# Course on Computer Communication and Networks

# Lecture 5
# Chapter 3; Transport Layer, Part B

EDA344/DIT 423, CTH/GU

**Based on the book Computer Networking: A Top Down Approach, Jim Kurose, Keith Ross, Addison-Wesley.**
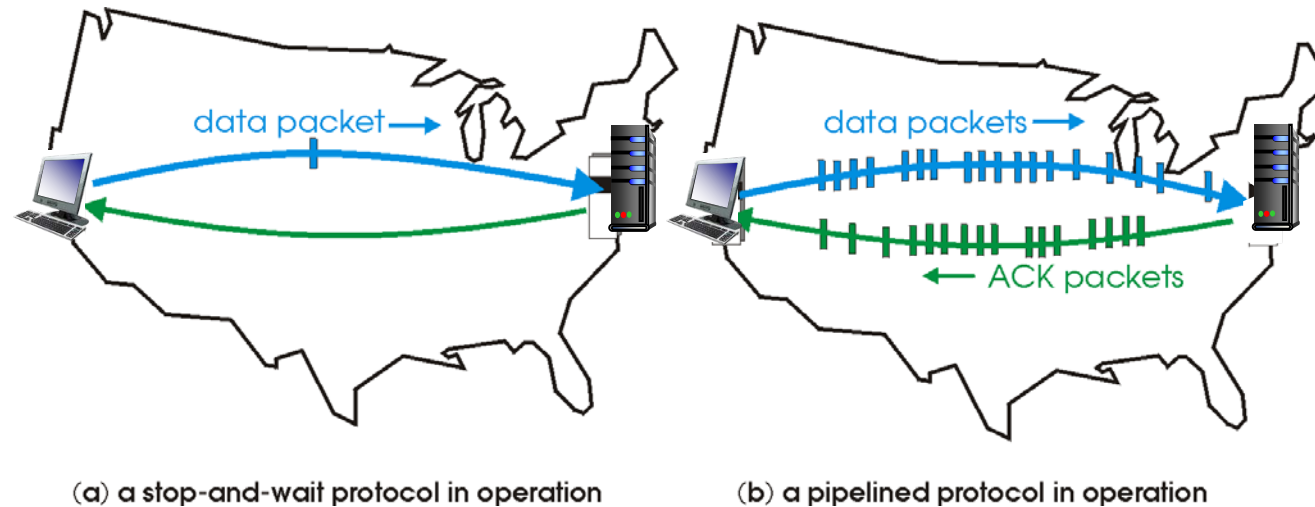
# Roadmap



- Transport layer services in Internet
- Addressing, multiplexing/demultiplexing
- Connectionless, unreliable transport: UDP
- principles of reliable data transfer
  - Efficiency perspective

- *Next lecture: connection-oriented transport: TCP*
  - *reliable transfer*
  - *flow control*
  - *connection management*
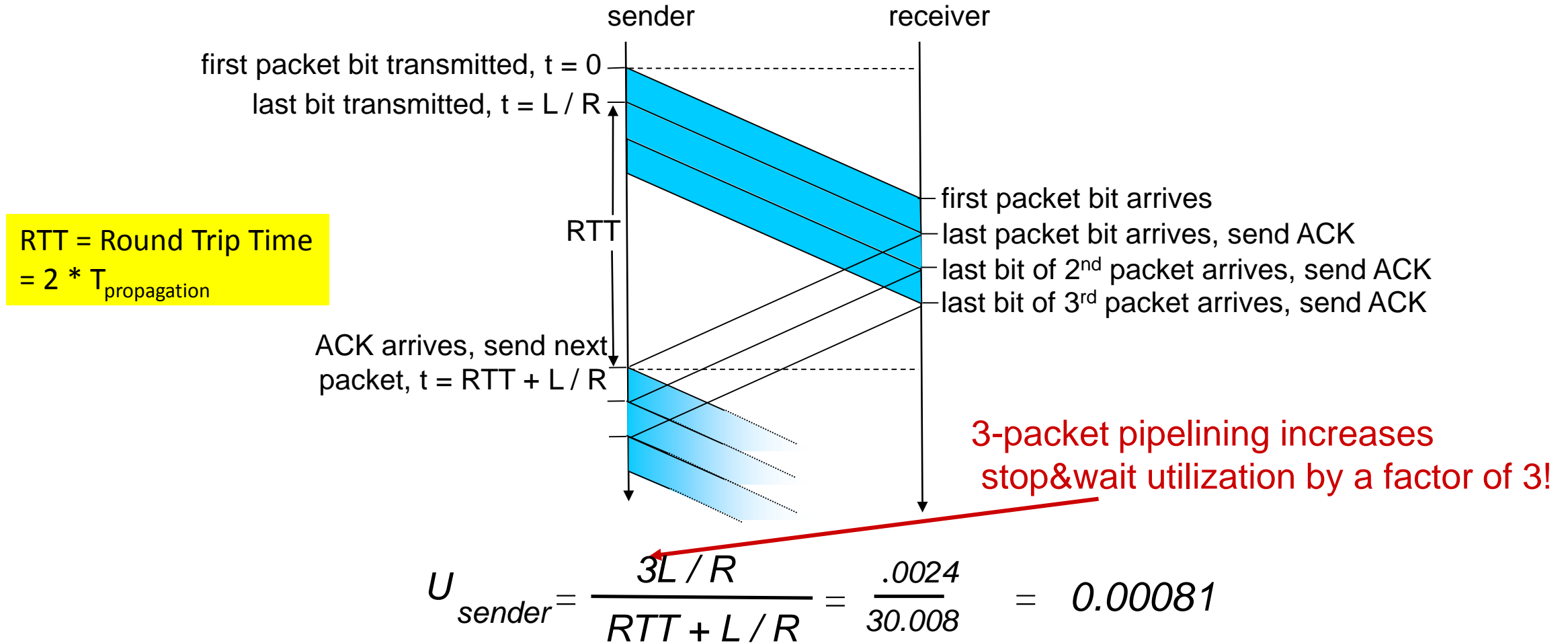  - *TCP congestion control*

# Pipelined ack-based error-control protocols

pipelining: sender allows multiple, "in-flight", yet-to-be-acknowledged pkts
- 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

# Pipelining: increased utilization



sender                    receiver

first packet bit transmitted, t = 0

last bit transmitted, t = L / R

RTT = Round Trip Time
= 2 * $T_{propagation}$

RTT

first packet bit arrives

last packet bit arrives, send ACK

last bit of $2^{nd}$ packet arrives, send ACK

last bit of $3^{rd}$ packet arrives, send ACK

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

3-packet pipelining increases
stop&wait utilization by a factor of 3!

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

# Pipelined protocols: ack-based error control

two generic forms (i.e. ack + book-keeping policies) to deal with lost data:
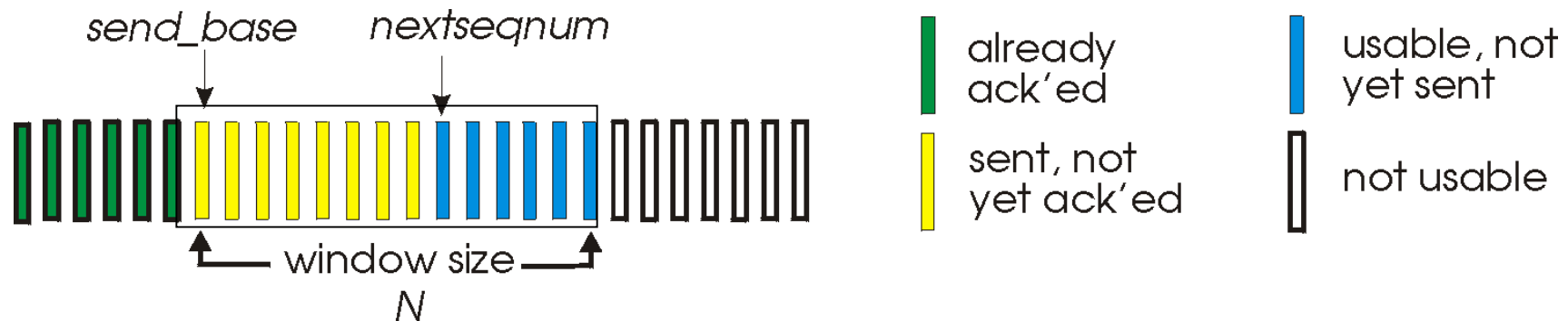
- *go-Back-n*
- *selective repeat*

# Go-Back-n

**sender**

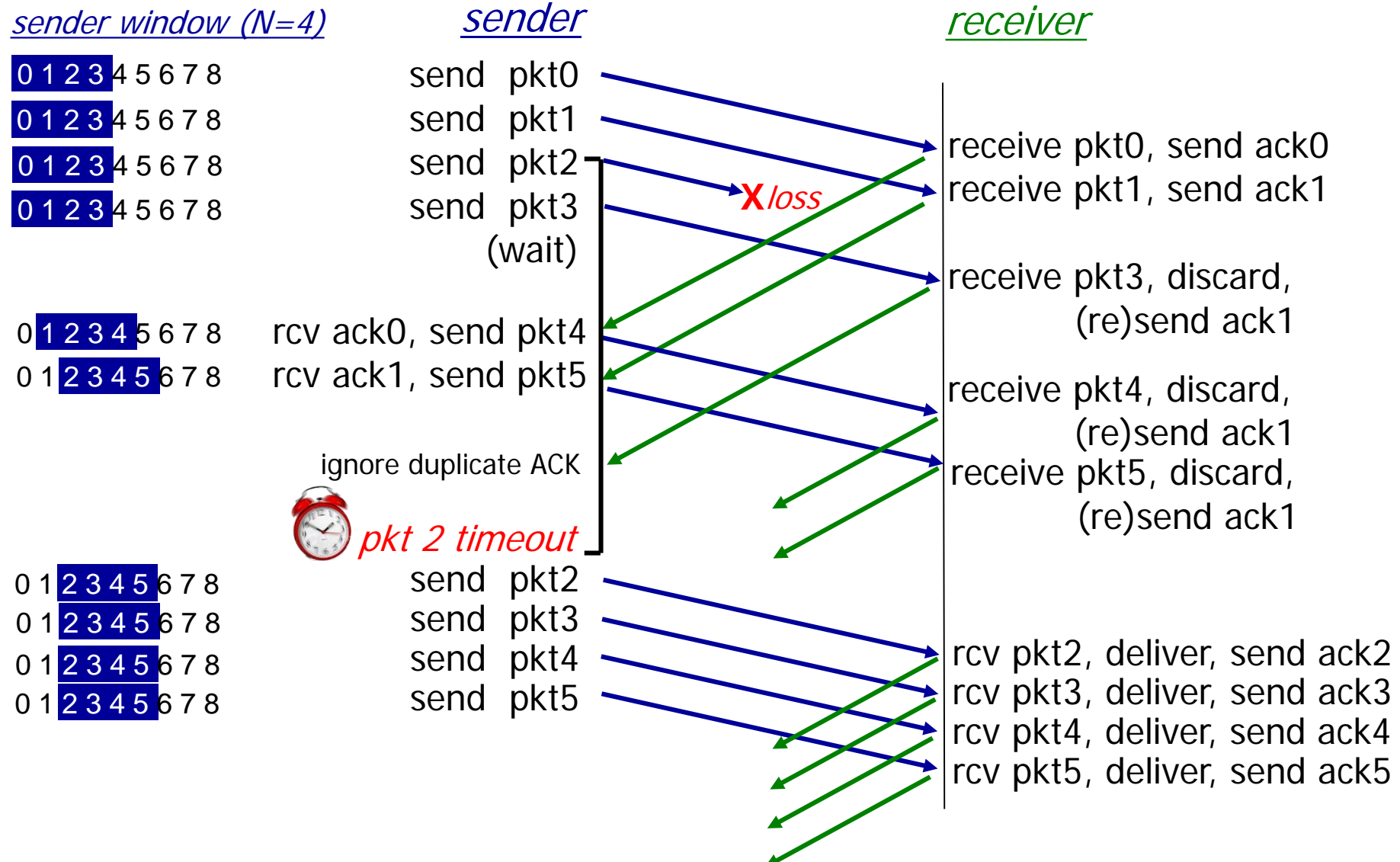- "window" of up to N, consecutive unack'ed pkts allowed

**receiver**

- Ack last correctly received pkt



- ACK(n): ACKs all pkts up to, including seq # n - *"cumulative ACK"*
  - may receive duplicate ACKs
- timer for oldest in-flight pkt
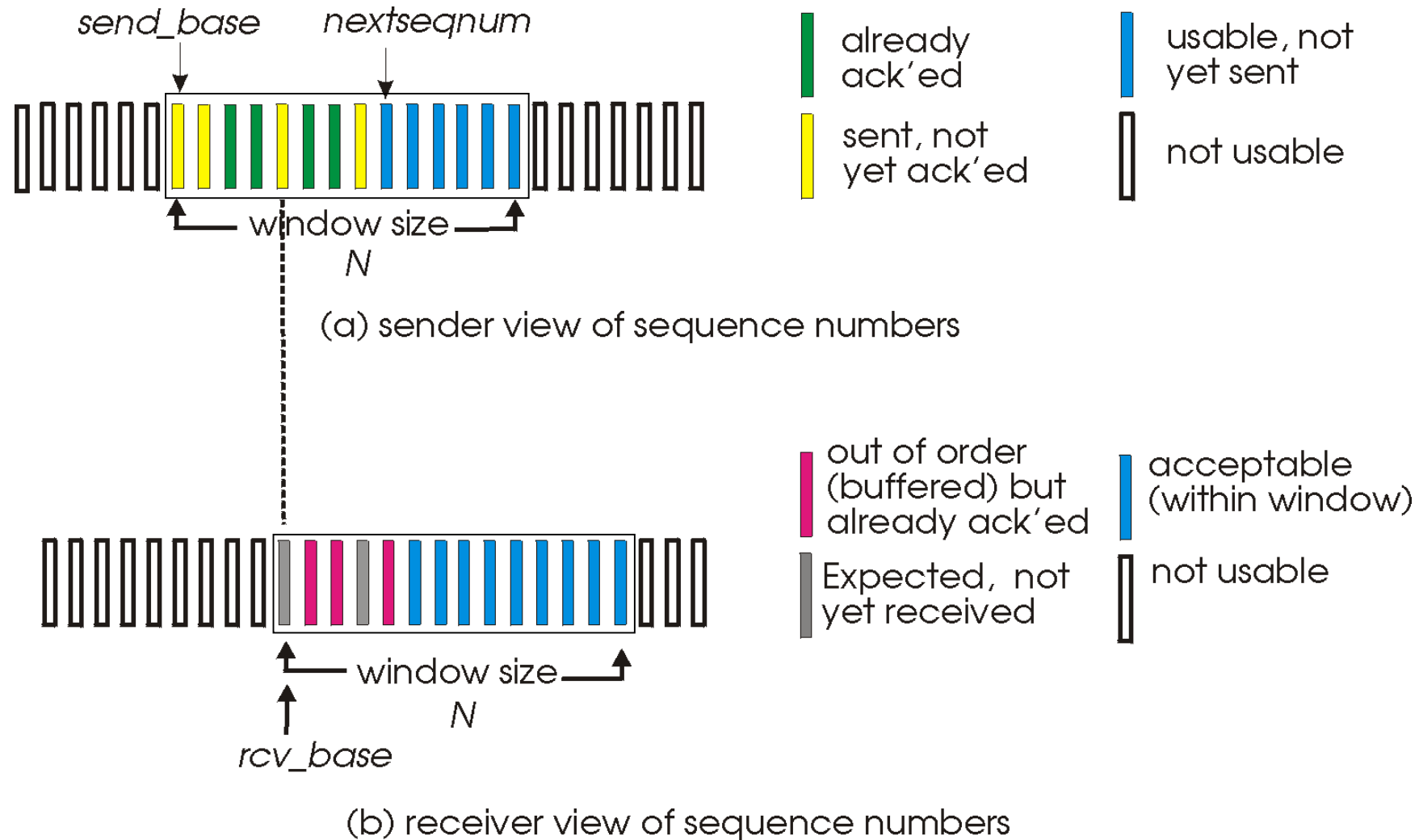- *timeout(n):* retransmit packet n and all higher seq # pkts in window

# GBn in action

_sender window (N=4)_          _sender_                          _receiver_

0 1 2 3 4 5 6 7 8          send  pkt0
0 1 2 3 4 5 6 7 8          send  pkt1
0 1 2 3 4 5 6 7 8          send  pkt2                            receive pkt0, send ack0
                                                      **X** _loss_  receive pkt1, send ack1
0 1 2 3 4 5 6 7 8          send  pkt3
                           (wait)
                                                                receive pkt3, discard,
                                                                      (re)send ack1
0 1 2 3 4 5 6 7 8          rcv ack0, send pkt4
0 1 2 3 4 5 6 7 8          rcv ack1, send pkt5
                                                                receive pkt4, discard,
                                                                      (re)send ack1
                           ignore duplicate ACK                 receive pkt5, discard,
                                                                      (re)send ack1

                           _pkt 2 timeout_

0 1 2 3 4 5 6 7 8          send  pkt2
0 1 2 3 4 5 6 7 8          send  pkt3
0 1 2 3 4 5 6 7 8          send  pkt4                            rcv pkt2, deliver, send ack2
0 1 2 3 4 5 6 7 8          send  pkt5                            rcv pkt3, deliver, send ack3
                                                                rcv pkt4, deliver, send ack4
                                                                rcv pkt5, deliver, send ack5
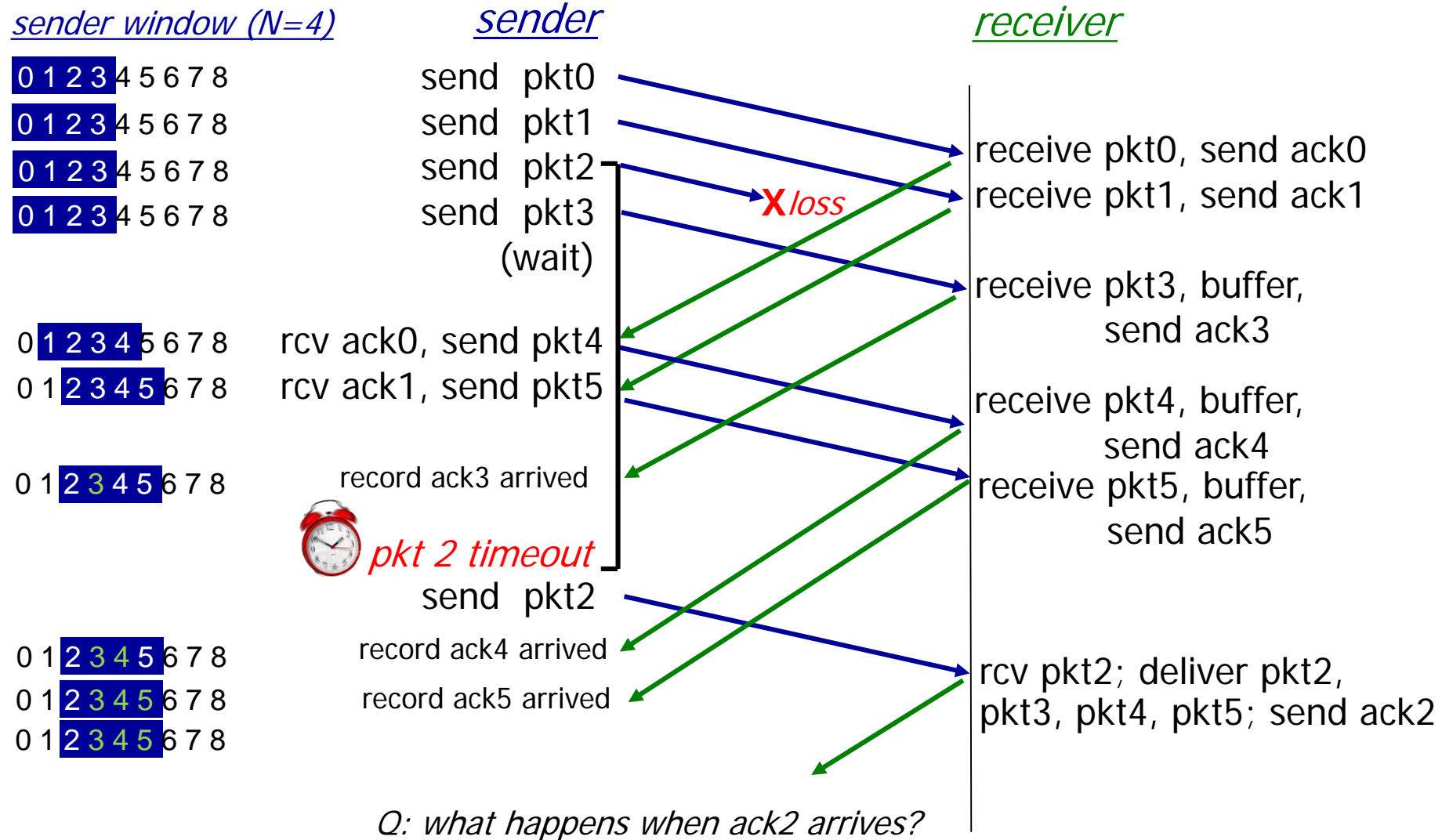
# Selective repeat: sender, receiver windows

- receiver *individually* acknowledges received pkts
  - buffers pkts for eventual in-order delivery to upper layer
- sender only resends pkts for which ACK not received
  - Requires timer for each unACKed pkt



send_base    nextseqnum

already ack'ed

usable, not yet sent

sent, not yet ack'ed

not usable

window size N

(a) sender view of sequence numbers

out of order (buffered) but already ack'ed

acceptable (within window)

Expected, not yet received

not usable

window size N

rcv_base

(b) receiver view of sequence numbers

# Selective repeat in action

*sender window (N=4)*          *sender*          *receiver*

`0 1 2 3` 4 5 6 7 8          send  pkt0

`0 1 2 3` 4 5 6 7 8          send  pkt1

`0 1 2 3` 4 5 6 7 8          send  pkt2          receive pkt0, send ack0

`0 1 2 3` 4 5 6 7 8          send  pkt3          receive pkt1, send ack1

**X** *loss*

(wait)

receive pkt3, buffer,
         send ack3

0 `1 2 3 4` 5 6 7 8          rcv ack0, send pkt4

0 1 `2 3 4 5` 6 7 8          rcv ack1, send pkt5          receive pkt4, buffer,
                                                            send ack4

0 1 `2 3 4 5` 6 7 8          record ack3 arrived          receive pkt5, buffer,
                                                            send ack5

*pkt 2 timeout*

send  pkt2

0 1 `2 3 4 5` 6 7 8          record ack4 arrived

0 1 `2 3 4 5` 6 7 8          record ack5 arrived          rcv pkt2; deliver pkt2,
                                                          pkt3, pkt4, pkt5; send ack2
0 1 `2 3 4 5` 6 7 8

*Q: what happens when ack2 arrives?*

# Roadmap

- Transport layer services in Internet
- Addressing, multiplexing/demultiplexing
- Connectionless, unreliable transport: UDP
- principles of reliable data transfer
  - Efficiency perspective: pipelined protocols & error control through go-back-n, selective-repeat
    - Sequence numbers

- *Next: connection-oriented transport: TCP*
  - *reliable transfer*
  - *flow control*
  - *connection management*
  - *TCP congestion control*

# Selective repeat:
# Sequence numbers

example:

- seq #'s: 0, 1, 2, 3
- window size=3

  ❖ duplicate data accepted as new in (b)

  Q: what relationship between seq # size and window size to avoid problem in (b)?

pkt0
pkt1
pkt2
pkt3
pkt0

0 1 2 3 0 1 2
0 1 2 3 0 1 2
0 1 2 3 0 1 2
0 1 2 3 0 1 2
0 1 2 3 0 1 2

0 1 2 3 0 1 2
0 1 2 3 0 1 2
0 1 2 3 0 1 2
0 1 2 3 0 1 2

X

*will accept packet with seq number 0*

(a) no problem

*receiver can't see sender side.*
*receiver behavior identical in both cases!*
*something's (very) wrong!*

pkt0
pkt1
pkt2

0 1 2 3 0 1 2
0 1 2 3 0 1 2
0 1 2 3 0 1 2

0 1 2 3 0 1 2
0 1 2 3 0 1 2
0 1 2 3 0 1 2

X
X
X

timeout
retransmit pkt0
0 1 2 3 0 1 2
pkt0

*will accept packet with seq number 0*

(b) oops!

# Question

## What do Ack's achieve besides reliability?

**Flow control**: receiver can ack its receiving capacity i.e. **avoid swamping the receiver**

**Flow control:** Sender/receiver (ie network edge) issue:  S cares to not overwhelm R

# Ack-based pipelining =>
# error-control & flow control at the same time!!!

sender                                    receiver

first packet bit transmitted, t = 0

last bit transmitted, t = L / R

RTT

first packet bit arrives

last packet bit arrives, send ACK

last bit of 2nd packet arrives, send ACK

last bit of 3rd packet arrives, send ACK

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

**Flow control:** Sender/receiver problem;  S cares to not overwhelm R

# Roadmap Transport Layer



- transport layer services
- multiplexing/demultiplexing
- connectionless transport: UDP
- principles of reliable data transfer
- **connection-oriented transport: TCP**
  - reliable transfer
    - Acknowledgements
    - Retransmissions
    - Connection management
    - Flow control and buffer space
  - Congestion control
    - Principles
    - TCP congestion control

# TCP: Overview  RFCs: 793,1122,1323, 2018, 5681

- ❖ point-to-point & full duplex data:
    - one sender, one receiver
    - bi-directional data flow in same connection
    - MSS: maximum segment size
- ❖ connection-oriented, reliable, in-order byte steam:
    - Needs handshaking (exchange of control msgs); inits sender & receiver state before data exchange
- ❖ Flow&error control: ack-based, pipelined:
- ❖ (+ extra) congestion control:
    - sender will not flood network

# TCP segment structure



URG: urgent data
(generally not used)

ACK: ACK # valid

PSH: push data now
(generally not used)

RST, SYN, FIN:
connection estab
(setup, teardown
commands)

Internet
checksum
(as in UDP)

32 bits

| source port # | dest port # |
| --- | --- |
| sequence number | |
| acknowledgement number | |
| head len / not used / U A P R S F | receive window |
| checksum | Urg data pointer |
| options (variable length) | |
| application data (variable length) | |

counting
by bytes
of data
(not segments!)

# bytes
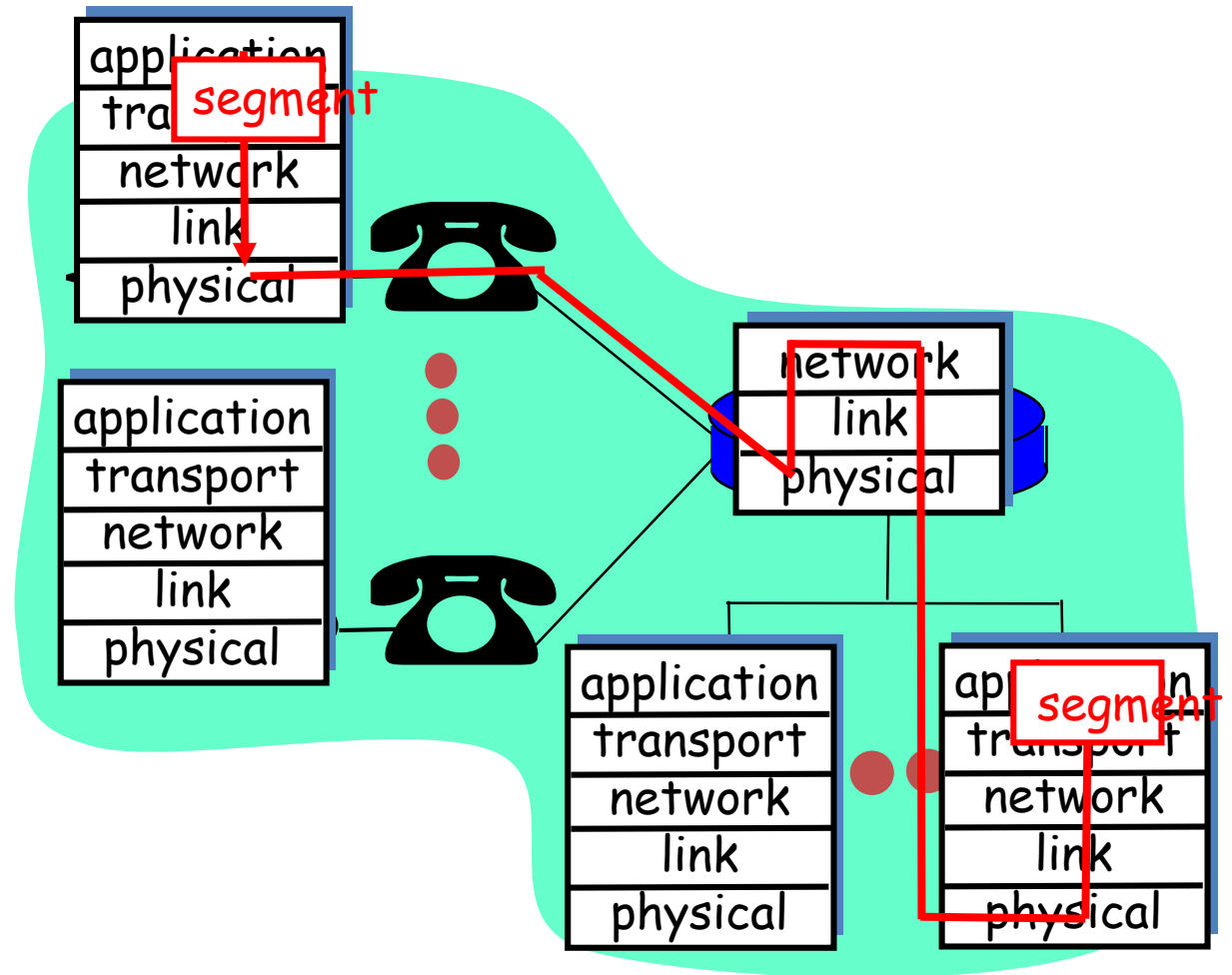rcvr willing
to buffer
*(flow control)*

# Roadmap Transport Layer

- transport layer services
- multiplexing/demultiplexing
- connectionless transport: UDP
- principles of reliable data transfer
- **connection-oriented transport: TCP**
  - reliable transfer
    - **Acknowledgements**
    - Retransmissions
    - Connection management
    - Flow control and buffer space
  - Congestion control
    - Principles
    - TCP congestion control

# TCP seq. numbers, ACKs
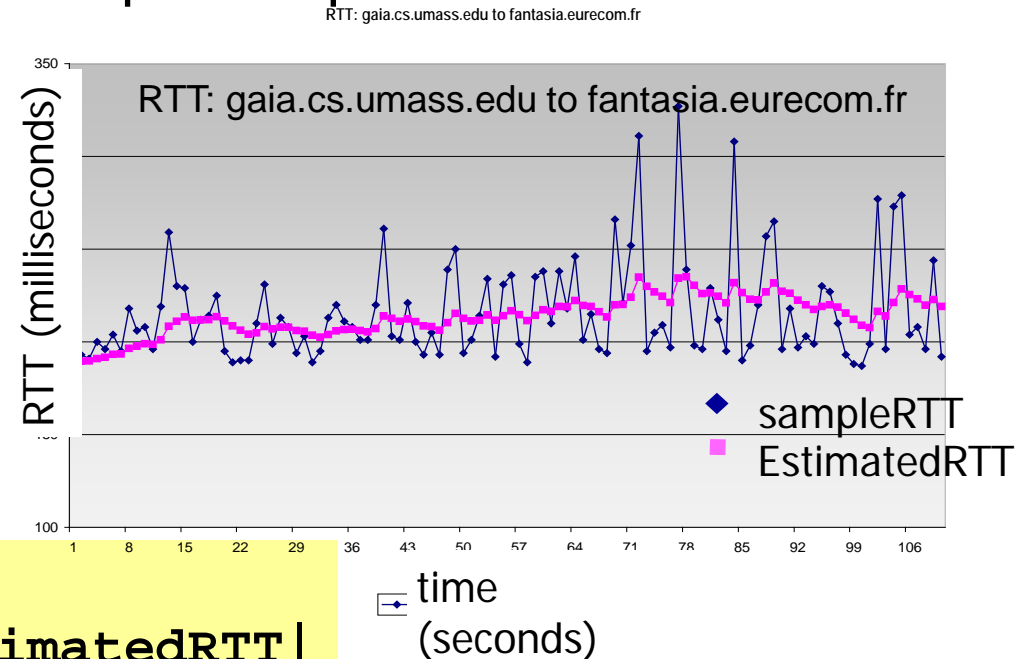
sequence numbers:

–"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 | |

window size
_N_

_sender sequence number space_

sent
ACKed

sent, not-
yet ACKed
("in-
flight")

usable
but not
yet sent

not
usable

incoming segment to sender

| source port # | dest port # |
|---|---|
| sequence number | |
| acknowledgement number | |
| A | rwnd |
| checksum | |

# TCP seq. numbers, ACKs

Always ack next in-order expected byte

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

Simple example scenario
Based on telnet msg exchange

Host A                    Host B

Seq=92, 8 bytes of data

ACK=100

X

Seq=92, 8 bytes of data

ACK=100

timeout

# TCP: cumulative Ack - retransmission scenarios

Host A       Host B       Host A       Host B

Seq=92, 8 bytes of data

Seq=100, 20 bytes of data

timeout

ACK=100

X

ACK=120

Seq=120,  15 bytes of data

**Cumulative ACK**

SendBase=92

Seq=92, 8 bytes of data

Seq=100, 20 bytes of data

timeout

ACK=100

ACK=120

SendBase=100

Seq=92,  8
bytes of data

SendBase=120

ACK=120

SendBase=120

**(Premature) timeout**

# TCP ACK generation [RFC 1122, RFC 5681]

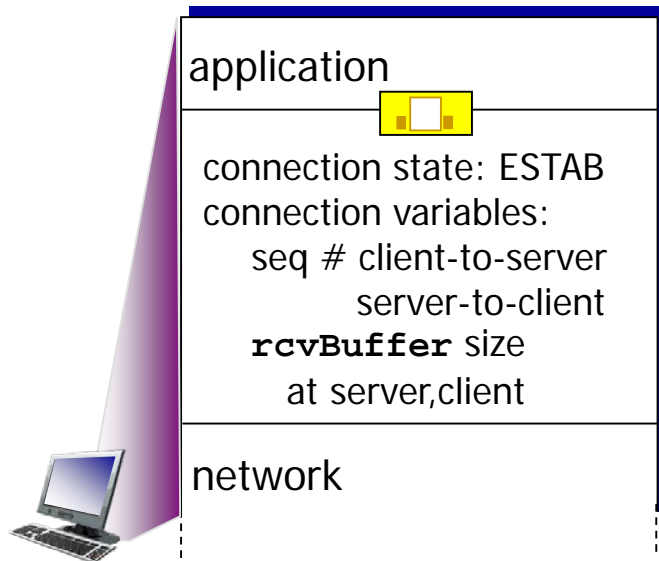| Event | TCP Receiver action |
|---|---|
| in-order segment arrival,<br>no gaps,<br>everything else already ACKed | **Delayed ACK**. Wait max 500ms for next segment then send ACK |
| in-order segment arrival,<br>no gaps,<br>one delayed ACK pending | immediately send single cumulative ACK |
| out-of-order segment arrival<br>higher-than-expect seq. #<br>gap detected | send (duplicate) ACK, indicating seq. # of next expected byte |

# Roadmap Transport Layer

- transport layer services
- multiplexing/demultiplexing
- connectionless transport: UDP
- principles of reliable data transfer
- **connection-oriented transport: TCP**
  - reliable transfer
    - Acknowledgements
    - **Retransmissions**
    - Connection management
    - Flow control and buffer space
  - Congestion control
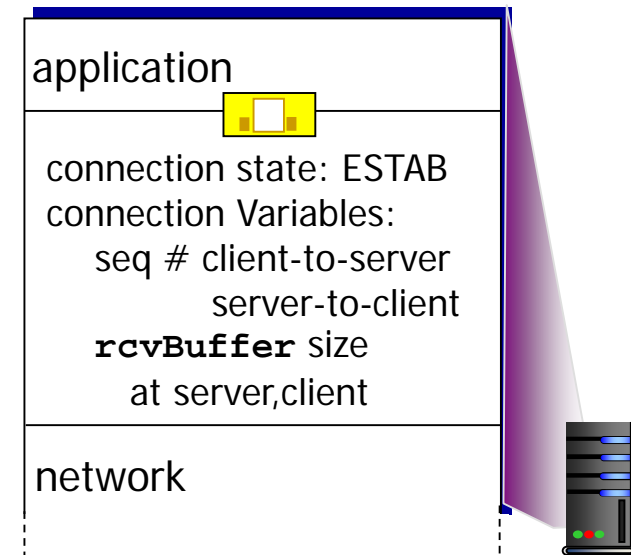    - Principles
    - TCP congestion control

# Q: how to set TCP timeout value?

❖ longer than end-to-end RTT
  ▪ but that varies!!!
❖ *too short* timeout*:*
  ❖ premature, unnecessary retransmissions
❖ *too long:*
  ❖ slow reaction to loss

# TCP round trip time, timeout estimation

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

❖ exponential weighted moving average: influence of past sample decreases exponentially fast

❖ typical value: $\alpha = 0.125$

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



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

◆ sampleRTT
■ EstimatedRTT

time (seconds)

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

(typically, $\beta = 0.25$)

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

"safety margin"

# TCP fast retransmit (RFC 5681)

❖ time-out can be long:
  ▪ long delay before resending lost packet

❖ IMPROVEMENT: detect lost segments via duplicate ACKs

*TCP fast retransmit*

if sender receives 3 duplicate ACKs for same data

• resend unacked segment with smallest seq #

  ▪ likely that unacked segment lost, so don't wait for timeout

Implicit NAK!
Q: Why need at least 3?

Host A                                    Host B

Seq=92, 8 bytes of data
Seq=100, 20 bytes of data
                              X

ACK=100
ACK=100
ACK=100

Seq=100, 20 bytes of data

timeout

# Q: Is TCP stateful or stateless?

# Is it possible to have reliable transfer over UDP?

# reliable transfer over UDP?

- add reliability at application layer

- top of UDP

ons and implementations for types of
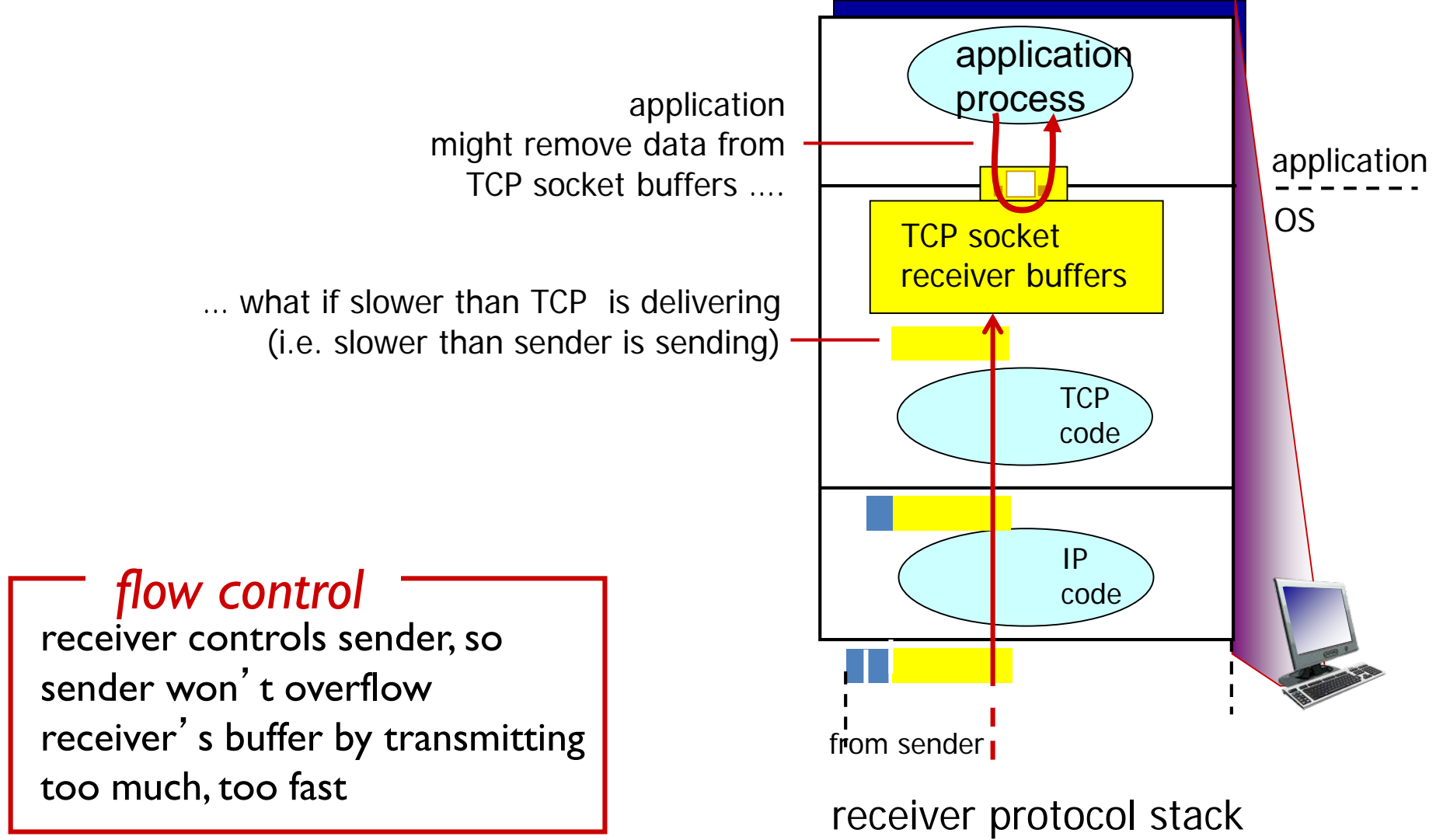ed and not need to be implemented
atch)

# Roadmap Transport Layer

- transport layer services
- multiplexing/demultiplexing
- connectionless transport: UDP
- principles of reliable data transfer
- **connection-oriented transport: TCP**
  - reliable transfer
    - Acknowledgements
    - Retransmissions
    - **Connection management**
    - Flow control and buffer space
  - Congestion control
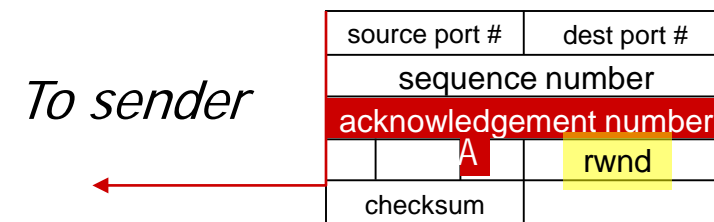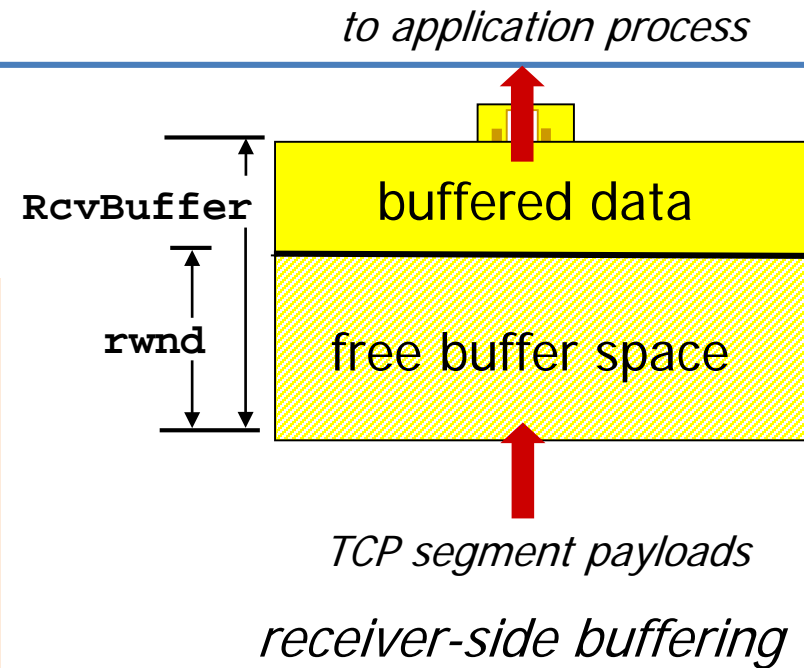    - Principles
    - TCP congestion control

# Connection Management

before exchanging data, sender/receiver "handshake":
- agree to establish connection + connection parameters

application

connection state: ESTAB
connection variables:
 seq # client-to-server
  server-to-client
 **rcvBuffer** size
 at server,client

network

application

connection state: ESTAB
connection Variables:
 seq # client-to-server
  server-to-client
 **rcvBuffer** size
 at server,client

network

```
Socket clientSocket =
    newSocket("hostname","port
    number");
```

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

# Setting up a connection: TCP 3-way handshake

*client state*

*server state*

LISTEN

choose init seq num, x
send TCP SYN msg

SYNSENT

SYN=1, Seq=x

choose init seq num, y
send TCP SYN/ACK
msg, acking SYN
Reserve buffer

SYN RCVD

SYN=1, Seq=y
ACK=1; ACKnum=x+1

received SYN/ACK(x)
server is live;
send ACK for SYN/ACK;
this segment may contain
client-to-server data

ESTAB

ACK=1, ACKnum=y+1

received ACK(y)
indicates client is live

ESTAB

# TCP: closing a connection

*client state*

*server state*

ESTAB

ESTAB

`clientSocket.close()`

FIN_WAIT_1

can no longer
send but can
receive data

FIN=1, seq=x

CLOSE_WAIT

ACK=1; ACKnum=x+1

can still
send data

FIN_WAIT_2

wait for server
close

FIN=1, seq=y

LAST_ACK

TIME_WAIT

can no longer
send data

ACK=1; ACKnum=y+1

timed wait
(typically 30s)

CLOSED

CLOSED

simultaneous FINs
can be handled

RST: alternative way to close connection
immediately, when **error** occurs

# Roadmap Transport Layer

- transport layer services
- multiplexing/demultiplexing
- connectionless transport: UDP
- principles of reliable data transfer
- **connection-oriented transport: TCP**
  - **reliable transfer**
    - Acknowledgements
    - Retransmissions
    - Connection management
    - **Flow control and buffer space**
  - Congestion control
    - Principles
    - TCP congestion control

# TCP flow control

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

**application
process**

application
- - - - - - - -
OS

TCP socket
receiver buffers

... what if slower than TCP is delivering
(i.e. slower than sender is sending)

TCP
code

*flow control*

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

IP
code

from sender

receiver protocol stack

# TCP flow control

- receiver "advertises" free buffer space through **rwnd** value in header
  - **RcvBuffer** size set via socket options (typical default 4 Kbytes)
  - OS can autoadjust **RcvBuffer**

- sender limits unacked ("in-flight") data to receiver's **rwnd** value
  - s.t. receiver's buffer will not overflow

*to application process*

RcvBuffer

rwnd

buffered data

free buffer space

*TCP segment payloads*

*receiver-side buffering*

*To sender*

| source port # | dest port # |
|---|---|
| sequence number | |
| acknowledgement number | |
| A | rwnd |
| checksum | |

# Roadmap Transport Layer

- transport layer services
- multiplexing/demultiplexing
- connectionless transport: UDP
- principles of reliable data transfer
- connection-oriented transport: TCP
  - reliable transfer
    - Acknowledgements
    - Retransmissions
    - Connection management
    - Flow control and buffer space
  - Congestion control
    - Principles
    - TCP congestion control

# Question:

## Difference between <u>congestion control</u> and <u>flow control</u>?

Congestion control = ***Avoid congesting the network***

**Congestion is network-core issue**

in contrast to

**flow-control, which is sender-receiver (i.e. network edge) issue**

# Principles of congestion control

*congestion*:

- informally: "many sources sending too much data too fast for *network* to handle"

- Manifestations?

  - lost packets (buffer overflow at routers)
  - long delays (queueing in router buffers)



http://photocarsonline.com/blog

# Distinction between flow control and congestion control



Transmission rate adjustment

Transmission network

Small-capacity receiver

(a)

Internal congestion

Large-capacity receiver

(b)

Fig. A. Tanenbaum
Computer Networks

Need for flow control

Need for congestion control

# Causes/costs of congestion

original data: $\lambda_{in}$

throughput: $\lambda_{out}$

Host A

output link capacity: R
link buffers

Host B

❖ Recall queueing behaviour + losses

❖ Losses => retransmissions => even higher load...

R/2

$\lambda_{out}$

$\lambda_{in}$  R/2

❖ Ideal per-connection throughput:
R/2 (if 2 connections)

delay

$\lambda_{in}$  R/2

C/2

$\lambda_{out}$

$\lambda_{in}'$ (incl. Retransmisions)

❖ reality 😦

# Approaches towards congestion control

approach taken by TCP

Not present in Internet's network layer protocols

**end-end congestion control:**

❖ no explicit feedback from network

❖ congestion inferred from end-system observed loss, delay

**network-assisted congestion control:**

❖ routers collaborate for optimal rates + provide feedback to end-systems eg.

- a single bit indicating congestion
- explicit rate for sender to send at

# Roadmap Transport Layer

- transport layer services
- multiplexing/demultiplexing
- connectionless transport: UDP
- principles of reliable data transfer
- connection-oriented transport: TCP
  - reliable transfer
    - Acknowledgements
    - Retransmissions
    - Connection management
    - Flow control and buffer space
  - Congestion control
    - Principles
    - **TCP congestion control**

# TCP congestion control:
## additive increase multiplicative decrease (AIMD)

❖ end-end control (no network assistance), sender limits transmission

How does sender perceive congestion?

- loss = timeout *or* 3 duplicate acks
- TCP sender reduces rate (**Congestion Window**) then

$$\text{rate} \approx \frac{\text{cwnd}}{\text{RTT}} \text{ bytes/sec}$$

- *Additive Increase:* increase **cwnd** by 1 MSS every RTT until loss detected

- *Multiplicative Decrease:* cut **cwnd** in half after loss

- **To start with: slow start**

AIMD saw tooth behavior: probing for bandwidth



additively increase window size ...

.... until loss occurs (then cut window in half)

**cwnd :** TCP sender congestion window size

time

# TCP Slow Start

❖ when connection begins, increase rate exponentially until first loss event:

- initially `cwnd` = 1 MSS
- double `cwnd` every ack of previous "batch"
- done by incrementing `cwnd` for every ACK received

❖ *summary:* initial rate is slow but ramps up exponentially fast

   ❖ then, saw-tooth

Host A                                    Host B

RTT

one segment

two segments

four segments

time

# TCP cwnd:
## from exponential to linear growth + reacting to loss



**Reno: loss indicated by timeout or 3 duplicate ACKs:** cwnd is cut in half; then grows linearly

**Implementation:**
- ❖ variable `ssthresh` (slow start threshold)
- ❖ on loss event, `ssthresh` = ½*`cwnd`

**Non-optimized: loss indicated by timeout:** cwnd set to 1 MSS; then window slow start to threshold, then grows linearly

# TCP's throughput (Fast recovery - Reno)

# 2 problems, joint solution: limit the rate of the sender!
## (or "How many windows does a TCP's sender maintain?")

rwnd

cwnd

cwnd, rwnd

Transmission rate adjustment

Transmission network

Internal congestion

Small-capacity receiver

Large-capacity receiver

(a)

(b)

Fig. A. Tanenbaum
Computer Networks

Need for flow control

Need for congestion control

# TCP combined flow-ctrl, congestion ctrl windows

*sender sequence number space*

Min{cwnd, rwnd}

last byte
ACKed

sent, not-
yet ACKed
("in-
flight")

last byte
sent

*TCP sending rate:*

❖ send min {cwnd, rwnd} bytes, wait for ACKS, then send more

sender limits transmission:

$$\text{LastByteSent} - \text{LastByteAcked} \leq \text{Min\{cwnd, rwnd\}}$$

❖ **cwnd** is dynamic, function of perceived network congestion,
❖ **rwnd** dymanically limited by receiver's buffer space

# TCP Fairness

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



TCP connection 1

TCP connection 2

bottleneck
router
capacity R

# Q: Can a TCP implementation deviate from the Congestion-Control standard?

# Roadmap Transport Layer

- transport layer services
- multiplexing/demultiplexing
- connectionless transport: UDP
- principles of reliable data transfer
- connection-oriented transport: TCP
  - reliable transfer
    - Acknowledgements
    - Retransmissions
    - Connection management
    - Flow control and buffer space
  - Congestion control
    - Principles
    - TCP congestion control

# Chapter 3: summary

❖ principles behind transport layer services:

- Addressing

- reliable data transfer

- flow control

- congestion control

❖ instantiation, implementation in the Internet

- UDP

- TCP

next:
- leaving the network "edge" (application, transport layers)
- into the network "core"

# Some more review questions on this part

- 
- rd ack and not a 2nd?
- le, method for detection of

- ndefinitely?
- ent?
- t and the end of connection?
- sfer if it uses UDP? How or why not?

# Reading instructions chapter 3

- ## KuroseRoss book

| Careful | Quick |
|---|---|
| 3.1, 3.2, 3.4-3.7 | 3.3 |

- ## Other resources (further  study)
  - Eddie Kohler, Mark Handley, and Sally Floyd. 2006. Designing DCCP: congestion control without reliability. *SIGCOMM Comput. Commun. Rev.* 36, 4 (August 2006), 27-38. DOI=10.1145/1151659.1159918 http://doi.acm.org/10.1145/1151659.1159918

  - http://research.microsoft.com/apps/video/default.aspx?id=104005

  - Exercise/throughput analysis TCP in extra

  -  slides

# Extra slides, for further study

# From RFC 1122: TCP Ack

- TCP SHOULD implement a delayed ACK, but an ACK should not be excessively delayed; in particular, the delay MUST be less than 0.5 seconds, and in a stream of full-sized segments there SHOULD be an ACK for at least every second segment.

- A delayed ACK gives the application an opportunity to update the window and perhaps **to send an immediate response**. In particular, in the case of character-mode remote login, a delayed ACK can reduce the number of segments sent by the server by a factor of 3 (ACK, window update, and echo character all combined in one segment).

- In addition, on some large multi-user hosts, a delayed ACK can substantially reduce protocol processing overhead by reducing the total number of packets to be processed.

- However, excessive delays on ACK's can disturb the round-trip timing and packet "clocking" algorithms.

- We also emphasize that this is a SHOULD, meaning that an implementor should indeed only deviate from this requirement after careful consideration of the implications.
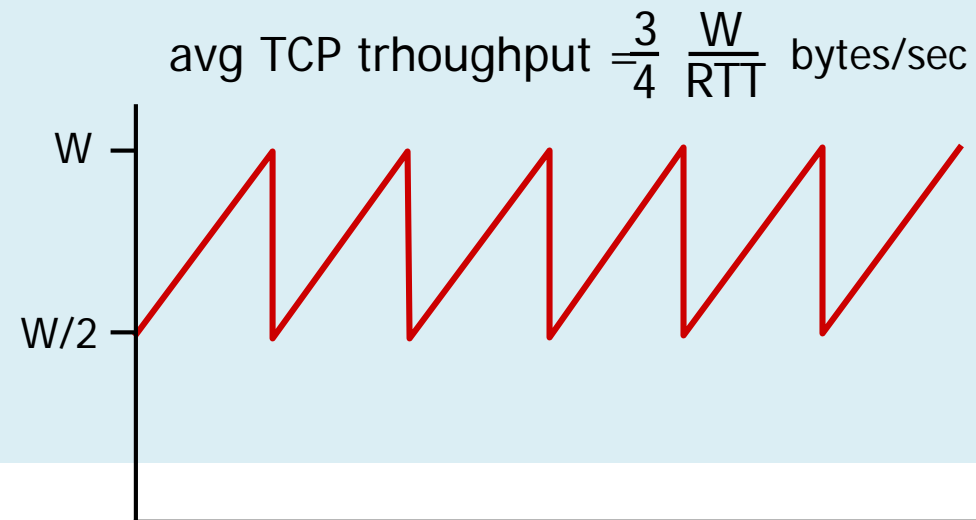
# TCP – Closing a connection: Reset

**RST**

- RST is used to signal an error condition and causes an immediate close of the connection on both sides

- RST packets are not supposed to carry data payload, except for an optional human-readable description of what was the reason for dropping this connection.

- Examples:
  - A TCP data segment when no session exists
  - Arrival of a segment with incorrect sequence number
  - Connection attempt to non-existing port
  - Etc.

# TCP throughput

- avg. TCP throughput as function of window size, RTT?
  - ignore slow start, assume always data to send

- W: window size (measured in bytes) where loss occurs
  - avg. window size (# in-flight bytes) is ¾ W
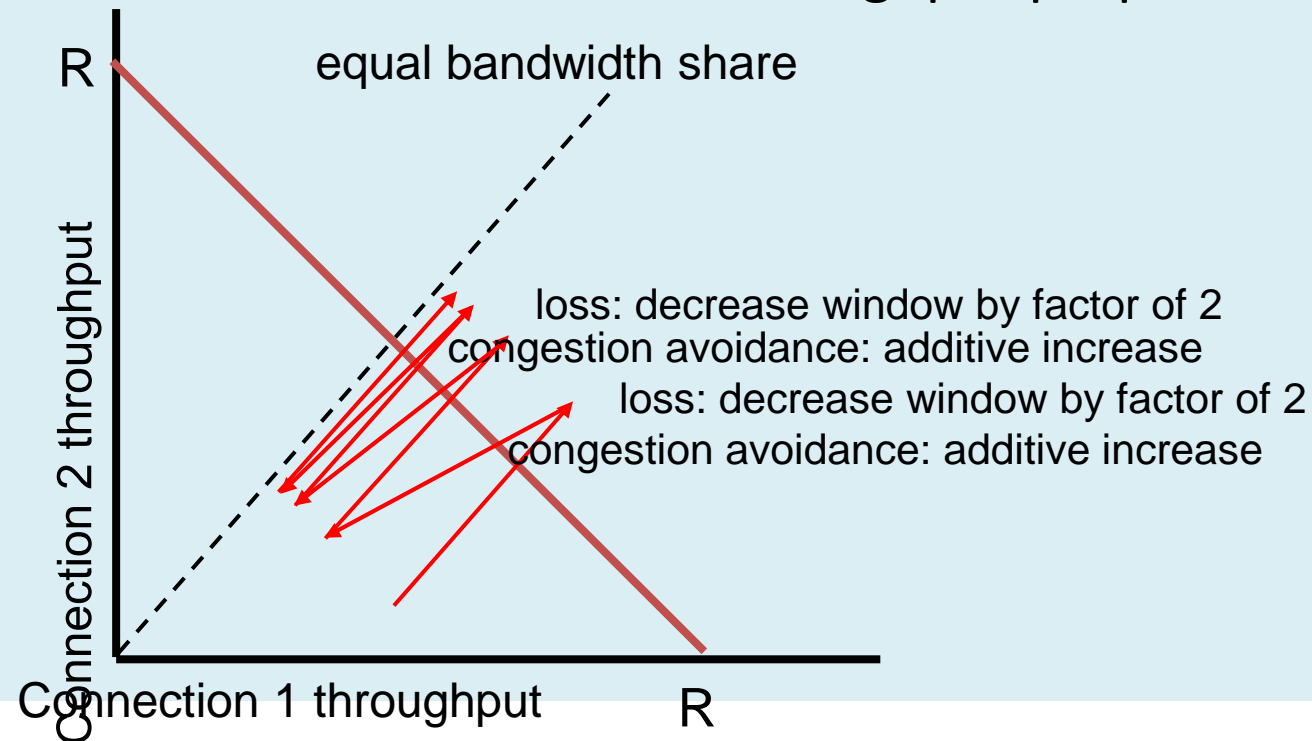  - avg. trhoughput is 3/4W per RTT

$$\text{avg TCP trhoughput} = \frac{3}{4} \frac{W}{RTT} \text{ bytes/sec}$$

**Marina Papatriantafilou – Transport layer part2: TCP**

# TCP Futures: TCP over "long, fat pipes"

- example: 1500 byte segments, 100ms RTT, want 10 Gbps throughput

- requires W = 83,333 in-flight segments

- throughput in terms of segment loss probability, L [Mathis 1997]:

$$\text{TCP throughput} = \frac{1.22 \cdot \text{MSS}}{\text{RTT} \sqrt{L}}$$

➔ to achieve 10 Gbps throughput, need a loss rate of L = $2 \cdot 10^{-10}$ *– a very small loss rate!*

- new versions of TCP for high-speed

# Why is TCP fair?

two competing sessions:

❖ additive increase gives slope of 1, as throughout increases

❖ multiplicative decrease decreases throughput proportionally



R

equal bandwidth share

loss: decrease window by factor of 2
congestion avoidance: additive increase

loss: decrease window by factor of 2
congestion avoidance: additive increase

Connection 2 throughput

Connection 1 throughput

R

# Fairness (more)

## Fairness and UDP

❖ multimedia apps often do not use TCP

- do not want rate throttled by congestion control

❖ instead use UDP:

- send audio/video at constant rate, tolerate packet loss

## Fairness, parallel TCP connections

❖ application can open multiple parallel connections between two hosts

❖ web browsers do this

❖ e.g., link of rate R with 9 existing connections:

- new app asks for 1 TCP, gets rate R/10
- new app asks for 11 TCPs, gets R/2

# TCP delay modeling (slow start – related)

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

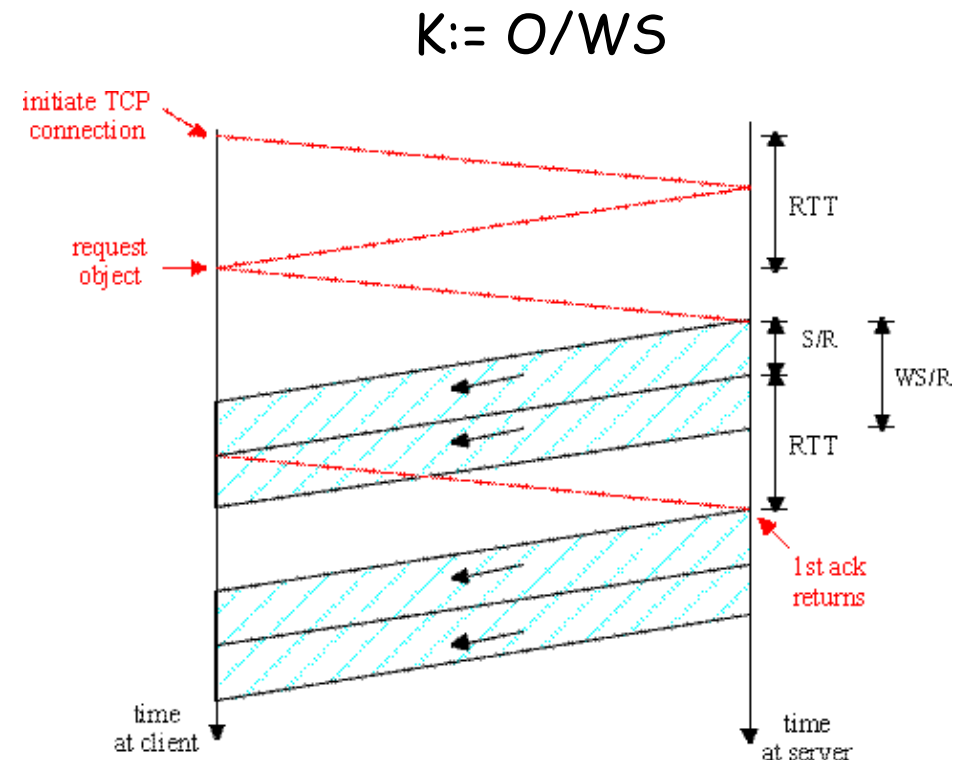- TCP connection establishment
- data transfer delay

Notation, assumptions:

- Assume one link between client and server of rate R
- Assume: fixed congestion window, W segments
- S: MSS (bits)
- O: object size (bits)
- no retransmissions (no loss, no corruption)
- Receiver has unbounded buffer

# TCP delay Modeling: simplified, fixed window

K:= O/WS



Case 1: WS/R > RTT + S/R:
ACK for first segment in window returns before window's worth of data nsent
delay = 2RTT + O/R

Case 2: WS/R < RTT + S/R:
wait for ACK after sending window's worth of data sent
delay = 2RTT + O/R
+ (K-1)[S/R + RTT - WS/R]

$$delay = \frac{O}{R} + 2RTT + \sum^{P} idleTime_p$$

# TCP Delay Modeling: Slow Start



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

Server idles:
 P = min{K-1,Q} times

where
- Q  = #times server stalls until cong. window is  larger than a "full-utilization" window (if the object were of unbounded size).

- K = #(incremental-sized) congestion-windows that "cover" the object.

Diagram labels:
- initiate TCP connection
- request object
- RTT
- first window = S/R
- second window = 2S/R
- third window = 4S/R
- fourth window = 8S/R
- object delivered
- complete transmission
- time at client
- time at server

**Example:**
- O/S  = 15 segments
- K = 4 windows
- Q = 2
- Server idles P = min{K-1,Q} = 2 times

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

$$\text{until server receives acknowledgement}$$

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

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

$$\text{delay} = \frac{O}{R} + 2RTT + \sum_{p=1}^{P} idleTime_p$$

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

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

initiate TCP connection

request object

RTT

first window = S/R

second window = 2S/R

third window = 4S/R

fourth window = 8S/R

object delivered

complete transmission

time at client

time at server

# TCP Delay Modeling

Recall K = number of windows that cover object

How do we calculate K ?

$$K = \min\{k : 2^0 S + 2^1 S + \mathsf{L} + 2^{k-1} S \geq O\}$$

$$= \min\{k : 2^0 + 2^1 + \mathsf{L} + 2^{k-1} \geq O/S\}$$

$$= \min\{k : 2^k - 1 \geq \frac{O}{S}\}$$

$$= \min\{k : k \geq \log_2(\frac{O}{S} + 1)\}$$

$$= \left\lceil \log_2(\frac{O}{S} + 1) \right\rceil$$

Calculation of Q, number of idles for infinite-size object, is similar.