# Lecture 15: Transport Layer Congestion Control

COMP 411, Fall 2022

Victoria Manfredi

WESLEYAN UNIVERSITY

# Today

1. Flow control

2. Congestion causes and costs

3. TCP congestion control

# TCP
# FLOW CONTROL

# What if sender overwhelms receiver?

**Problem**

Application may remove data from TCP socket buffers ….

… slower than TCP receiver is delivering (sender is sending)

Receiver protocol stack

Application process

App

OS

TCP socket receiver buffers

TCP code

IP code

from sender

5

# TCP flow control

## Receiver provides feedback to sender

- so sender doesn't overflow receiver's buffer
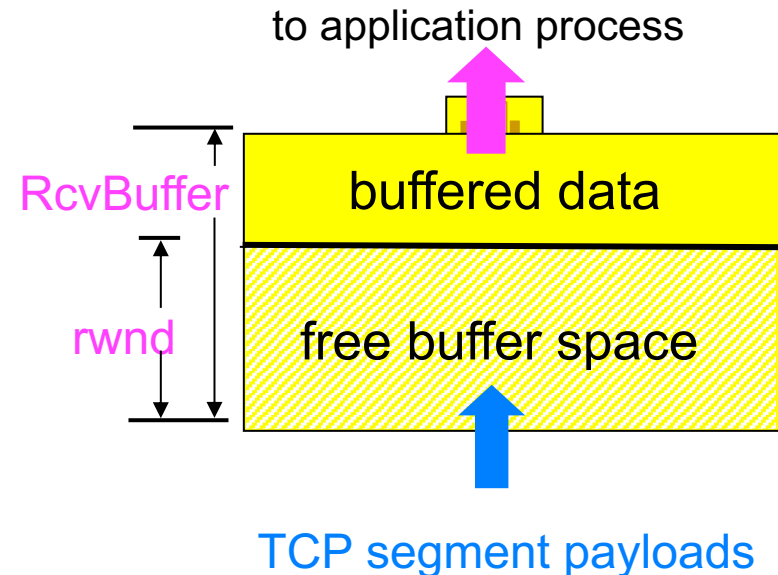- sender and receiver each maintain window

## Receiver

- rwnd: free space in RcvBuffer
- puts rwnd in TCP header of receiver-to-sender segments

## Sender

- limits unacked data to rwnd
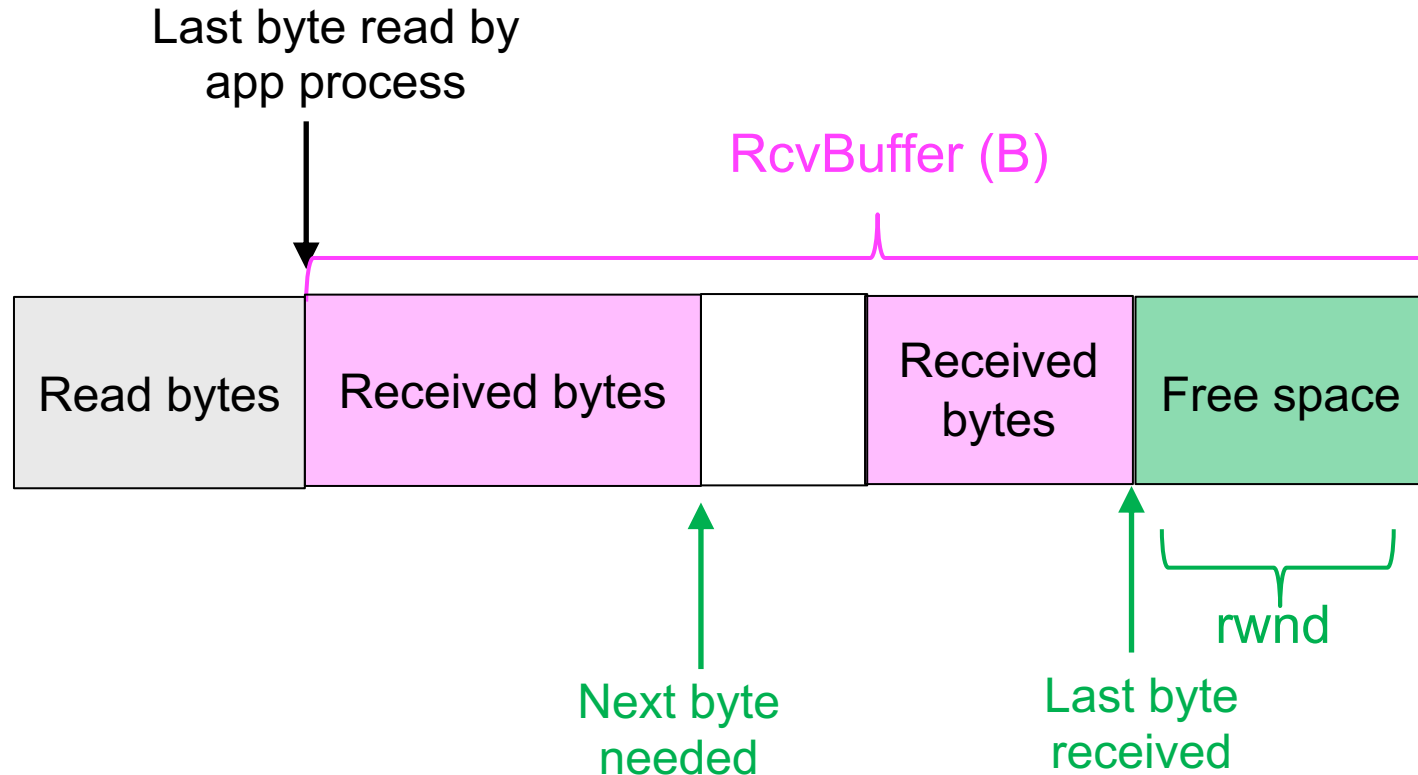- ensures RcvBuffer will not overflow

## Receiver-side buffering



to application process

RcvBuffer

buffered data

rwnd

free buffer space

TCP segment payloads

# Receive window (rwnd)

# Receiver use of receive window (rwnd)

Keeps track of available space in its RcvBuffer

Last byte read by
app process

RcvBuffer (B)

| Read bytes | Received bytes | | Received bytes | Free space |

Next byte
needed

Last byte
received

rwnd

rwnd = B − (last byte received − last byte read)

8

# Sender use of receive window (rwnd)

Limits # of in-flight segments of sender

SendBuffer (B)

| ACK'd bytes | Sent bytes | | No data |
|---|---|---|---|

rwnd

1st unACK'd byte

Last byte can send (= last byte written by app)

Sending rate limited to: rwnd bytes/RTT seconds

# Sender use of receive window (rwnd)

Problem: if rwnd = 0, what happens?



SendBuffer (B)

| ACK'd bytes | | | No data |

rwnd

1st unACK'd byte

Last byte can send (= last byte written by app)

No ACKs sent: receiver has no way to let sender know rwnd increased

Solution: send segments with 1 byte of data, which receiver ACKs

# Congestion
# CAUSES AND COSTS

# What if sender overwhelms network?

Receive buffer is not only resource limitation
- every packet travels through path of routers
- routers may be congested, have long queues …

Causes of network congestion
- many senders compete for network resources
- senders lack knowledge
  - amount of resources available (bandwidth)
  - # of other senders competing

# Costs of network congestion

As queues in bottleneck link fill up: large packet delays, dropped packets

Bad feedback loop!

As timeouts expire at sender due to delays/drops: packets retransmitted

Problem
- retransmission treats symptoms but not underlying problem

Q: how to solve underlying problem of congestion?
- reduce sending rate … but what should sending rate be?
  - depends on available bandwidth
  - sender increases/decreases sending rate based on congestion level

# Recall link and network resources are shared

1. Hosts: divide data to send into fixed-length packets

Host 1

2. Routers: interleave packets from different hosts on links

Host 2

www.google.com

# Scenario 1

No retransmission, 2 senders, 2 receivers

Original data: $\lambda_{in}$  **Host A**

Throughput: $\lambda_{out}$

**Host B**

Infinite buffers:
unlimited shared
output link buffers

Output link capacity: R

No loss

Max per-connection throughput: R/2

Large delays as arrival rate, $\lambda_{in}$, approaches capacity

Even though high throughput when close to capacity, also high delay!

Q: Why R/2?

15

# Scenario 2: retransmission

Sender retransmits timed-out packet

- $\lambda_{in} = \lambda_{out}$: app-layer input equals app-layer output
- $\lambda'_{in} \geq \lambda_{in}$: transport-layer input includes retransmissions



Original data: $\lambda_{in}$

**Host A**

Retransmitted + original data: $\lambda'_{in}$

**Host B**

Throughput: $\lambda_{out}$

Finite buffers: limited shared output link buffers

Output link capacity: R

Loss

Performance depends on how retransmission performed…

16

# Scenario 2: retransmission + perfect knowledge

## Idealization: perfect knowledge

– sender sends only when router buffers available

– no loss occurs, so $\lambda'_{in} = \lambda_{in}$



λ'in  λin

λout

Copy

Finite buffers

Output link
capacity: R

Free buffer
space

# Scenario 2: retransmission only when lost

## Idealization: known loss

– packets can be lost, dropped at router due to full buffers

– only resend packet known to be lost



when sending at R/2, some packets are retransmissions but asymptotic goodput is still R/2 (why?)

$\lambda'_{in}$

$\lambda_{in}$

Copy

$\lambda_{out}$

Finite buffers

Output link capacity: R

Free buffer space

# Scenario 2: retransmission causing duplicates

## Realistic: duplicates

- packets can be lost, dropped at router due to full buffers

- sender times out prematurely
  - sends 2 copies, both delivered



when sending at R/2, some packets are retransmissions but asymptotic goodput is still R/2 (why?)

$\lambda_{out}$

$\lambda'_{in}$

$\lambda_{in}$

$\lambda'_{in}$

timeout

Copy

$\lambda_{out}$

Finite buffers

Output link capacity: R

Free buffer space

# TCP
# CONGESTION CONTROL

# Goals of TCP congestion control

1.  ## Discover available bandwidth
    –  how much bandwidth can be used without causing congestion
        •  will vary over time
    –  estimate starting from no information

2.  ## Correctly set sending rate
    –  should not exceed available bandwidth

3.  ## Fairness
    –  no user gets all of the bandwidth

# TCP Congestion Control

## Sender limits transmission

`LastByteSent-LastByteAcked ≤` **`cwnd`**

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

*sender sequence number space*



last byte ACKed

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

last byte sent

## TCP sending rate

– roughly
  - send cwnd bytes
  - wait RTT for ACKs
  - send more bytes

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

## Q: How does sender estimate cwnd?

# To estimate cwnd

Detect congestion
- **delays**
  - large RTTs: too variable to be used in practice

- **duplicate ACKs**
  - isolated loss

Use to adjust cwnd,
affecting sending rate

- **timer expired**
  - multiple losses

How to intuitively adjust cwnd
- ACK received: increase cwnd
- loss detected: decrease cwnd

# 3 states in TCP finite state machine

Goal: send segments, adjust cwnd as needed

## 1. Slow start
– determine available bandwidth starting from no info

## 2. Congestion avoidance
– deal with fluctuations in bandwidth

## 3. Fast recovery
– quickly recover from isolated lost packets

We'll first look at different states, then full FSM

# Slow start: initialization

## Initial rate is "slow"

- relative to original TCP which had no congestion control
- initially cwnd = 1 MSS

## Ramp up exponentially fast

- every time ACK received
  - cwnd = cwnd + MSS
- essentially doubles cwnd every RTT



**Host A**

**Host B**

RTT

one segment

two segments

four segments

time

# Congestion avoidance

## Additive Increase Multiplicative Decrease (AIMD)

– probe cautiously for usable bandwidth

– additive increase

  • **cautious:** increase cwnd by 1 MSS every RTT until loss detected

– multiplicative decrease

  • **aggressive:** cut cwnd in half after loss

additively increase window size …

… until loss occurs, then cut window in half

cwnd

**AIMD saw tooth behavior:** probing for bandwidth

time

# Finite state machine

**Slow start**

$\Lambda$
cwnd = 1 MSS
ssthresh = 64 KB
dupACKcount = 0

dup ACK
dupACKcount++

new ACK
cwnd = cwnd+MSS
dupACKcount = 0
*txmt new segment(s), as allowed*

new ACK
cwnd = cwnd + MSS · (MSS/cwnd)
dupACKcount = 0
*txmt new segment(s), as allowed*

**Congestion avoidance**

cwnd > ssthresh
$\Lambda$

dup ACK
dupACKcount++

timeout
ssthresh = cwnd/2
cwnd = 1 MSS
dupACKcount = 0
*retxmt missing segment*

timeout
ssthresh = cwnd/2
cwnd = 1 MSS
dupACKcount = 0
*retxmt missing segment*

timeout
ssthresh = cwnd/2
cwnd = 1
dupACKcount = 0
*retxmt missing segment*

dupACKcount == 3
ssthresh= cwnd/2
cwnd = ssthresh + 3
*retxmt missing segment*

dupACKcount == 3
ssthresh= cwnd/2
cwnd = ssthresh+3MSS
*retxmt missing segment*

New ACK
cwnd = ssthresh
dupACKcount = 0

**Fast recovery**

dup ACK
cwnd = cwnd + MSS
*txmt new segment(s), as allowed*

27