# Network Transport Layer: TCP

**Qiao Xiang**, Congming Gao, Qiang Su
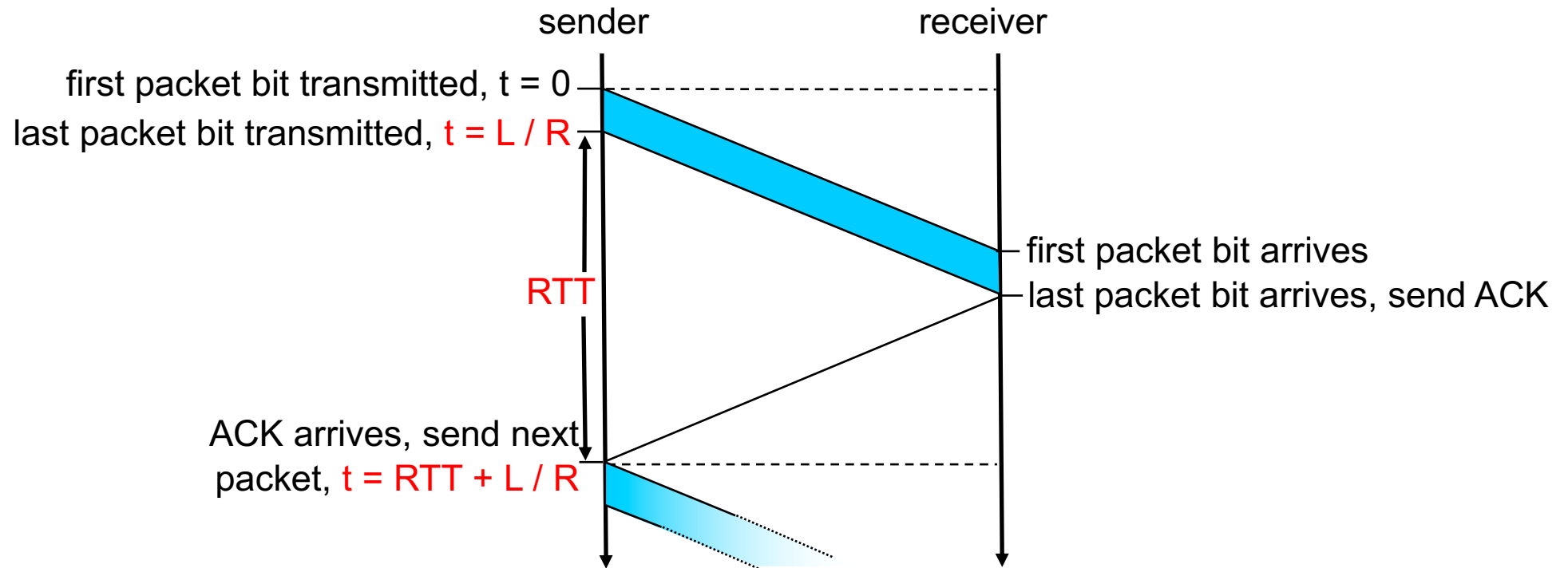
https://sngroup.org.cn/courses/cnns-xmuf25/index.shtml

11/04/2025

# Outline

- ❑ Admin and recap
- ❑ TCP

# rdt3.0: Stop-and-Wait Performance

first packet bit transmitted, t = 0

last packet bit transmitted, t = L / R

RTT

first packet bit arrives

last packet bit arrives, send ACK
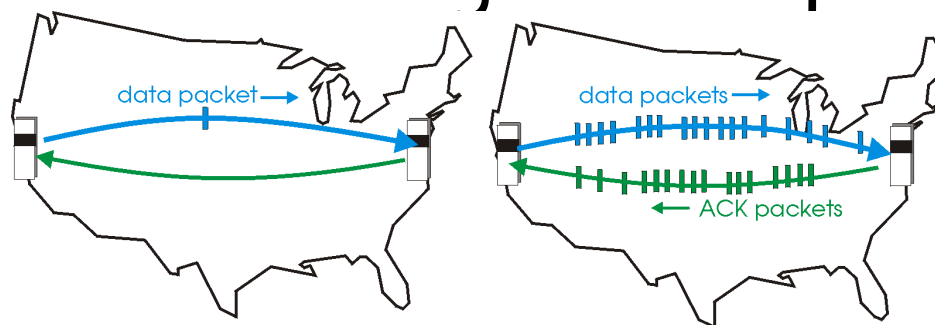
ACK arrives, send next
packet, t = RTT + L / R

What is $U_{sender}$: utilization – fraction of time link busy sending?

Assume: 1 Gbps link, 15 ms e-e prop. delay, 1KB packet

# 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 ≥ 2W |
| Buffer size at receiver | 1 | W |
| Complexity | Simpler | More complex |

# TCP Reliable Data Transfer
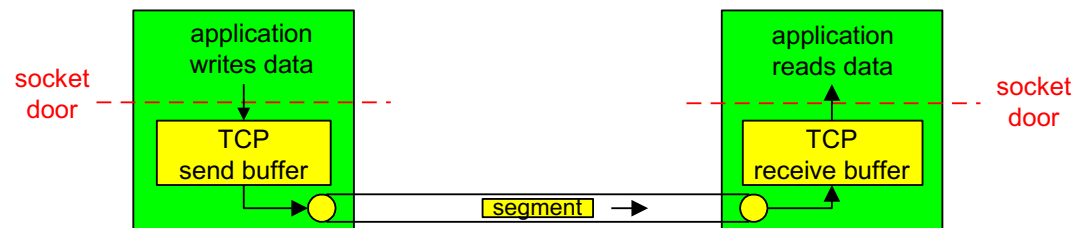
❑ **Connection-oriented:**

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

❑ **Full duplex data:**

　o bi-directional data flow in same connection
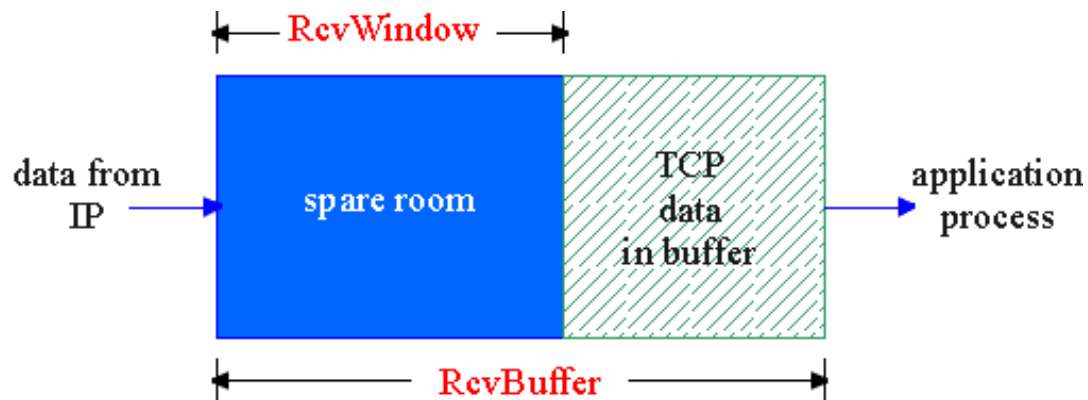
❑ **A sliding window protocol**

　o 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

application writes data

socket door

TCP send buffer

application reads data

socket door

TCP receive buffer

segment

# Flow Control

❑ **receive side of a connection has a receive buffer:**

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



❑ **app process may be slow at reading from buffer**

❑ **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`

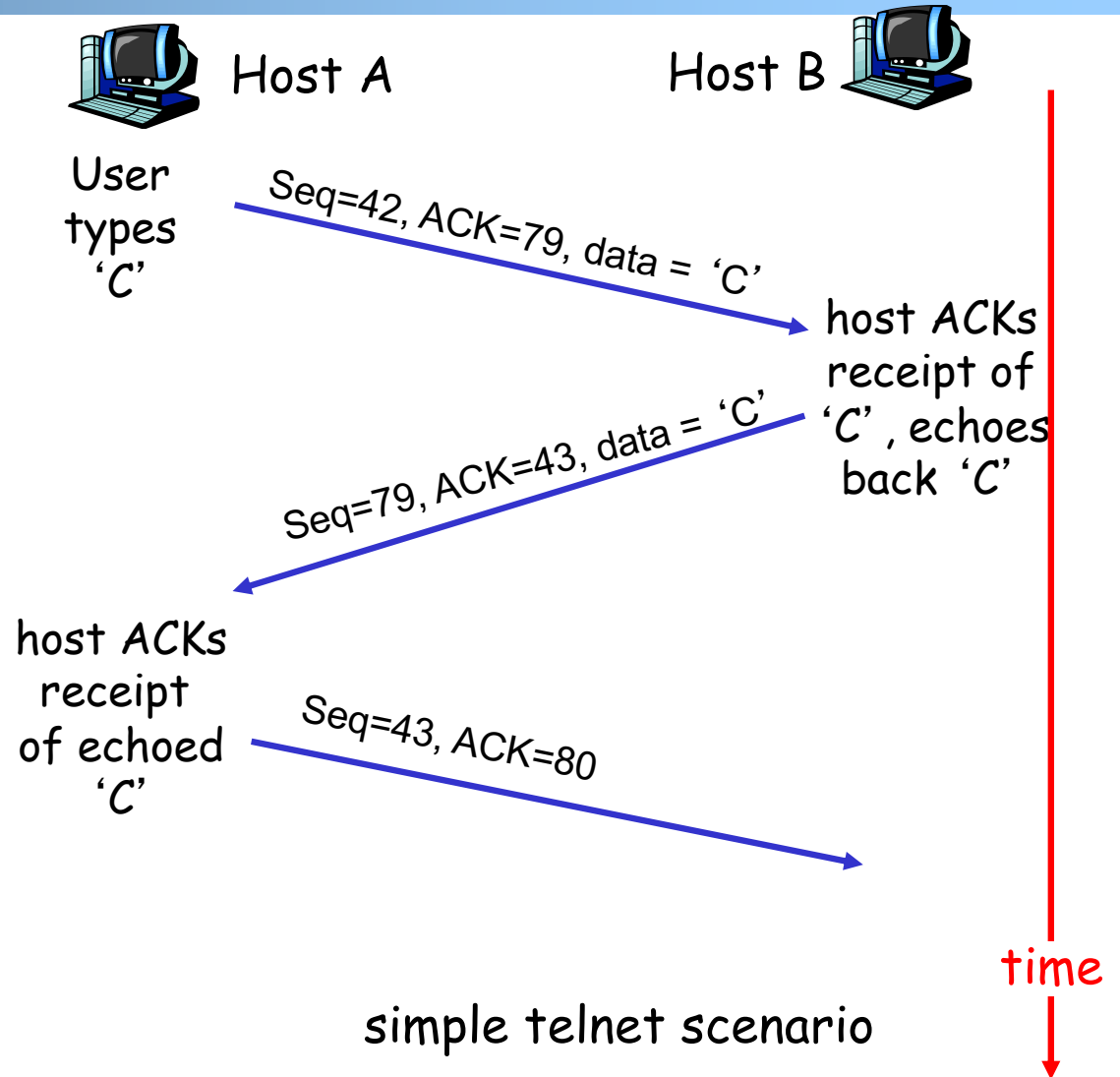| source port # | | | | | | dest port # | |
|---|---|---|---|---|---|---|---|
| sequence number | | | | | | | |
| acknowledgement number | | | | | | | |
| head len | not used | U | A | P | R | S | F | rcvr window size |
| checksum | | | | | | ptr urgent data | |
| Options (variable length) | | | | |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

Host A           Host B

User types 'C'

Seq=42, ACK=79, data = 'C'

host ACKs receipt of 'C', echoes back 'C'

Seq=79, ACK=43, data = 'C'

host ACKs receipt of echoed 'C'

Seq=43, ACK=80

time

simple telnet scenario

# 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-expect 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

   o perfect channel

   o channel with bit errors

   o channel with bit errors and losses

   o sliding window: reliability with throughput

❑ TCP reliability

   o 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 RTT

- Network collapse in the mid-1980s
  - UCB ←→ 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

  o premature timeout

  o unnecessary retransmissions; many duplicates
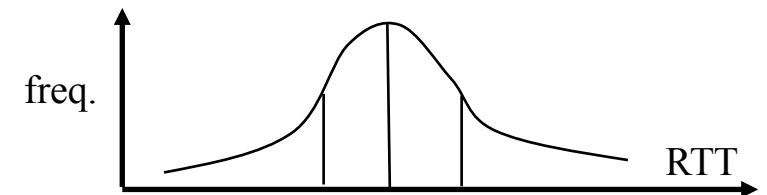
❑ too long

  o 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

freq.

RTT

**Solution:**
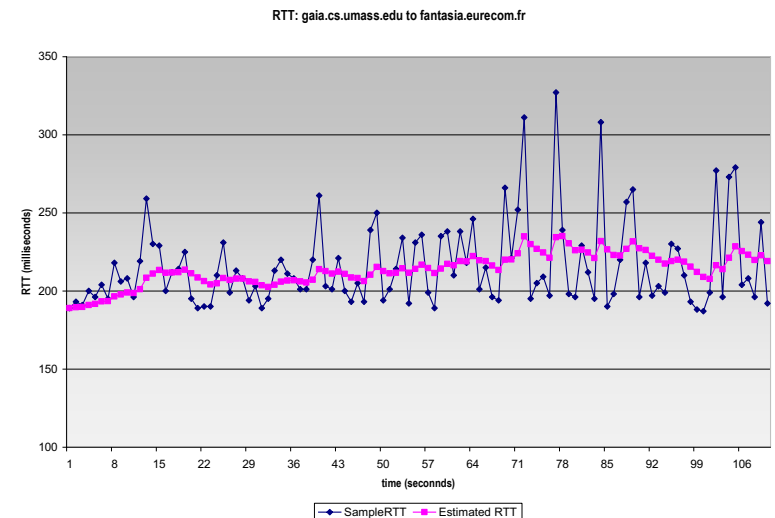- `Set Timeout RTO = avg +` "safety margin" based on variation

TCP approach:

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

# Compute EstRTT and DevRTT

❑ Exponential weighted moving average (EWMA)

    o   influence of past sample decreases exponentially fast

$$\texttt{EstRTT = (1-alpha)*EstRTT + alpha*SampleRTT}$$
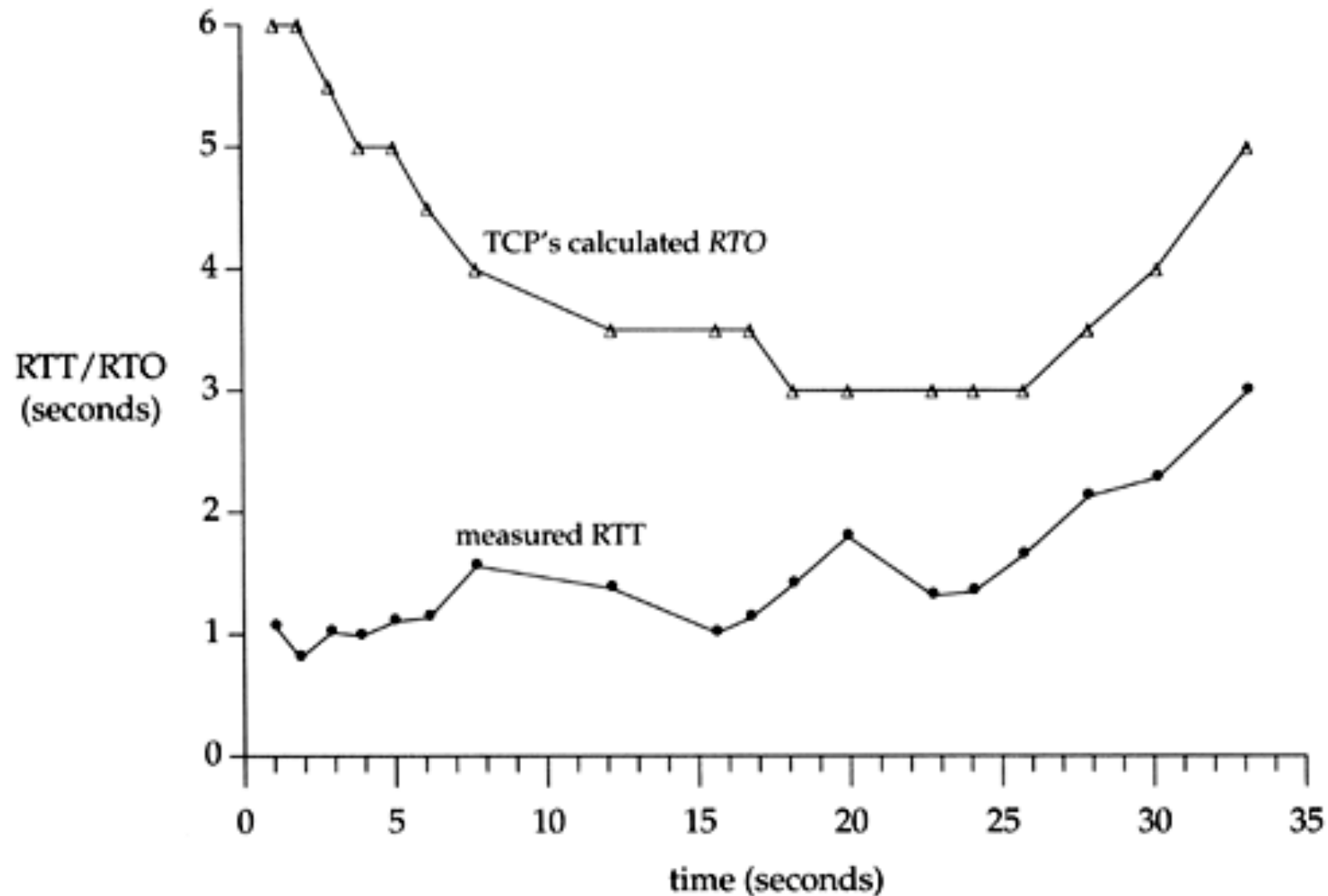
RTT: gaia.cs.umass.edu to fantasia.eurecom.fr

– **SampleRTT**: measured time
from segment transmission
until ACK receipt

- typical value: `alpha` = 0.125



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

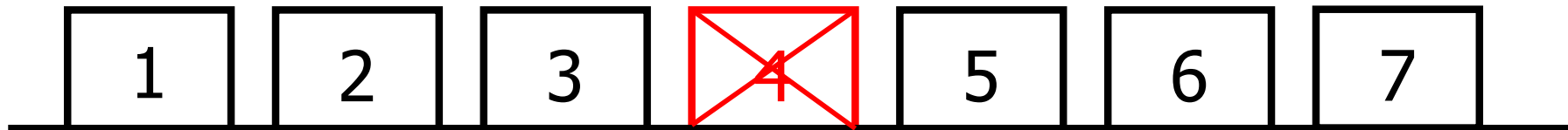$$\texttt{(typically, beta = 0.25)}$$

# 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

| 1 | 2 | 3 | ☒ | 5 | 6 | 7 |

Acknowledgements (waiting seq#)

| 2 | 3 | 4 | | 4 | 4 | 4 |

# 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

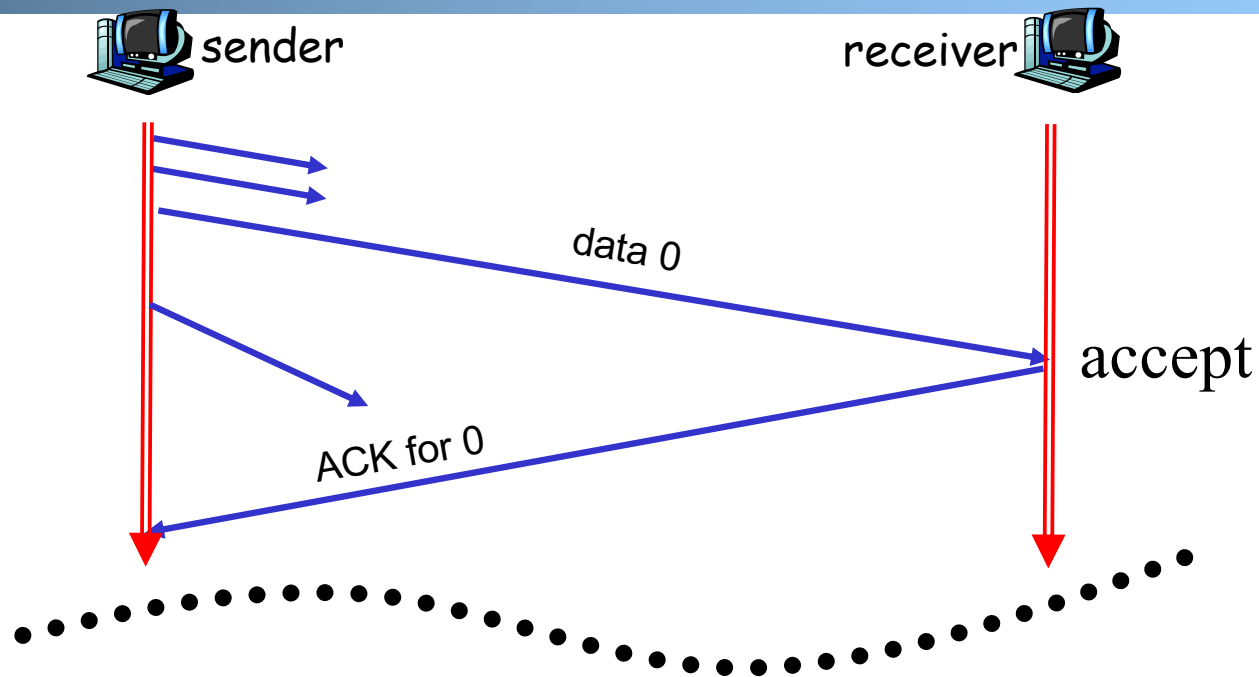# TCP: reliable data transfer

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)
            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 */
```
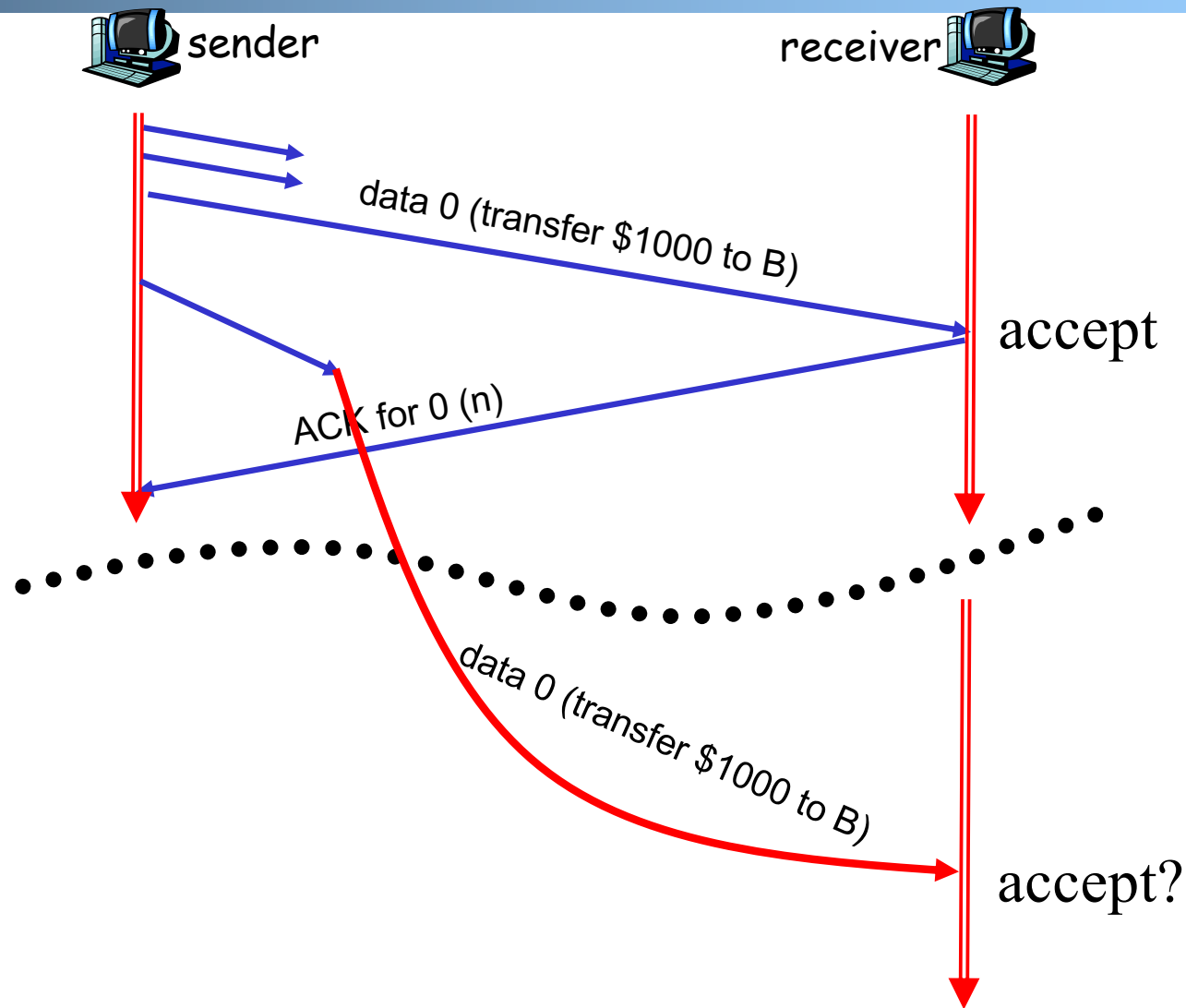
# Outline

❑ Admin and Recap

❑ Reliable data transfer

    o    perfect channel

    o    channel with bit errors

    o    channel with bit errors and losses

    o    sliding window: reliability with throughput

❑ TCP reliability

    o  data seq#, ack, buffering

    o  timeout realization

    ➢ *connection management*

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



sender

receiver

data 0

accept

ACK for 0

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

sender                                            receiver

data 0 (transfer $1000 to B)

accept

ACK for 0 (n)

data 0 (transfer $1000 to B)

accept?

# 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



sender          receiver

I have data to send

deliver

Challenge (nonce)

Demonstrate knowing nonce; data

# Three Way Handshake (TWH) [Tomlinson 1975]

Host A                          Host B

SYN(seq=x)

notify initial seq#. Accept?

ACK(seq=x), SYN(seq=y)

think of y as a challenge

ACK(seq=y)

accept data only after
verified y is bounced back
x is the init. seq

DATA(seq=x+1)

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?**
  - o so that each side can release resource and remove state about the connection (do not want dangling socket)
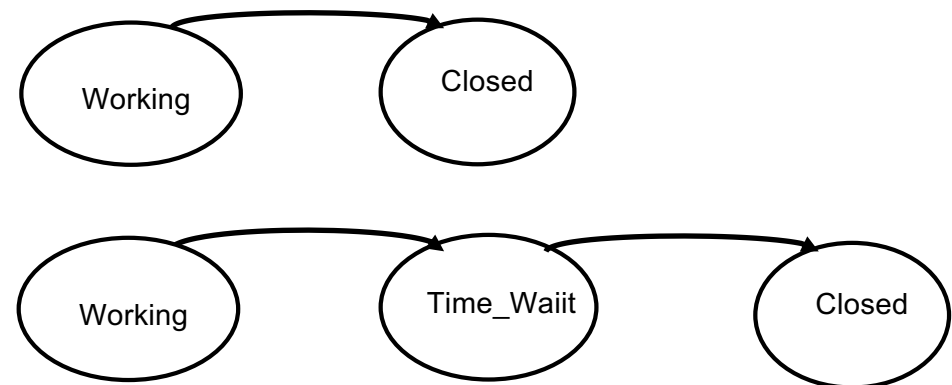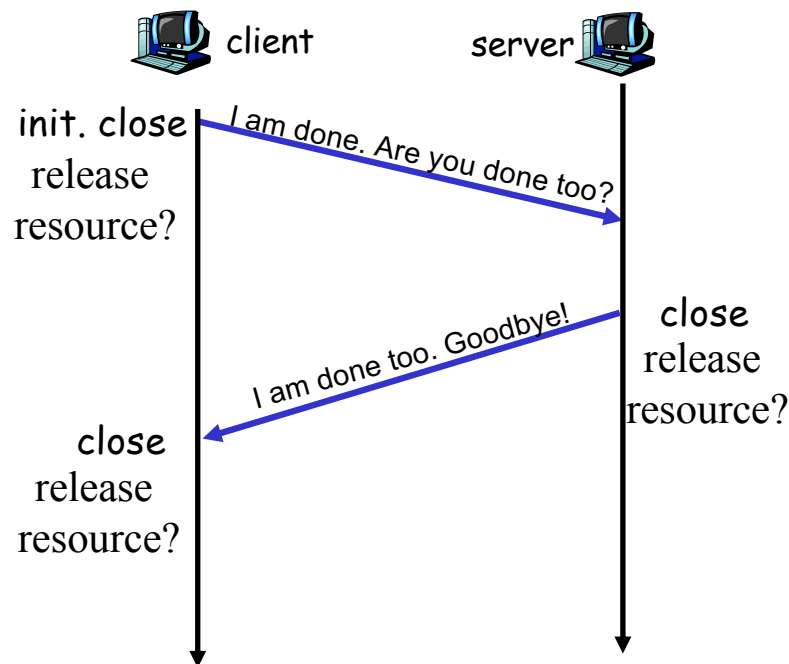
client      server

init. close
release
resource?

*I am done. Are you done too?*

*I am done too. Goodbye!*

close
release
resource?

close
release
resource?

# 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

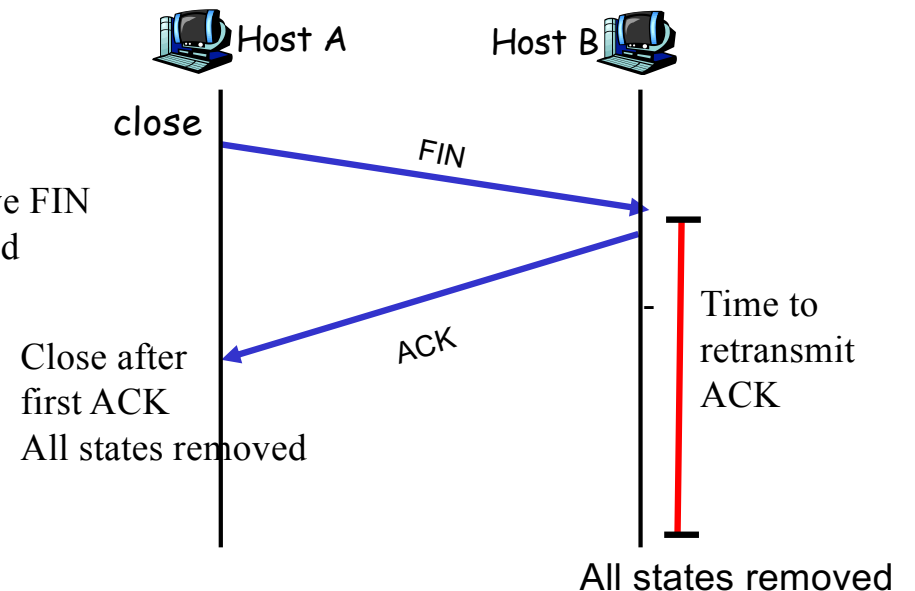○ Instead of message-driven state transition, use a timeout based transition; use timeout to handle error cases
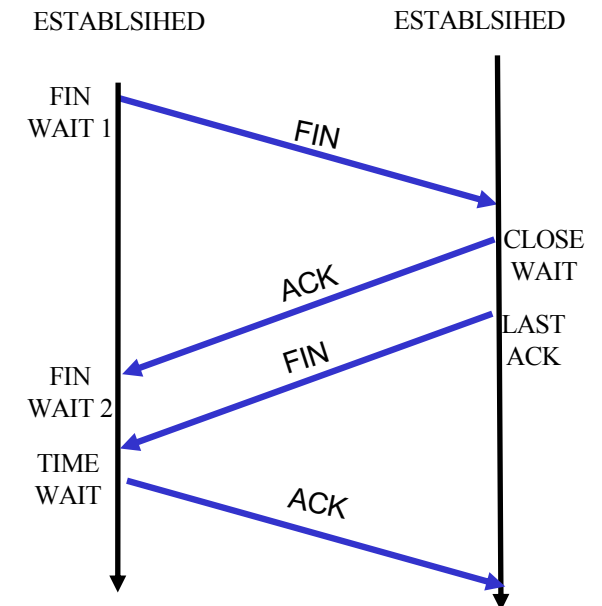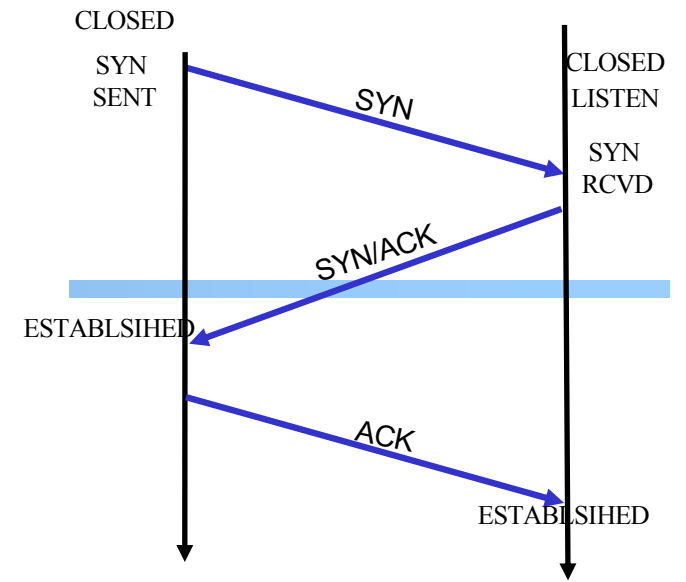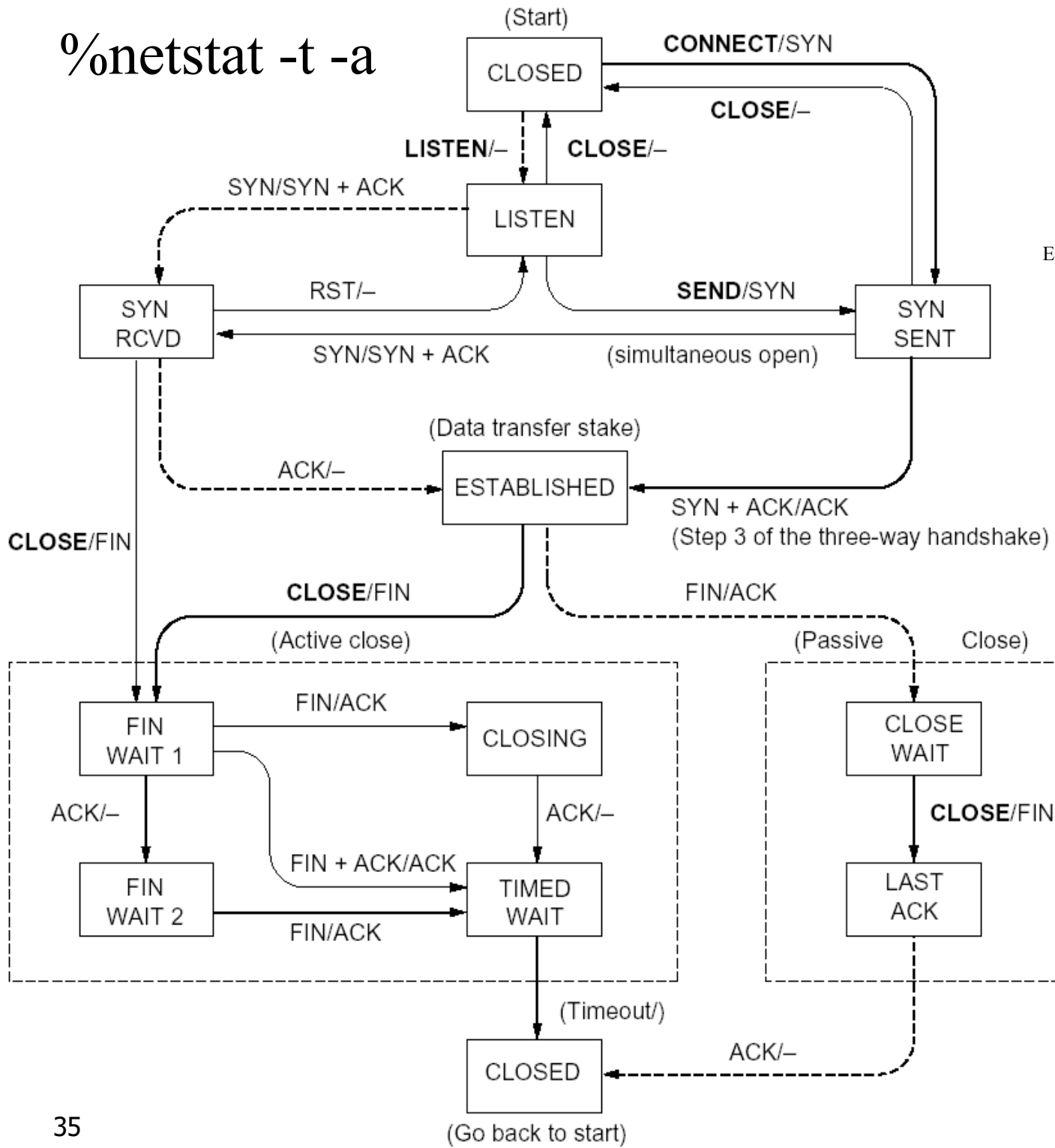
# Time_Wait Design Options



Design 1 (initiator time wait)

Host A    Host B

close

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

FIN

Close after receive FIN
All states removed

ACK

All states removed

Design 2 (receiver time wait)

Host A    Host B

close

FIN

ACK

Close after first ACK
All states removed

- Time to retransmit ACK

All states removed
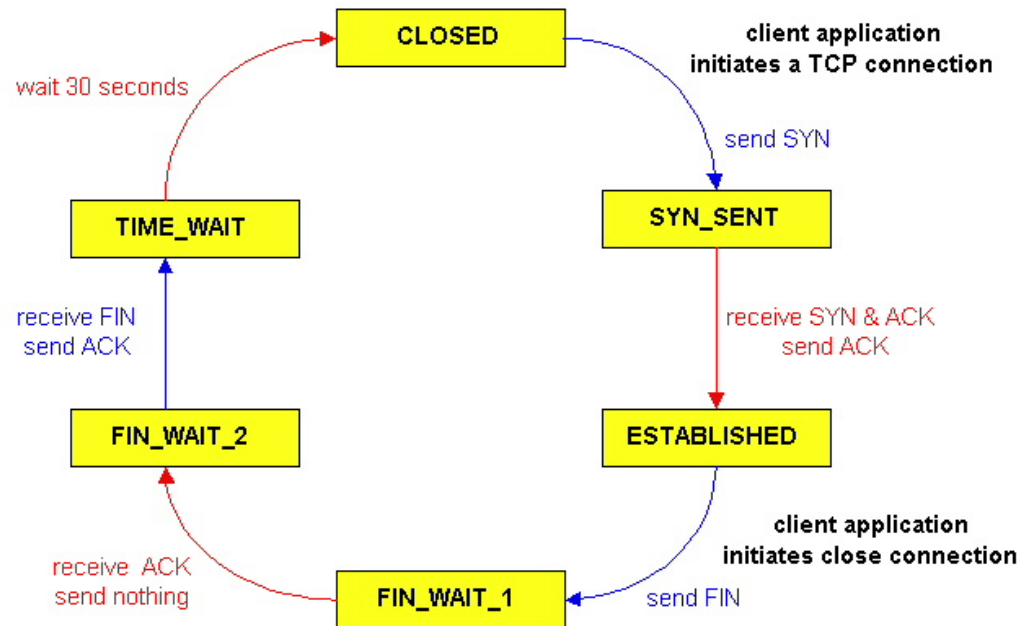
%netstat -t -a

# TCP Connection Management

## TCP lifecycle: init SYN/FIN

# TCP Connection Management

TCP lifecycle: wait for SYN/FIN



CLOSED

SYN SENT
CLOSED

SYN

SYN RCVD

SYN/ACK

ESTABLSIHED

ACK

ESTABLSIHED                    ESTABLSIHED

FIN WAIT 1

FIN

CLOSE WAIT

ACK

LAST ACK

FIN WAIT 2

FIN

TIME WAIT

ACK

CLOSED

server application creates a listen socket

receive ACK send nothing

LISTEN

LAST_ACK

receive SYN send SYN & ACK

send FIN

CLOSE_WAIT

SYN_RCVD

receive ACK send nothing

receive FIN send ACK

ESTABLISHED

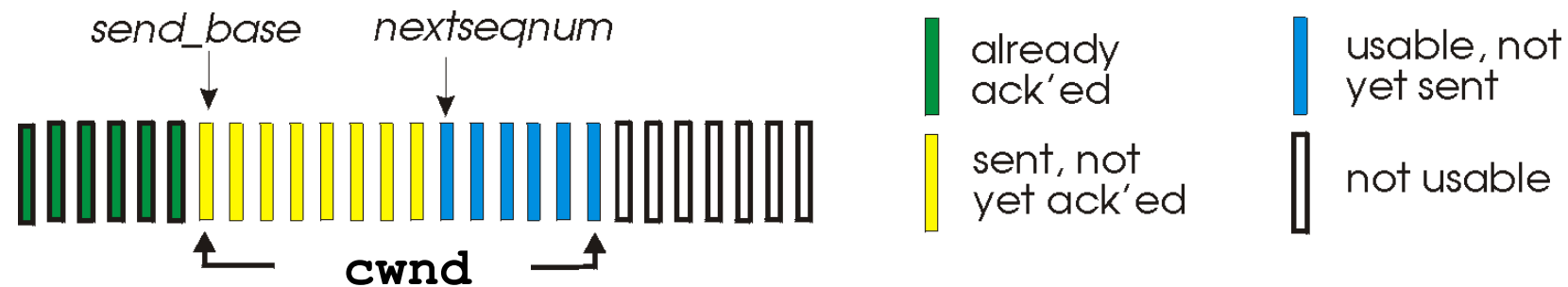# A Summary of Questions

❑ Basic structure: sliding window protocols

❑ How to determine the "right" parameters?

   ✓ timeout: mean + variation

   ○ sliding window size?

# Sliding Window Size Function: Rate Control

❑ **Transmission rate determined by congestion window size, `cwnd`, over segments:**



❑ **cwnd segments, each with MSS bytes sent in one RTT:**

$$\text{Rate} = \frac{\text{cwnd} * \text{MSS}}{\text{RTT}} \text{Bytes/sec}$$

Assume W is small enough. Ignore small details. MSS: Minimum Segment Size

# Some General Questions

Big picture question:

❑ How to determine a flow's sending rate?

For better understanding, we need to look at a few basic questions:

❑ What is congestion (cost of congestion)?

❑ Why are desired properties of congestion control?
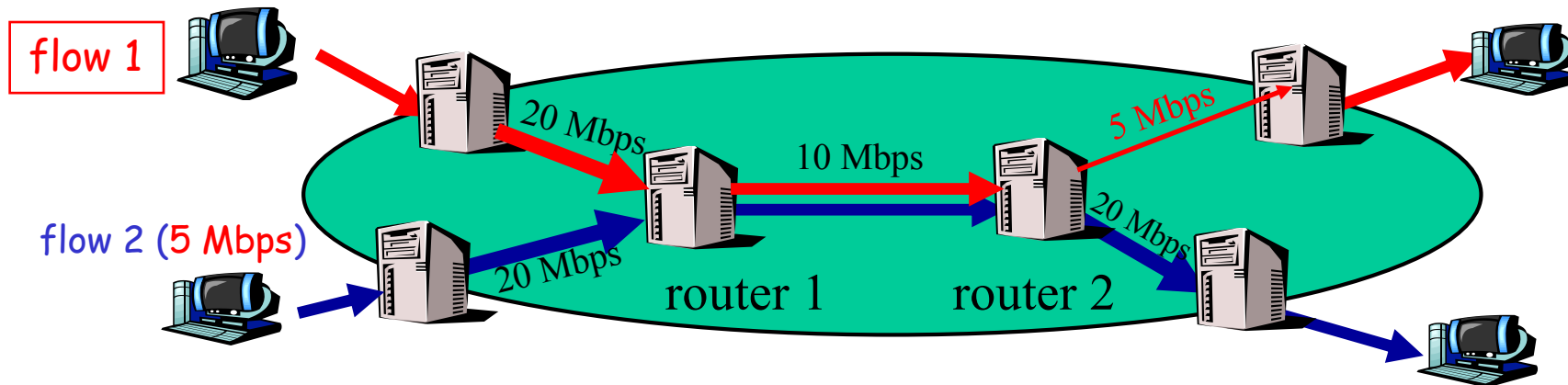
# Roadmap

- ❑ What is congestion
- ❑ The basic CC alg
- ❑ TCP/reno CC
- ❑ TCP/Vegas
- ❑ A unifying view of TCP/Reno and TCP/Vegas
- ❑ Network wide resource allocation
  - ○ Framework
  - ○ Axiom derivation of network-wide objective function
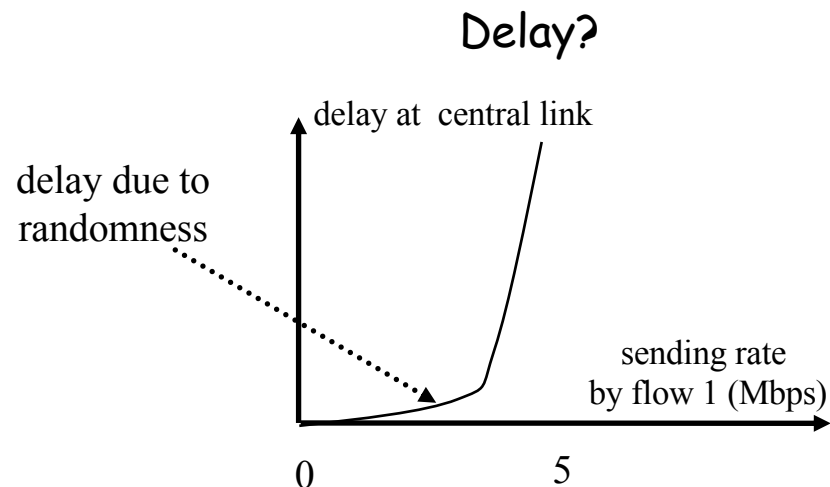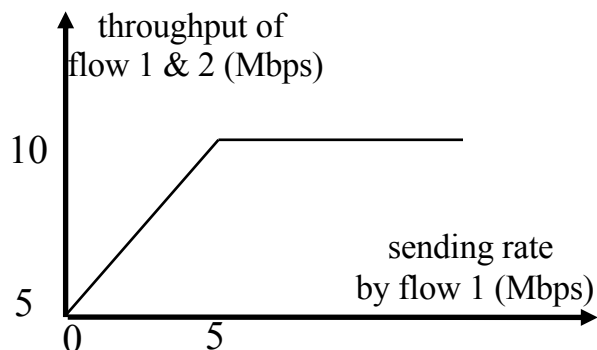  - ○ Derive distributed algorithm

# Outline

- ❏ Admin and recap
- ❏ TCP Reliability
- ❏ Transport congestion control
  - ➢ *what is congestion (cost of congestion)*

# Cause/Cost of Congestion: Single Bottleneck



flow 1

20 Mbps

10 Mbps

5 Mbps

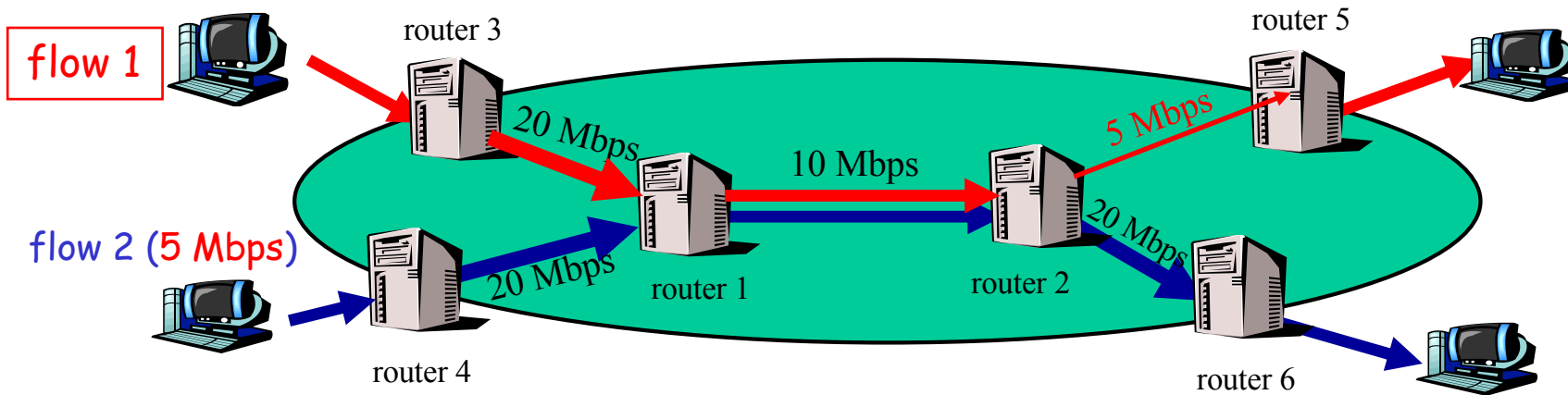flow 2 (5 Mbps)

20 Mbps

20 Mbps

router 1          router 2

- Flow 2 has a fixed sending rate of 5 Mbps
- We vary the sending rate of flow 1 from 0 to 20 Mbps
- Assume
  - no retransmission; link from router 1 to router 2 has infinite buffer

throughput: e2e packets
delivered in unit time

Delay?
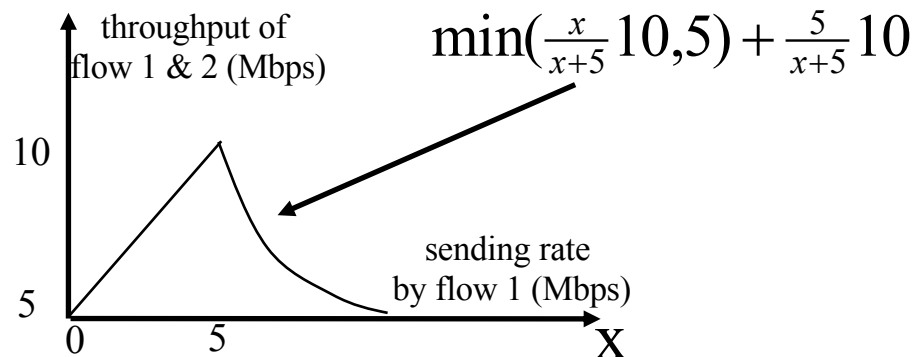
throughput of
flow 1 & 2 (Mbps)

delay at central link

delay due to
randomness

10

sending rate
by flow 1 (Mbps)

sending rate
by flow 1 (Mbps)

5

0          5

0                    5

# Cause/Cost of Congestion: Single Bottleneck



❑ Assume
- o   no retransmission
- o   the link from router 1 to router 2 has finite buffer
- o   throughput: e2e packets delivered in unit time

$$\min(\tfrac{x}{x+5}10,5) + \tfrac{5}{x+5}10$$

throughput of flow 1 & 2 (Mbps)

sending rate by flow 1 (Mbps)

X

10

5

0

5

❑ Zombie packet: a packet dropped at the link from router 2 to router 5; the upstream transmission from router 1 to router 2 used for that packet was wasted!

44

# Summary: The Cost of Congestion

When sources sending rate too high for the *network* to handle":

❑ Packet loss =>

- o wasted upstream bandwidth when a pkt is discarded at downstream

- o wasted bandwidth due to retransmission (a pkt goes through a link multiple times)

❑ High delay