# Information Network

Lecture 8 : Reliable Data Transfer

Holger Thies

京都大学

KYOTO UNIVERSITY

# Chapter 3 outline

京都大学

# Principles of reliable data transfer



reliable service *abstraction*

# Principles of reliable data transfer



reliable service *abstraction*

reliable service *implementation*

# Principles of reliable data transfer



Complexity of reliable data transfer protocol  will depend (strongly) on characteristics of unreliable channel (lose, corrupt, reorder data?)

# Principles of reliable data transfer

Sender, receiver do *not* know the "state" of each other, e.g., was a message received?
- unless communicated via a message



sending process

application
transport
data

sender-side of reliable data transfer protocol

transport
network

reliable service *implementation*

# Sequence numbers

- Packets are sent one by one and are numbered (sequence numbers).

- When a packet is sent, its sequence number is contained in the header information.

- When the receiver receives a packet and it does not contain any errors an acknowledgment message (ACK) is sent for this packet and the receiver remembers the sequence number.

- The sender waits for the acknowledgment packet until it sends the next packet.

- If the sender does not receive an acknowledgment for some time, the same packet is sent again.

- If the receiver receives a packet with a sequence number other than expected or the packet has errors, the packet is discarded, and an acknowledgment message for the last valid packet is send.

京都大学

# Reliable data transfer using ack messages



(a) no loss

(b) packet loss

# Reliable data transfer using ack messages

**sender**

send pkt0
→ pkt0 →
rcv pkt0
send ack0

rcv ack0
← ack0 ←
send pkt1
→ pkt1 →
rcv pkt1
send ack1
← ack1 ✗ loss

⏰ *timeout*
resend pkt1
→ pkt1 →
rcv pkt1
(detect duplicate)
send ack1

rcv ack1
← ack1 ←
send pkt2
→ pkt2 →
rcv pkt2
send ack2
← ack2 ←

**receiver**

(c) ACK loss

---

**sender**

send pkt0
→ pkt0 →
rcv pkt0
send ack0

rcv ack0
← ack0 ←
send pkt1
→ pkt1 →
rcv pkt1
send ack1
← ack1

⏰ *timeout*
resend pkt1
→ pkt1 →
rcv pkt1
(detect duplicate)
send ack1

rcv ack1
send pkt2
→ pkt2 →
rcv pkt2
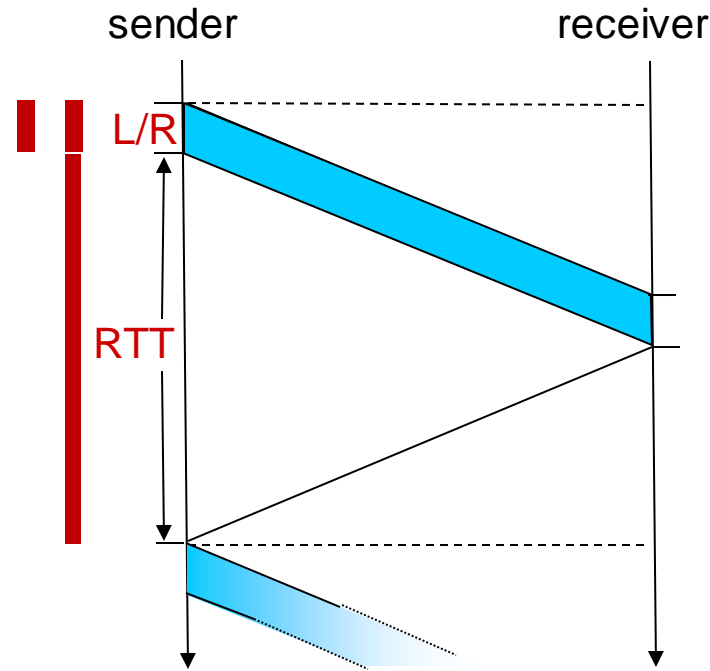send ack2

rcv ack1
(ignore)
← ack1 ←
← ack2 ←
→ pkt3 →

**receiver**

(d) premature timeout/ delayed ACK

---

Note: 1 bit sequence number (0 or 1) suffices

# Stop and wait protocol

$$U_{sender} = \frac{L/R}{RTT + D}$$

$$= \frac{.008}{30.008}$$

$$= 0.027\%$$



- D is the transmission delay (the amount of time required to push all the packet's bits into the wire)
- if RTT=30 msec, 1KB packet every 30 msec: 33kB/sec throughput over 1 Gbps link
- Protocol limits performance of underlying infrastructure (channel)

# Pipelined protocols

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



(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

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
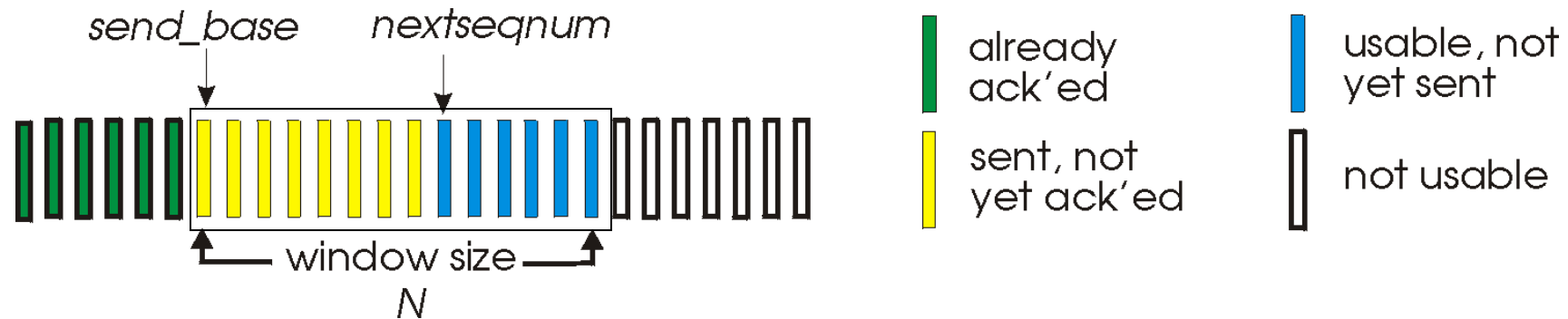
RTT

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

3-packet pipelining increases
utilization by a factor of 3!

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

京都大学

# Go-Back-N: sender

- k-bit sequence number in packet header

- "window" of up to N, consecutive unacknowledged packets allowed



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

# Go-Back-N: receiver

- ACK-only: always send ACK for correctly-received packet so far, with highest *in-order* seq number
  - may generate duplicate ACKs
  - need only remember `rcv_base`

- on receipt of out-of-order packet:
  - can discard (don't buffer) or buffer: an implementation decision
  - re-ACK packet with highest in-order sequence number

Receiver view of sequence number space:



`rcv_base`

received and ACKed

Out-of-order: received but not ACKed

Not received

# Go-Back-N 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 ─── **X** *loss*
`0 1 2 3` 4 5 6 7 8    send  pkt3
                      (wait)

                                              receive pkt0, send ack0
                                              receive pkt1, send ack1

                                              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
                                              receive pkt5, discard,
                      ignore duplicate ACK         (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
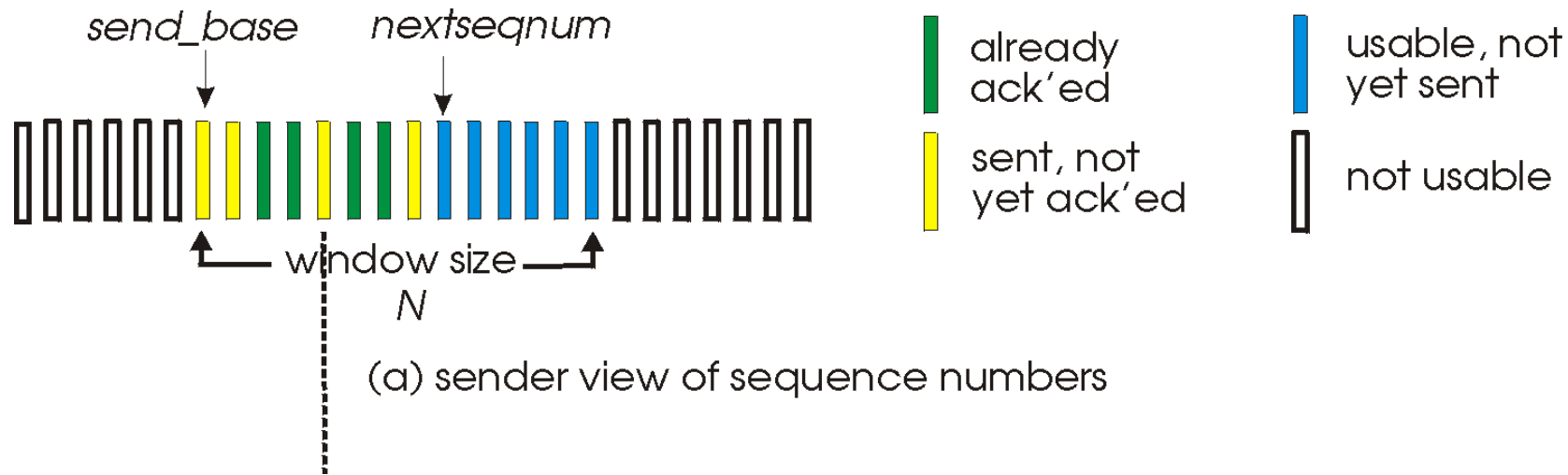                                         rcv pkt5, deliver, send ack5

京都大学

# Selective repeat: the approach

▪ *pipelining*:  *multiple* packets in flight

▪ *receiver individually ACKs* all correctly received packets
  - buffers packets, as needed, for in-order delivery to upper layer

▪ sender:
  - maintains (conceptually) a timer for each unACKed pkt
    - timeout: retransmits single unACKed packet  associated with timeout
  - maintains (conceptually) "window" over  *N* consecutive seq #s
    - limits pipelined, "in flight" packets to be within this window

# Selective repeat: sender, receiver windows



(a) sender view of sequence numbers

send_base     nextseqnum

window size N

- already ack'ed (green)
- sent, not yet ack'ed (yellow)
- usable, not yet sent (blue)
- not usable (white)

# Example

Host A sends a message to Host B consisting of 10 packets using a pipelined reliable data transfer protocol with the Go-Back-N strategy.

Assume the window size is 4 (that is, N=4) and that every 6th packet that A sends to B is lost (but no acknowledgment message from B to A is lost and there are no premature timeouts).

How many packets will A have to send in total to transmit the message to B?

# Chapter 3 outline

3.1 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

- segment structure
- reliable data transfer
- flow control
- connection management
- congestion control

# TCP: Transmission Control Protocol

RFC: 793

TRANSMISSION CONTROL PROTOCOL

DARPA INTERNET PROGRAM
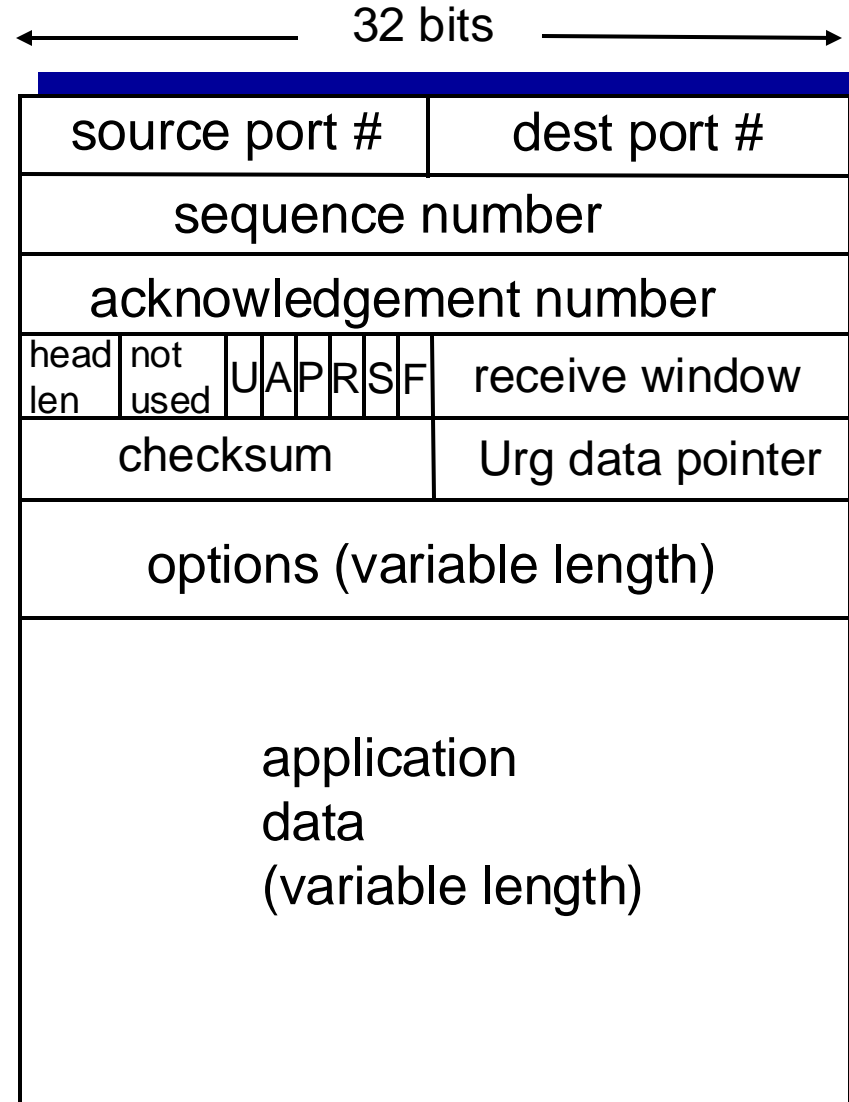
PROTOCOL SPECIFICATION

September 1981

1.  INTRODUCTION

The Transmission Control Protocol (TCP) is intended for use as a highly
reliable host-to-host protocol between hosts in packet-switched computer
communication networks, and in interconnected systems of such networks.

This document describes the functions to be performed by the
Transmission Control Protocol, the program that implements it, and its
interface to programs or users that require its services.
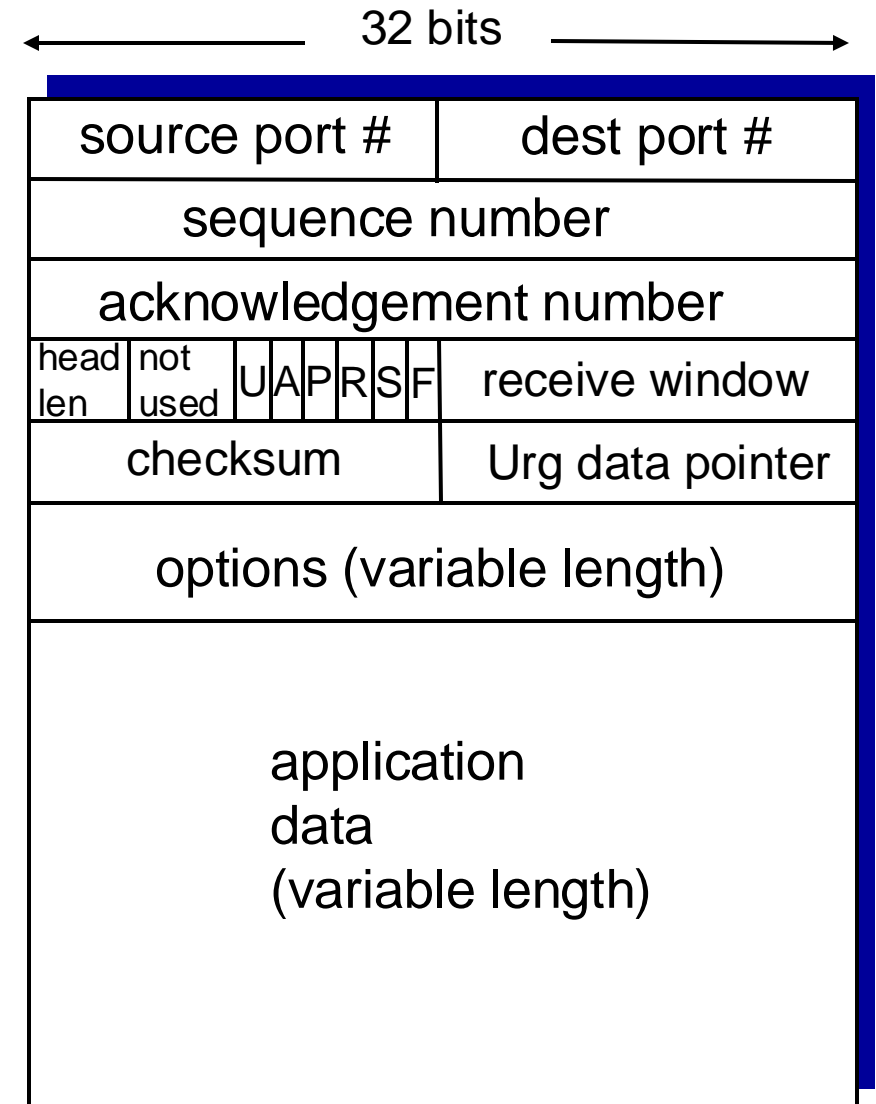
京都大学

# TCP: Overview

- **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
- **connection-oriented:**
  - handshaking (exchange of control mesages) to initialize sender and receiver state before data exchange
- **flow controlled:**
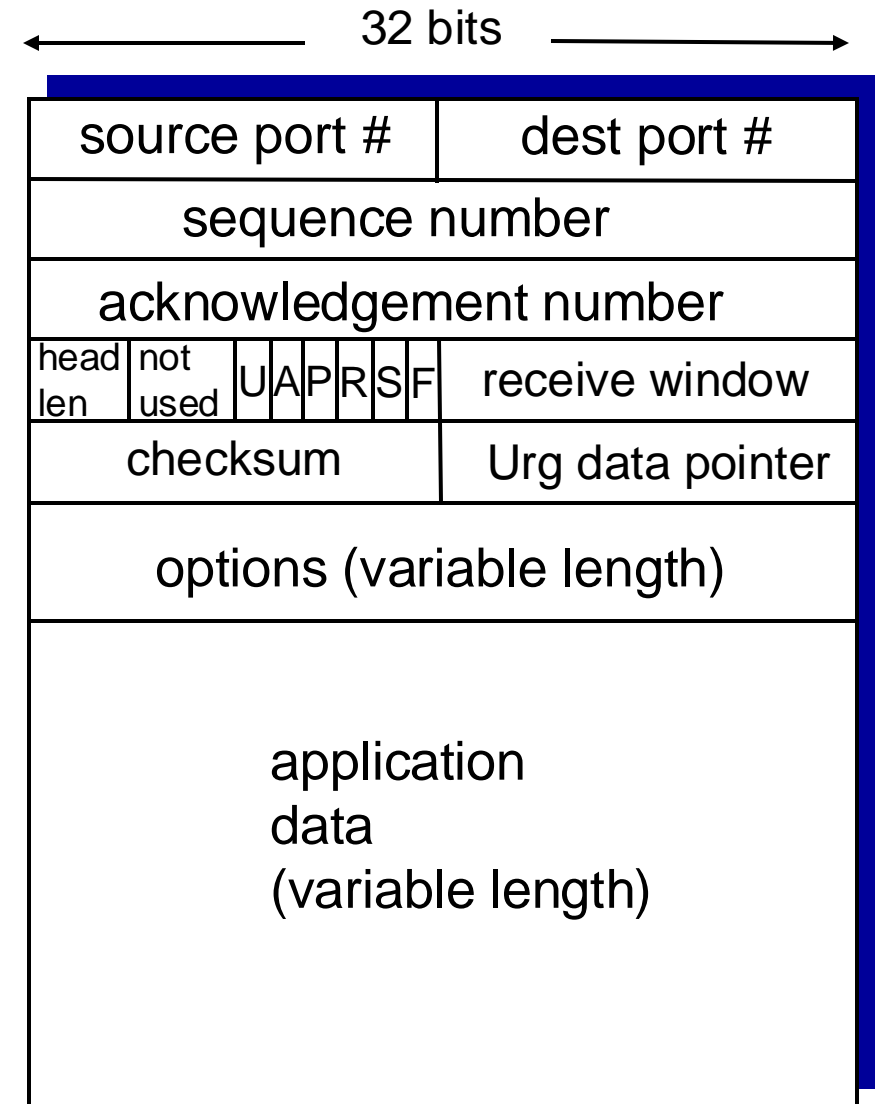  - sender will not overwhelm receiver

京都大学

# TCP segment structure

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)

京都大学

# TCP Segment Header

- **Source Port, Destination Port**

- **Sequence Number**
  - At the transport layer application data is split into several smaller segments.
  - Sequence Number is used to keep track of the position of the current segment in the sequence.

- **Acknowledgment Number: Number of the next expected segment**
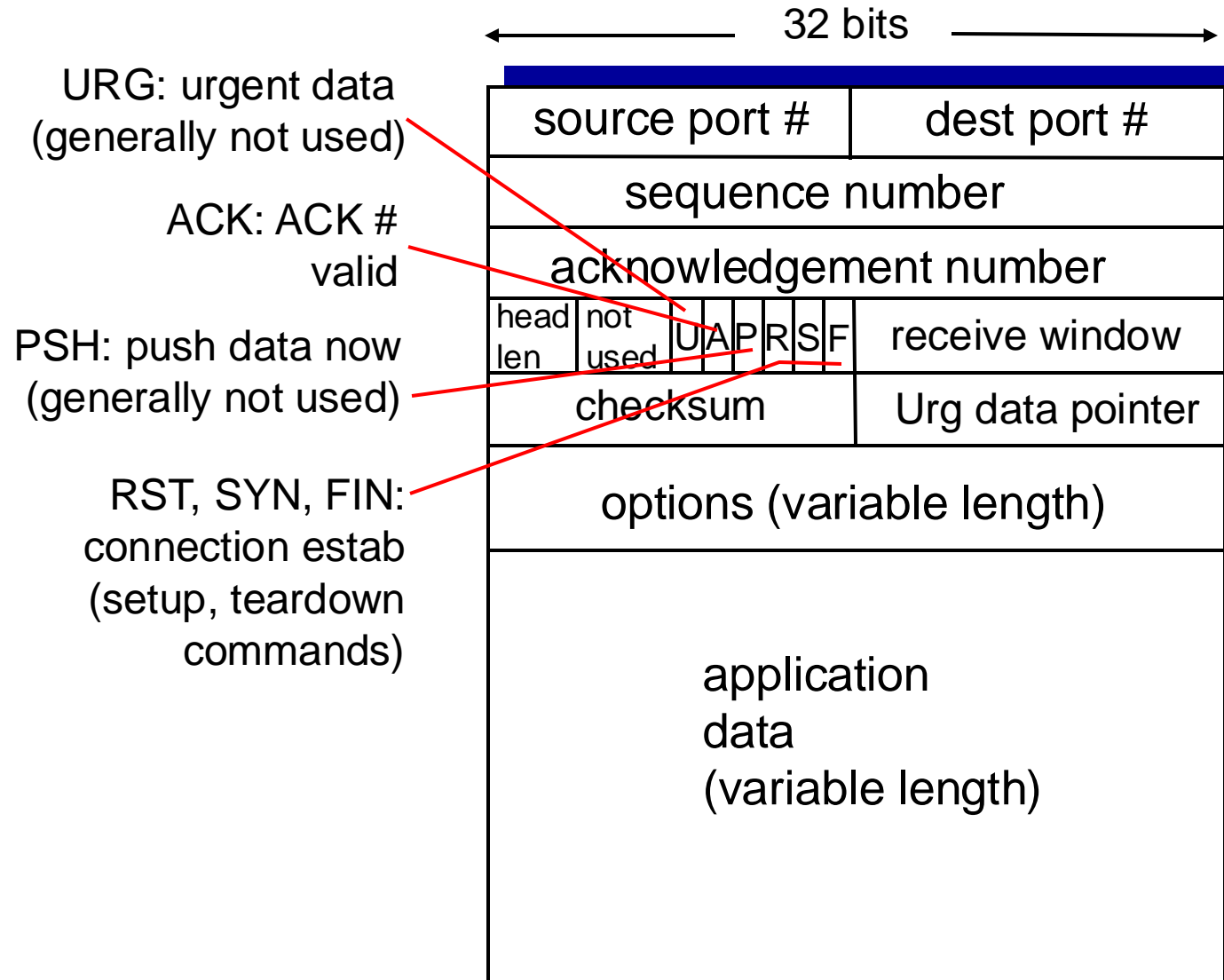  - Used for reliable data transfer.

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)

京都大学

# TCP Segment Header

- head len: Length of Header
  - For the receiver to know where the header ends and the application data begins
- 6 TCP control flags
- receive window
  - Number of bytes the receiver is willing to accept at a time
- Checksum
  - Same as in UDP
- Urg data pointer
  - Point out segments that are urgent
  - Rarely used
- Options
  - Additional options such as more complicated flow control
  - rarely used

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)

京都大学

# TCP control flags

URG: urgent data
(generally not used)

ACK: ACK #
valid

PSH: push data now
(generally not used)

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

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)

京都大学

3-27

# TCP seq. numbers, ACKs

**sequence numbers:**

- byte stream "number" of first byte in segment's data

**acknowledgements:**

- Sequence number of next byte expected from other side

- cumulative acknowledgments

- TCP does not specify how the receiver handles out-of-order segments, it is up to implementor

outgoing segment from sender

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

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 | urg pointer |

京都大学

# TCP sequence numbers, ACKs

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 scenario

京都大学

# Chapter 3 outline

京都大学

# TCP reliable data transfer

Reliability:

The TCP must recover from data that is damaged, lost, duplicated, or delivered out of order by the internet communication system.  This is achieved by assigning a sequence number to each octet transmitted, and requiring a positive acknowledgment (ACK) from the receiving TCP.  If the ACK is not received within a timeout interval, the data is retransmitted.  At the receiver, the sequence numbers are used to correctly order segments that may be received out of order and to eliminate duplicates.  Damage is handled by adding a checksum to each segment transmitted, checking it at the receiver, and discarding damaged segments.

As long as the TCPs continue to function properly and the internet system does not become completely partitioned, no transmission errors will affect the users.  TCP recovers from internet communication system errors.

京都大学

# TCP reliable data transfer

- TCP creates reliable data transfer service on top of IP's unreliable service
  - pipelined segments
  - cumulative acks
  - single retransmission timer
- retransmissions triggered by:
  - timeout events
  - duplicate acks

京都大学

# TCP Sender (simplified)

event: data received from application

- create segment with seq number

- seq number is byte-stream number of first data byte in segment
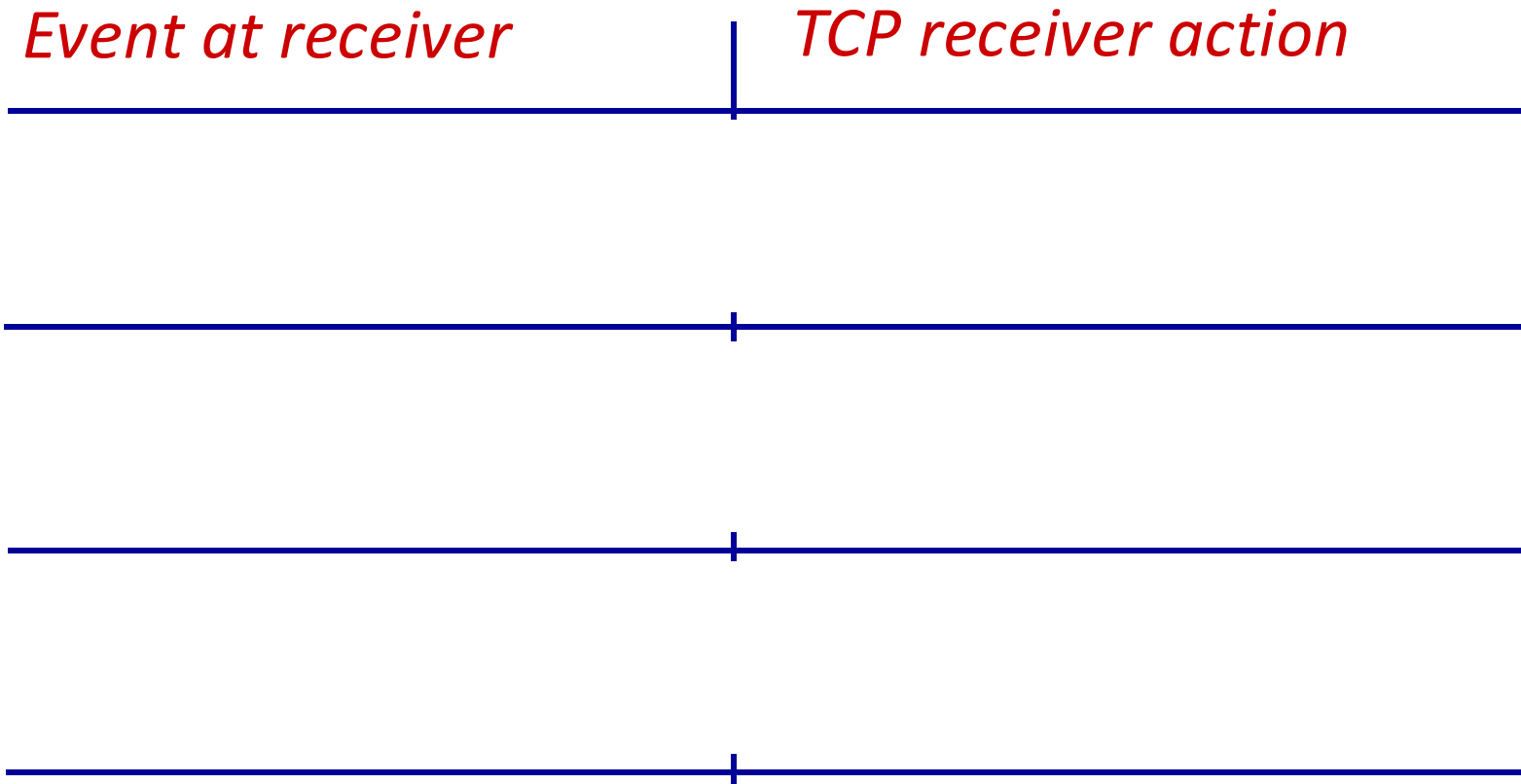
- start timer if not already running
  - Timer is for oldest unacknowledged segment

*event: timeout*

- retransmit segment that caused timeout
- restart timer

*event: ACK received*

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

京都大学

# TCP Receiver: ACK generation [RFC 5681]

| Event at receiver | TCP receiver action |
| --- | --- |
| | |
| | |
| | |

# TCP: retransmission scenarios



lost ACK scenario

premature timeout

# TCP: retransmission scenarios



Host A                    Host B

timeout

Seq=92, 8 bytes of data

Seq=100, 20 bytes of data

ACK=100

**X**

ACK=120

Seq=120,  15 bytes of data

cumulative ACK

京都大学

# 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 number
- likely that unacked segment lost, so don't wait for timeout

京都大学

# TCP fast retransmit

Host A                                    Host B

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

ACK=100

ACK=100

ACK=100

ACK=100

Seq=100, 20 bytes of data

timeout

fast retransmit after sender
receipt of triple duplicate ACK

# Chapter 3 outline

京都大学

# Flow control and congestion control

- Both flow control and congestion control are about slowing down the sender, when more data is sent than can be handled.

- **Flow control** is about slowing down the sender when more data is sent than <u>the receiver</u> can handle.

- **Congestion control** is about slowing down the sender when more data is sent than <u>the network</u> can handle.

- In general, congestion control is more complex than flow control as it needs to operate across the entire network, involving multiple devices etc.
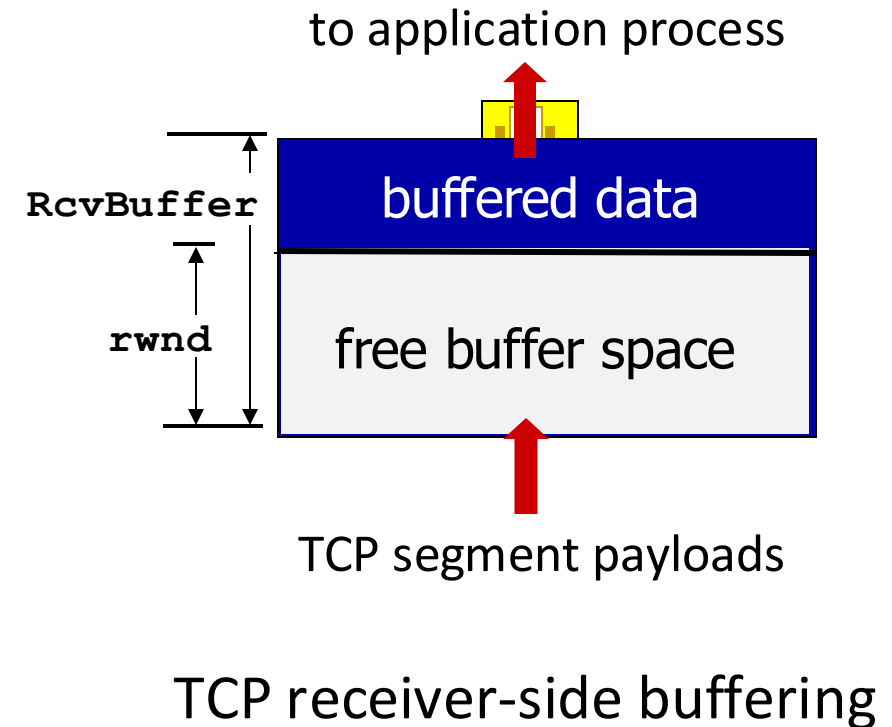
京都大学

# 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

application process

application

OS

TCP socket receiver buffers

TCP code

IP code

from sender

receiver protocol stack

京都大学

# TCP flow control

- TCP receiver "advertises" free buffer space in **rwnd** field in TCP header
  - **RcvBuffer** size set via socket options (typical default is 4096 bytes)
  - many operating systems autoadjust **RcvBuffer**

- sender limits amount of unacknowledged ("in-flight") data to received **rwnd**

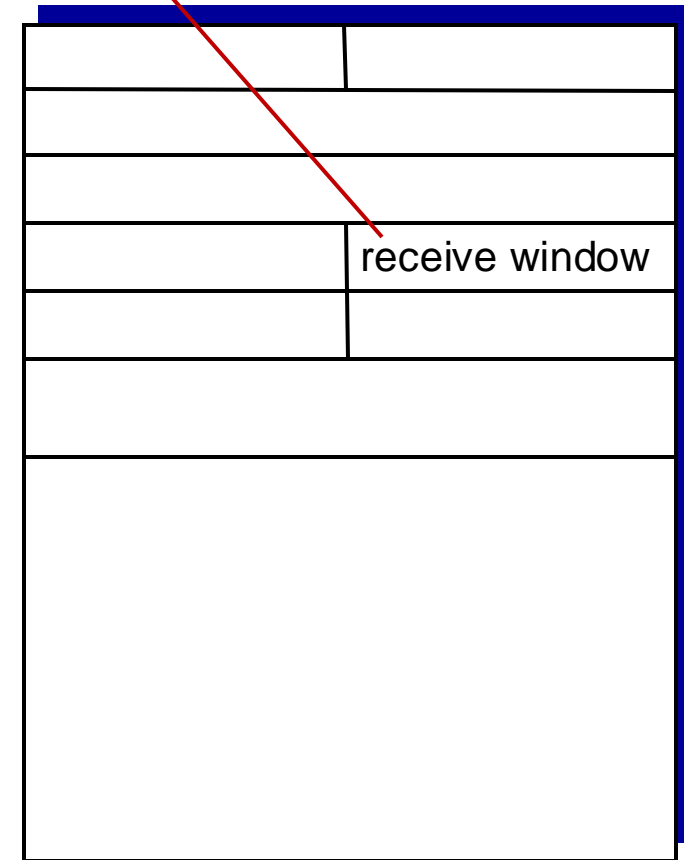- guarantees receive buffer will not overflow

to application process

RcvBuffer

buffered data

rwnd

free buffer space

TCP segment payloads

TCP receiver-side buffering

# TCP flow control

- TCP receiver "advertises" free buffer space in **rwnd** field in TCP header
  - **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 received **rwnd**

- guarantees receive buffer will not overflow

flow control: # bytes receiver willing to accept

receive window

TCP segment format

# Chapter 3 outline
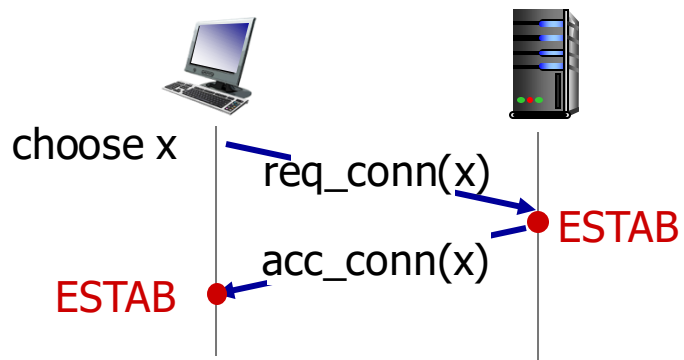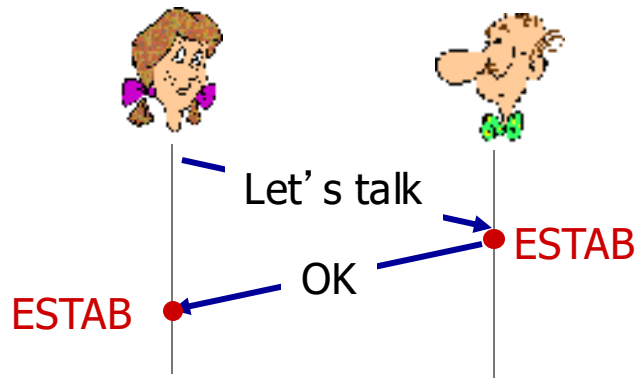
京都大学

# Connection Management

before exchanging data, sender/receiver "handshake":

- agree to establish connection (each knowing the other willing to establish connection)

- agree on 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

京都大学

# Agreeing to establish a connection

2-way handshake:



*Q:* will 2-way handshake always work in network?
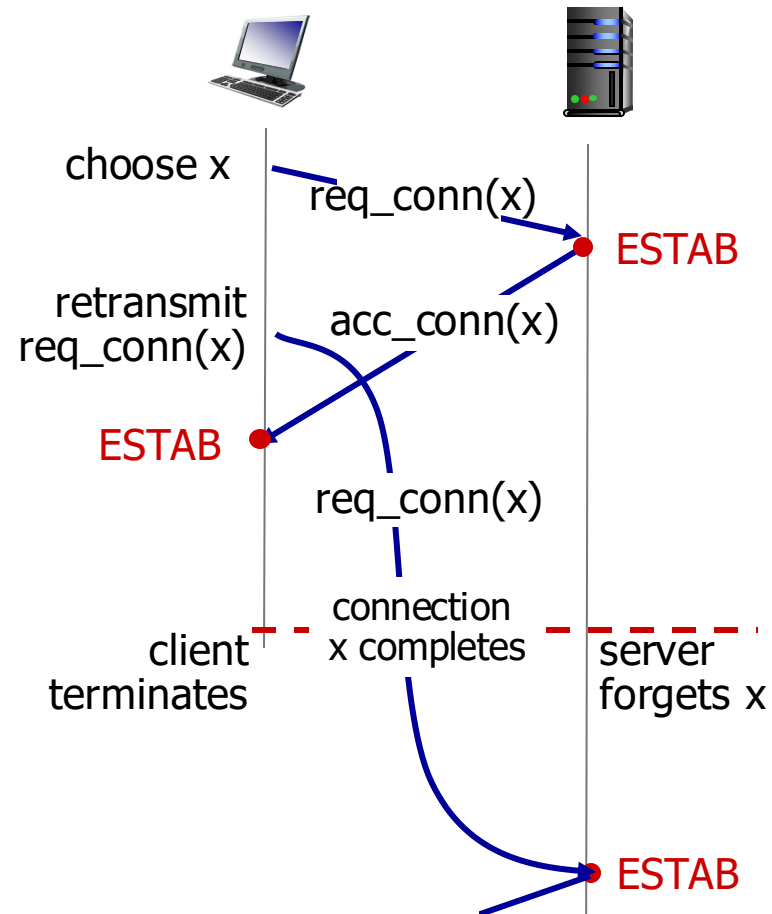
- variable delays

- retransmitted messages (e.g. req_conn(x)) due to message loss

- message reordering

- can't "see" other side

# 2-way handshake scenarios



choose x

req_conn(x)

ESTAB

acc_conn(x)

ESTAB

data(x+1)

accept
data(x+1)

ACK(x+1)

connection
x completes

No problem!

# 2-way handshake scenarios



choose x

req_conn(x)

ESTAB

retransmit
req_conn(x)

acc_conn(x)

ESTAB

req_conn(x)

connection
x completes

client
terminates

server
forgets x

ESTAB

❌ Problem: half open connection! (no client)

# 2-way handshake scenarios

choose x

req_conn(x)

ESTAB

retransmit
req_conn(x)

acc_conn(x)

ESTAB

data(x+1)

accept
data(x+1)

retransmit
data(x+1)

connection
x completes

server
forgets x

client
terminates

req_conn(x)

ESTAB

data(x+1)

accept
data(x+1)

❌ Problem: duplicate
data accepted!

# TCP 3-way handshake

*client state*

*server state*

LISTEN

LISTEN

choose init seq num, x
send TCP SYN msg

SYNSENT

SYNbit=1, Seq=x

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

SYN RCVD

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

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

ESTAB

ACKbit=1, ACKnum=y+1

received ACK(y)
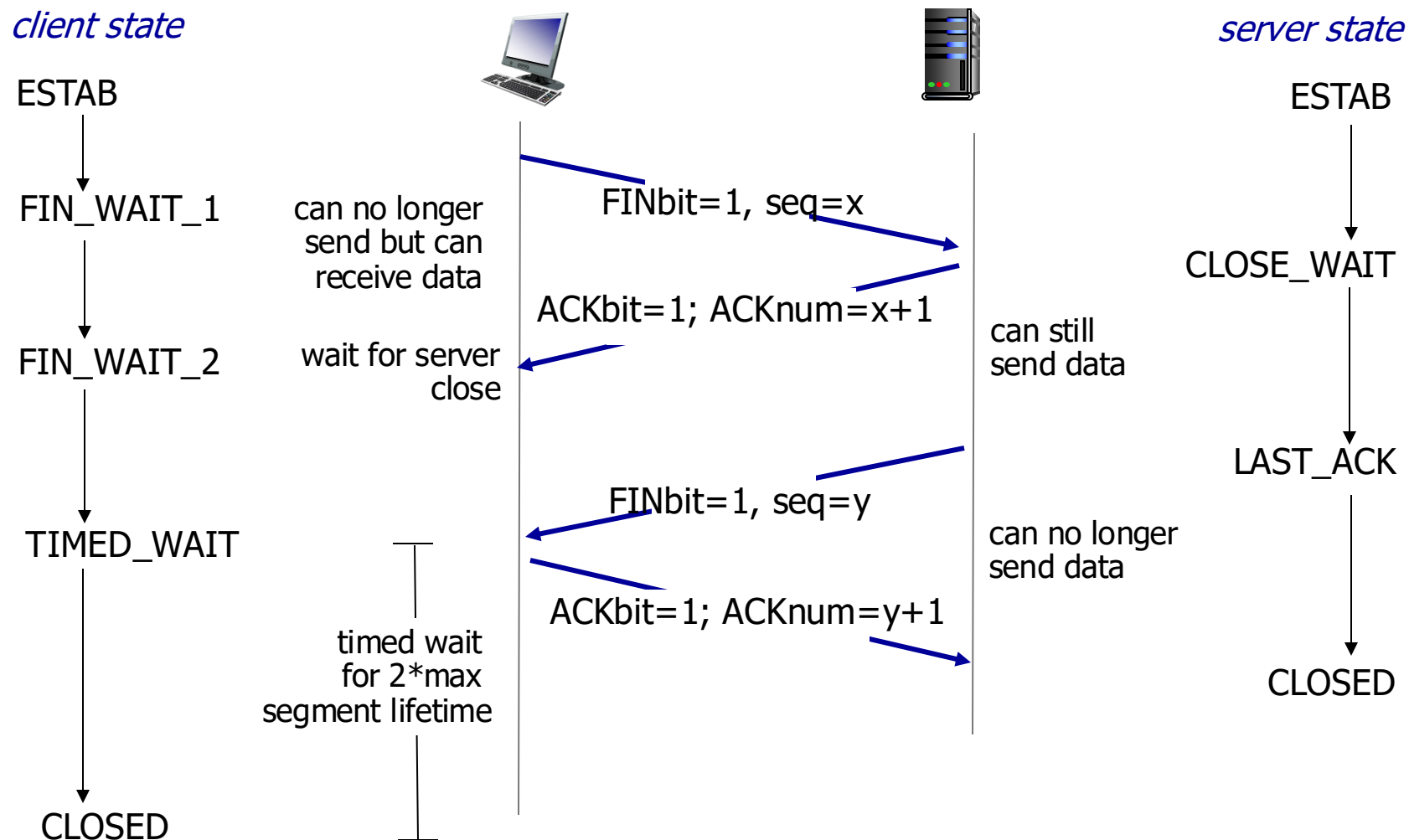indicates client is live

ESTAB

京都大学

# 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

京都大学

# TCP: closing a connection

client state

server state

ESTAB

ESTAB

FIN_WAIT_1

can no longer
send but can
receive data

FINbit=1, seq=x

CLOSE_WAIT

ACKbit=1; ACKnum=x+1

FIN_WAIT_2

wait for server
close

can still
send data

FINbit=1, seq=y

LAST_ACK

TIMED_WAIT

can no longer
send data

ACKbit=1; ACKnum=y+1

timed wait
for 2*max
segment lifetime

CLOSED

CLOSED

# Chapter 3 outline

3.1 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

- segment structure
- reliable data transfer
- flow control
- connection management
- congestion control

京都大学

# Principles of congestion control

*congestion*:

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

- different from flow control!

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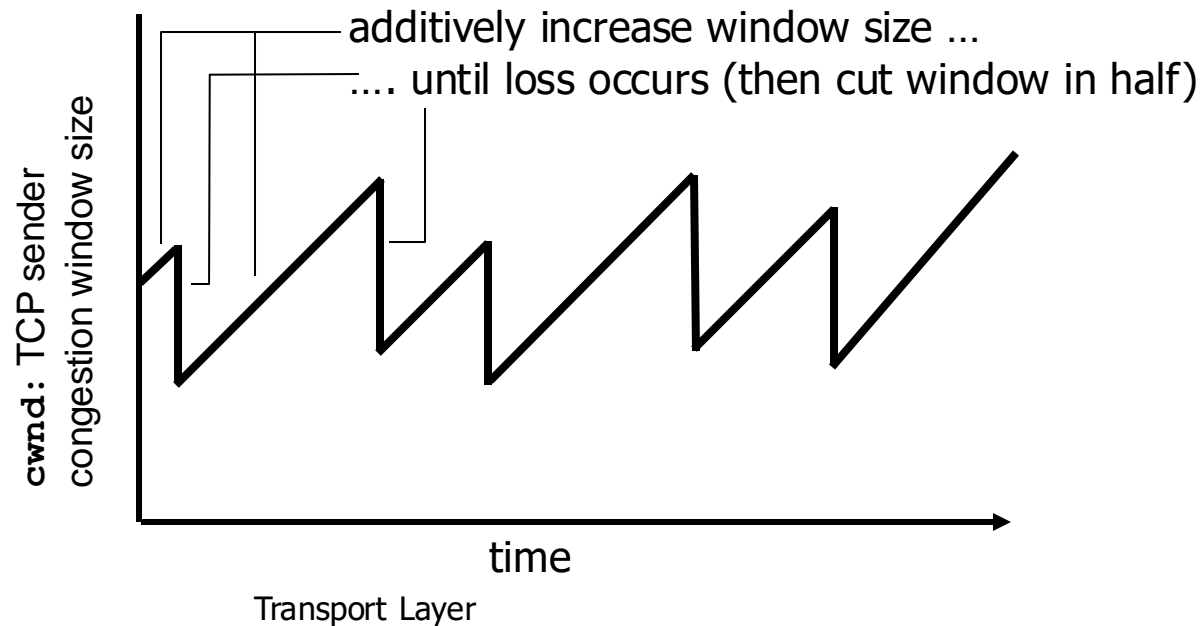  - long delays (queueing in router buffers)

# TCP congestion control

- TCP limits the send rate when the network is congested.
- How does TCP know when the network is congested?
  - If the network is congested, packets get lost.
  - Whenever TCP thinks that a packet got lost, it assumes that the network is congested.
- How does TCP slow down the send rate?
  - By limiting the size of the sender window.
- ACK received -> no congestion, increase window size
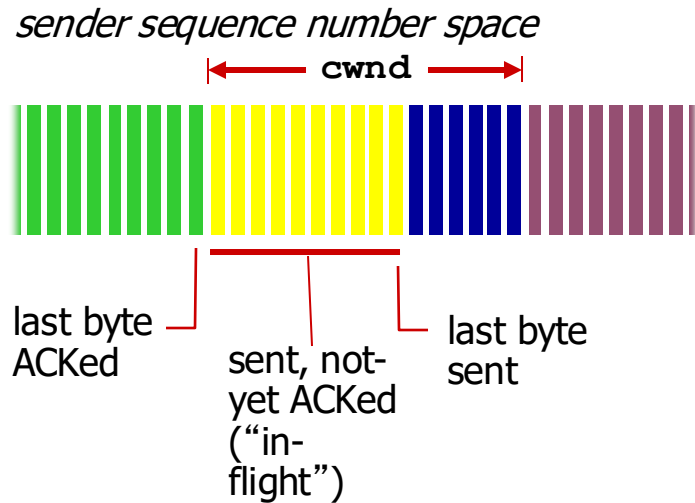- ACK not received -> congestion, decrease window size

# TCP congestion control: additive increase multiplicative decrease

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

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

京都大学

# TCP Congestion Control: details

*sender sequence number space*



last byte
ACKed

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

last byte
sent

- sender limits transmission:

$$\text{LastByteSent-LastByteAcked} \leq \text{cwnd}$$

- **cwnd** is dynamic, function of perceived network congestion

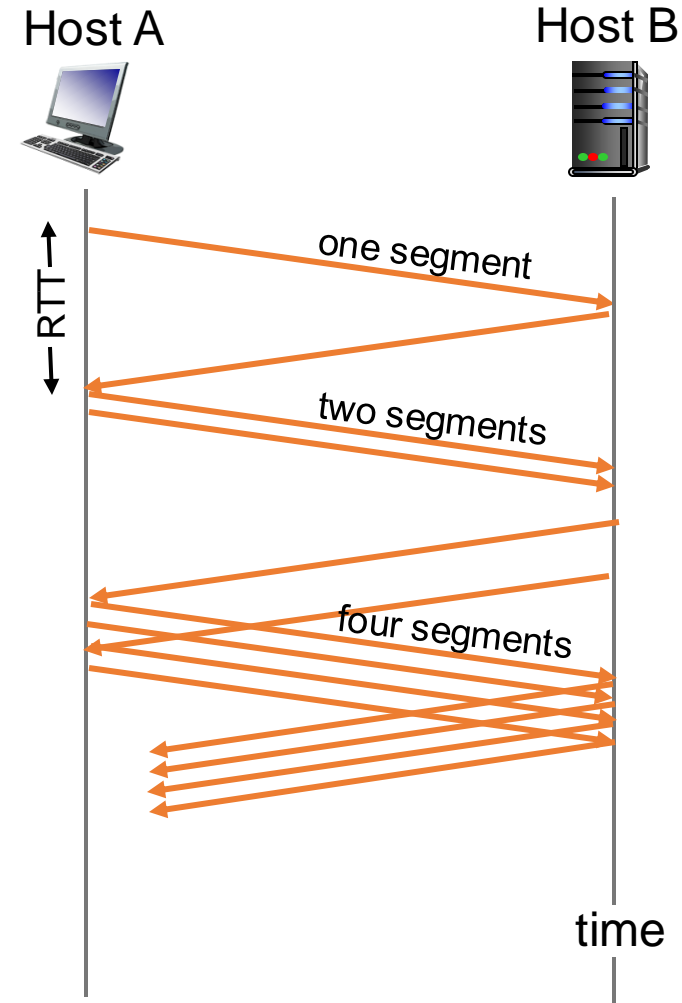- The actual window size is limited by min(rwnd, cwnd)

*TCP sending rate:*

- *roughly:* send cwnd bytes, wait RTT for ACKS, then send more bytes

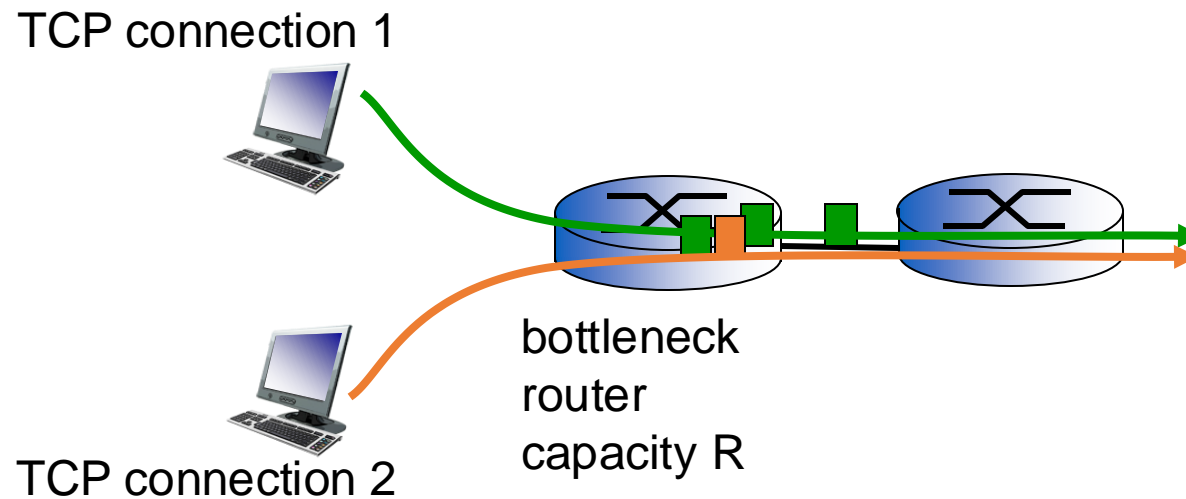$$\text{rate} \approx \frac{\text{cwnd}}{\text{RTT}} \text{ bytes/sec}$$

# TCP Slow Start

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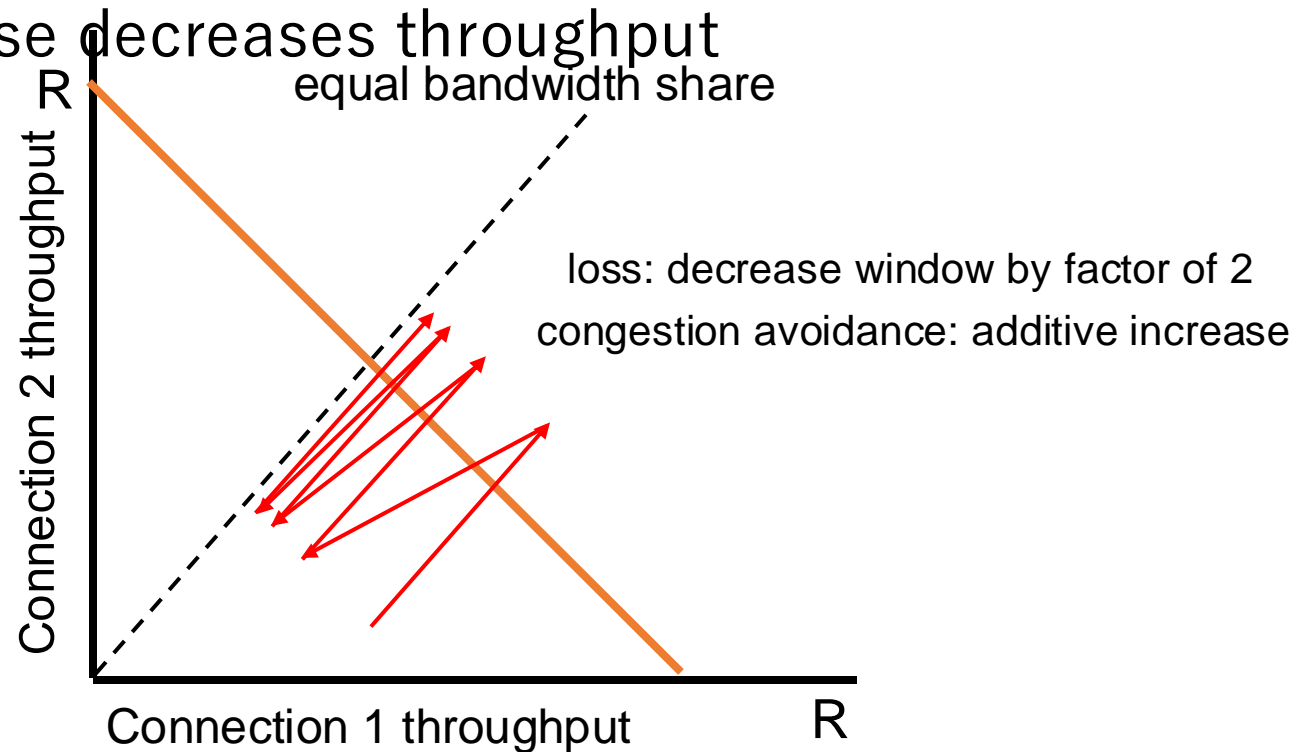- *summary:* initial rate is slow but ramps up exponentially fast

Host A                    Host B

RTT

one segment

two segments

four segments

time

京都大学

3-60

# 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

# Why is TCP fair?

two competing sessions:

- additive increase gives slope of 1, as throughout increases

- multiplicative decrease decreases throughput proportionally



equal bandwidth share

loss: decrease window by factor of 2

congestion avoidance: additive increase

Connection 2 throughput

Connection 1 throughput

# 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

京都大学