

---

# Network Transport Layer: TCP

**Qiao Xiang, Congming Gao**

<https://sngroup.org.cn/courses/cnns-xmuf23/index.shtml>

**11/14/2023**

# Outline

---

- ❑ Admin and recap
- ❑ TCP Reliability

# Admin

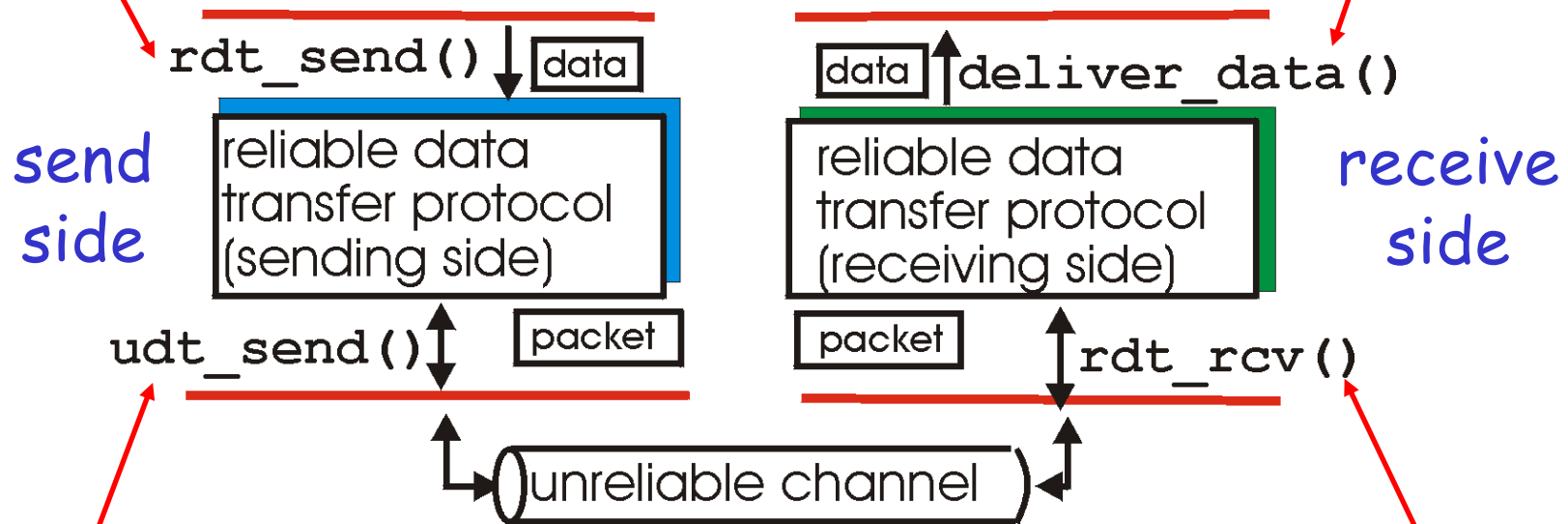
---

- Lab 3 due on Nov. 19
- Lab 4 to be posted this week

# Recap: Reliable Data Transfer Context

**rdt\_send()** : called from above,  
(e.g., by app.)

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

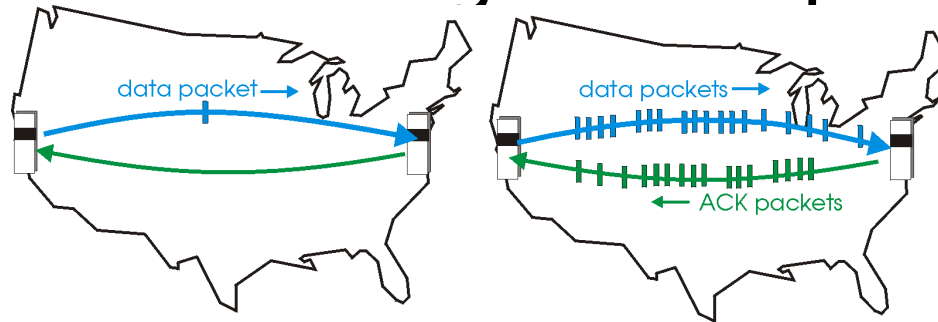


**udt\_send()** : called by rdt,  
to transfer packet over  
unreliable channel to receiver

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

# Recap: Reliable Transport

## Basic structure: sliding window protocols



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

(b) a pipelined protocol in operation

General technique: pipelining.

## Realization: GBN or SR

	Go-back-n	Selective Repeat
data bandwidth: sender to receiver (avg. number of times a pkt is transmitted)	Less efficient $\frac{1-p+pw}{1-p}$	More efficient $\frac{1}{1-p}$
ACK bandwidth (receiver to sender)	More efficient	Less efficient
Relationship between M (the number of seq#) and W (window size)	$M > W$	$M \geq 2W$
Buffer size at receiver	1	W
Complexity	Simpler	More complex

# Outline

---

- Admin and Recap
  - *TCP reliability*

# TCP Reliable Data Transfer

## □ Connection-oriented:

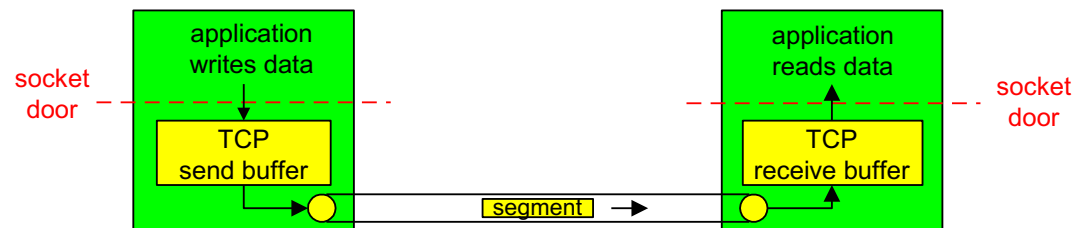
- connection management
  - setup (exchange of control msgs) init's sender, receiver state before data exchange
  - close

## □ Full duplex data:

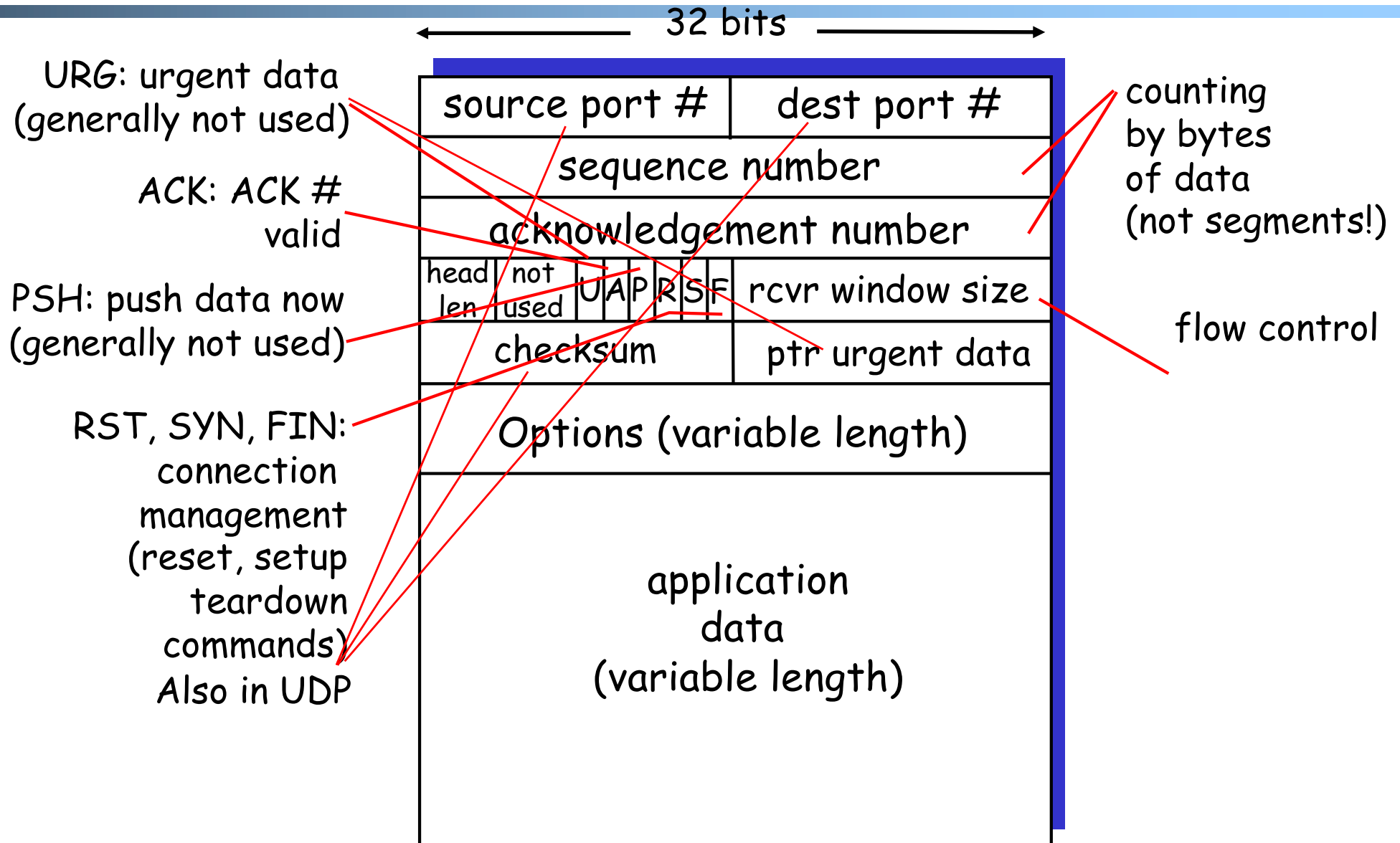
- bi-directional data flow in same connection

## □ A sliding window protocol

- a combination of go-back-n and selective repeat:
  - send & receive buffers
  - cumulative acks
  - TCP uses a single retransmission timer
  - do not retransmit all packets upon timeout



# TCP Segment Structure





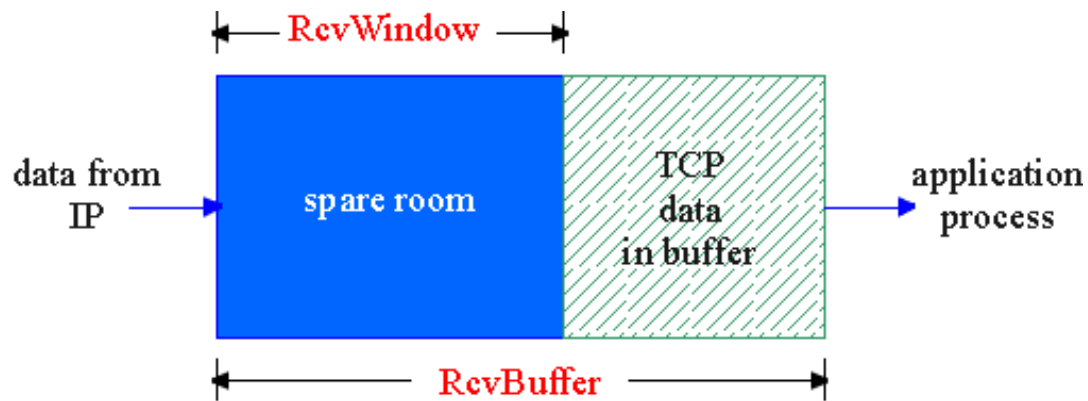
# Outline

---

- ❑ Admin and Recap
- ❑ Reliable data transfer
  - perfect channel
  - channel with bit errors
  - channel with bit errors and losses
  - sliding window: reliability with throughput
- ❑ TCP reliability
  - *data seq#, ack, buffering*

# Flow Control

- receive side of a connection has a receive buffer:



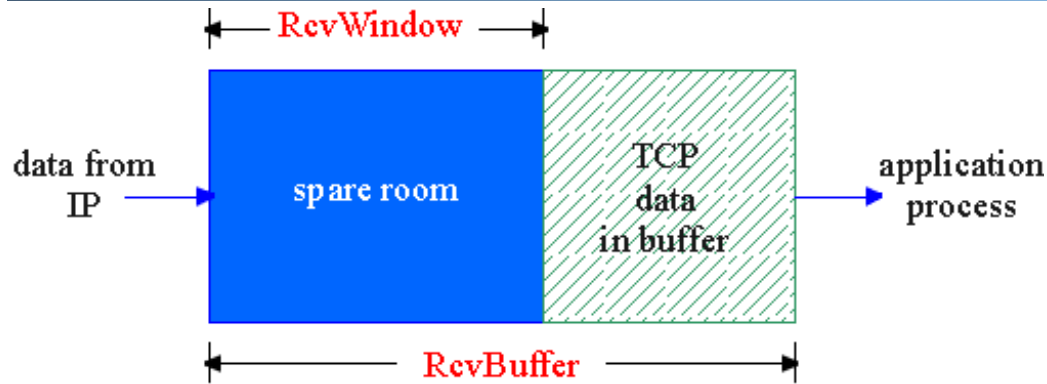
- app process may be slow at reading from buffer

## 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



□ spare room in buffer  
= **RcvWindow**

source port #		dest port #	
sequence number			
acknowledgement number			
head len	not used	U	A
P	R	S	F
checksum		ptr urgent data	
Options (variable length)			
application data (variable length)			

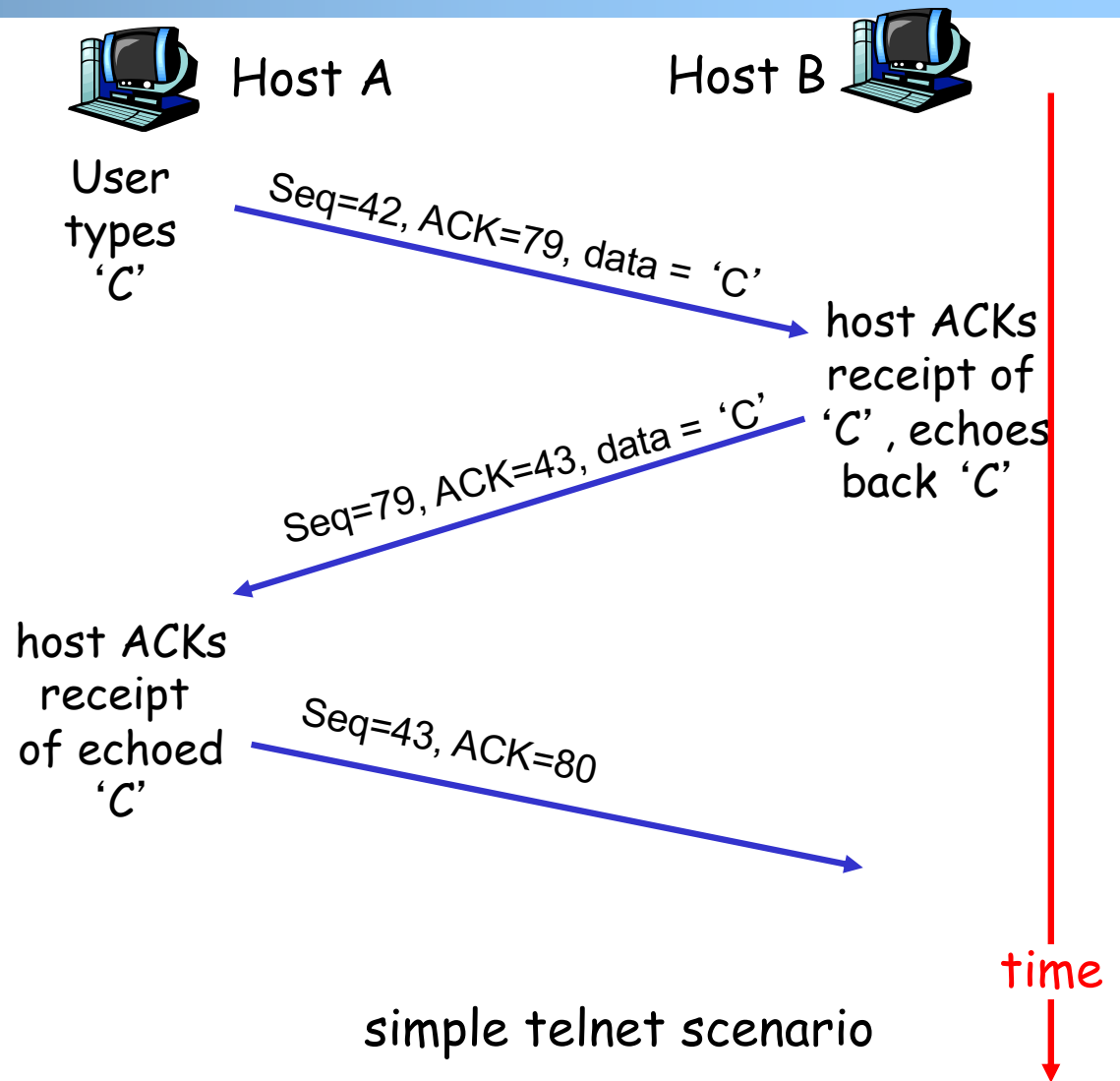
# TCP Seq. #'s and ACKs

## Seq. #'s:

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

## ACKs:

- seq # of next byte **expected** from other side
- **cumulative** ACK in standard header
- selective ACK in options



# TCP Send/Ack Optimizations

---

- ❑ TCP includes many tune/optimizations, e.g.,
  - the “small-packet problem”: sender sends a lot of small packets (e.g., telnet one char at a time)
    - Nagle’s algorithm: do not send data if there is small amount of data in send buffer and there is an unack’d segment
  - the “ack inefficiency” problem: receiver sends too many ACKs, no chance of combing ACK with data
    - Delayed ack to reduce # of ACKs/combine ACK with reply

# TCP Receiver 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 segment has ACK pending

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

# Outline

---

- ❑ Admin and Recap
- ❑ Reliable data transfer
  - perfect channel
  - channel with bit errors
  - channel with bit errors and losses
  - sliding window: reliability with throughput
- ❑ TCP reliability
  - data seq#, ack, buffering
  - *timeout realization*

# TCP Reliable Data Transfer

---

- ❑ Basic structure: sliding window protocol
- ❑ Remaining issue: How to determine the “right” parameters?
  - timeout value?
  - sliding window size?



# History

---

- ❑ Key parameters for TCP in mid-1980s
  - fixed window size  $W$
  - timeout value =  $2 \text{ RTT}$
  
- ❑ Network collapse in the mid-1980s
  - UCB  $\leftrightarrow$  LBL throughput dropped by 1000X !
  
- ❑ The intuition was that the collapse was caused by wrong parameters...

# Timeout: Cost of Timeout Param

---

Why is good timeout value important?

❑ too short

- premature timeout
- unnecessary retransmissions; many duplicates

❑ too long

- slow reaction to segment loss

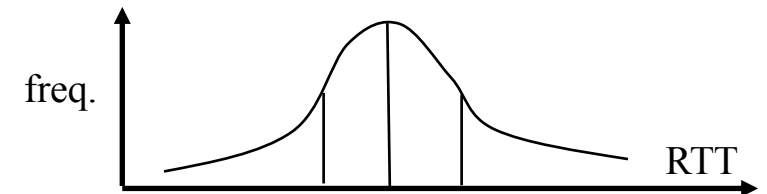
Q: Is it possible to set Timeout as a constant?

Q: Any problem w/ the early approach: Timeout = 2 RTT

# Setting Timeout

## Problem:

- ❑ Ideally, we set timeout = RTT, but RTT is not a fixed value  
=> using the average of RTT will generate many timeouts due to network variations
- ❑ Possibility: using the average/median of RTT
- ❑ Issue: this will generate many timeouts due to network variations



## Solution:

- ❑ Set Timeout RTO = avg + “safety margin” based on variation

## TCP approach:

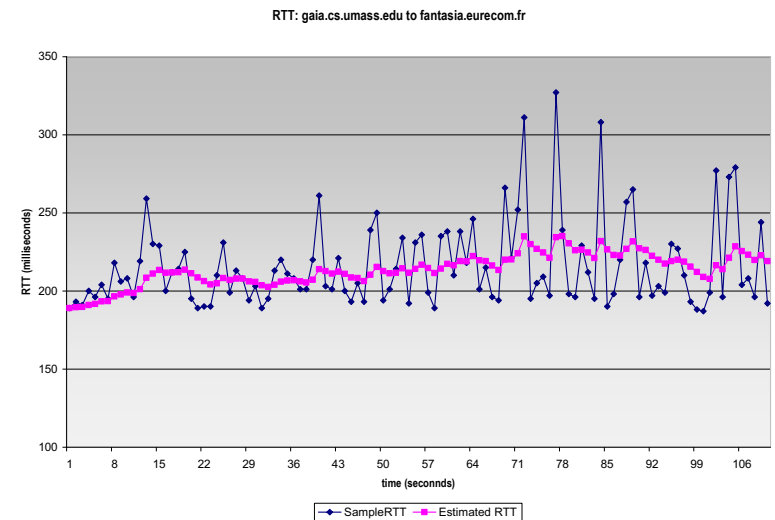
$$\text{Timeout} = \text{EstRTT} + 4 * \text{DevRTT}$$

# Compute EstRTT and DevRTT

- Exponential weighted moving average (EWMA)
  - influence of past sample decreases exponentially fast

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

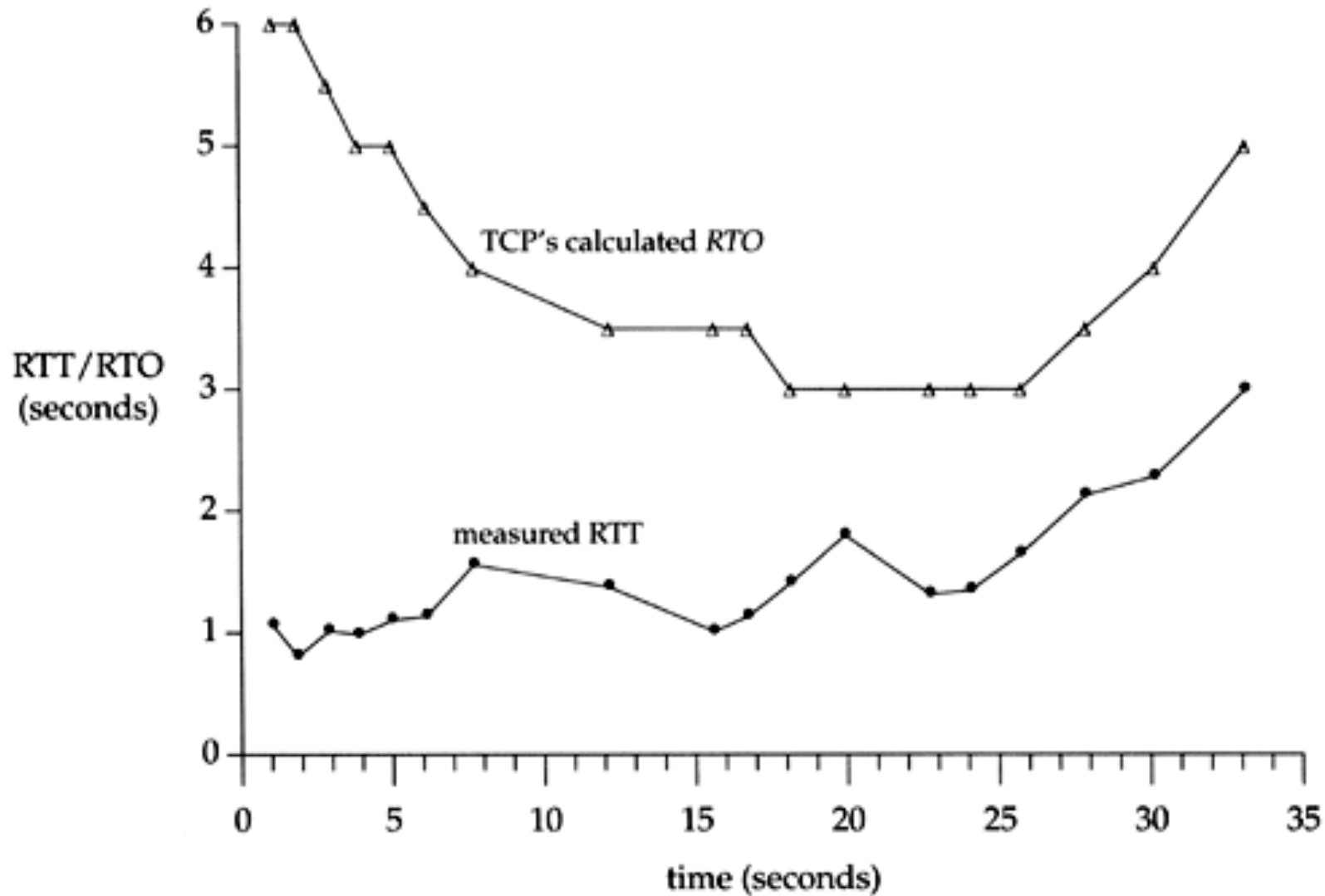
- **SampleRTT**: measured time from segment transmission until ACK receipt
- typical value:  $\alpha = 0.125$



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

(typically,  $\beta = 0.25$ )

# An Example TCP Session



# Fast Retransmit

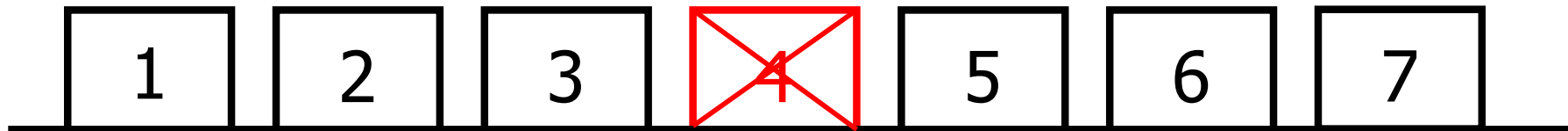
---

- ❑ Issue: Timeout period often relatively long:
  - long delay before resending lost packet
- ❑ Question: Can we detect loss faster than RTT?
  
- ❑ 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:
  - resend segment before timer expires

# Triple Duplicate Ack

---

Packets



Acknowledgements (waiting seq#)



# Fast Retransmit:

```
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



# TCP: reliable data transfer

Simplified  
TCP  
sender

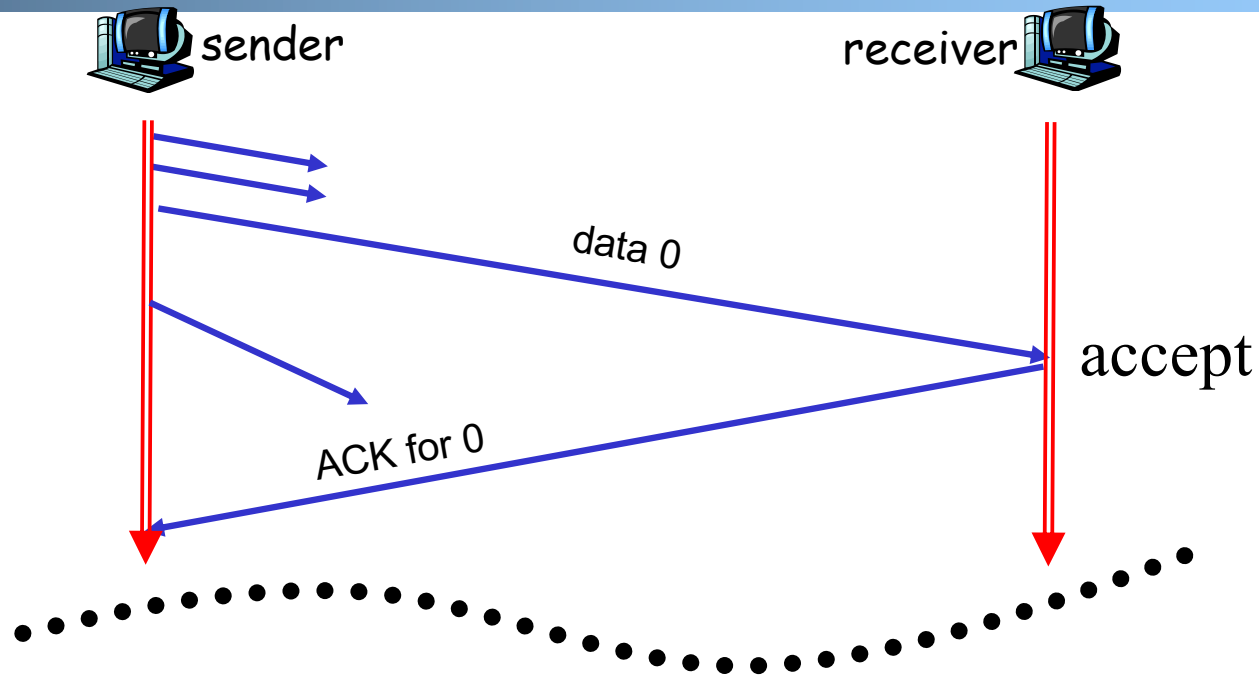
```
00 sendbase = initial_sequence number agreed by TWH
01 nextseqnum = initial_sequence number by TWH
02 loop (forever) {
03     switch(event)
04     event: data received from application above
05         if (window allows send)
06             create TCP segment with sequence number nextseqnum
06             if (no timer) start timer
07             pass segment to IP
08             nextseqnum = nextseqnum + length(data)
09             else put packet in buffer
09     event: timer timeout for sendbase
10         retransmit segment
11         compute new timeout interval
12         restart timer
13     event: ACK received, with ACK field value of y
14         if (y > sendbase) { /* cumulative ACK of all data up to y */
15             cancel the timer for sendbase
16             sendbase = y
17             if (no timer and packet pending) start timer for new sendbase
17             while (there are segments and window allow)
18                 sent a segment;
18         }
19         else { /* y==sendbase, duplicate ACK for already ACKed segment */
20             increment number of duplicate ACKs received for y
21             if (number of duplicate ACKS received for y == 3) {
22                 /* TCP fast retransmit */
23                 resend segment with sequence number y
24                 restart timer for segment y
25             }
26     } /* end of loop forever */
```

# Outline

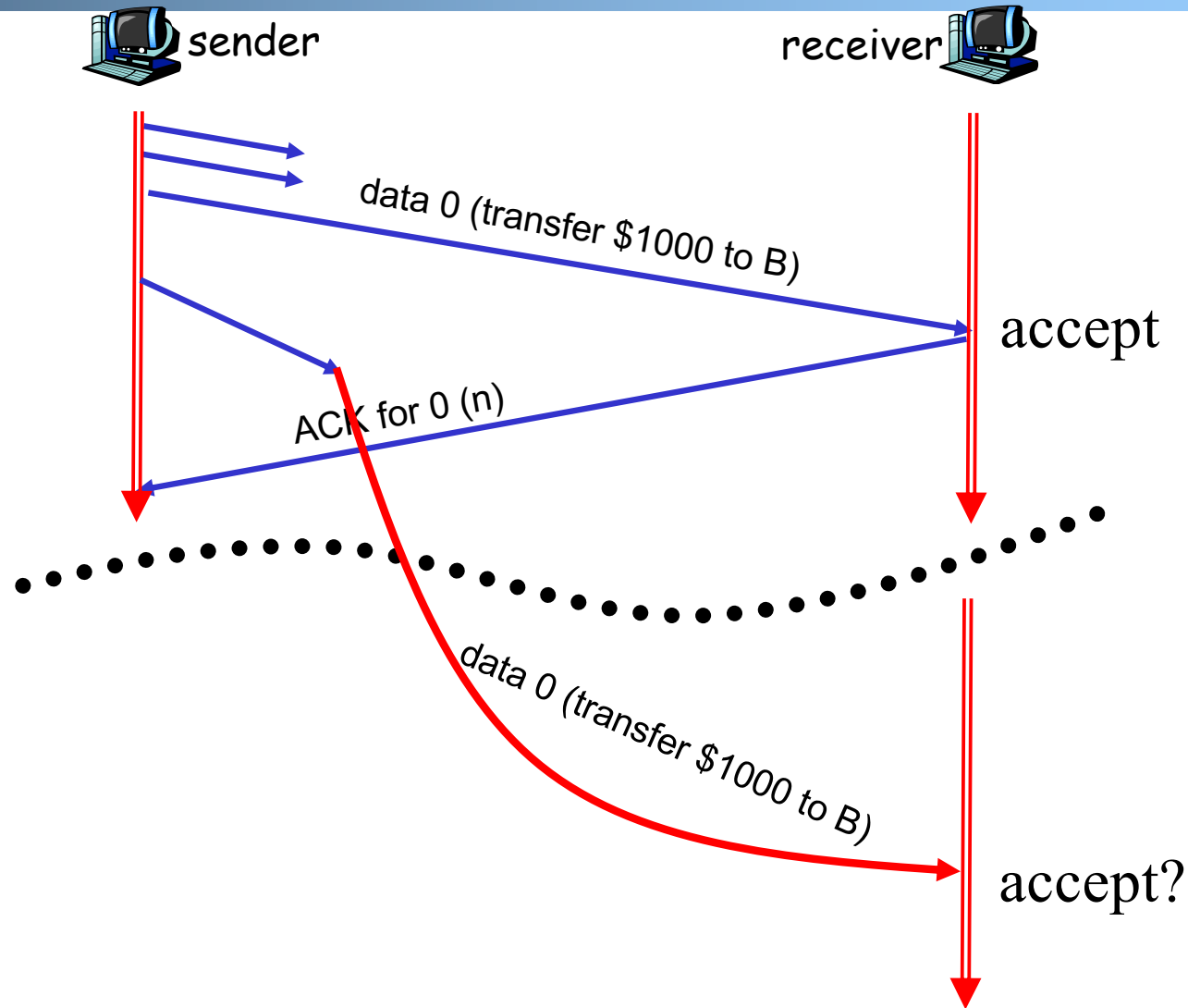
---

- ❑ Admin and Recap
- ❑ Reliable data transfer
  - perfect channel
  - channel with bit errors
  - channel with bit errors and losses
  - sliding window: reliability with throughput
- ❑ TCP reliability
  - data seq#, ack, buffering
  - timeout realization
  - *connection management*

# Why Connection Setup/When to Accept (Safely Deliver) First Packet?



# Why Connection Setup/When to Accept (Safely Deliver) First Packet?



# Transport "Safe-Setup" Principle

---

- A general safety principle for a receiver R to accept a message from a sender S is the general "authentication" principle, which consists of two conditions:

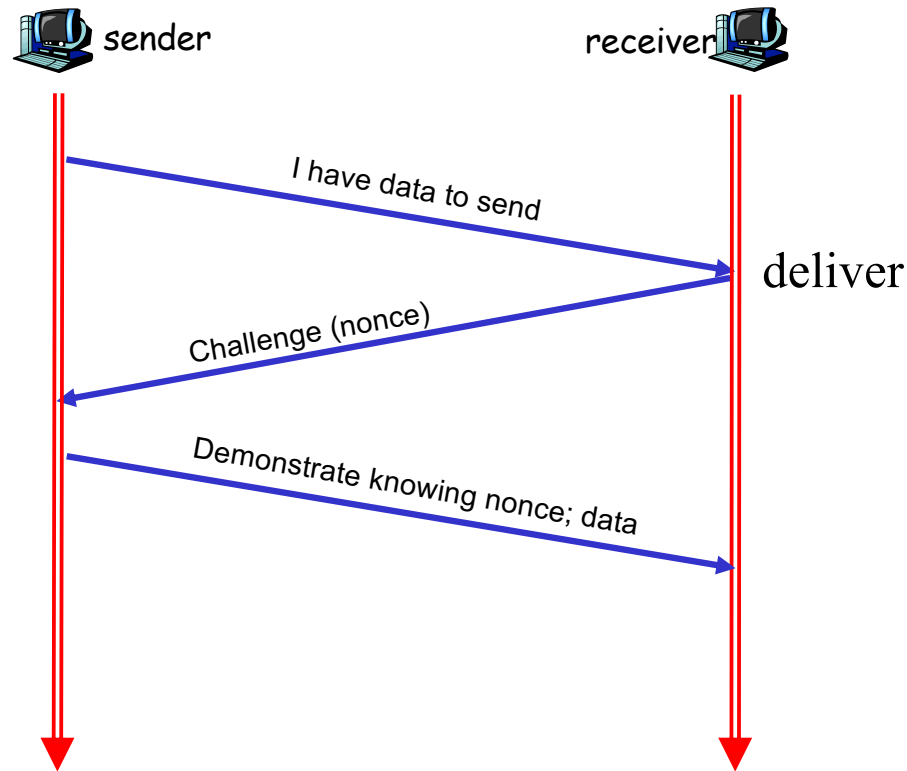
Transport authentication principle:

- [p1] Receiver can be sure that what Sender says is **fresh**
- [p2] Receiver receives something that **only Sender can say**

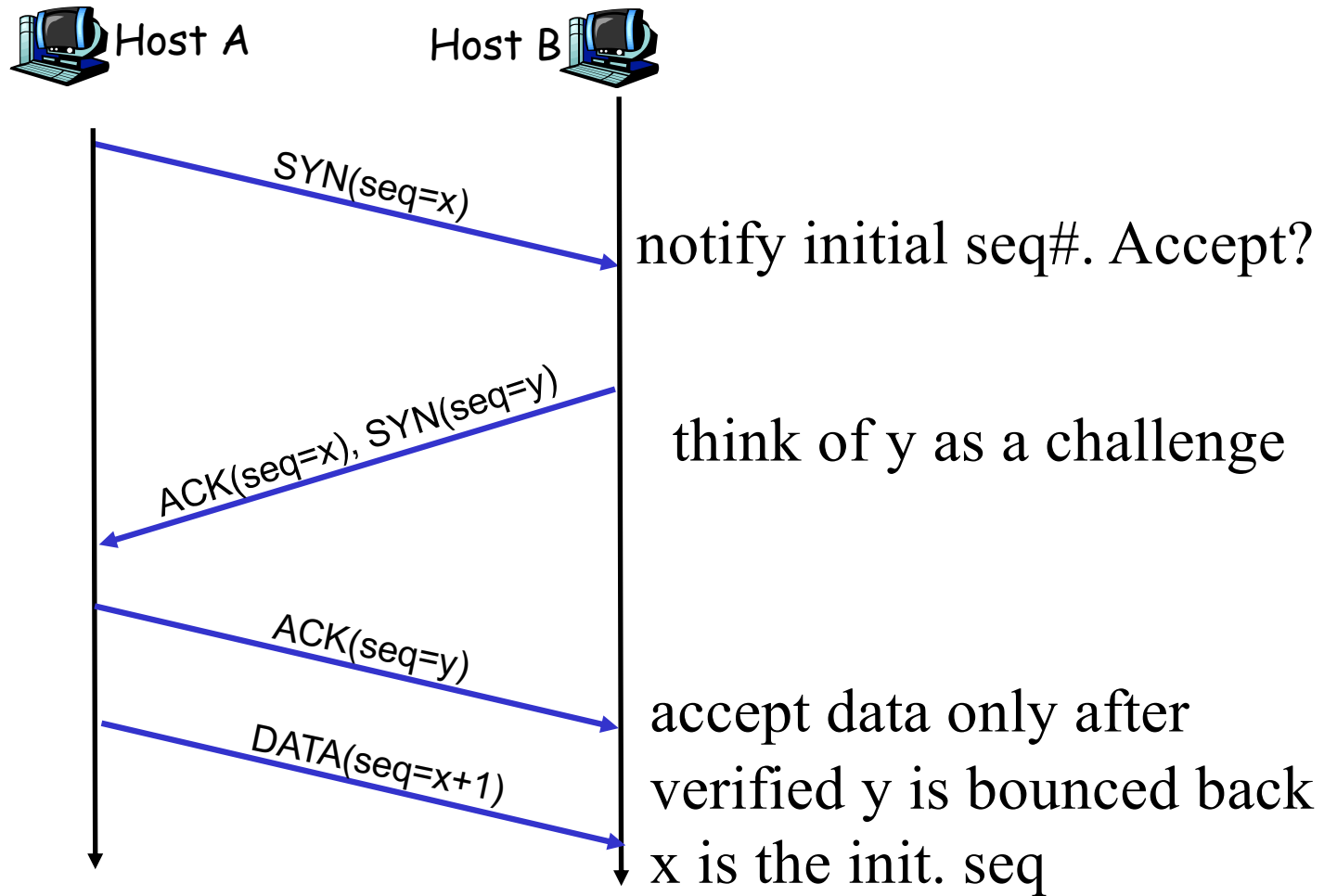
We first assume a secure setting: no malicious attacks.

Exercise: Techniques to allow a receiver to check for freshness (e.g., add a time stamp)?

# Generic Challenge-Response Structure Checking Freshness



# Three Way Handshake (TWH) [Tomlinson 1975]



SYN: indicates connection setup

# Make "Challenge y" Robust

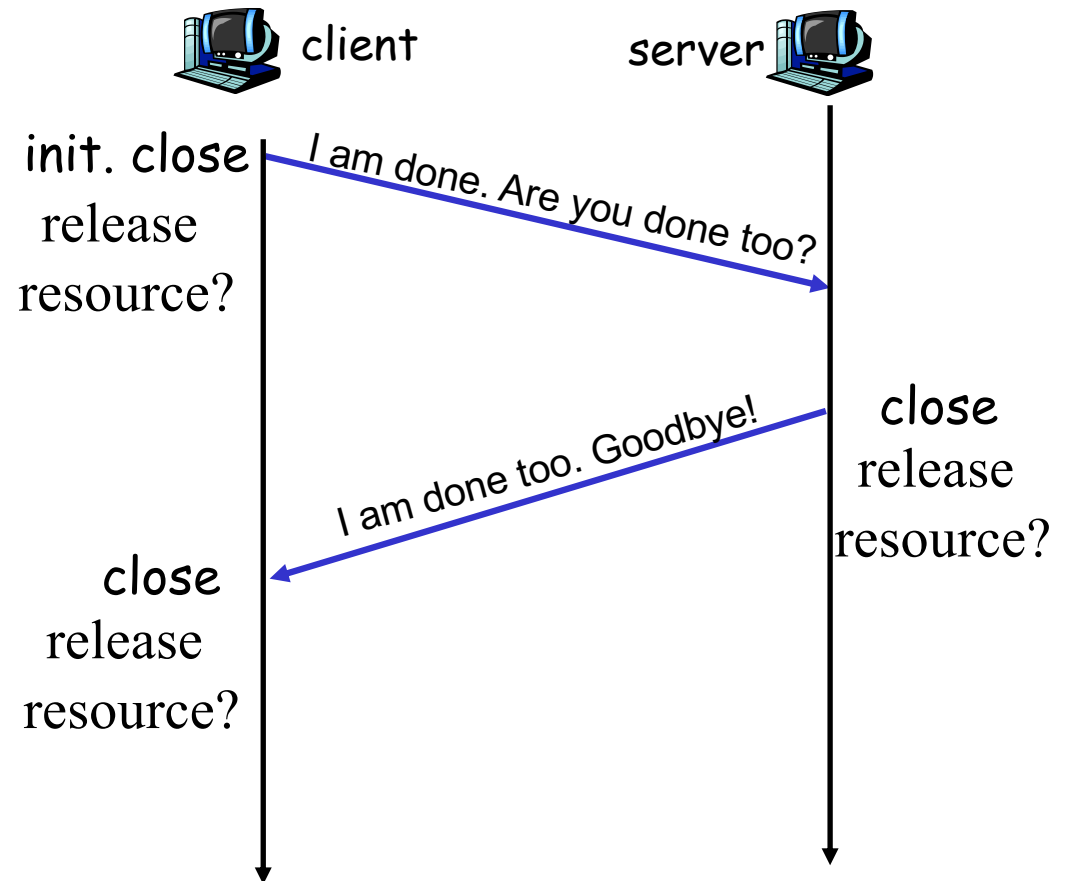
---

- ❑ To avoid that "SYNC ACK y" comes from reordering and duplication
  - for each connection (sender-receiver pair), ensuring that two identically numbered packets are never outstanding at the same time
    - network bounds the life time of each packet
    - a sender will not reuse a seq# before it is sure that all packets with the seq# are purged from the network
    - seq. number space should be large enough to not limit transmission rate
  
- ❑ Increasingly move to cryptographic challenge and response

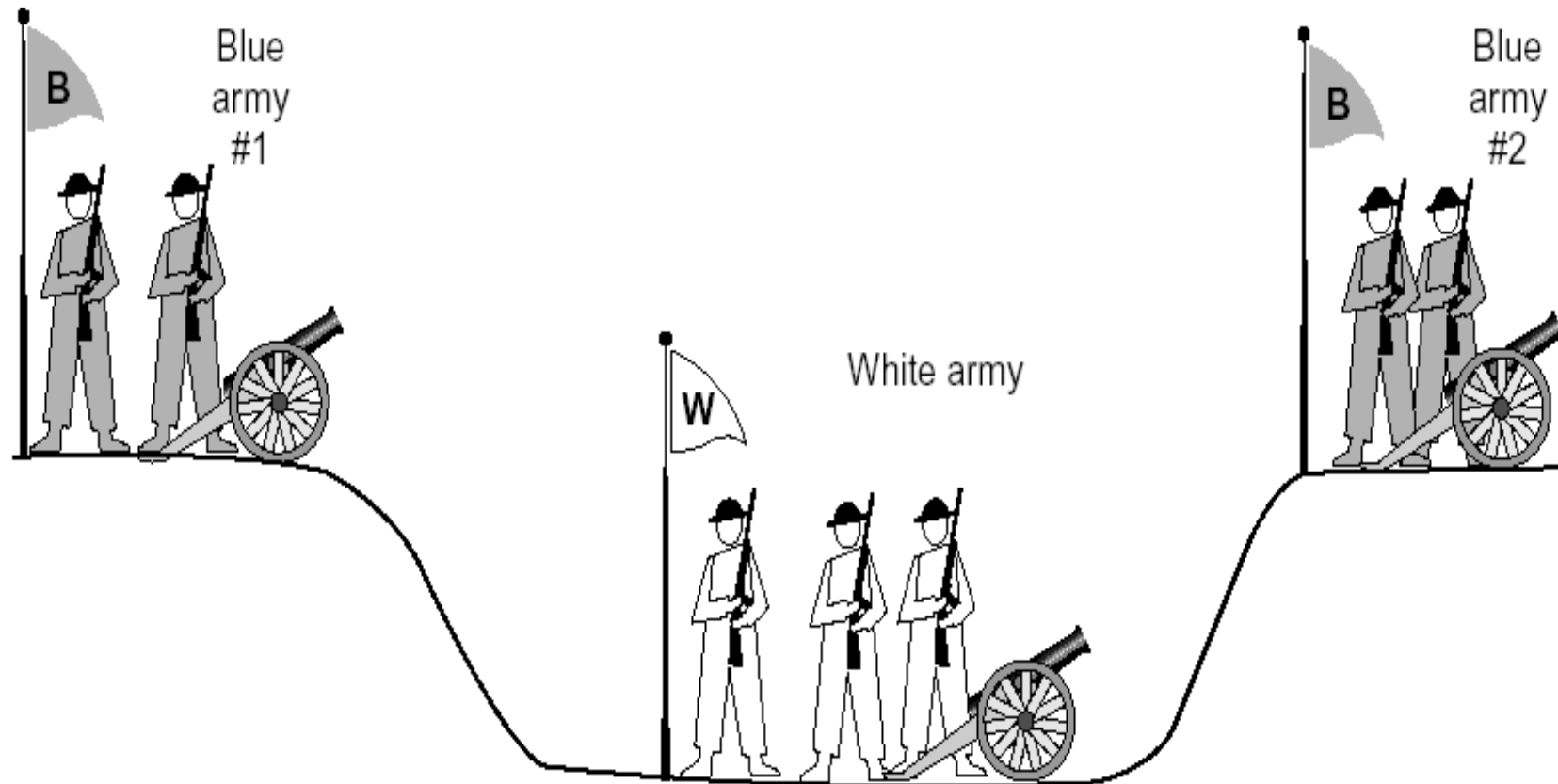


# Connection Close

- ❑ Why connection close?
  - so that each side can release resource and remove state about the connection (do not want dangling socket)



# General Case: The Two-Army Problem

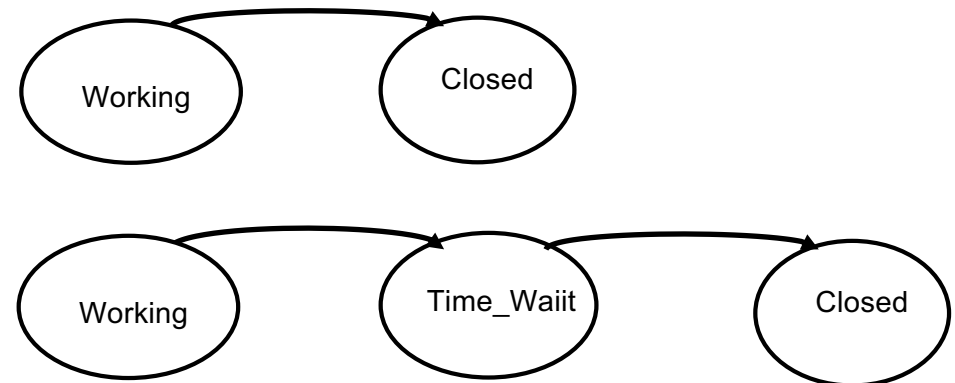
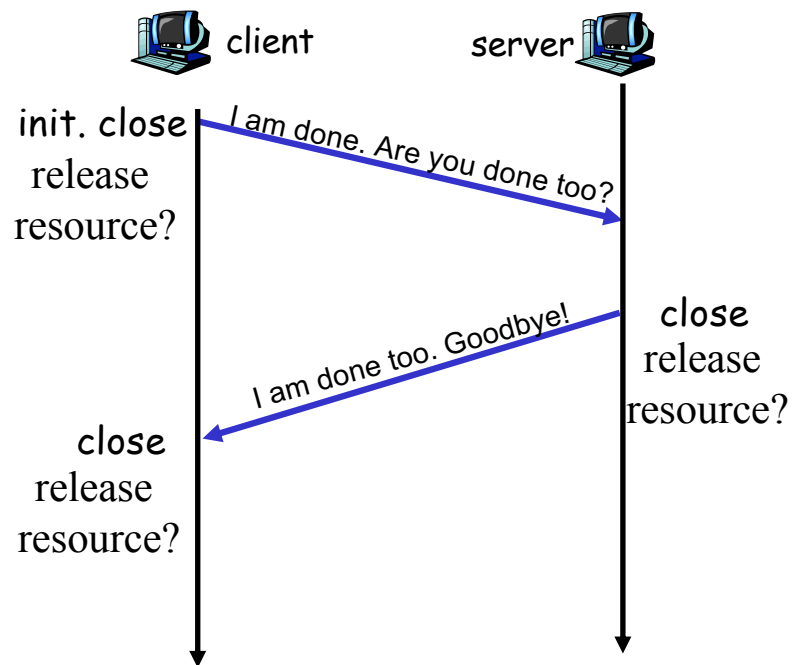


The gray (blue) armies need to agree on whether or not they will attack the white army. They achieve agreement by sending messengers to the other side. If they both agree, attack; otherwise, no. Note that a messenger can be captured!

# Time Wait

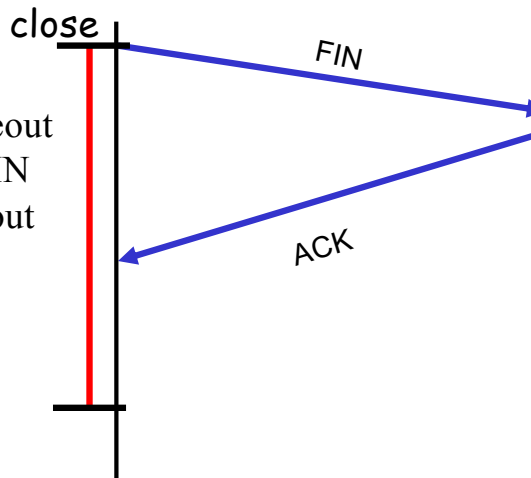
## □ Generic technique: Timeout to "solve" infeasible problem

- Instead of message-driven state transition, use a timeout based transition; use timeout to handle error cases



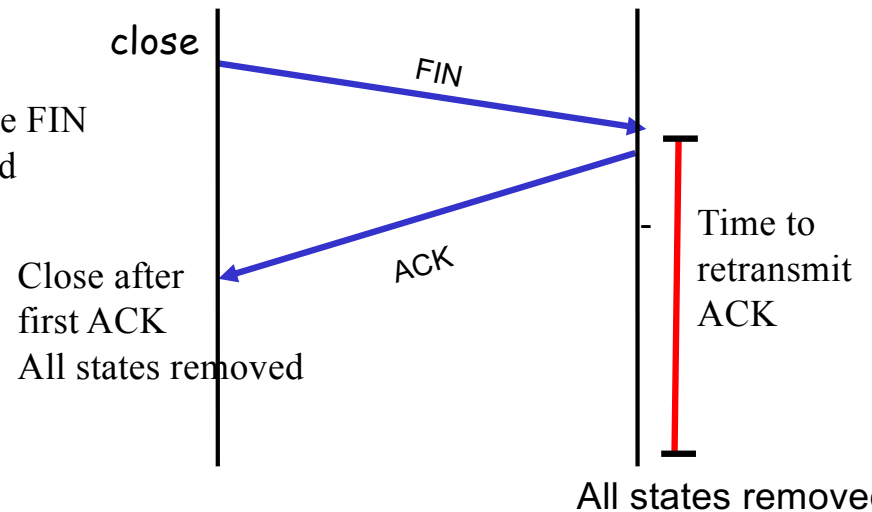
# Time\_Wait Design Options

Design 1 (initiator time wait)

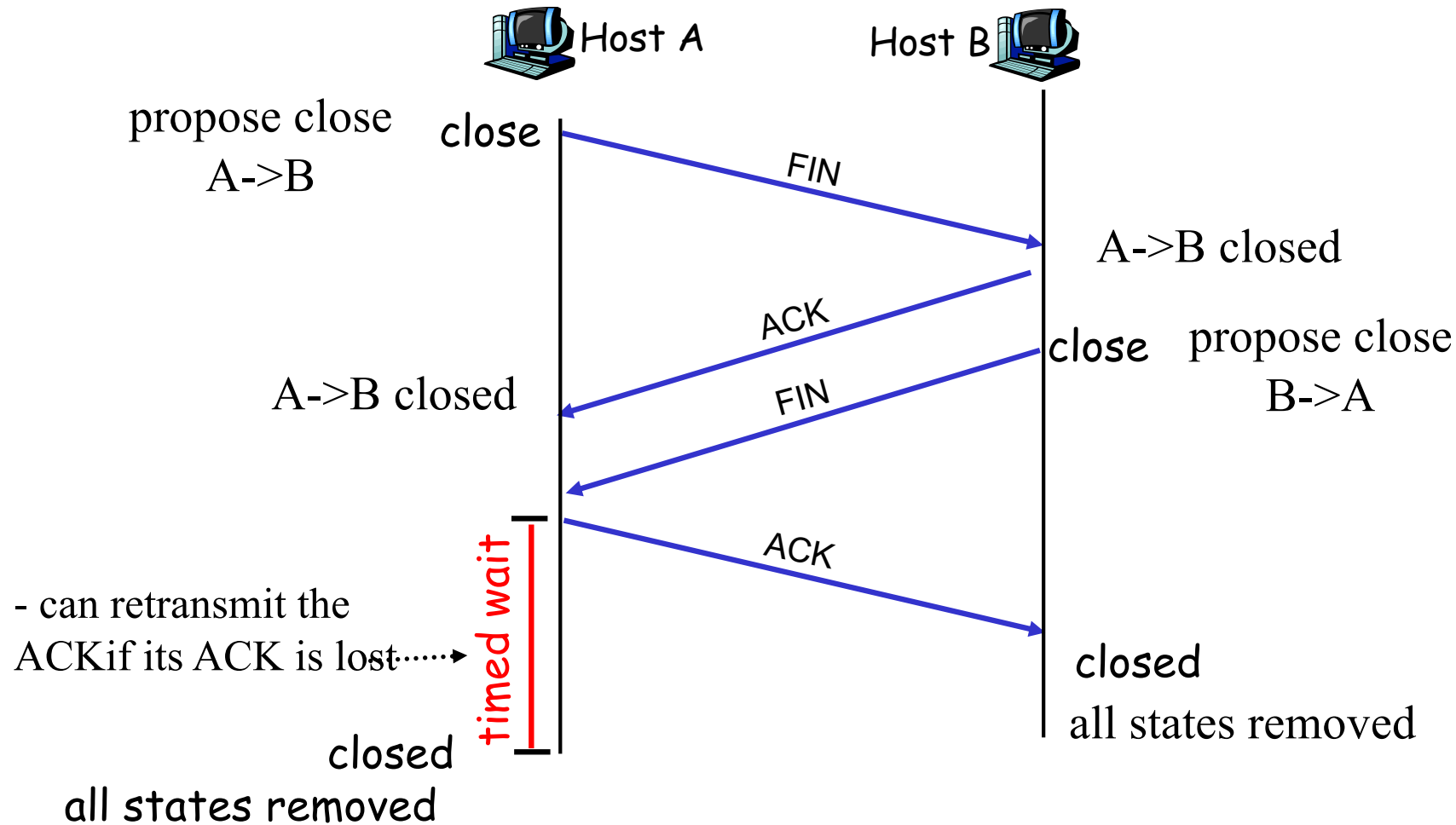


- Time = n x timeout
- Time to retry FIN after each timeout

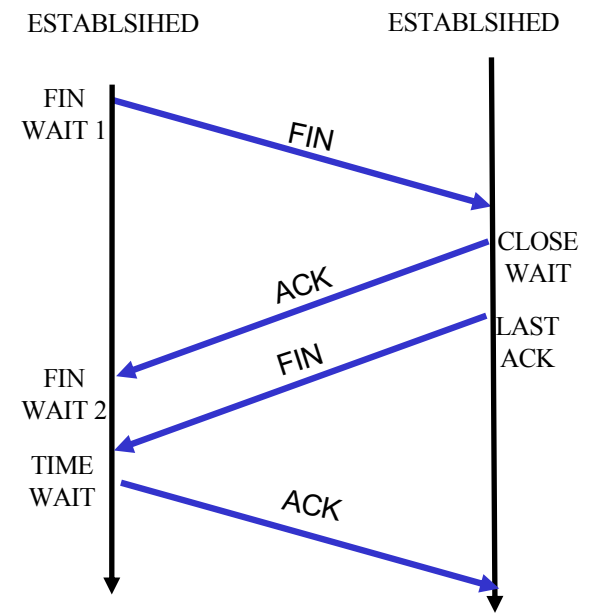
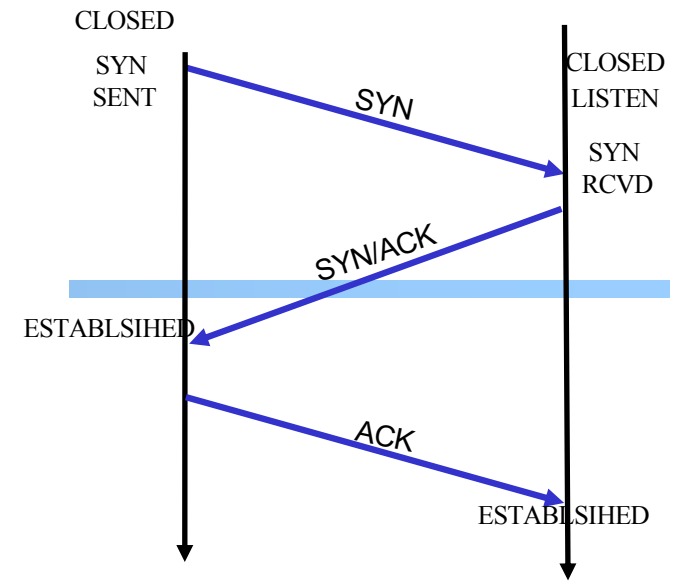
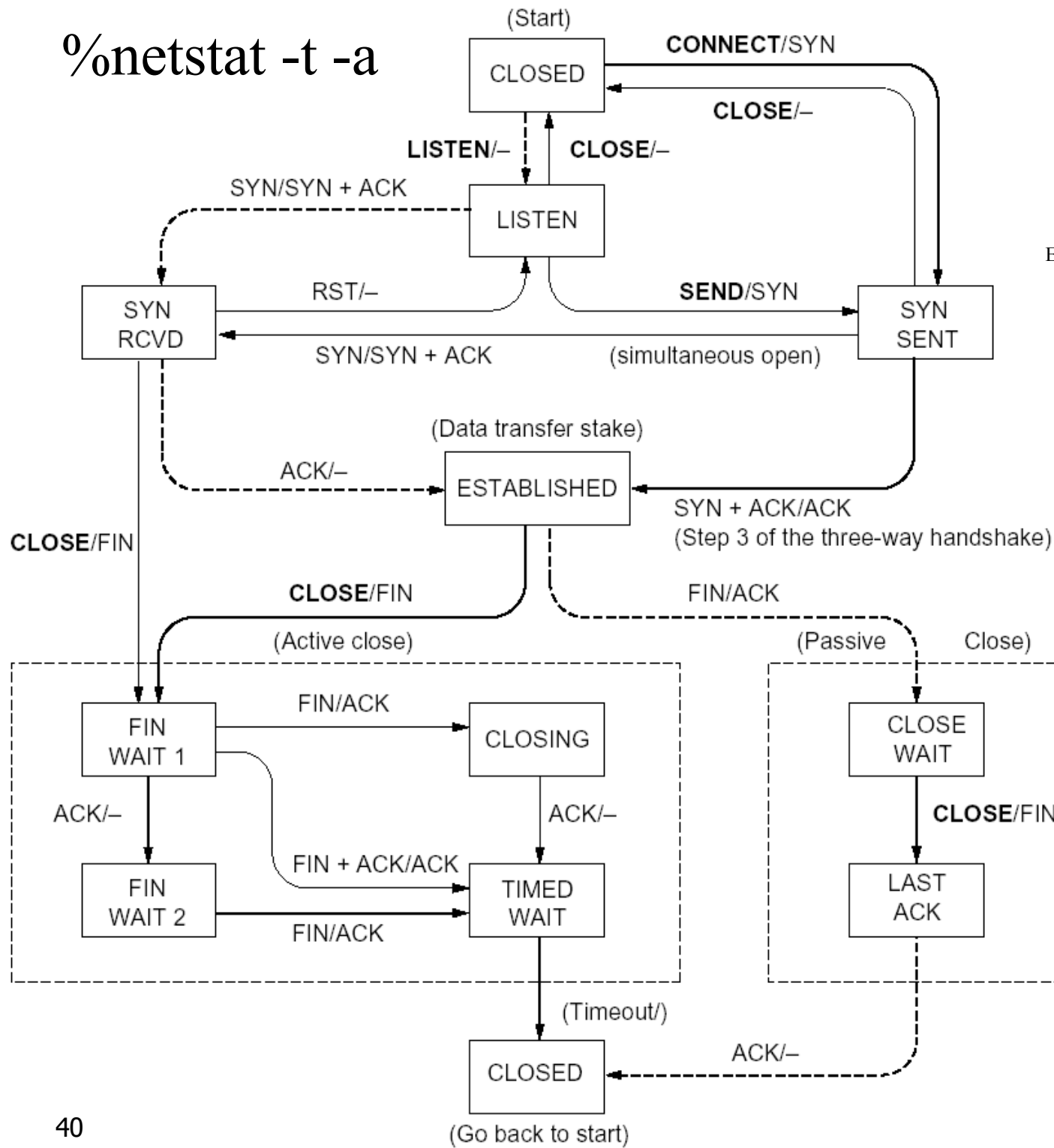
Design 2 (receiver time wait)



# TCP Four Way Teardown (For Bi-Directional Transport)

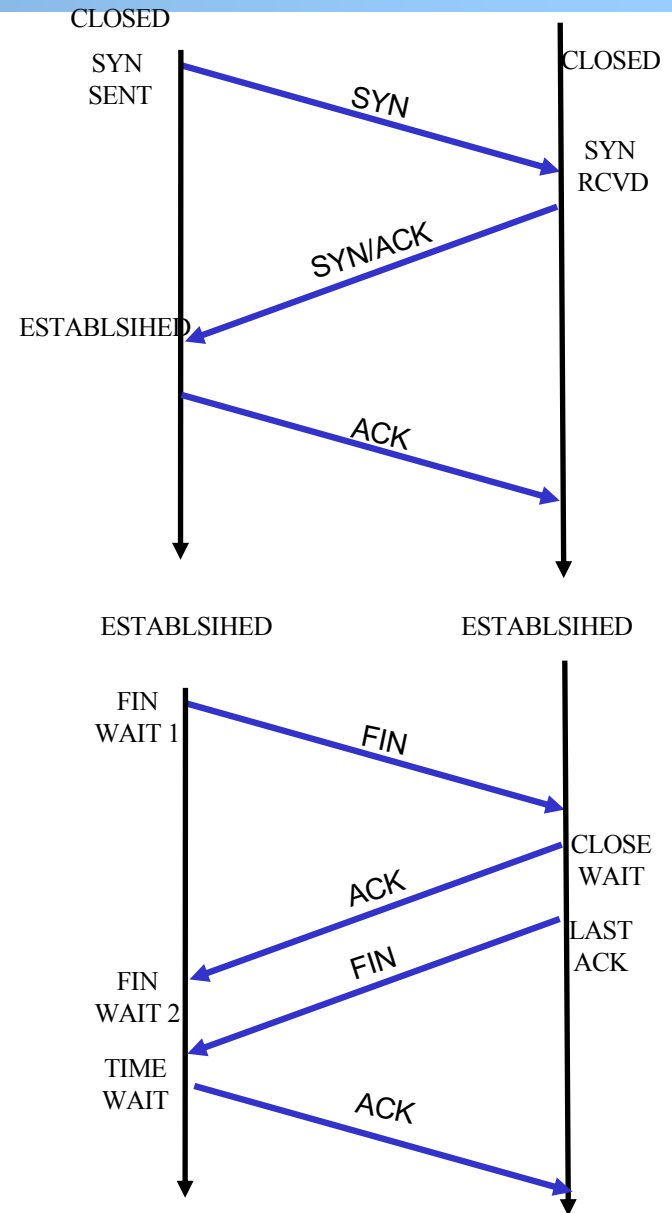
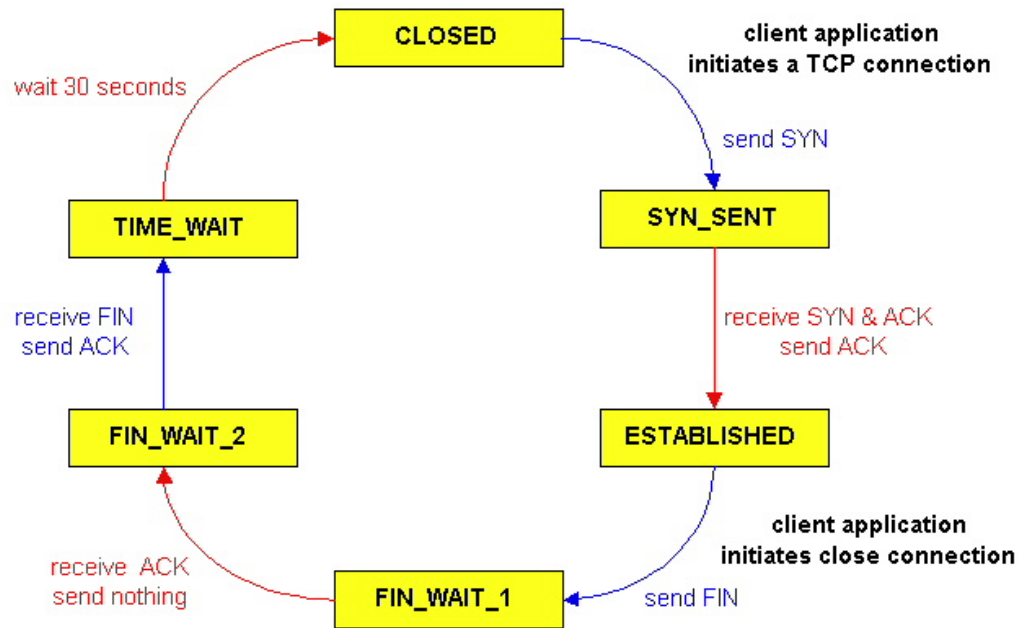


# %netstat -t -a



# TCP Connection Management

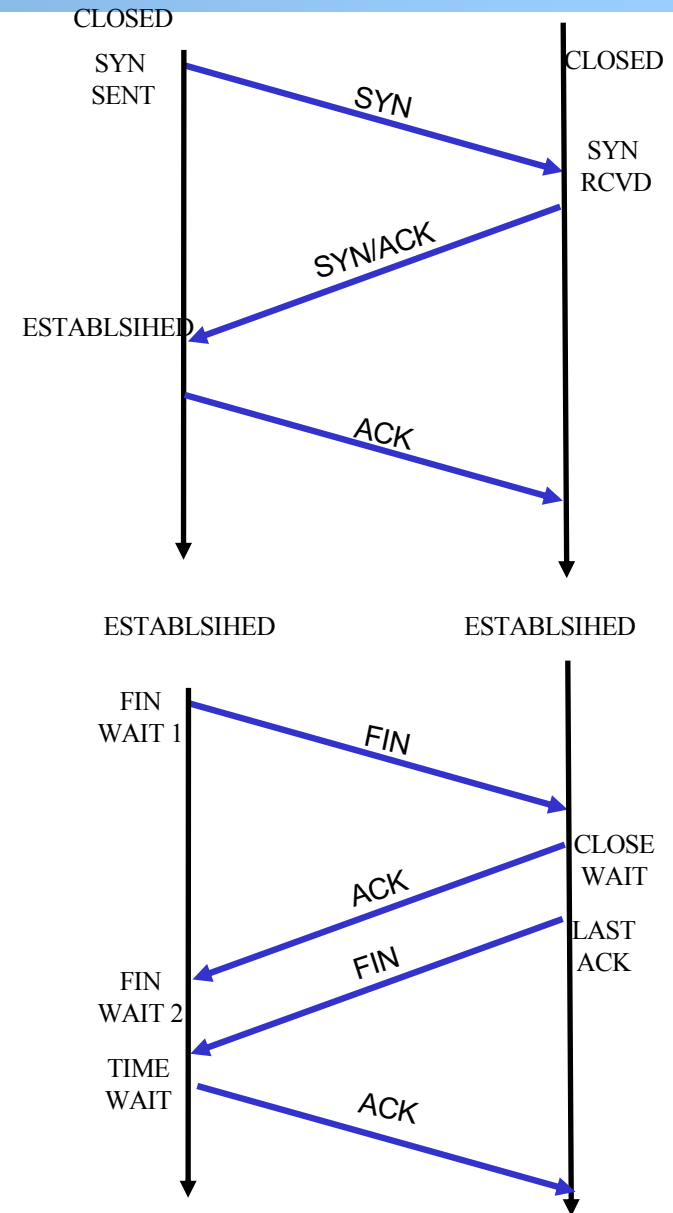
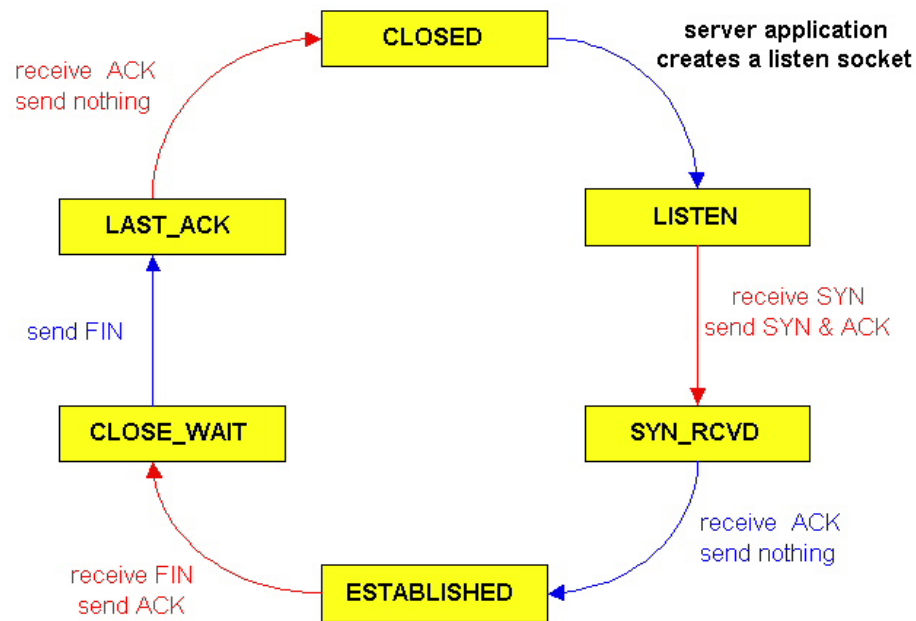
## TCP lifecycle: init SYN/FIN



<http://dsd.lbl.gov/TCP-tuning/ip-sysctl-2.6.txt>

# TCP Connection Management

TCP lifecycle: wait for SYN/FIN





# A Summary of Questions

---

- ❑ Basic structure: sliding window protocols
- ❑ How to determine the “right” parameters?
  - ✓ timeout: mean + variation
  - sliding window size?