# Information Network

Lecture 9: TCP

Holger Thies

京都大学

KYOTO UNIVERSITY

# Reminder

- Class grade: Homework and Short Test (40%), Final report (60%)

- In-class short test in lecture 14 (1/22).

- The test must be taken in class.

- If you can not attend the class (for a good reason) then you must contact the lecturer <u>before the test</u> to schedule a replacement test.
    - Email: thies.holger.5c@kyoto-u.ac.jp

- If you do not attend class 14 without priorly contacting the lecturer, the test will be counted as 0 points.

- Final report needs to be submitted until 1/24.

# Final report

The final report of this lecture must be submitted on PandA until the deadline 1/24 **in word or pdf format.**

Your report should be in a single document and address **both** of the two topics below:

1. Choose either the transport layer or the network layer of the Internet Protocol Stack
   - Summarize the main purpose of the layer and explain what services it can provide.
   - Describe some of the main protocols that operate on this layer in detail
   - Explain how the layer interacts with and uses protocols and services of the upper and lower layers.
2. Describe a network application or protocol (on any layer) **that was not covered in class.**
   - Summarize the protocol's purpose, functionality, and its role within the Internet Protocol Stack.
   - Make some suggestions on improving the selected network protocol's efficiency, security, or scalability.

京都大学

# Final report

- Your report should be 5-7 pages long in total (including headings, references, figures, etc.) and has to be written in English.

- The font size of the main text should not be larger than 12pt.

- Otherwise, the format is free,  as long as the report is written clearly and demonstrates understanding of the topic.

- Each student must write their own report.

- All external references must be cited properly.

- The deadline for submissions is January 24 (Friday).

- Late submissions are not allowed. Please submit your report early to avoid potential technical difficulties.

- **Any form of plagiarism, <u>including the use of generative AI (ChatGPT etc.)</u> is strictly prohibited.**

- **In case plagiarism is discovered there can be severe consequences such as automatically failing all courses in the academic year.**

  - See also Kyoto University Guideline for Liberal Arts Courses: https://www.z.k.kyoto-u.ac.jp/zenkyo/lesson

京都大学

# Chapter 3 outline

京都大学

# 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.
```
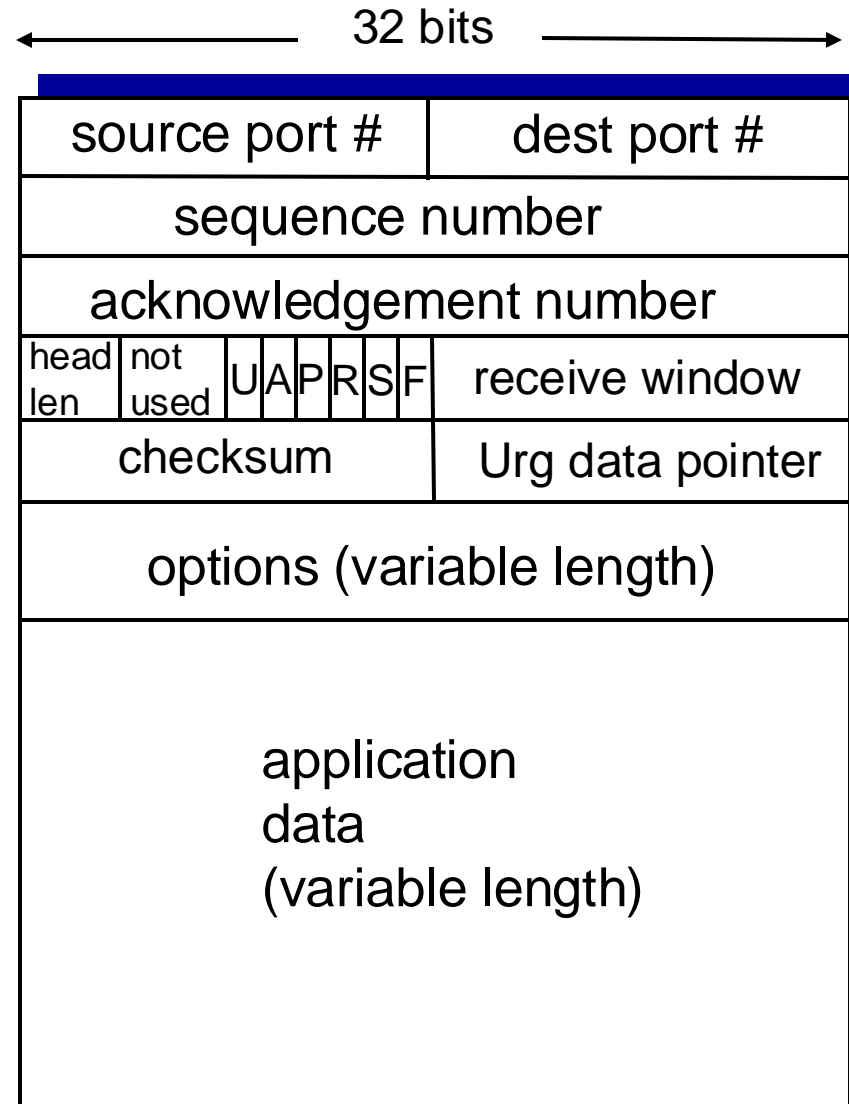
京都大学

# 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
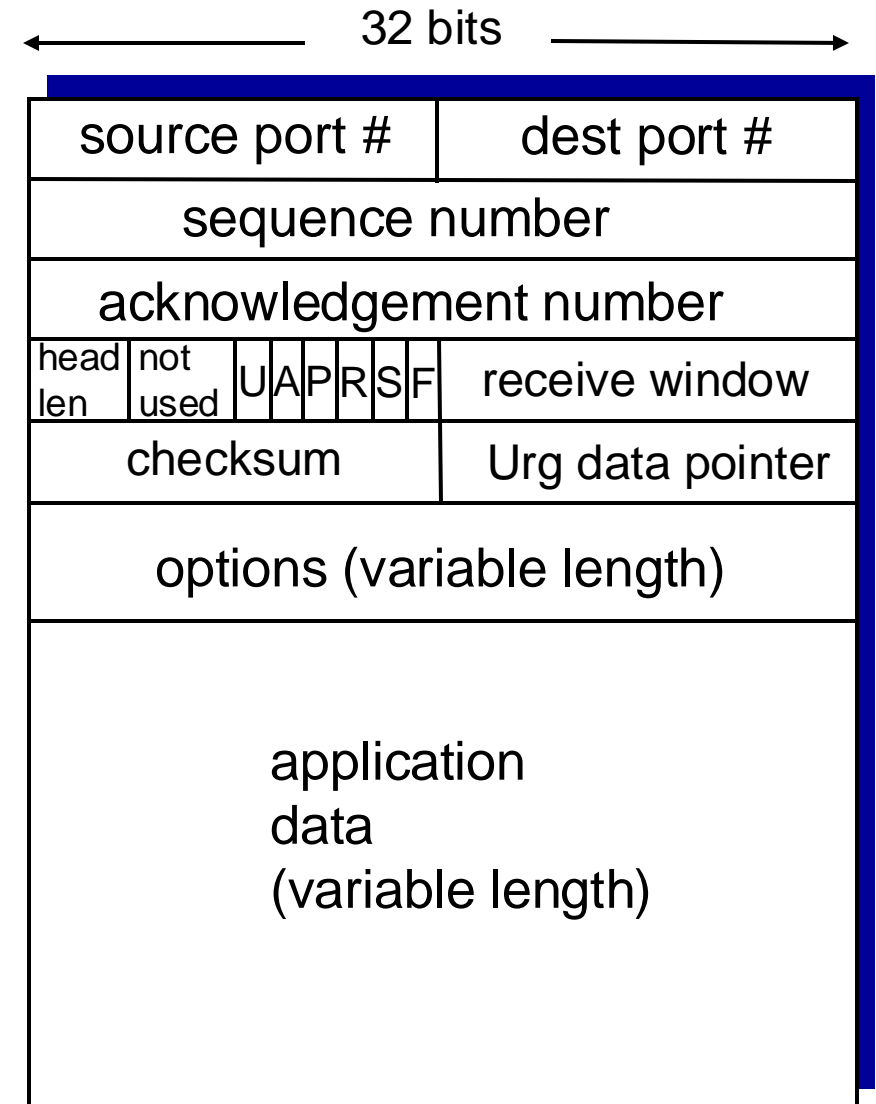
- flow controlled:
  - sender will not overwhelm receiver
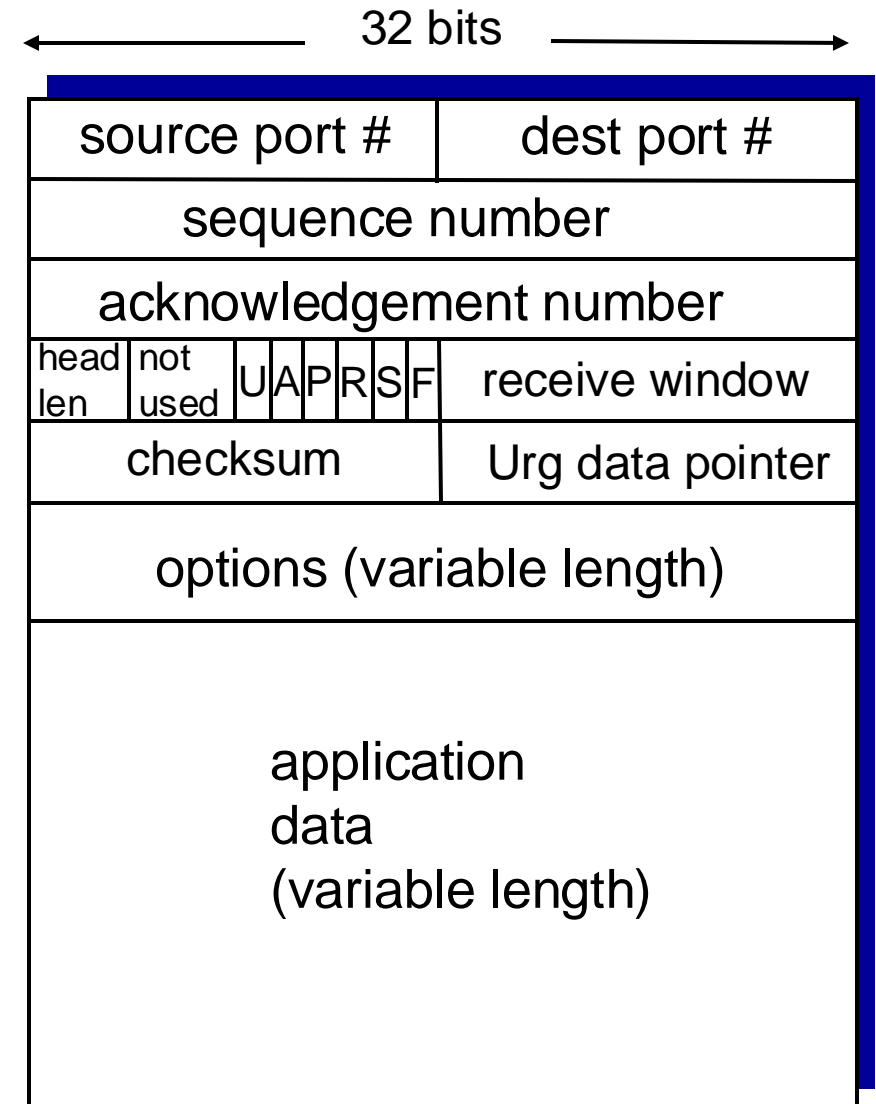
京都大学

# 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<br>data<br>(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 |
|:---:|:---:|

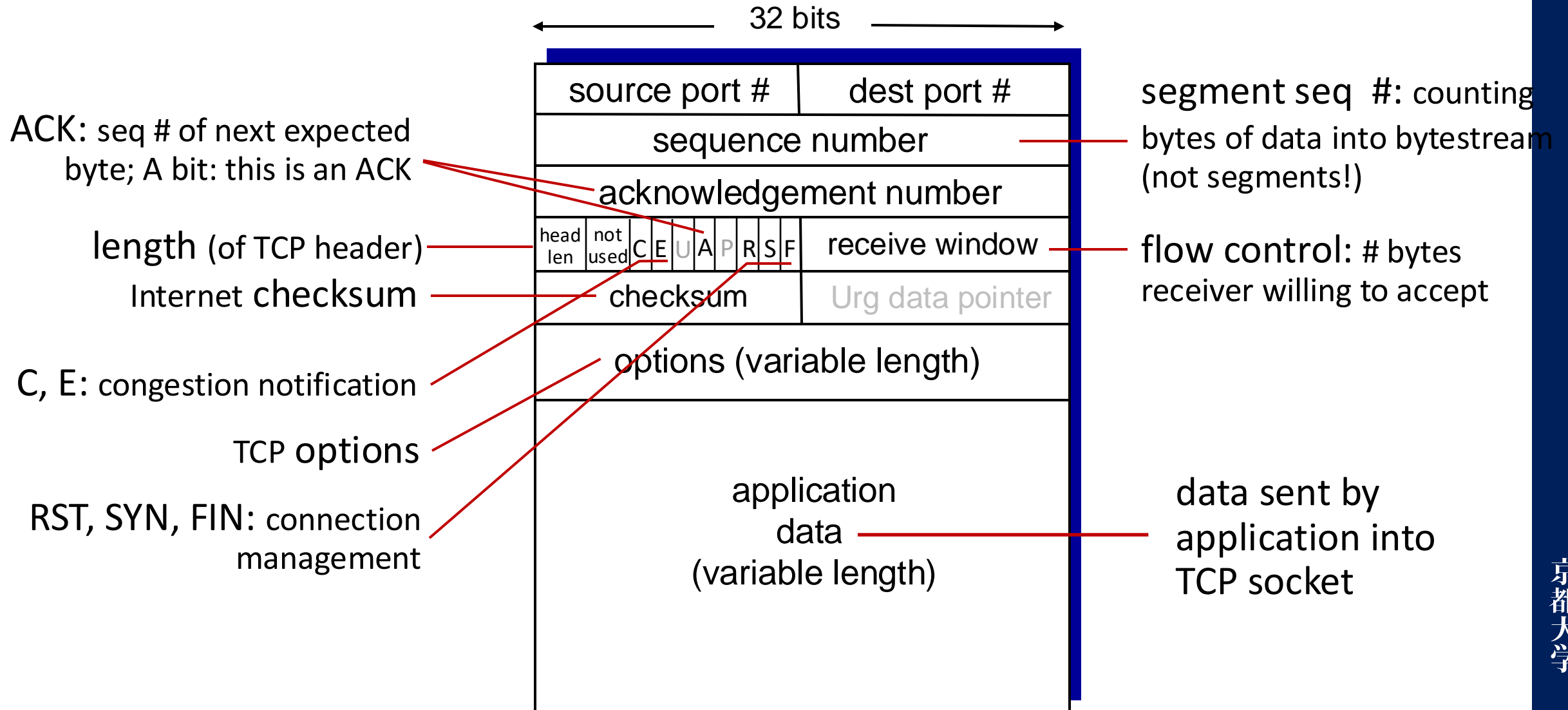options (variable length)

application
data
(variable length)

# 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

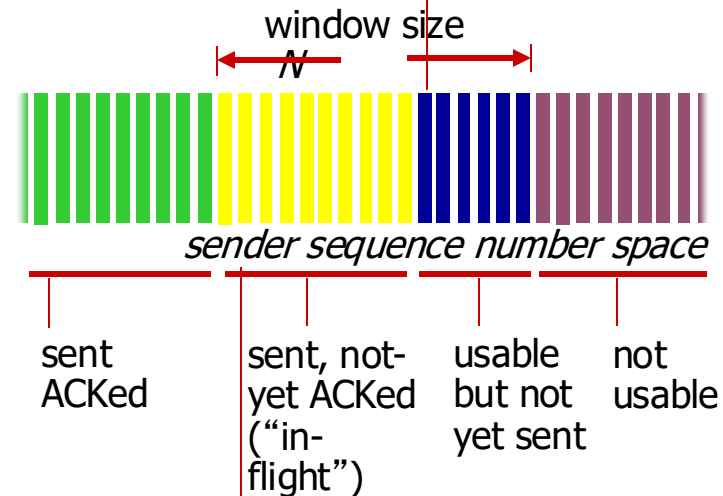# TCP segment structure

# 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
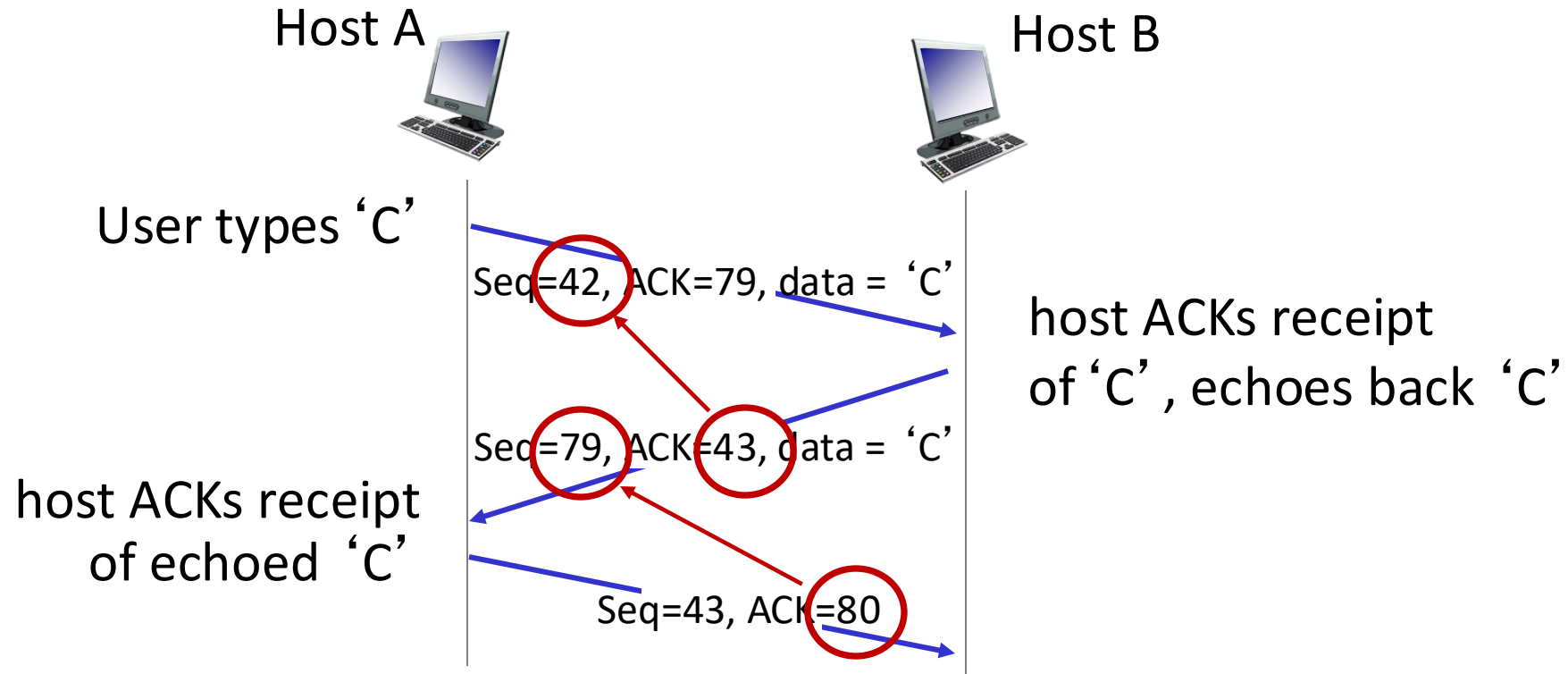
outgoing segment from sender

| 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 |

# 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

# Another example

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
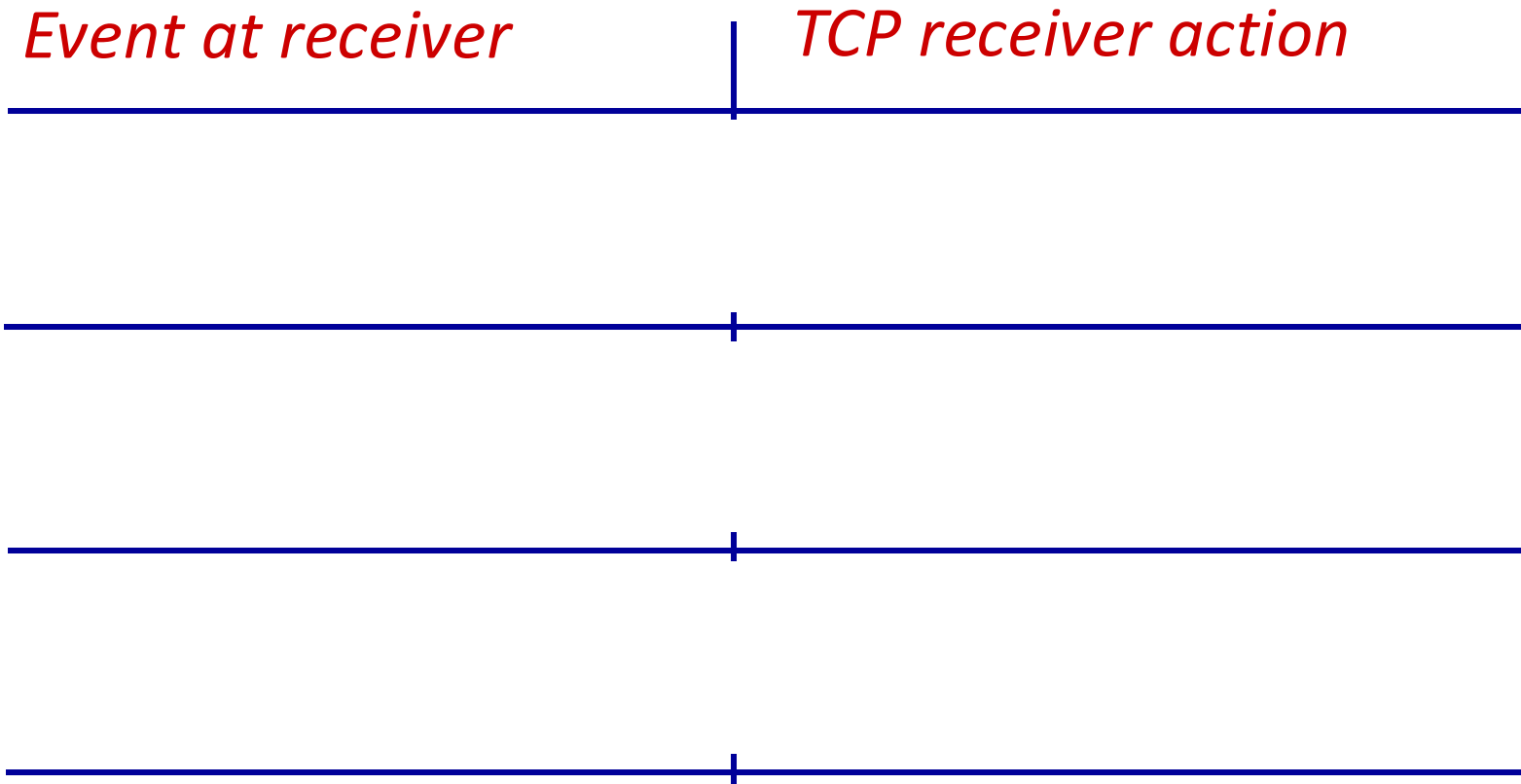  - Timer is for oldest unacknowledged segment

*event: timeout*

- retransmit segment that caused timeout
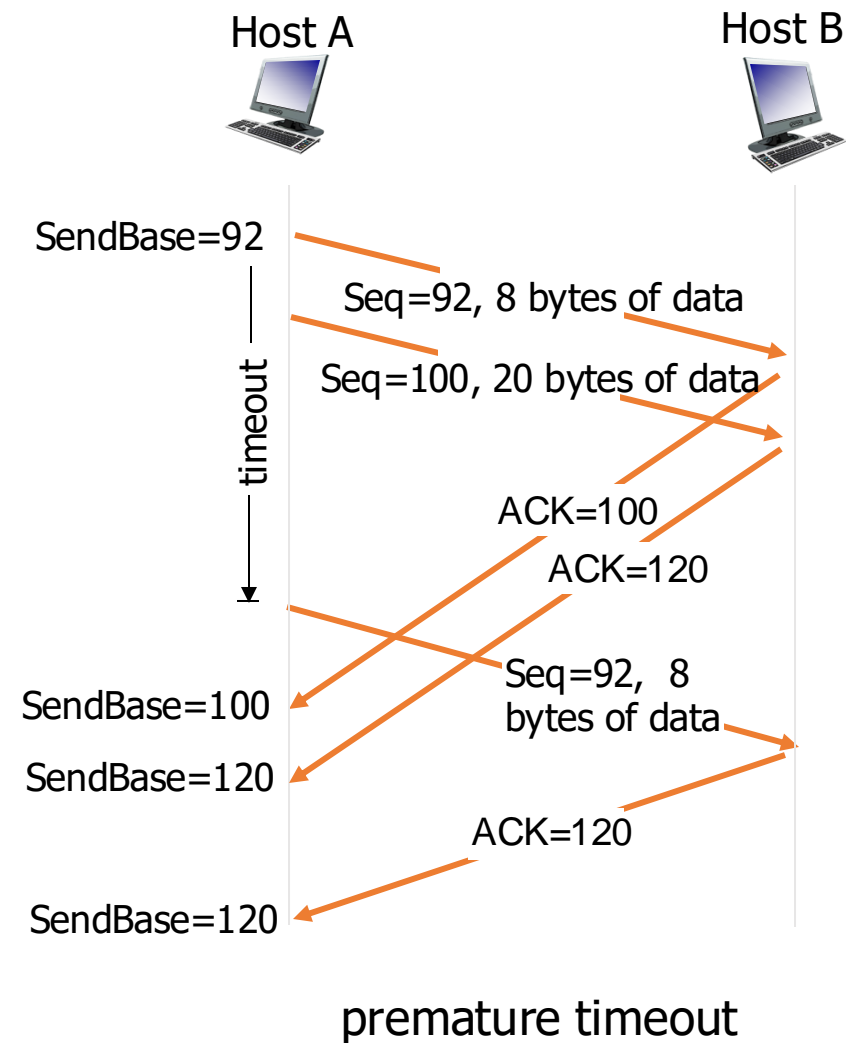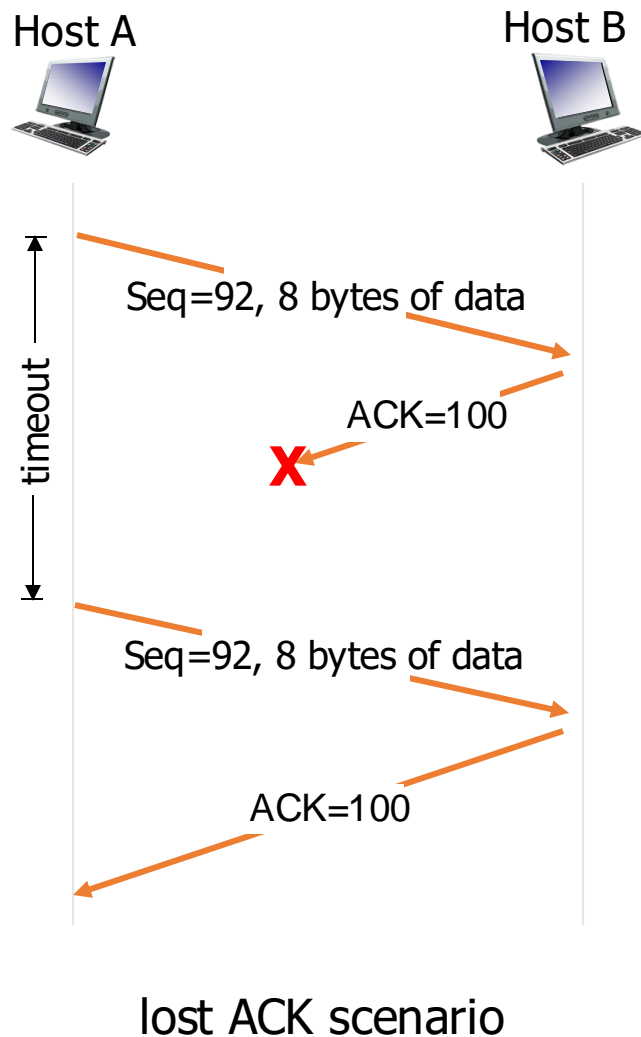- restart timer

*event: ACK received*

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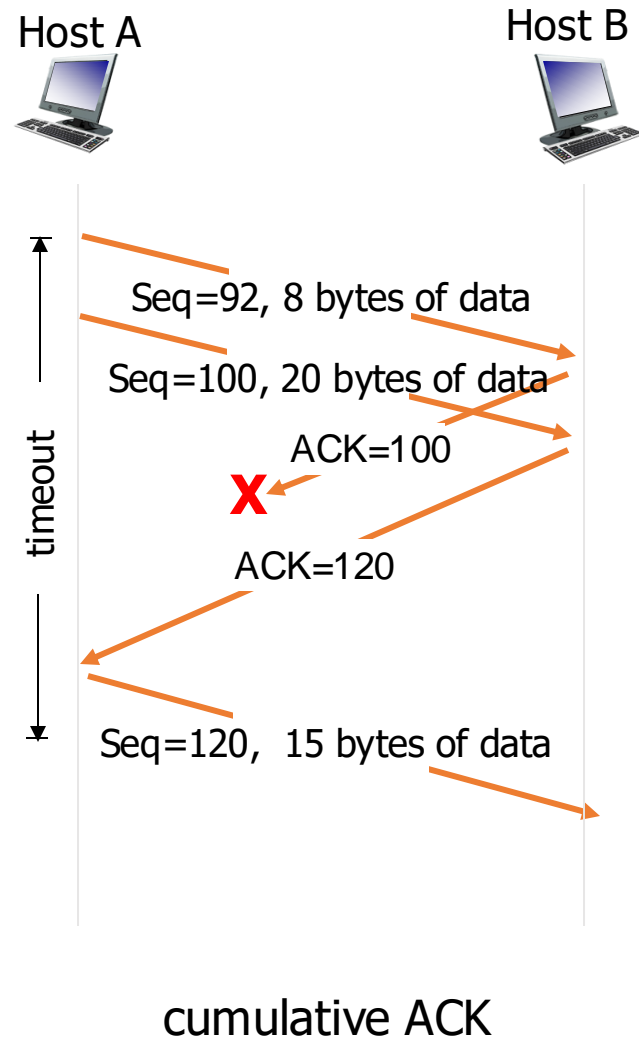  - start timer if there are still unACKed segments

京都大学

# TCP Receiver: ACK generation [RFC 5681]

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

京都大学

# TCP: retransmission scenarios



Host A          Host B

SendBase=92

Seq=92, 8 bytes of data

timeout

Seq=92, 8 bytes of data

ACK=100

**X**

Seq=92, 8 bytes of data

ACK=100

lost ACK scenario

Host A          Host B

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

京都大学

3-20

# TCP: retransmission scenarios

Host A                    Host B

Seq=92, 8 bytes of data

Seq=100, 20 bytes of data

ACK=100

X

ACK=120

Seq=120,  15 bytes of data

timeout

cumulative ACK

京都大学
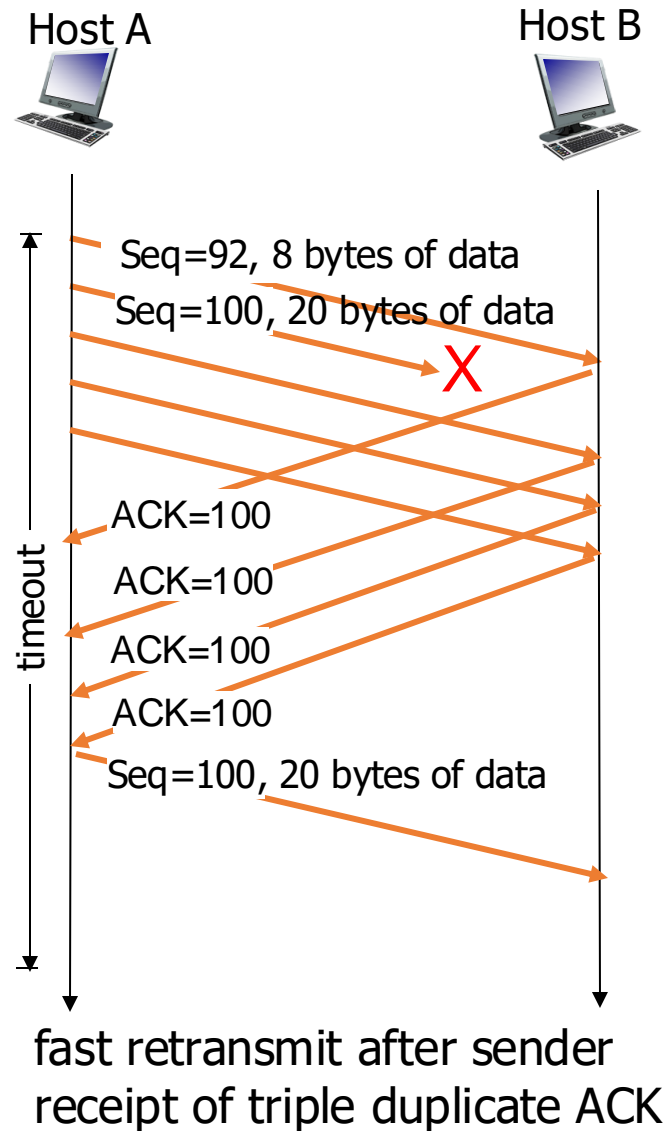
# TCP Example

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

# TCP fast retransmit

Host A                    Host B

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

ACK=100

timeout

ACK=100

ACK=100

ACK=100

Seq=100, 20 bytes of data

fast retransmit after sender
receipt of triple duplicate ACK

京都大学

# Chapter 3 outline
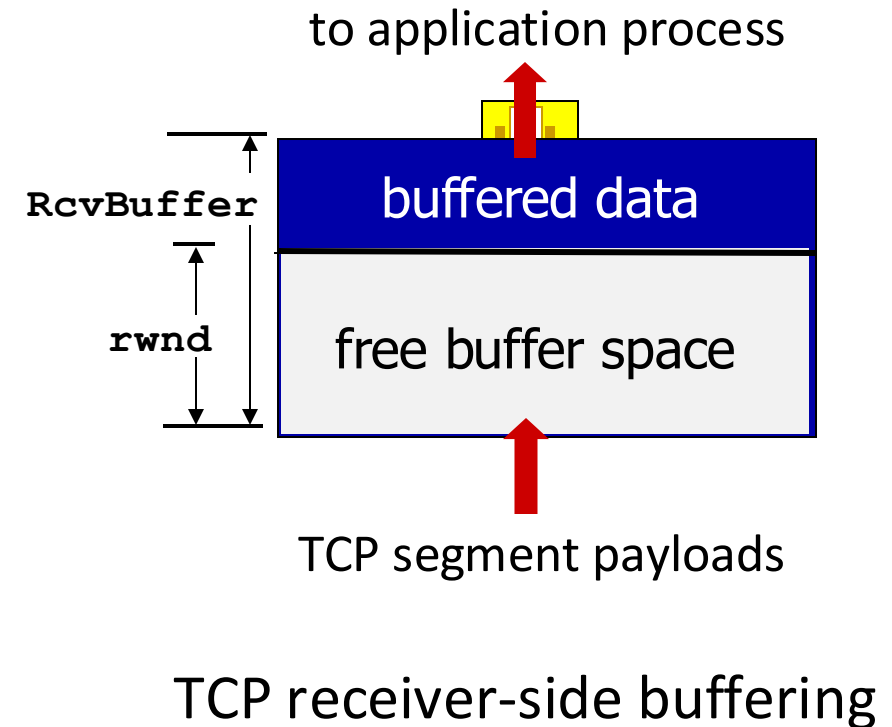
京都大学

# 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
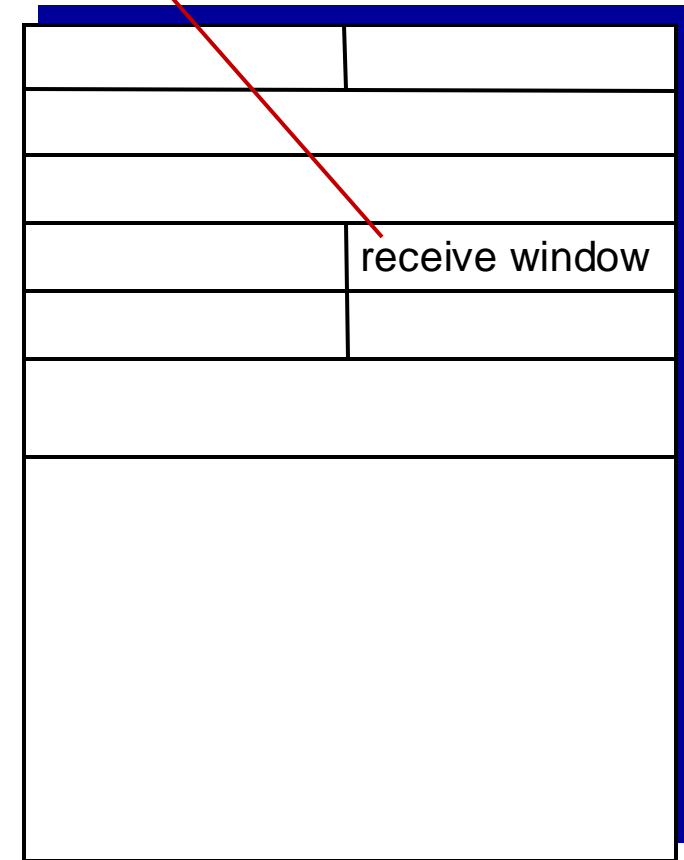
京都大学

# 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



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

3.1 transport-layer services

3.2 multiplexing and demultiplexing

3.3 connectionless transport: UDP
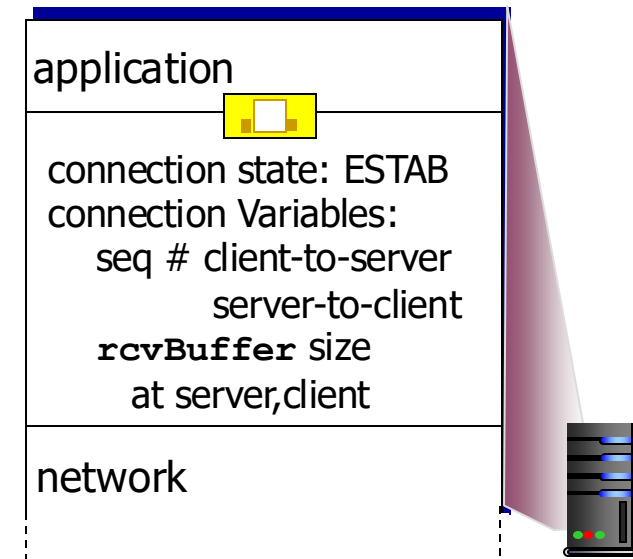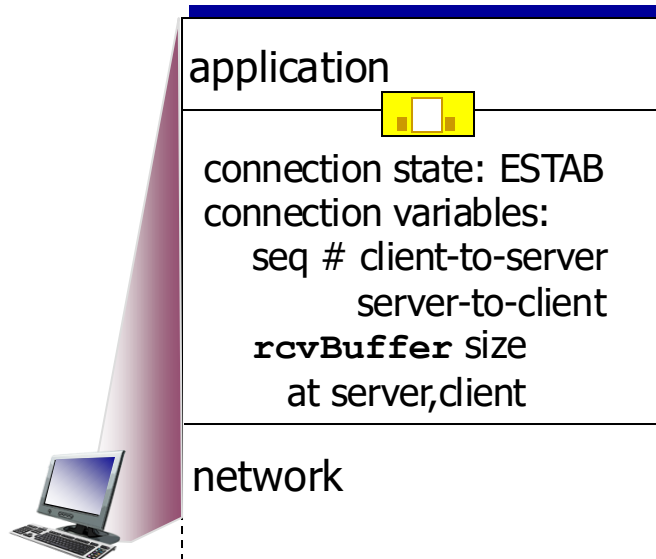
3.4 principles of reliable data transfer

3.5 connection-oriented transport: TCP
- segment structure
- reliable data transfer
- flow control
- connection management
- congestion control
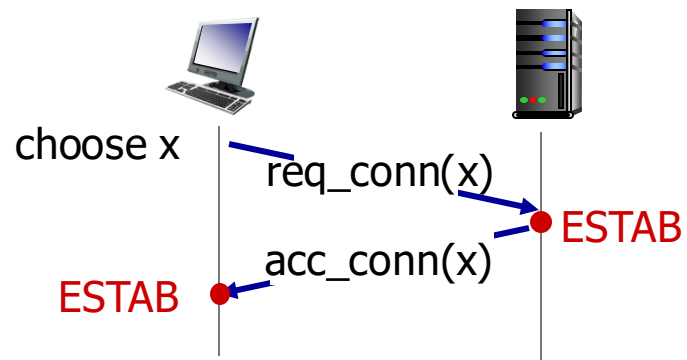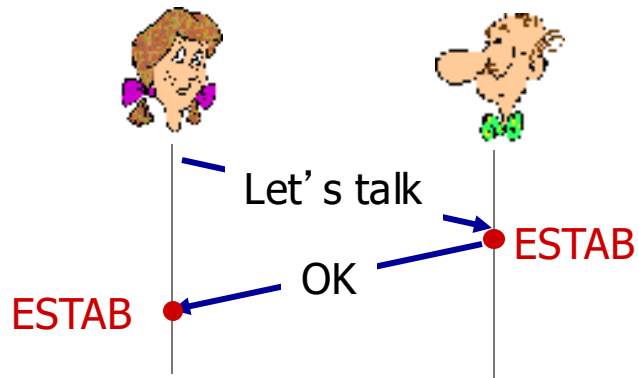
# Connection Management

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
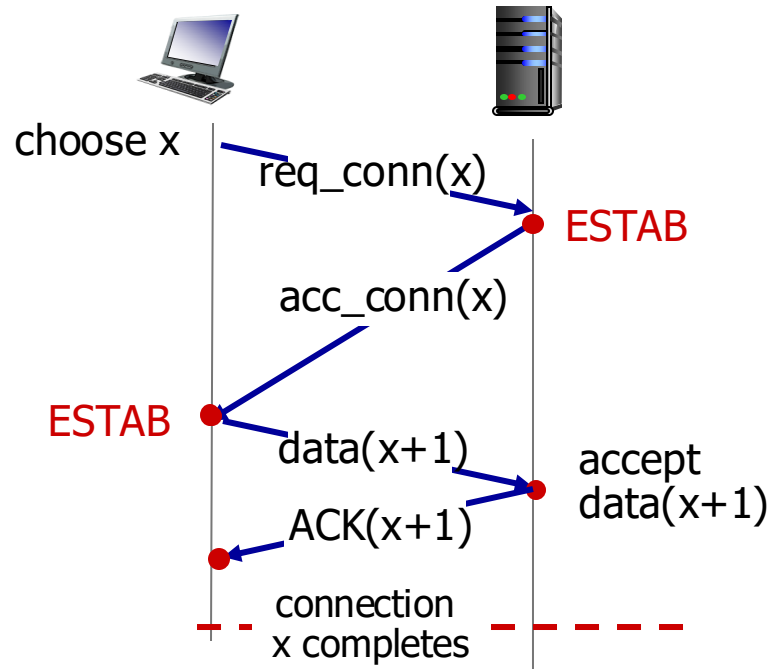
## 2-way handshake:



*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!

京都大学

# 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

client
terminates

server
forgets x

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
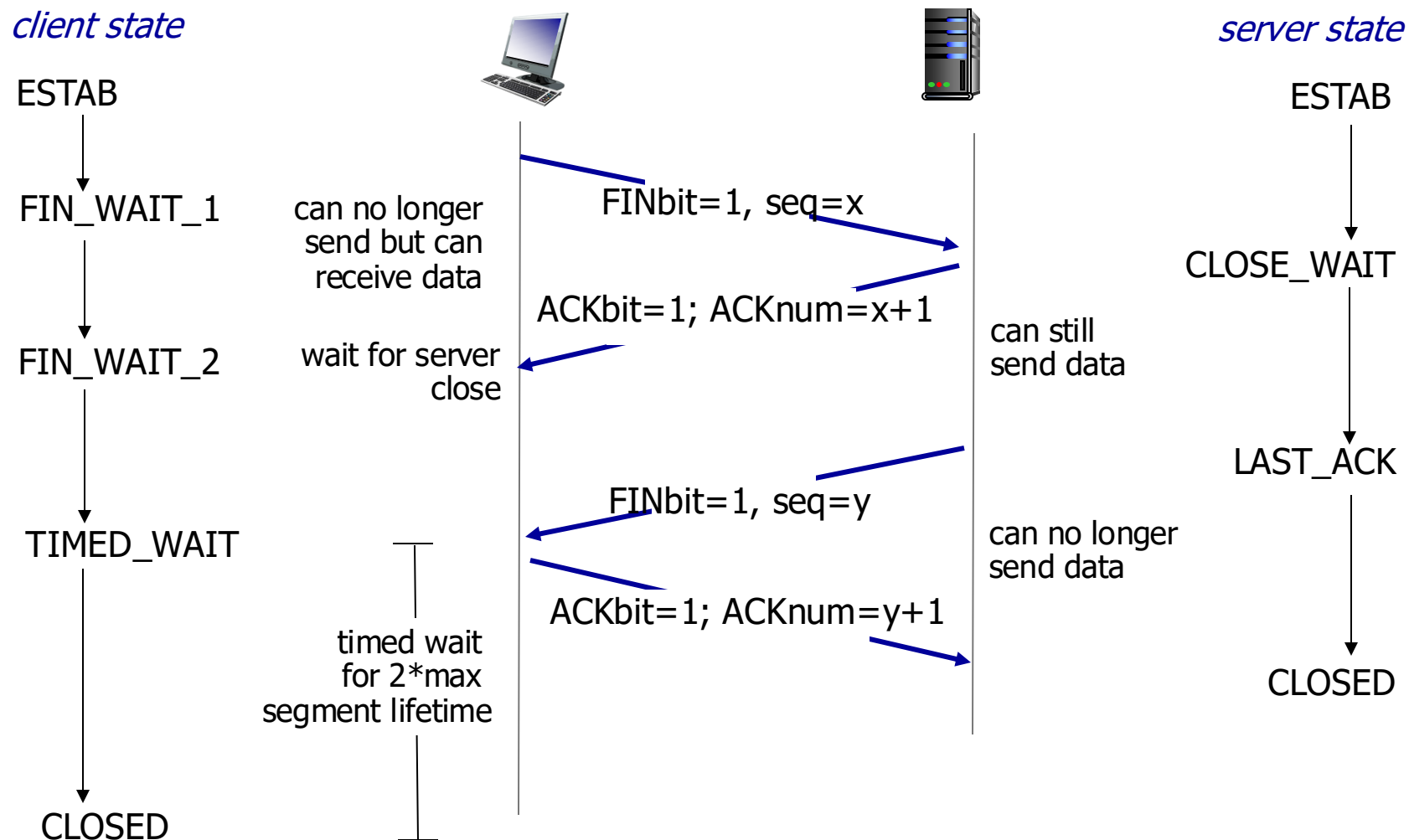
# TCP: closing a connection

- 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*

ESTAB

FIN_WAIT_1

can no longer
send but can
receive data

FIN_WAIT_2

wait for server
close

TIMED_WAIT

timed wait
for 2*max
segment lifetime

CLOSED

FINbit=1, seq=x

ACKbit=1; ACKnum=x+1

FINbit=1, seq=y

ACKbit=1; ACKnum=y+1

*server state*

ESTAB

CLOSE_WAIT

can still
send data

LAST_ACK

can no longer
send data

CLOSED

京都大学

# Chapter 3 outline

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)
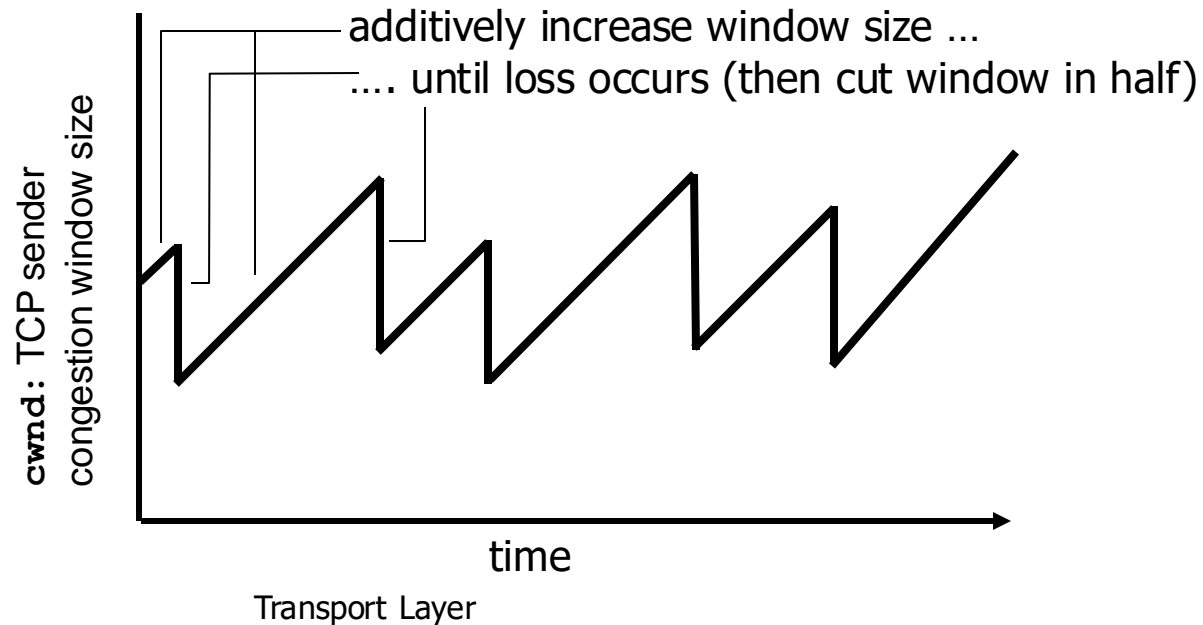  - long delays (queueing in router buffers)

# 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
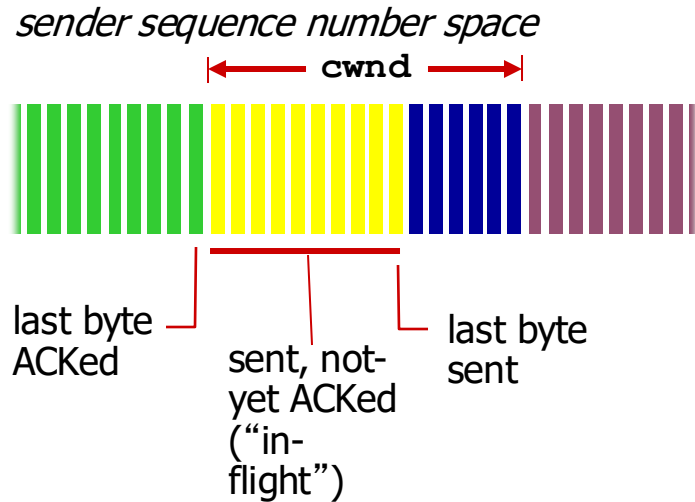- ACK not received -> congestion, decrease window size

京都大学

# 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

Transport Layer

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

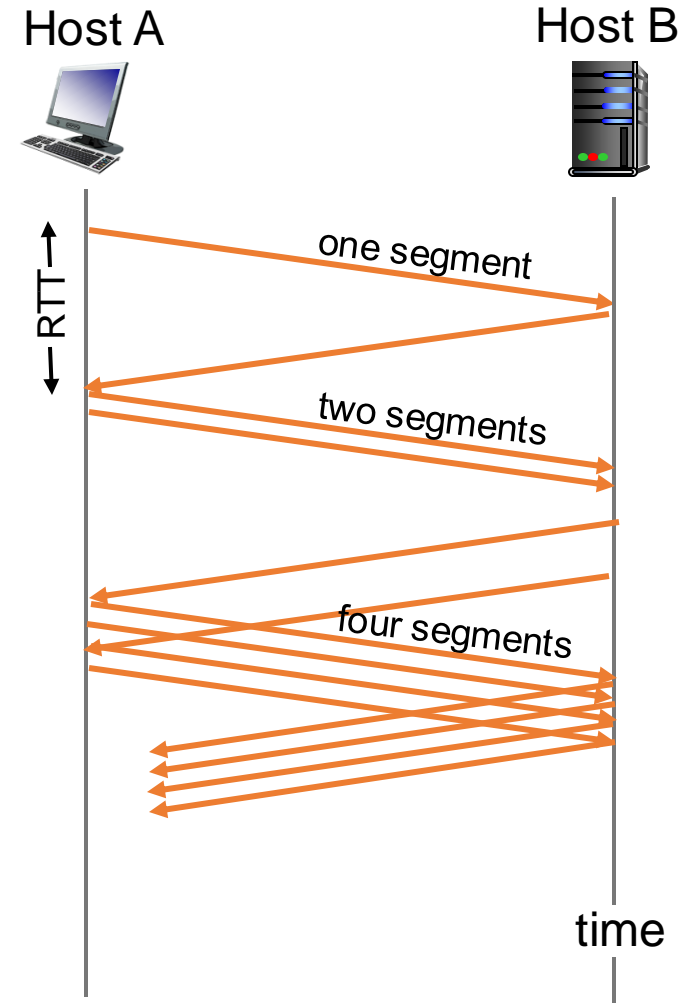- The actual window size is limited by min(rwnd, cwnd)

*TCP sending rate:*

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

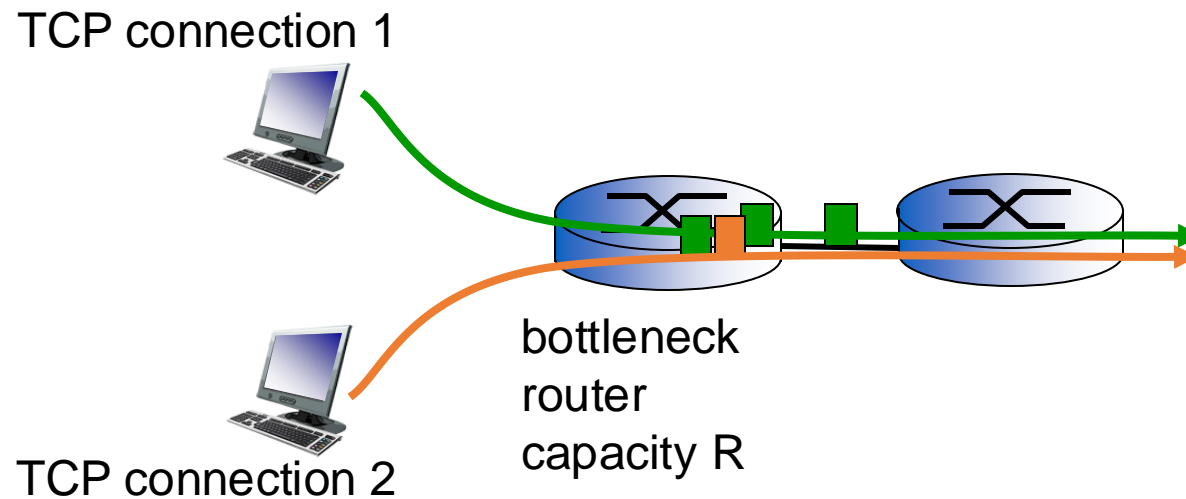$$\text{rate} \approx \frac{\text{cwnd}}{\text{RTT}} \text{ bytes/sec}$$

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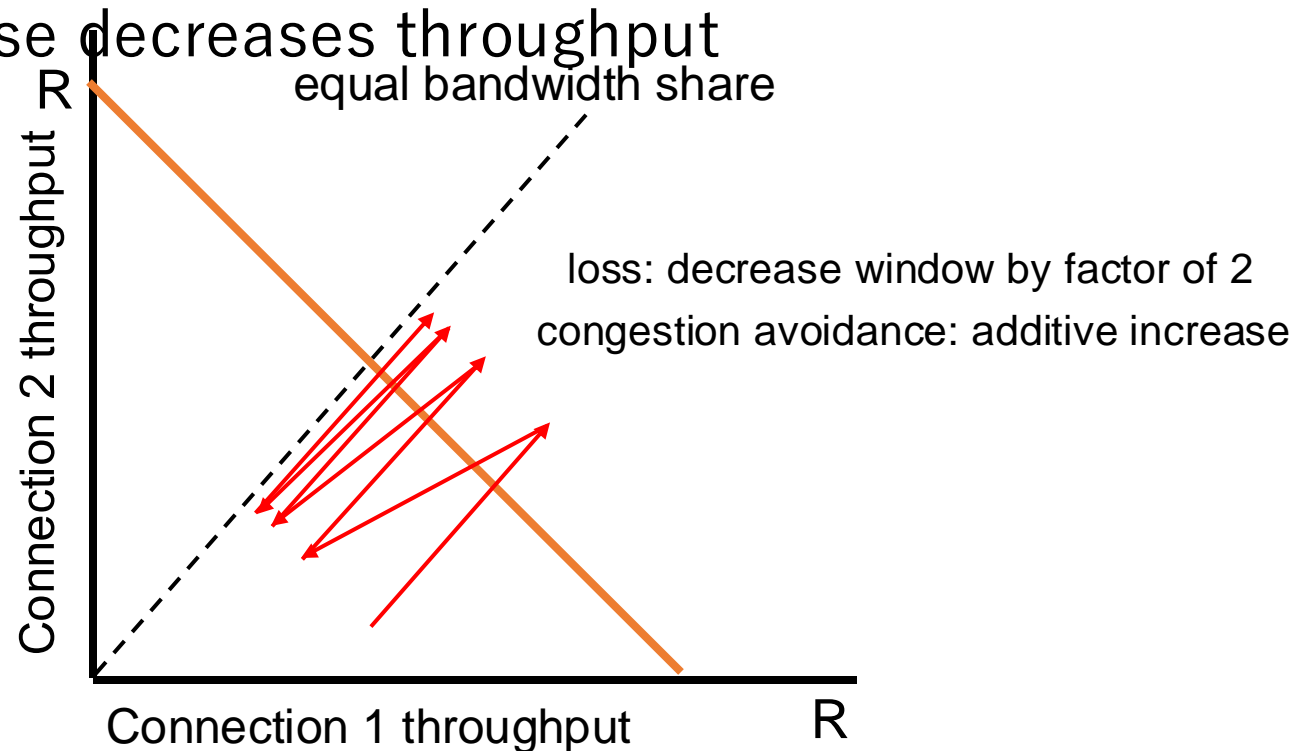- *summary:* initial rate is slow but ramps up exponentially fast

Host A                    Host B

RTT

one segment

two segments

four segments

time

京都大学

# 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

京都大学

# 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

R

R

京都大学

# 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

# Transport layer: summary

- principles behind transport layer services:
    - multiplexing, demultiplexing
    - 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"

京都大学