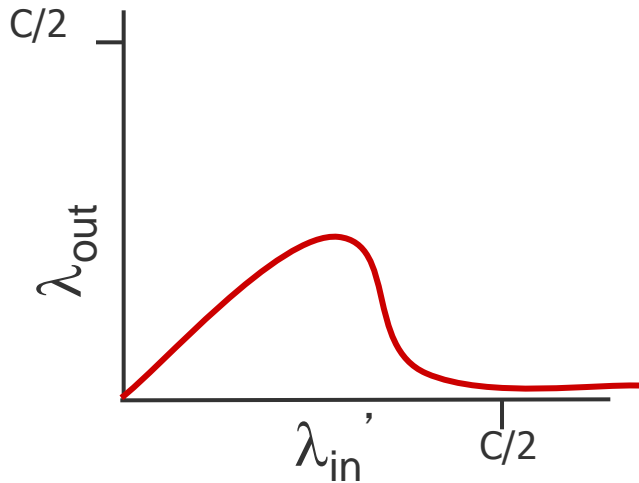


- ❖ four senders
- ❖ multihop paths
- ❖ timeout/retransmit

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

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

# Congestion: Causes/Costs 3



- another “cost” of congestion:
- when packet dropped, any “upstream transmission capacity used for that packet was wasted!

# Congestion: Control Approaches

## end-end congestion control:

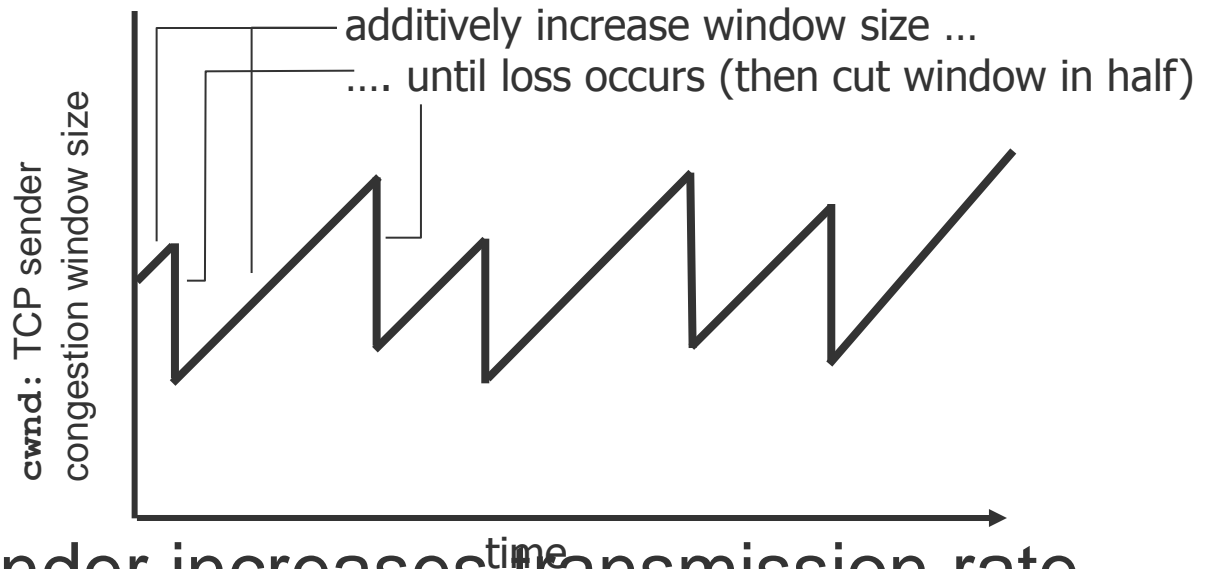
- ❖ no explicit feedback from network
- ❖ congestion inferred from end-system observed loss, delay
- ❖ approach taken by TCP

## network-assisted congestion control:

- ❖ routers provide feedback to end systems
  - single bit indicating congestion (SNA, DECbit, TCP/IP ECN, ATM)
  - explicit rate for sender to send at

# Congestion: TCP Control

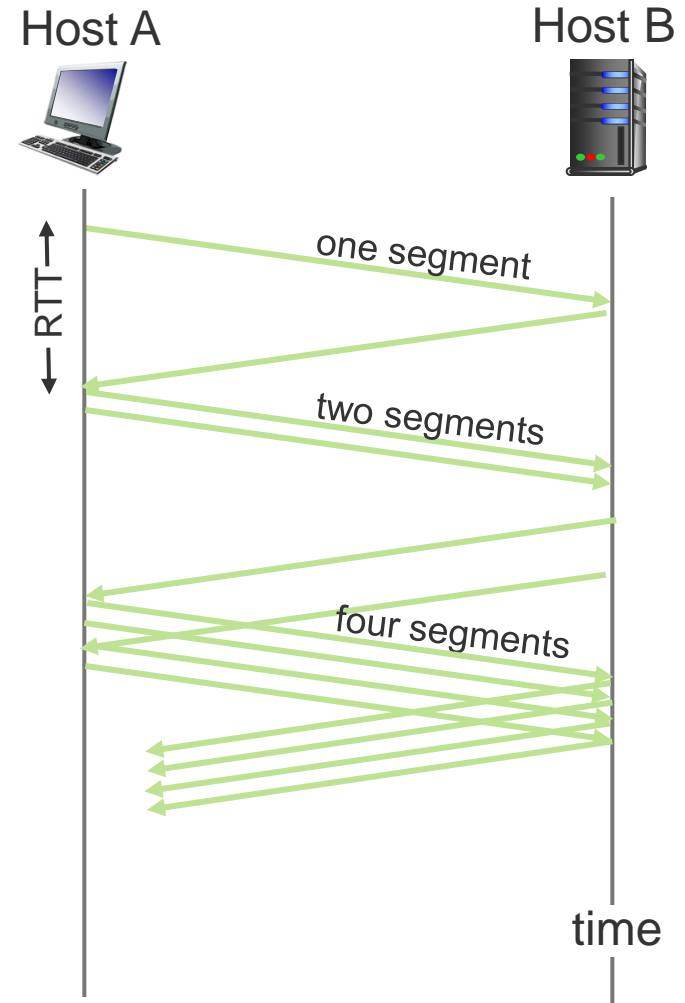
AIMD saw tooth behavior: probing for bandwidth



- ❖ *approach*: sender increases transmission rate (window size), probing for usable bandwidth, until loss occurs
  - *additive increase*: increase **cwnd** by 1 MSS every RTT until loss detected
  - *multiplicative decrease*: cut **cwnd** in half after loss

# TCP Slow Start

- ❖ when connection begins, increase rate exponentially until first loss event:
  - initially **cwnd** = 1 MSS
  - double **cwnd** every RTT
  - done by incrementing **cwnd** for every ACK received
- ❖ summary: initial rate is slow but ramps up exponentially fast



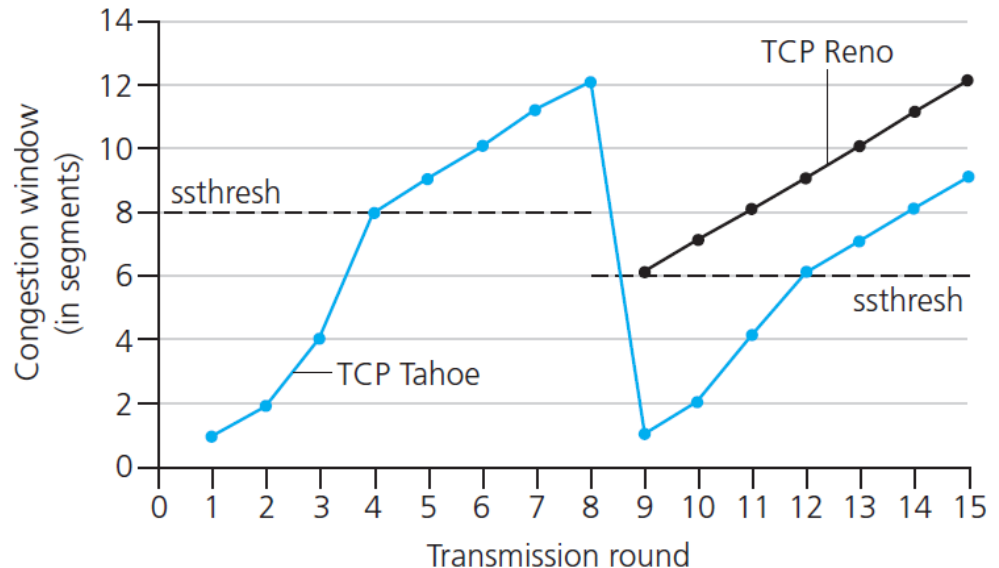
# TCP: switching from slow start to CA

**Q:** when should the exponential increase switch to linear?

**A:** when **cwnd** gets to 1/2 its value before timeout

## Implementation:

- ❖ variable **ssthresh**
- ❖ on loss event, **ssthresh** is set to 1/2 of **cwnd** just before loss event



# Fairness: TCP

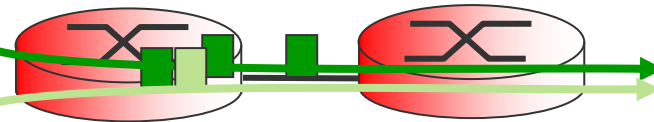
TCP connection 1



TCP connection 2



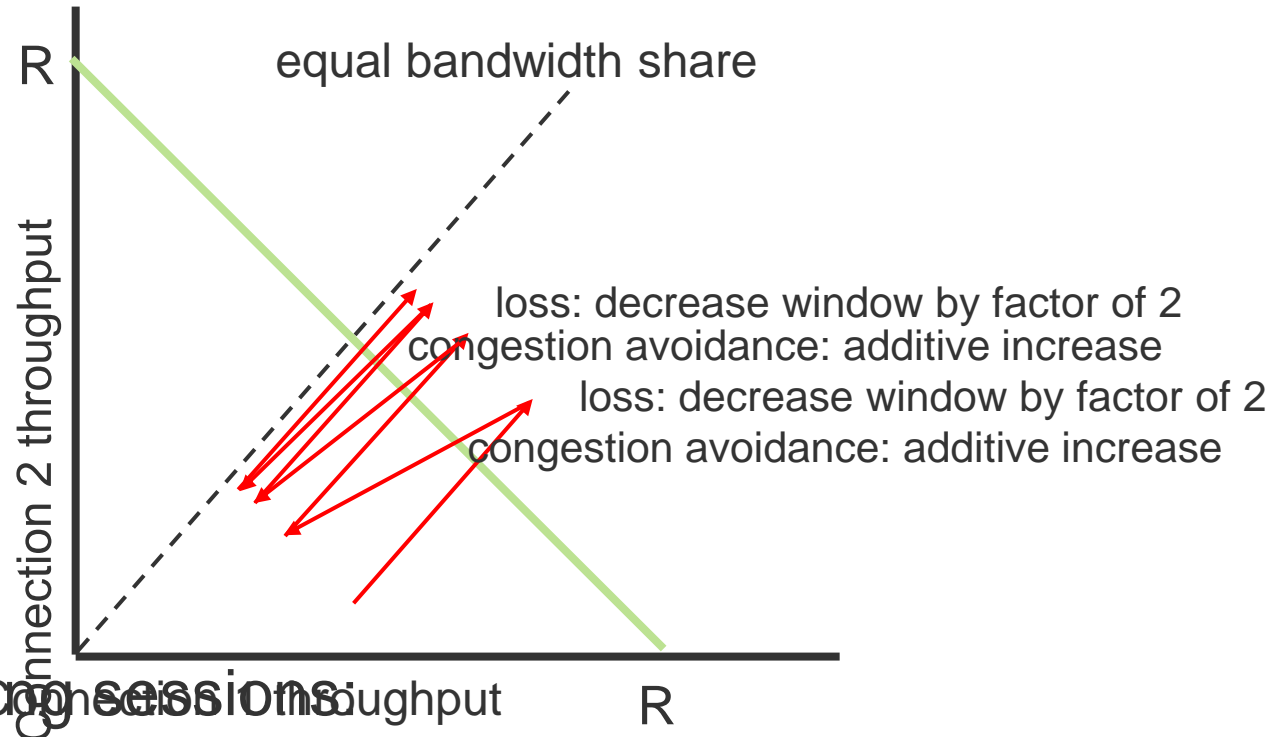
bottleneck  
router  
capacity R



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



# Fairness: TCP



two competing sessions

- ❖ additive increase gives slope of 1, as throughput increases
- ❖ multiplicative decrease decreases throughput proportionally

# Summary

- Examined end-to-end issues
- End-point identification, multiplexing/demultiplexing
  - all that UDP does
- Reliable data transfer
  - one important aspect is adaptive timeout
- Flow control
- Congestion control
- TCP use all to simplify applications