



# Chapter 3:Transport Layer

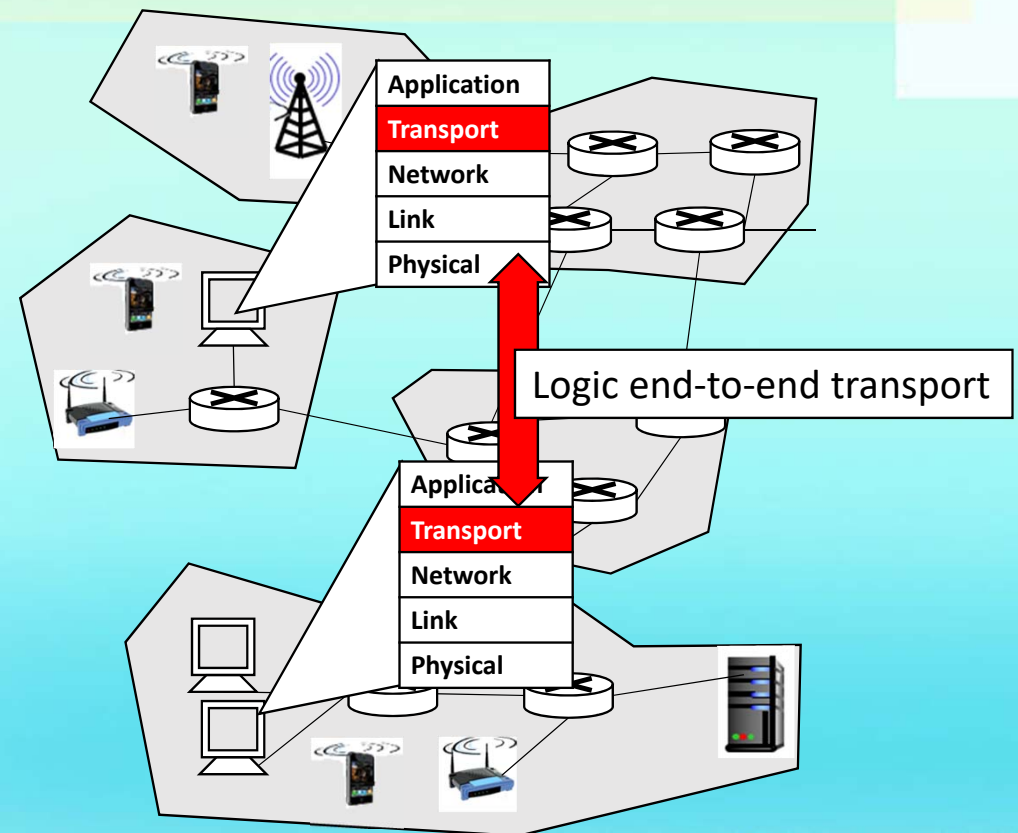
# Chapter 3: road map

- 3.1 Introduction and 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
- 3.6 Principles of Congestion Control
- 3.7 TCP Congestion Control

# 3.1 Introduction and Transport-Layer Services

## 3.1.1 Relationship Between Transport Layer and Network Layer

- 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**



### 3.1.1 Relationship Between Transport Layer and Network Layer

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

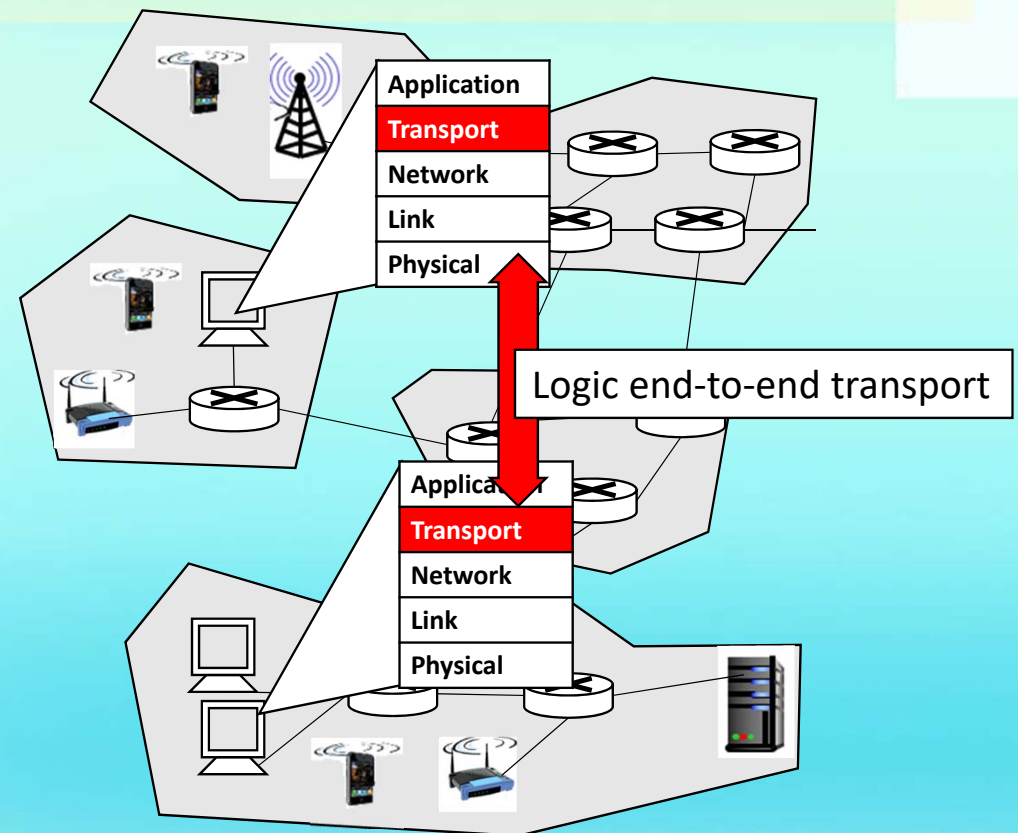
#### household analogy:

*12 kids in Ann's house sending letters to 12 kids in Bill's house:*

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

## 3.1.2 Overview of the Transport Layer in the Internet

- 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



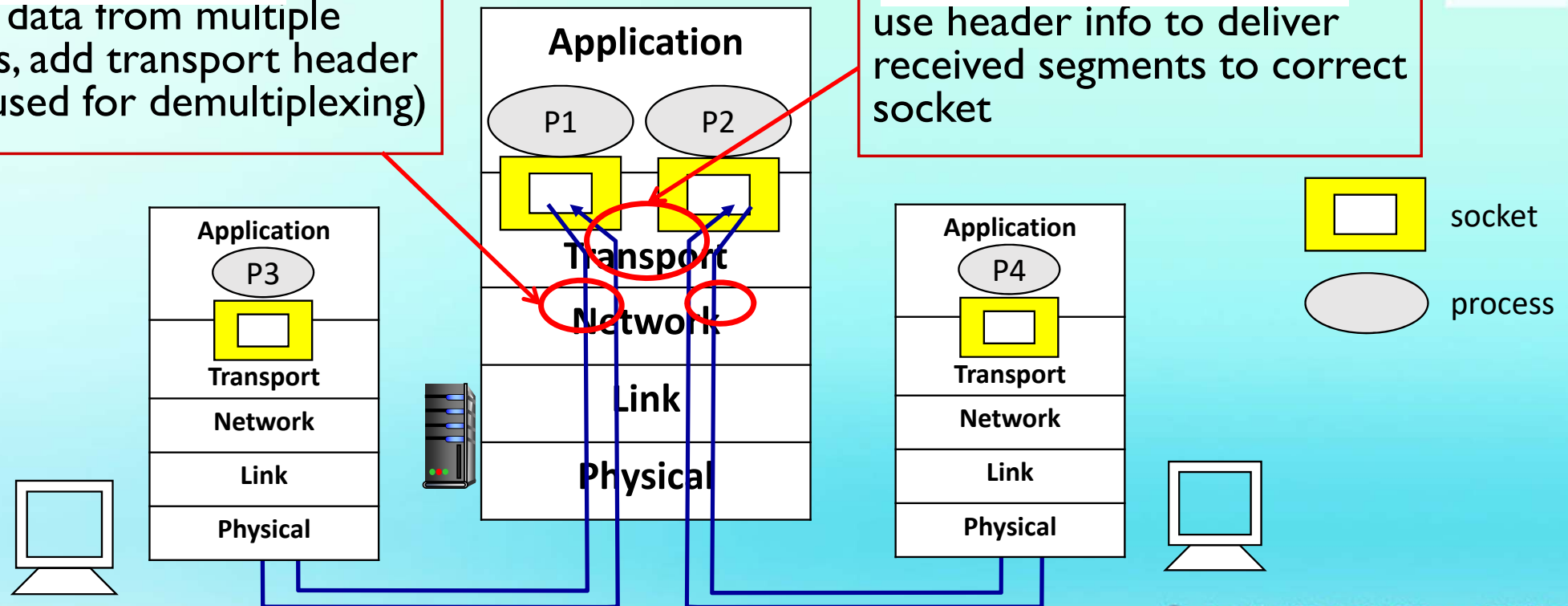
## 3.2 Multiplexing and Demultiplexing

### *multiplexing at sender:*

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

### *demultiplexing at receiver:*

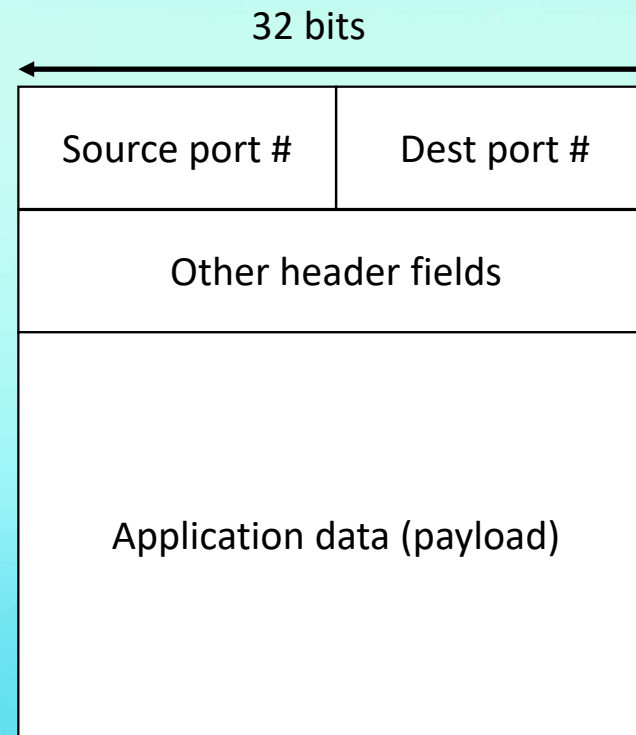
use header info to deliver received segments to correct socket



## 3.2 Multiplexing and Demultiplexing

### *How demultiplexing works*

- host receives IP datagrams
  - each datagram has source IP address, destination IP address
  - each datagram carries one transport-layer segment
  - each segment has source, destination port number
- host uses *IP addresses & port numbers* to direct segment to appropriate socket



TCP/UDP segment format



## 3.2 Multiplexing and Demultiplexing

### *Connectionless demultiplexing*

- *recall*: created socket has host-local port #:

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

- *recall*: when creating datagram to send into UDP socket, must specify

- destination IP address
- destination port #

- when host receives UDP segment:

- checks destination port # in segment
- directs UDP segment to socket with that port #

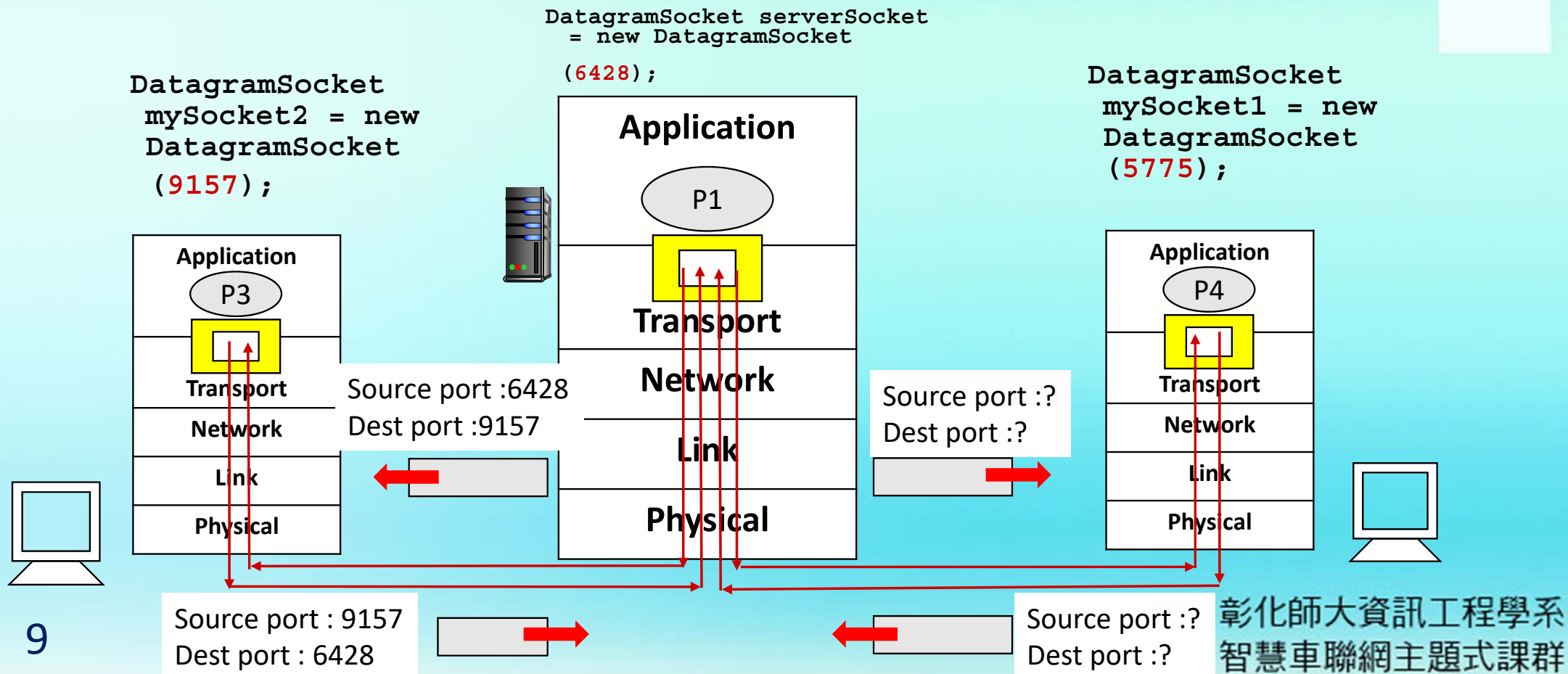


IP datagrams with *same dest. port #*, but different source IP addresses and/or source port numbers will be directed to *same socket* at dest



## 3.2 Multiplexing and Demultiplexing

### *Connectionless demultiplexing : example*



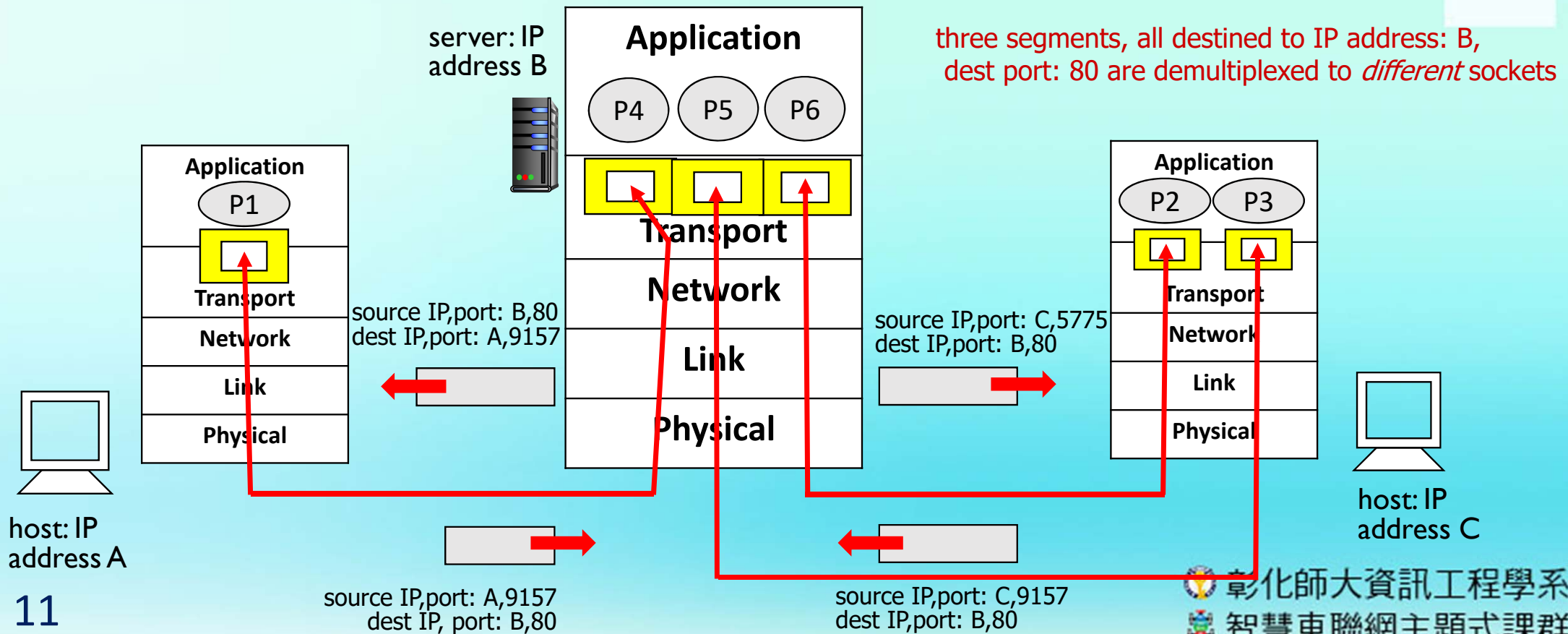
## 3.2 Multiplexing and Demultiplexing

### *Connection-oriented demux*

- TCP socket identified by 4-tuple:
  - source IP address
  - source port number
  - dest IP address
  - dest port number
- demux: receiver 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

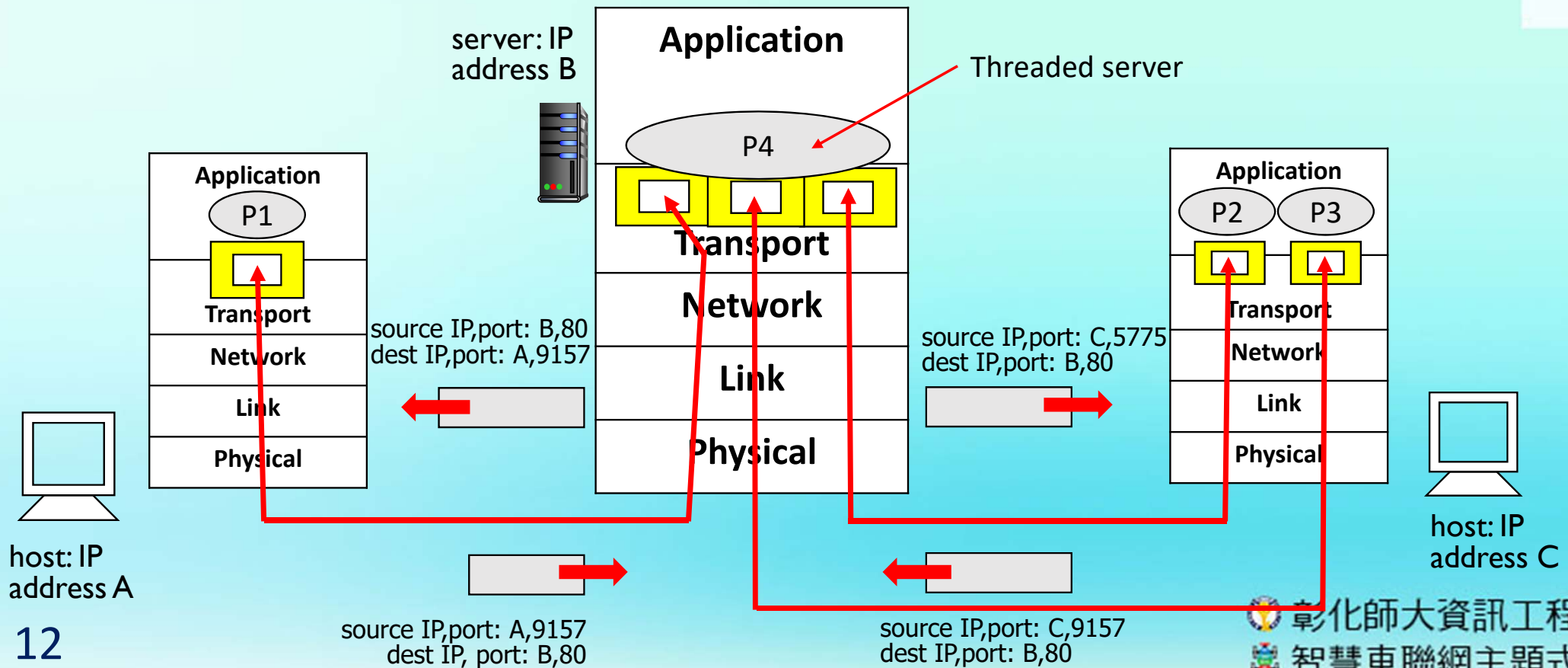
## 3.2 Multiplexing and Demultiplexing

### *Connection-oriented demux : example*



## 3.2 Multiplexing and Demultiplexing

### *Connection-oriented demux : example*



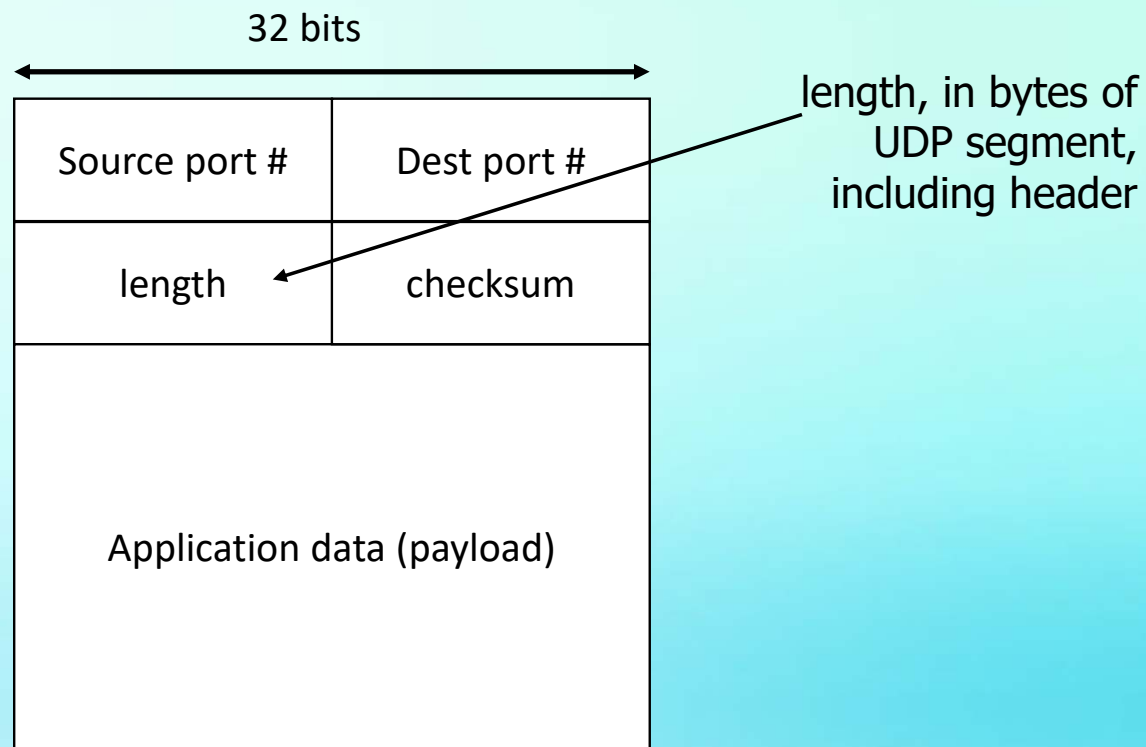
## 3.3 Connectionless Transport : UDP

### *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
- UDP use:
  - streaming multimedia apps (loss tolerant, rate sensitive)
  - DNS
  - SNMP
- reliable transfer over UDP:
  - add reliability at application layer
  - application-specific error recovery!

## 3.3 Connectionless Transport : UDP

### *UDP Segment Structure*



UDP segment format

#### why is there a UDP?

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



## 3.3 Connectionless Transport : UDP

### UDP Checksum

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

#### sender:

- treat segment contents, including header fields, as sequence of 16-bit integers
- checksum: addition (one'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 ....*



## 3.3 Connectionless Transport : UDP

### *UDP Checksum - example*

example: add two 16-bit integers

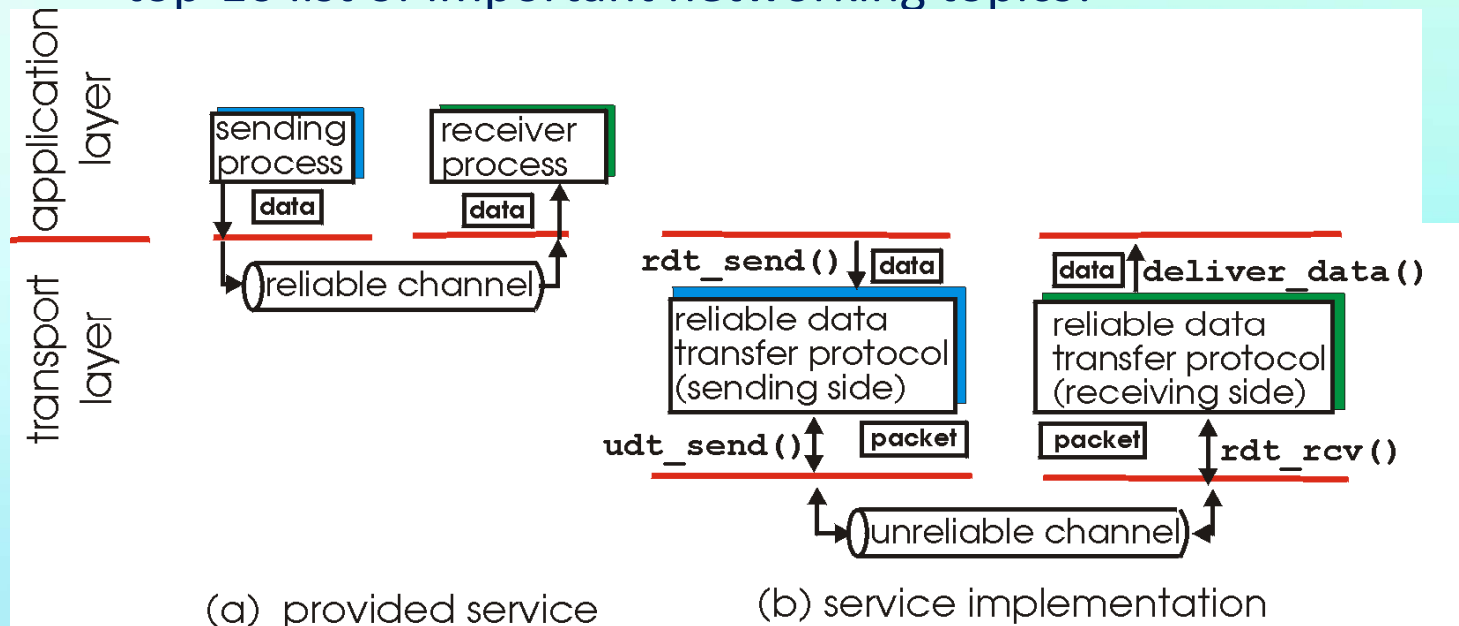
	1	1	1	1	0	0	1	1	0	0	1	1	0	0	1	1	0
	1	1	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1
<hr/>																	
wraparound	1	1	0	1	1	1	0	1	1	1	0	1	1	1	0	1	1
<hr/>																	
sum	1	1	0	1	1	1	0	1	1	1	0	1	1	1	1	0	0
checksum	1	0	1	0	0	0	1	0	0	0	1	0	0	0	0	1	1

Note: when adding numbers, a carryout from the most significant bit needs to be added to the result

## 3.4 Principles of Reliable Data Transfer

### *Building a Reliable Data Transfer Protocol*

- ❖ important in application, transport, link layers
  - top-10 list of important networking topics!



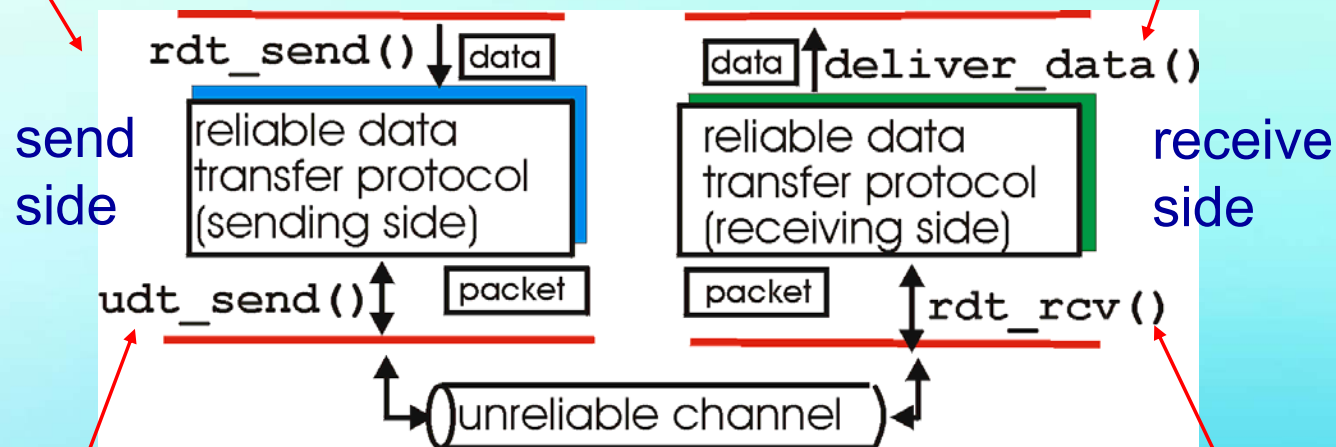
- ❖ characteristics of unreliable channel will determine complexity of reliable data transfer protocol (rdt)

## 3.4 Principles of Reliable Data Transfer

### *Building a Reliable Data Transfer Protocol*

**rdt\_send()** : called from above, (e.g., by app.). Passed data to deliver to receiver upper layer

**deliver\_data()** : called by **rdt** to deliver data to upper



**udt\_send()** : called by rdt, to transfer packet over unreliable channel to receiver

**rdt\_rcv()** : called when packet arrives on rcv-side of channel

## 3.4 Principles of Reliable Data Transfer

*rdt1.0: reliable transfer over a reliable channel*

- underlying channel perfectly reliable
  - no bit errors
  - no loss of packets

## 3.4 Principles of Reliable Data Transfer

### *rdt2.0: channel with bit errors*

- ❖ underlying channel may flip bits in packet
  - checksum to detect bit errors
- ❖ *the question*: how to recover from errors:

*How do humans recover from “errors”  
during conversation?*

## 3.4 Principles of Reliable Data Transfer

### *rdt2.0: channel with bit errors*

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

## 3.4 Principles of Reliable Data Transfer

### *rdt2.0: fatal flaw*

#### what happens if ACK/NAK corrupted?

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

#### handling duplicates:

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

#### stop and wait

sender sends one packet,  
then waits for receiver  
response



## 3.4 Principles of Reliable Data Transfer

### *rdt2.1: discussion*

#### sender:

- seq # added to pkt
- two seq. #'s (0,1) will suffice. Why?
- must check if received ACK/NAK corrupted
- twice as many states
  - state must “remember” whether “expected” pkt should have seq # of 0 or 1

#### 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

## 3.4 Principles of Reliable Data Transfer

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

#### new assumption:

underlying channel can also lose packets (data, ACKs)

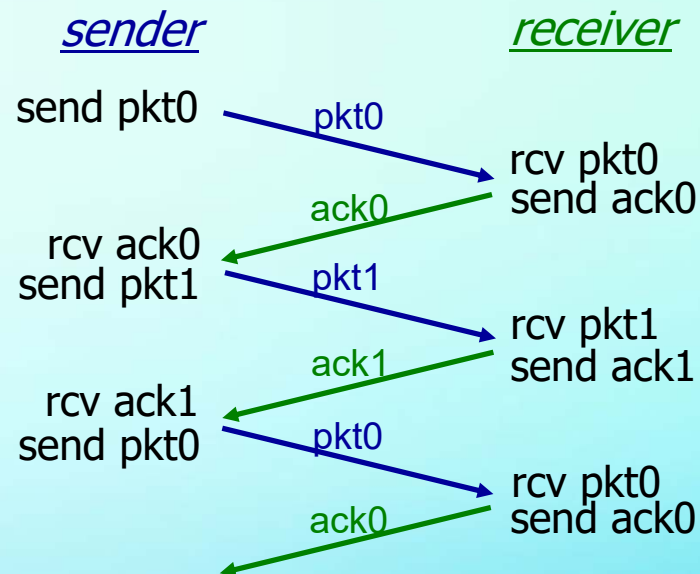
- checksum, seq. #, ACKs, retransmissions will be of help ... but not enough

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

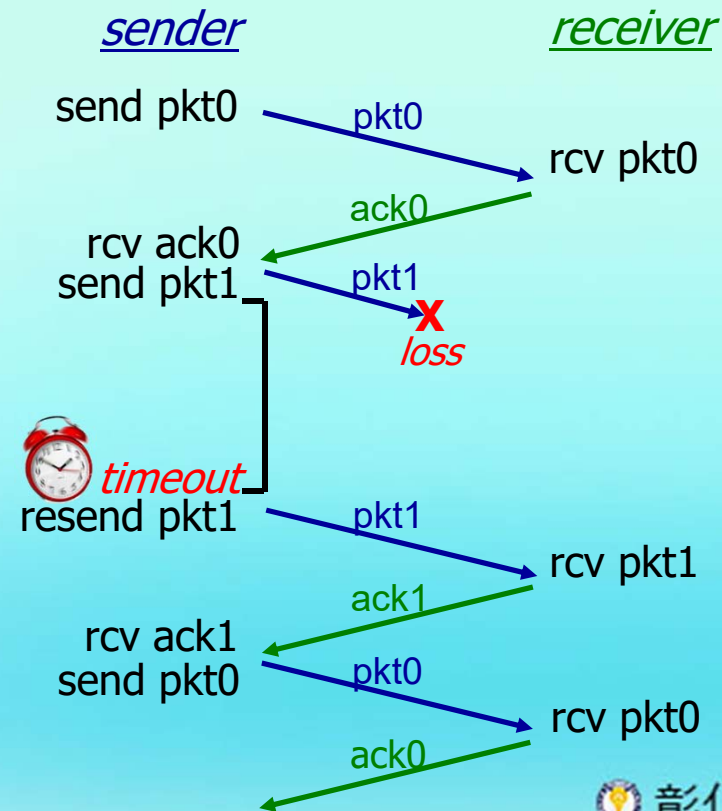
- retransmits if no ACK received in this time
- if pkt (or ACK) just delayed (not lost):
  - retransmission will be duplicate, but seq. #'s already handles this
  - receiver must specify seq # of pkt being ACKed
- requires countdown timer

## 3.4 Principles of Reliable Data Transfer

### *rdt3.0: in action*



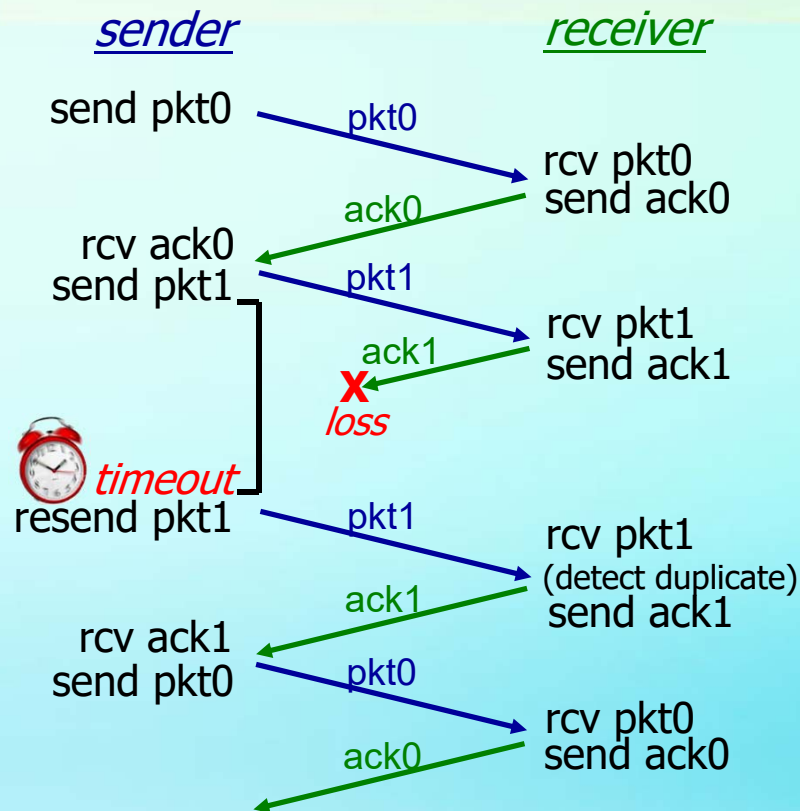
(a) no loss



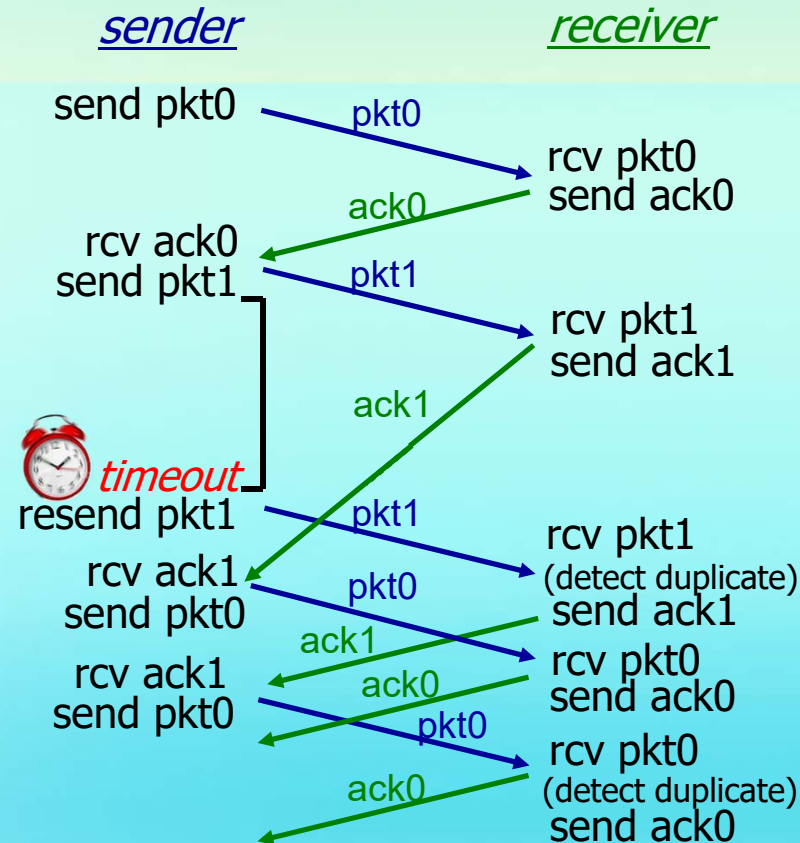
(b) packet loss

# 3.4 Principles of Reliable Data Transfer

## rdt3.0: in action



(c) ACK loss



(d) premature timeout/ delayed ACK

## 3.4 Principles of Reliable Data Transfer

### *rdt3.0: in action*

- rdt3.0 is correct, but performance stinks
- e.g.: 1 Gbps link, 15 ms prop. delay, 8000 bit packet:

$$D_{trans} = \frac{L}{R} = \frac{8000 \text{ bits}}{10^9 \text{ bits/sec}} = 8 \text{ microsecs}$$

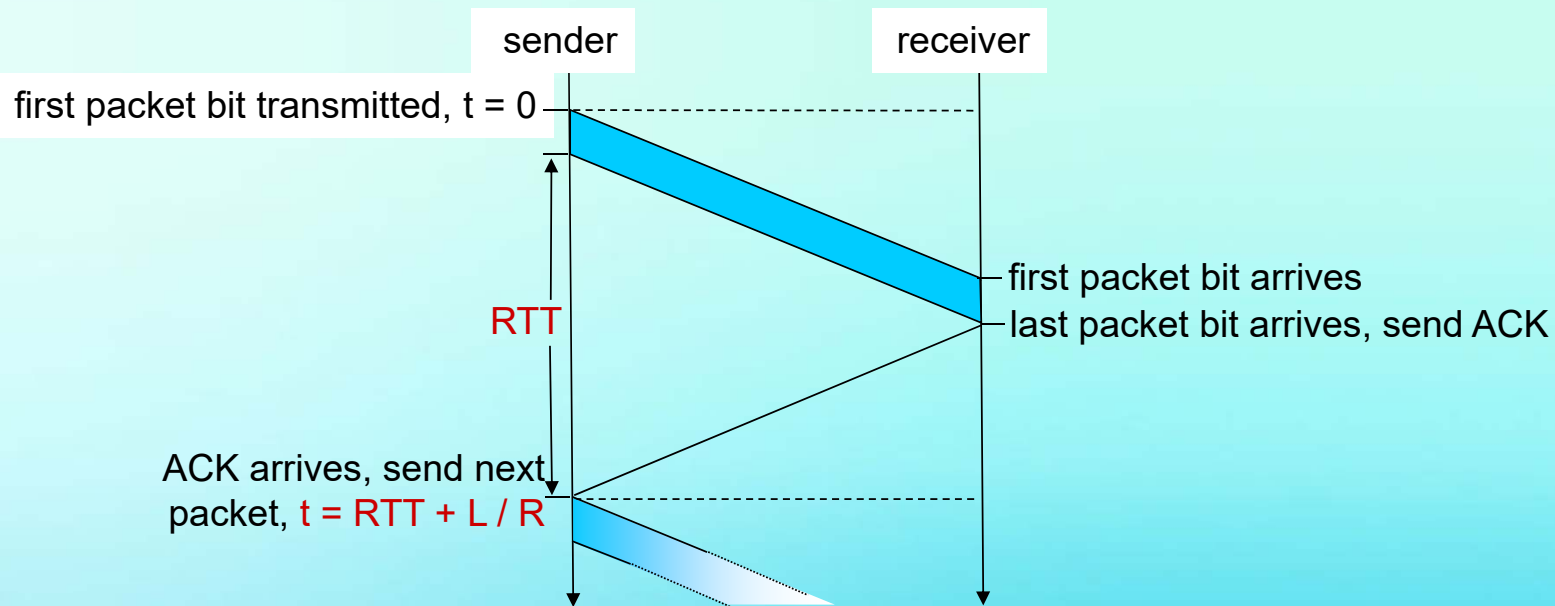
- $U_{sender}$ : **utilization** – fraction of time sender busy sending

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

- if RTT=30 msec, 1KB pkt every 30 msec: 33kB/sec thruput over 1 Gbps link
- network protocol limits use of physical resources!

## 3.4 Principles of Reliable Data Transfer

### *rdt3.0:stop-wait operation*



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

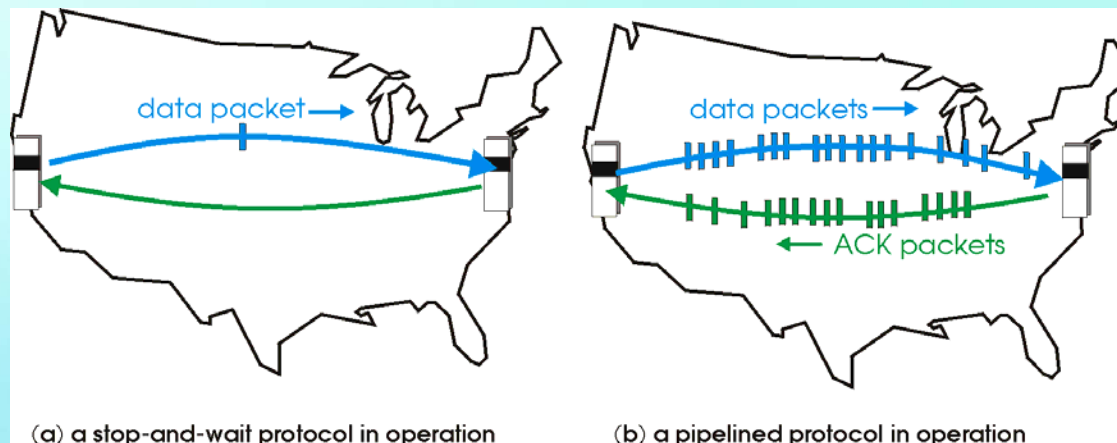


## 3.4 Principles of Reliable Data Transfer

### *Pipelined of Reliable Data Transfer*

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

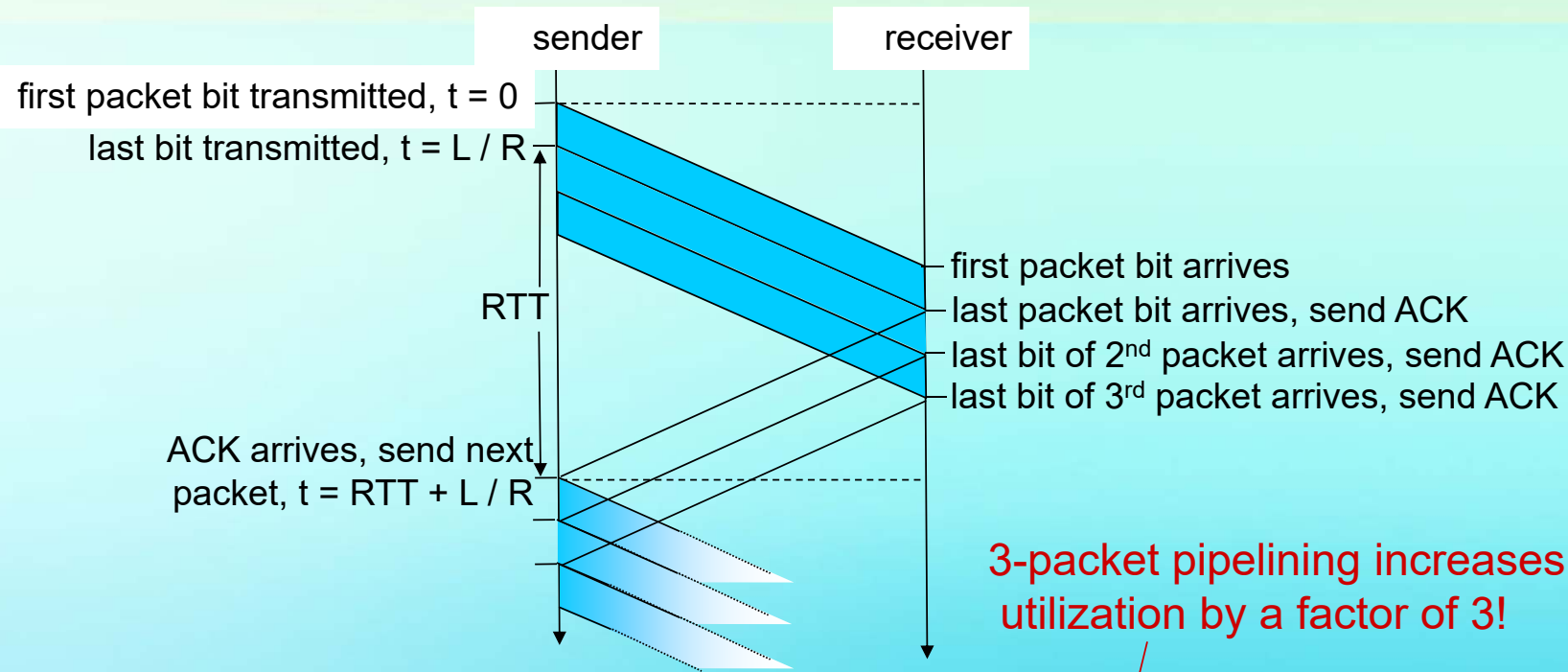
- range of sequence numbers must be increased
- buffering at sender and/or receiver



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



# Pipelined of Reliable Data Transfer : increased utilization



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

## 3.4 Principles of Reliable Data Transfer

### *Protocol overview*

#### Go-back-N:

- sender can have up to N unacked packets in pipeline
- receiver only sends *cumulative ack*
  - doesn't ack packet if there's a gap
- sender has timer for oldest unacked packet
  - when timer expires, retransmit *all* unacked packets

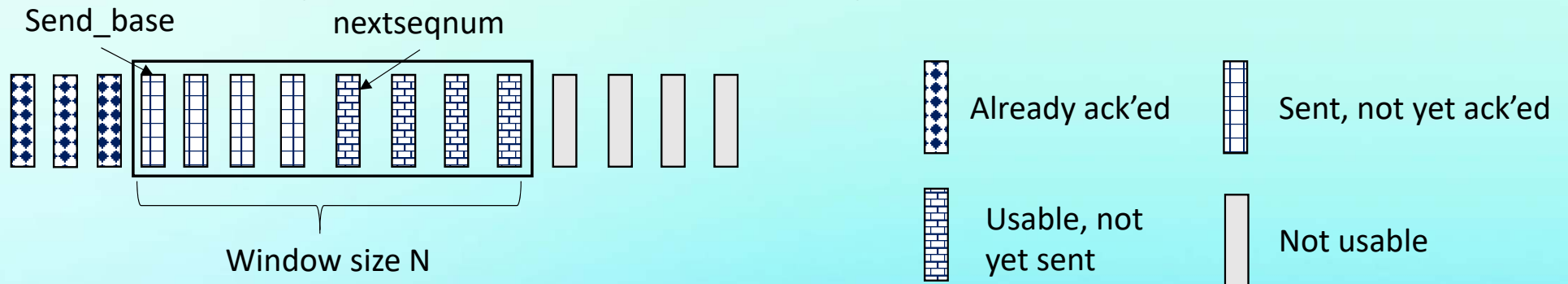
#### Selective Repeat:

- sender can have up to N unack'ed packets in pipeline
- rcvr sends *individual ack* for each packet
- sender maintains timer for each unacked packet
  - when timer expires, retransmit only that unacked packet

## 3.4 Principles of Reliable Data Transfer

### *Go-Back-N: sender*

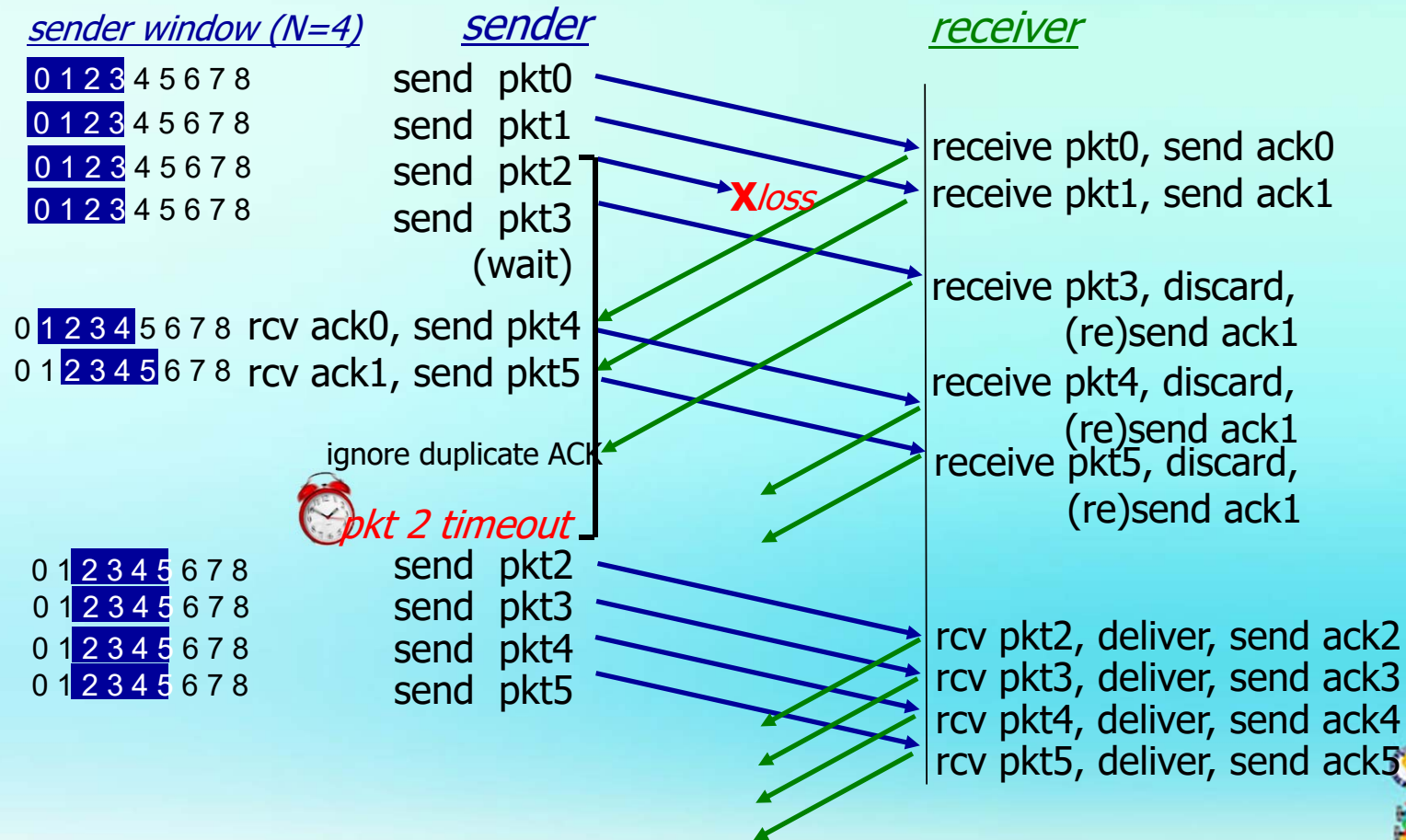
- k-bit seq # in pkt header
- “window” of up to N, consecutive unack’ed pkts allowed



- 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

# 3.4 Principles of Reliable Data Transfer

## Go-Back-N: in action



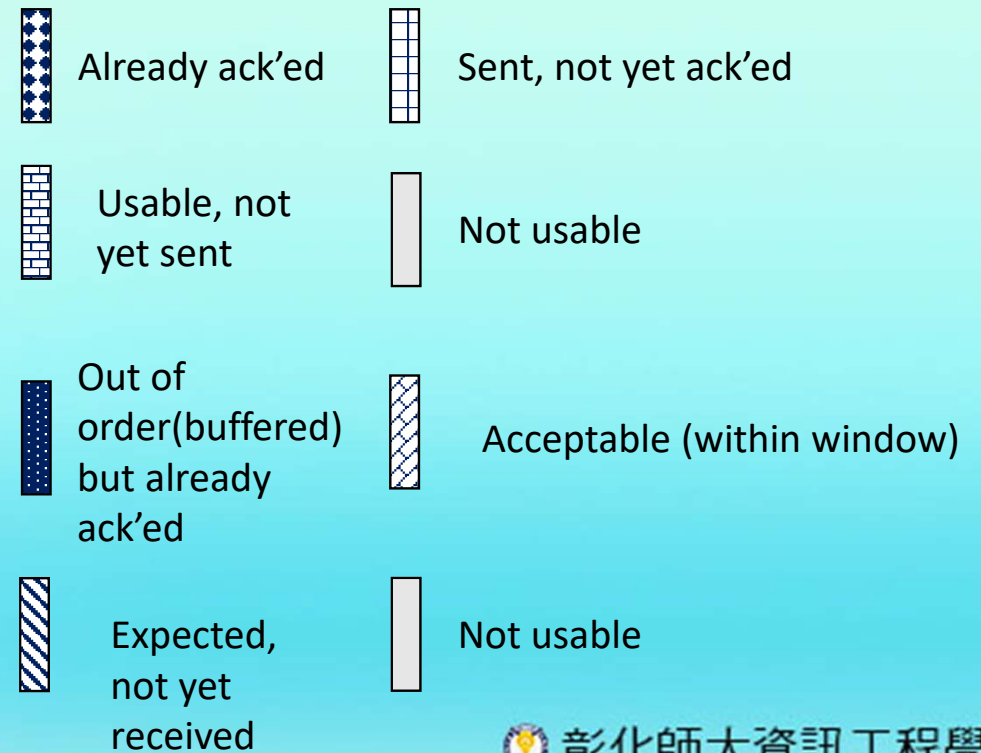
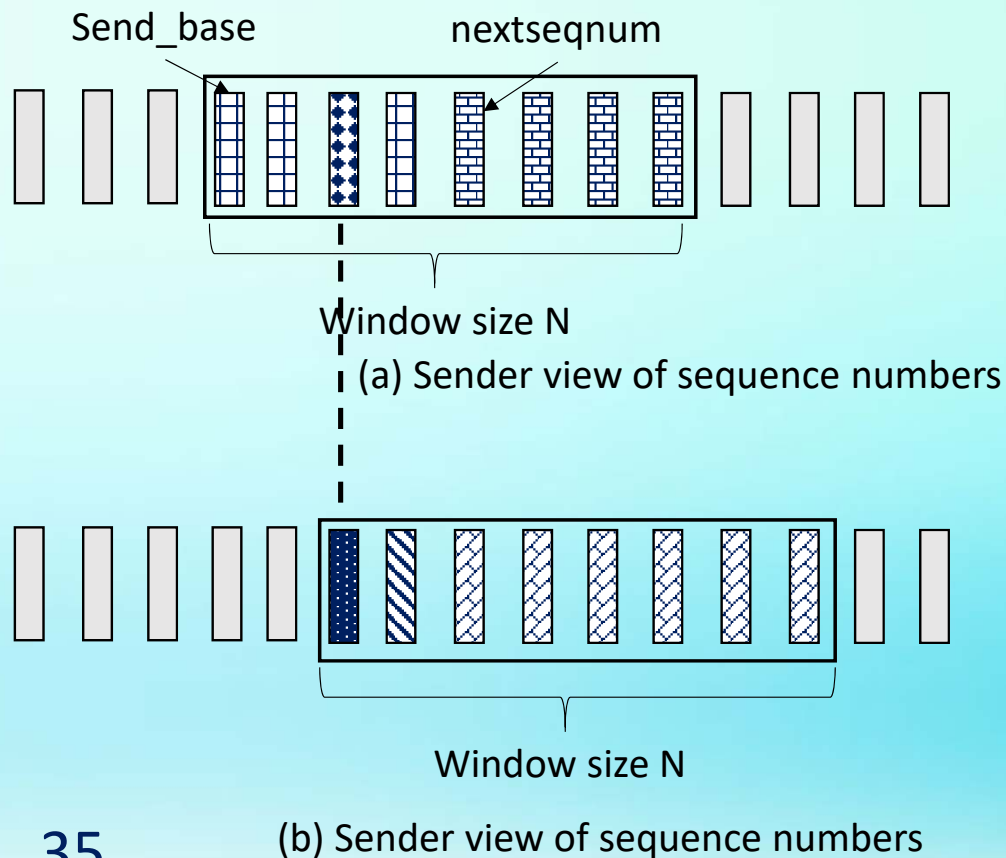
## 3.4 Principles of Reliable Data Transfer

### *Selective Repeat*

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

## 3.4 Principles of Reliable Data Transfer

### *Selective Repeat: sender, receiver windows*





## 3.4 Principles of Reliable Data Transfer

### *Selective Repeat*

#### sender

data from above:

- if next available seq # in window, send pkt

timeout(n):

- resend pkt n, restart timer

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

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

#### receiver

pkt n in [rcvbase, rcvbase+N-1]

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

pkt n in [rcvbase-N, rcvbase-1]

- ACK(n)

otherwise:

- ignore



## 3.4 Principles of Reliable Data Transfer

### *Selective Repeat : in action*

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 ack4 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?*



彰化師大資訊工程學系



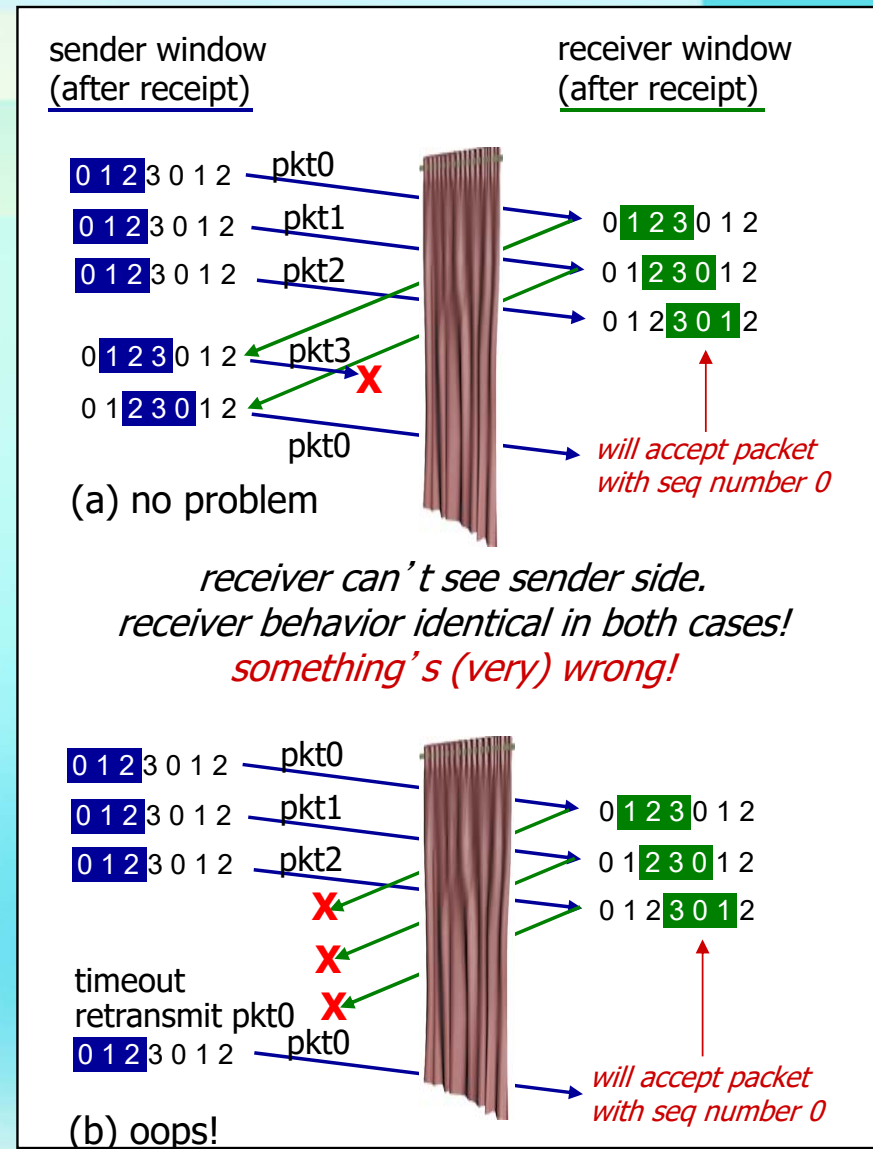
智慧車聯網主題式課群

## 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)?



## 3.5 Connection-oriented transport : TCP

### *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

#### ➤ full duplex data:

- bi-directional data flow in same connection
- MSS: maximum segment size

#### ➤ connection-oriented:

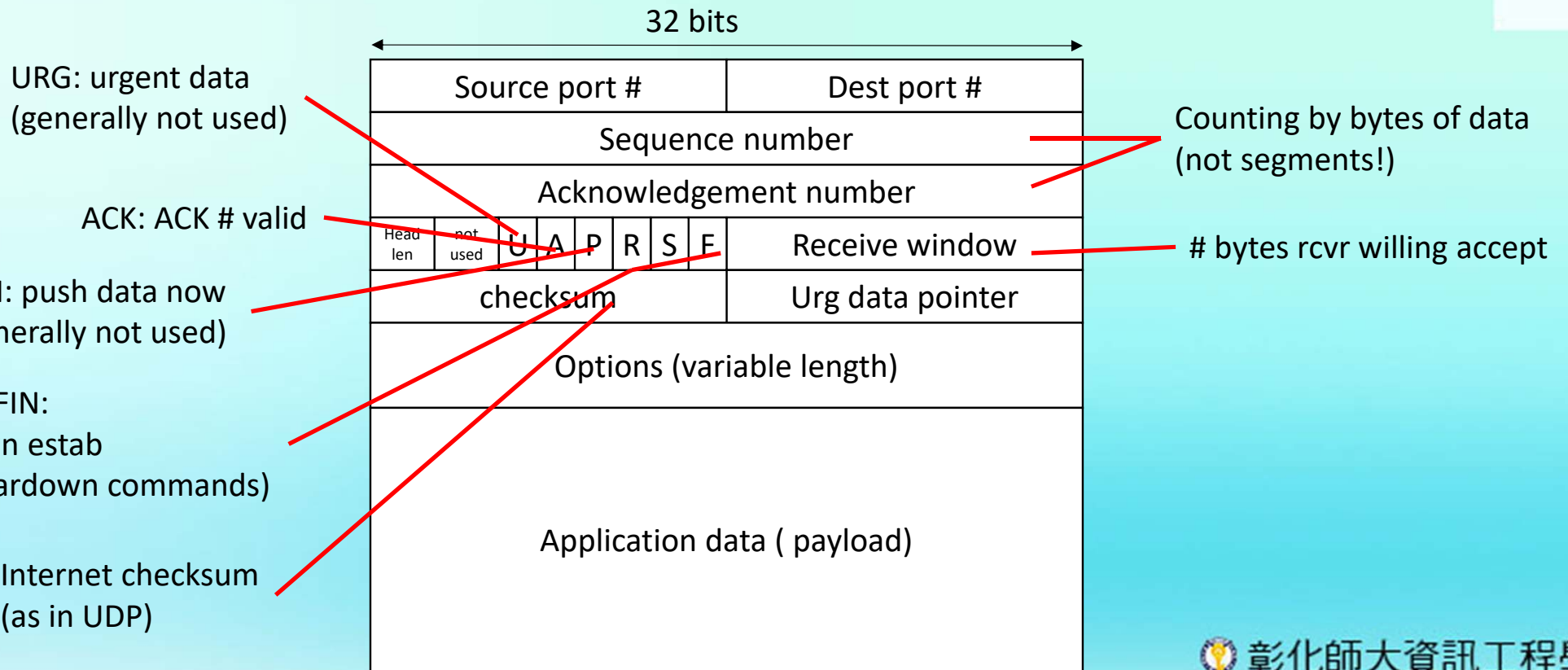
- handshaking (exchange of control msgs) initiates sender, receiver state before data exchange

#### ➤ flow controlled:

- sender will not overwhelm receiver

## 3.5 Connection-oriented transport : TCP

### *TCP: segment structure*



## 3.5 Connection-oriented transport : TCP

### *TCP: seq. numbers, ACKs*

#### sequence numbers:

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

#### acknowledgements:

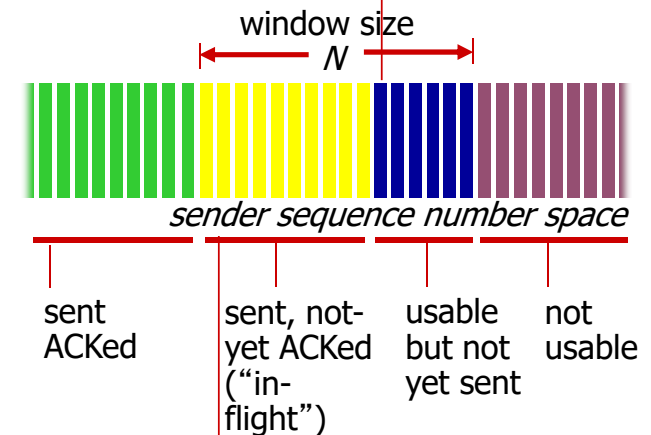
- seq # of next byte expected from other side
- cumulative ACK

**Q:** how receiver handles out-of-order segments

- **A:** TCP spec doesn’t say, - up to implementor

outgoing segment from sender

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

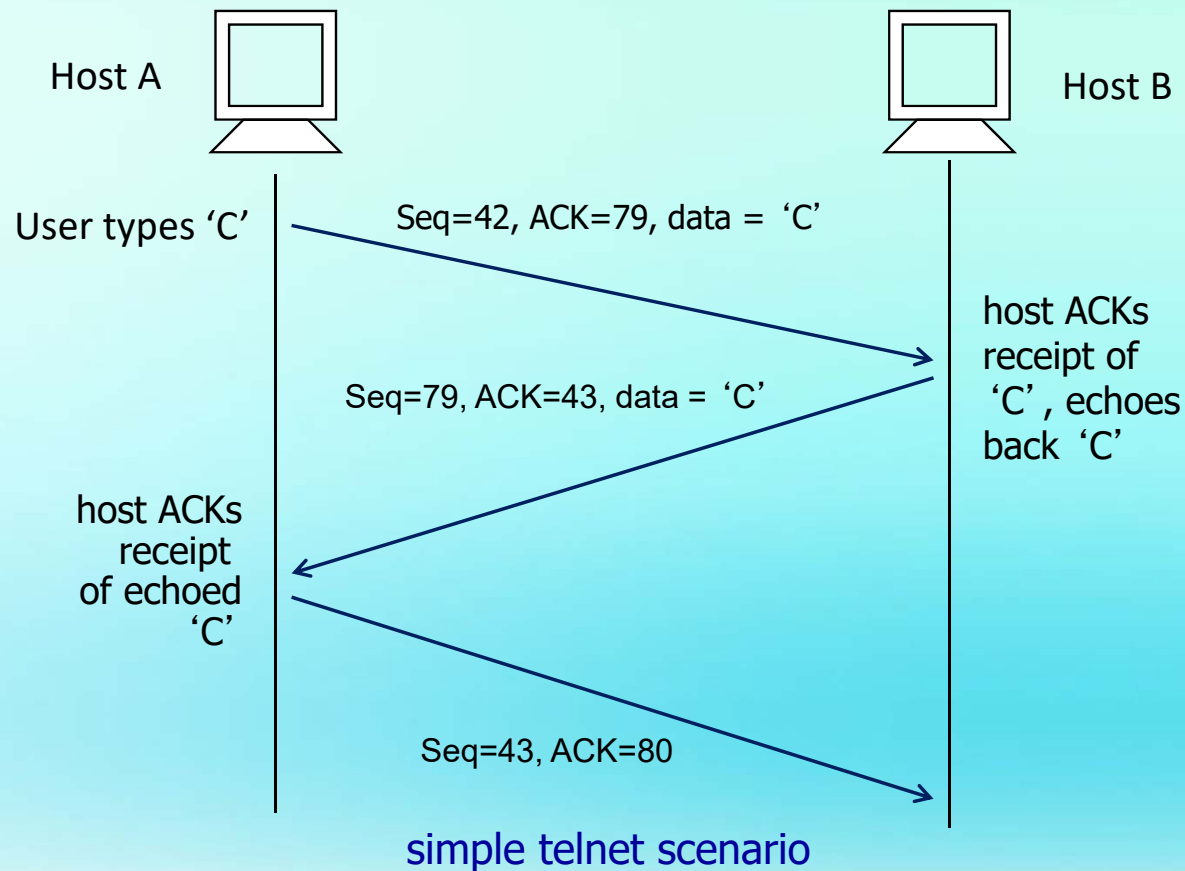


incoming segment to sender

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

## 3.5 Connection-oriented transport : TCP

*TCP: seq. numbers, ACKs*





## 3.5 Connection-oriented transport : TCP

*TCP: round trip time, 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
  - ignore retransmissions
- **SampleRTT** will vary, want estimated RTT “smoother”
  - average several *recent* measurements, not just current **SampleRTT**

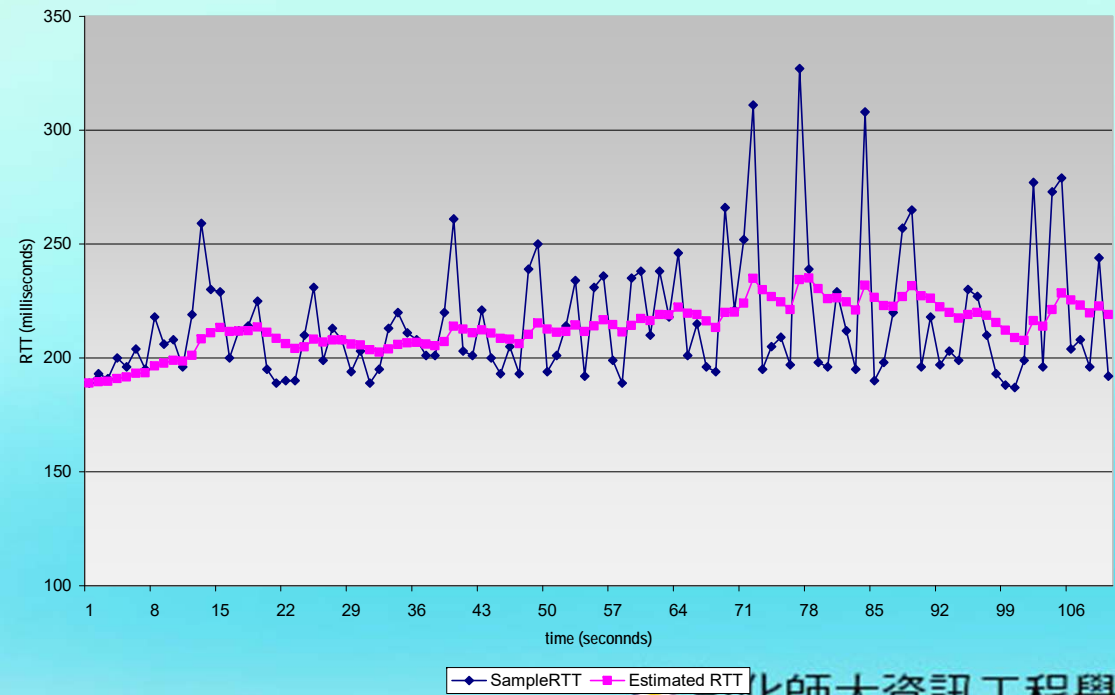
## 3.5 Connection-oriented transport : TCP

*TCP: round trip time, timeout*

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

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

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



## 3.5 Connection-oriented transport : TCP

*TCP: round trip time, timeout*

➤ timeout interval: **EstimatedRTT** plus “safety margin”

- large variation in **EstimatedRTT** -> larger safety margin

➤ estimate SampleRTT deviation from EstimatedRTT:

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

(typically,  $\beta = 0.25$ )

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



↑  
estimated RTT

↑  
“safety margin”

## 3.5 Connection-oriented transport : TCP

### *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

## 3.5 Connection-oriented transport : TCP

### *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:*

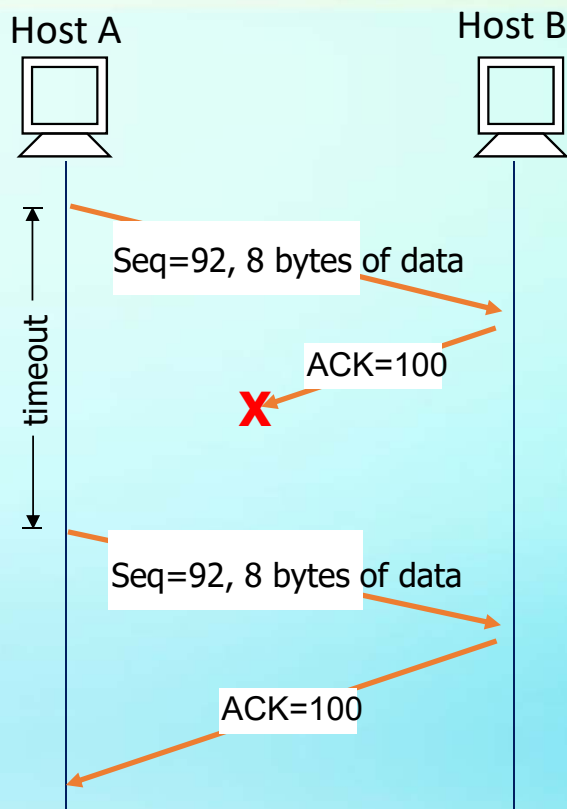
- retransmit segment that caused timeout
- restart timer

#### *ack rcvd:*

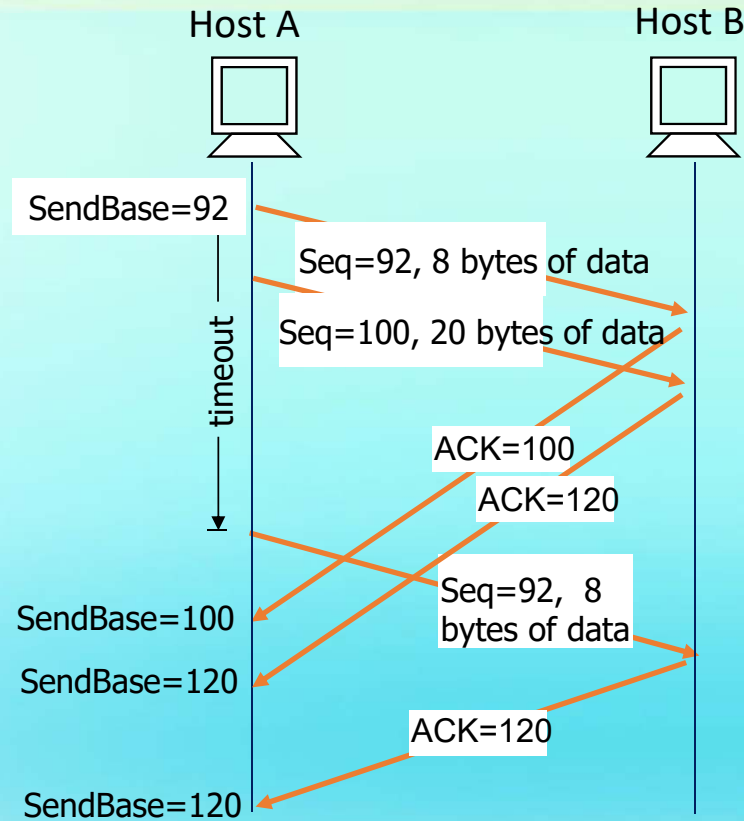
- if ack acknowledges previously unacked segments
  - update what is known to be ACKed
  - start timer if there are still unacked segments

## 3.5 Connection-oriented transport : TCP

### *TCP retransmission scenarios*



lost ACK scenario

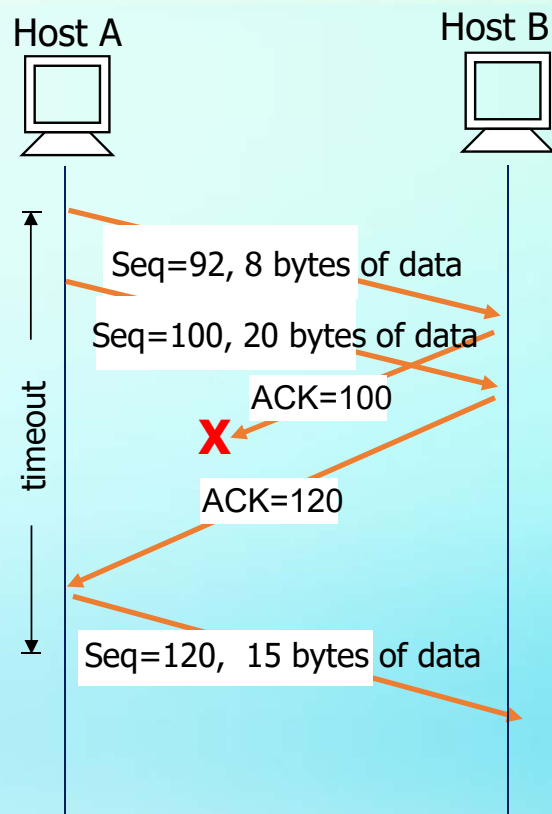


premature timeout



## 3.5 Connection-oriented transport : TCP

### *TCP retransmission scenarios*



cumulative ACK

## 3.5 Connection-oriented transport : TCP

### *TCP ACK generation [RFC 1122, RFC 2581]*

<i>event at receiver</i>	<i>TCP receiver action</i>
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 segment has ACK pending	immediately send single cumulative ACK, ACKing both in-order segments
arrival of out-of-order segment higher-than-expect seq. # . Gap detected	immediately send <i>duplicate ACK</i> , 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

## 3.5 Connection-oriented transport : TCP

### *TCP fast retransmit*

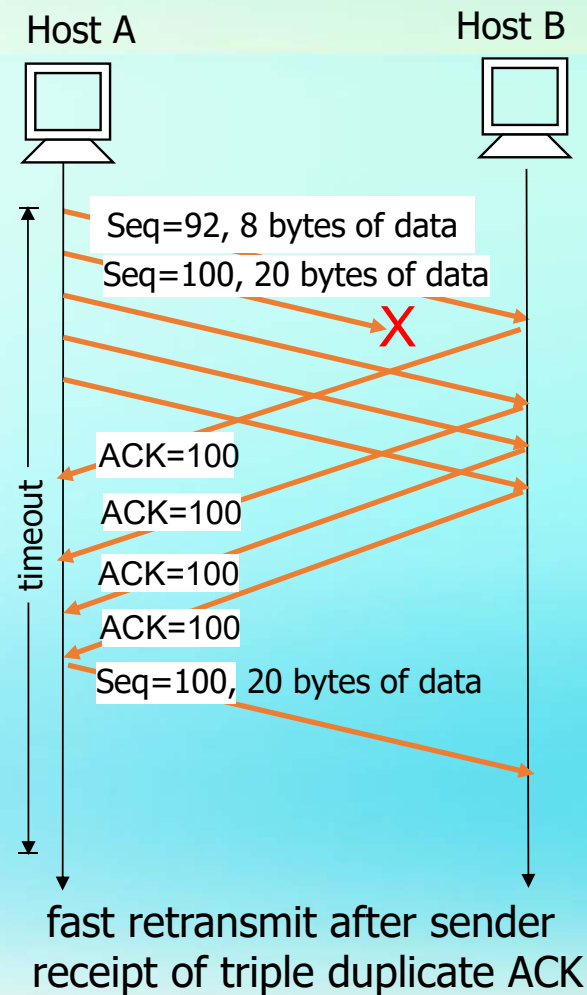
- time-out period often relatively long:
  - long delay before resending lost packet
- detect lost segments via duplicate ACKs.
  - sender often sends many segments back-to-back
  - if segment is lost, there will likely be many duplicate ACKs.

#### *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

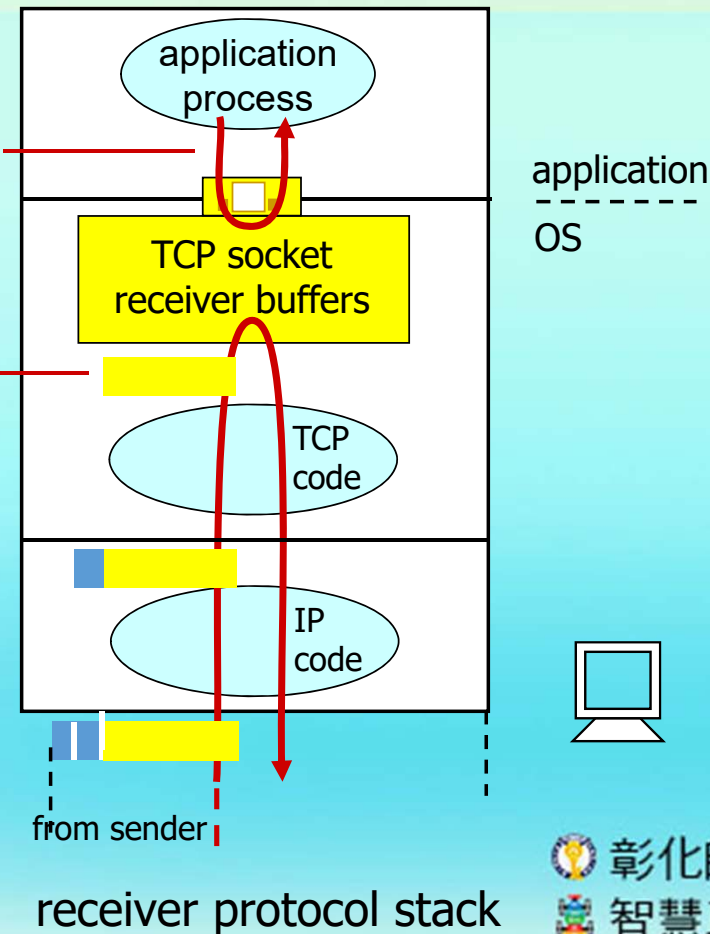


# TCP flow control

application may  
remove data from  
TCP socket buffers ....

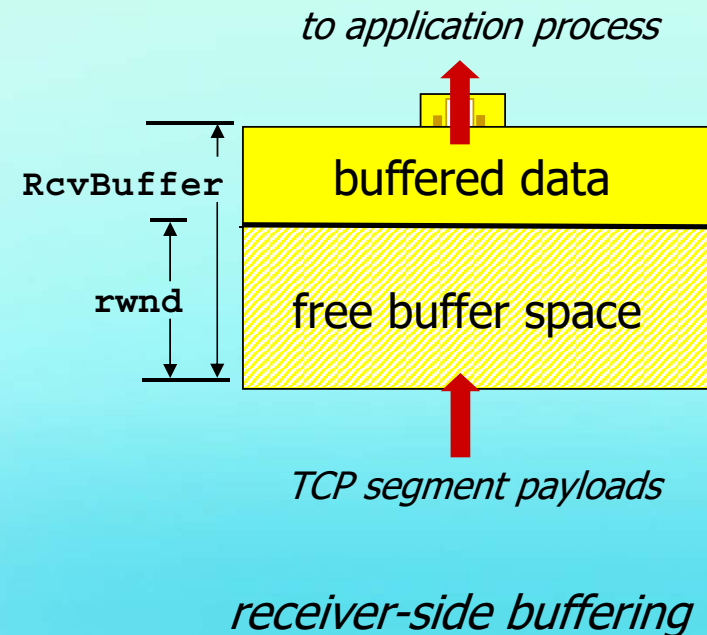
... slower than TCP  
receiver is delivering  
(sender is sending)

**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

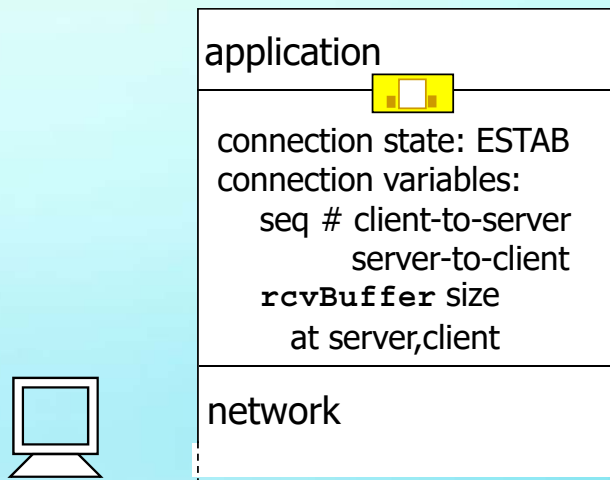




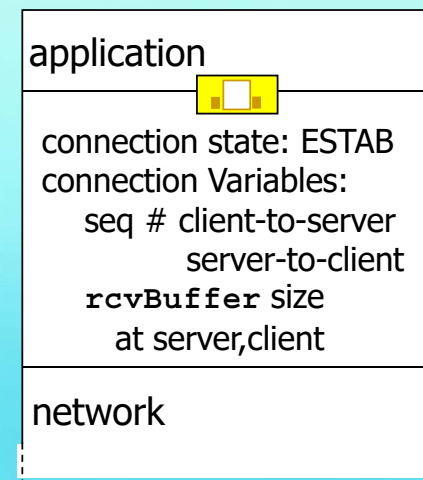
# connection management

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");
```



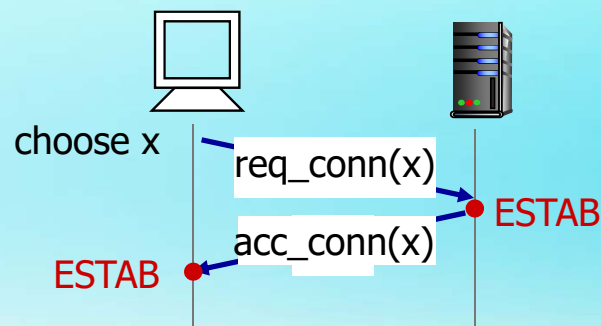
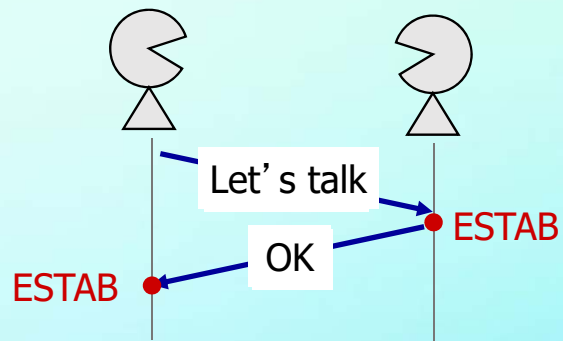
```
Socket connectionSocket =  
    welcomeSocket.accept();
```



# connection management

## agreeing to establish a connection

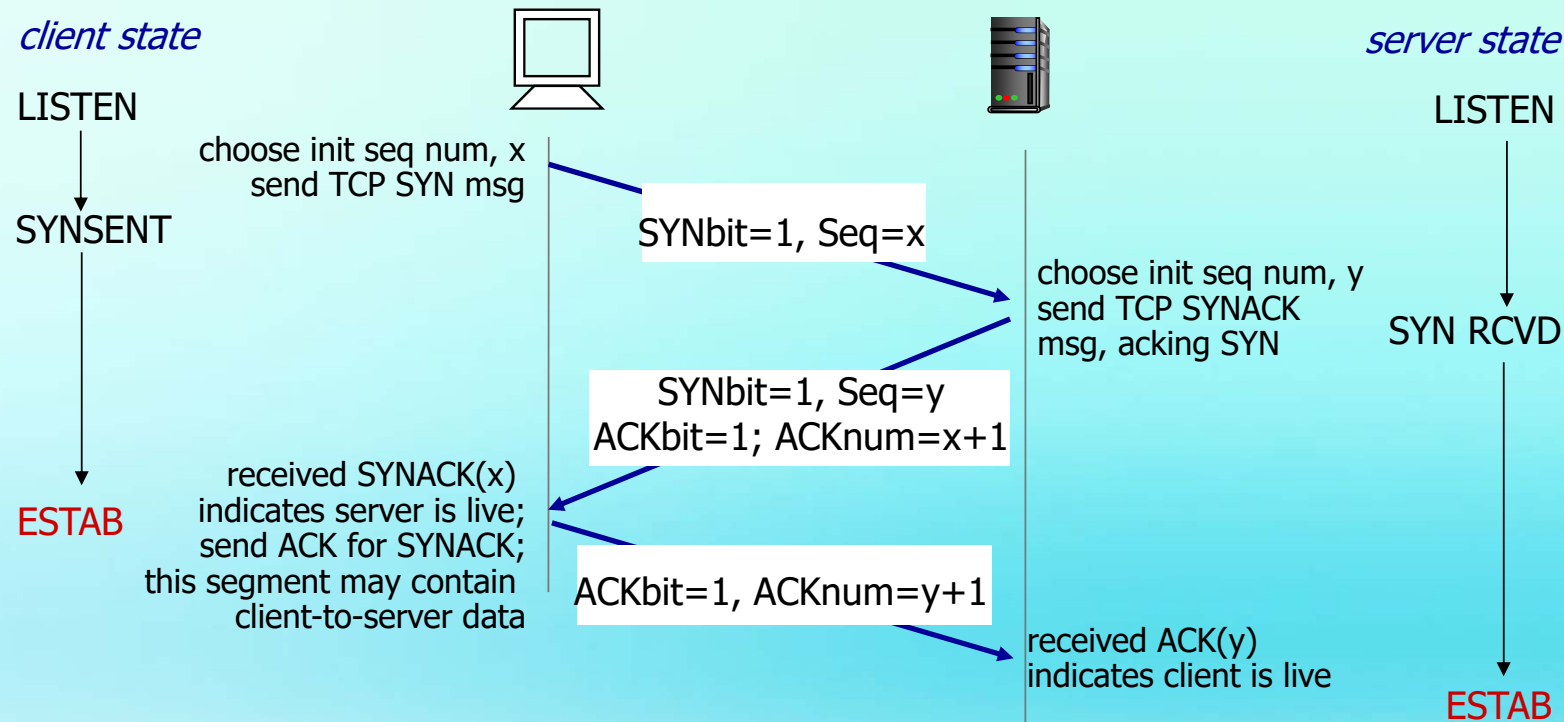
2-way handshake:



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

# TCP 3-way handshake



## *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

