# Chapter 3
# Transport Layer

The course notes are adapted for Bucknell's CSCI 363
Xiannong Meng
Spring 2016

*Computer
Networking: A Top
Down Approach*
6th edition
Jim Kurose, Keith Ross
Addison-Wesley
March 2012

Transport Layer 3-1
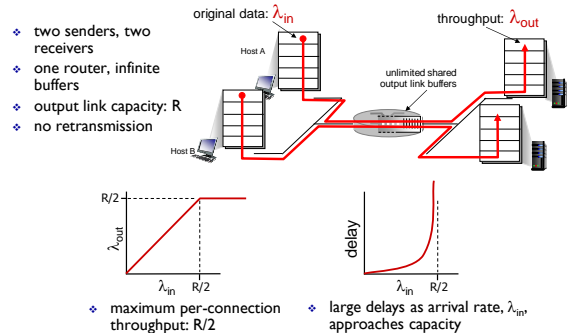
---

# Chapter 3 outline

Transport Layer 3-2

---

# Principles of congestion control

*congestion*:
- informally: "too many sources sending too much
  data too fast for *network* to handle"
- different from flow control!
  - flow control: between hosts
  - congestion control: hosts and network
- manifestations:
  - lost packets (buffer overflow at routers)
  - long delays (queueing in router buffers)
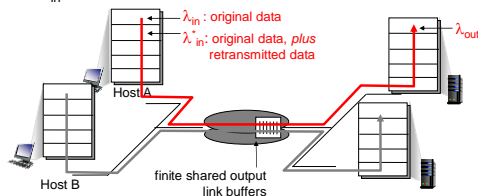- a top-10 problem!

Transport Layer 3-3

---

# Causes/costs of congestion: scenario 1

- two senders, two
  receivers
- one router, infinite
  buffers
- output link capacity: R
- no retransmission



original data: $\lambda_{in}$
throughput: $\lambda_{out}$
Host A
unlimited shared
output link buffers
Host B

- maximum per-connection
  throughput: R/2
- large delays as arrival rate, $\lambda_{in}$,
  approaches capacity

Transport Layer 3-4

---

# Causes/costs of congestion: scenario 2

- one router, *finite* buffers
- sender retransmission of timed-out packet
  - application-layer input = application-layer output: $\lambda_{in}$ =
    $\lambda_{out}$
  - transport-layer input includes *retransmissions* : $\lambda^{*}_{in}$ >=
    $\lambda_{in}$



$\lambda_{in}$ : original data
$\lambda^{*}_{in}$: original data, *plus*
retransmitted data
$\lambda_{out}$
Host A
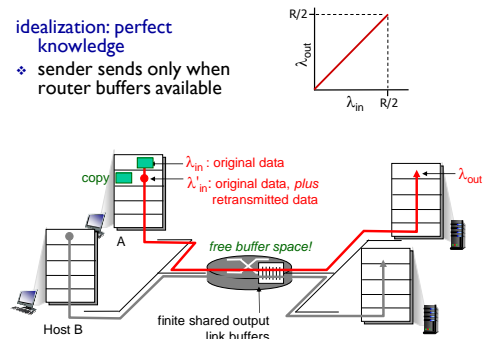Host B
finite shared output
link buffers

Transport Layer 3-5

---

# Causes/costs of congestion: scenario 2

idealization: perfect
knowledge
- sender sends only when
  router buffers available



copy
$\lambda_{in}$ : original data
$\lambda^{'}_{in}$: original data, *plus*
retransmitted data
$\lambda_{out}$
A
*free buffer space!*
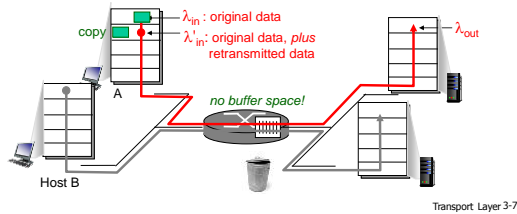Host B
finite shared output
link buffers

Transport Layer 3-6

1

## Causes/costs of congestion: scenario 2

*Idealization: known loss*

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

❖ sender only resends if packet *known* to be lost



λ_in : original data
λ'_in : original data, *plus* retransmitted data
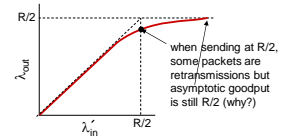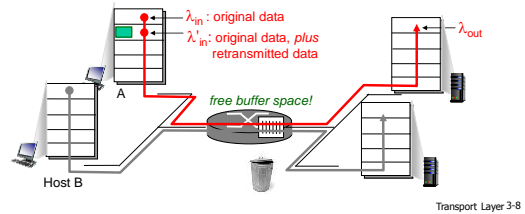
copy

A

no buffer space!

λ_out

Host B

---

## Causes/costs of congestion: scenario 2

*Idealization: known loss*

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

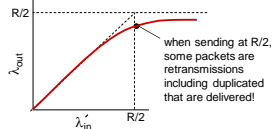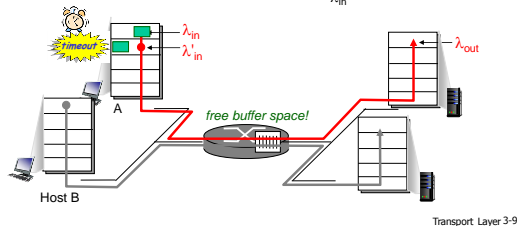❖ sender only resends if packet *known* to be lost



R/2

λ_out

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

λ'_in     R/2

λ_in : original data
λ'_in : original data, *plus* retransmitted data

A

free buffer space!

λ_out

Host B

---

## Causes/costs of congestion: scenario 2

*Realistic: duplicates*

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

❖ sender times out prematurely, sending *two* copies, both of which are delivered



R/2

λ_out

when sending at R/2, some packets are retransmissions including duplicated that are delivered!

λ'_in     R/2

timeout

λ_in
λ'_in

A

free buffer space!

λ_out

Host B

---

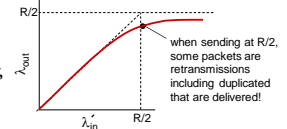## Causes/costs of congestion: scenario 2

*Realistic: duplicates*

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

❖ sender times out prematurely, sending *two* copies, both of which are delivered



R/2

λ_out

when sending at R/2, some packets are retransmissions including duplicated that are delivered!

λ'_in     R/2

**"costs" of congestion:**
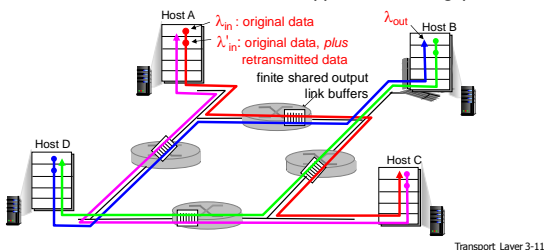
❖ more work (retrans) for given "goodput"

❖ unneeded retransmissions: link carries multiple copies of pkt
  ▪ decreasing goodput

---

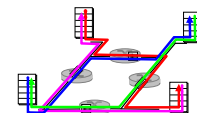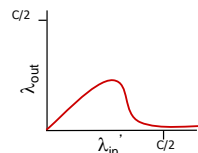## Causes/costs of congestion: scenario 3

❖ four senders
❖ multihop paths
❖ timeout/retransmit

Q: what happens as λ_in and λ'_in increase ?

A: as red λ'_in increases, all arriving blue pkts at upper queue are dropped, blue throughput → 0



Host A

λ_in : original data
λ'_in : original data, *plus* retransmitted data

λ_out     Host B

finite shared output link buffers

Host D

Host C

---

## Causes/costs of congestion: scenario 3



C/2

λ_out

λ'_in     C/2

another "cost" of congestion:

❖ when packet dropped, any "upstream" transmission capacity used for that packet was wasted!

2

## Approaches towards congestion control

two broad approaches towards congestion control:

| end-end congestion control: | network-assisted congestion control: |
|---|---|
| ❖ no explicit feedback from network | ❖ routers provide feedback to end systems |
| ❖ congestion inferred from end-system observed loss, delay | ▪ single bit indicating congestion (SNA, DECbit, TCP/IP ECN, ATM) |
| ❖ approach taken by TCP | ▪ explicit rate for sender to send at |

## Chapter 3 outline

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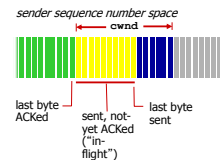…. until loss occurs (then cut window in half)

**cwnd**: TCP sender congestion window size

time

## TCP Congestion Control: details

*sender sequence number space*

cwnd

last byte ACKed | sent, not-yet ACKed ("in-flight") | last byte sent

- sender limits transmission:

```
LastByteSent-
LastByteAcked   ≤  cwnd
```

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

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

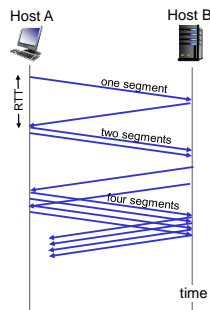$$\text{rate} \approx \frac{\text{cwnd}}{\text{RTT}} \text{ bytes/sec}$$

## 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

## TCP: detecting, reacting to loss

- loss indicated by timeout:
  - **cwnd** set to 1 MSS;
  - window then grows exponentially (as in slow start) to threshold, then grows linearly
- loss indicated by 3 duplicate ACKs: TCP RENO
  - recv'd ACKs indicate network capable of delivering some segments
  - **cwnd** is cut in half window then grows linearly
- TCP Tahoe (Van Jacobson 1988) always sets **cwnd** to 1 (timeout or 3 duplicate acks)
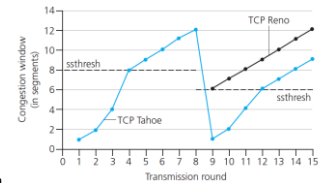
## Tahoe, Reno, and Vegas

- TCP Tahoe (~1988 Van Jacobson): BSD Unix 4.3, a.k.a. BSD Network Release 1.0 (BNR1), additive increase and multiplicative decrease, slow start, no fast retransmission

- TCP Reno (~1990?): BNR2, BNR1 plus fast retransmission, header prediction (fast path for pure ACKs and in-order packets), delayed ACKs

- TCP Vegas (~1994 Brakmo, O'Malley, and Peterson): varying congestion window size $w$ between $a$ and $b$, based on $diff$ = (expected – sample) rate of transmission. If $diff < a$ (more capacity available), increase $w$ by one, if $diff > b$ (showing congestion), decrease $w$ by one

## TCP: switching from slow start to CA

Q: when should the exponential increase switch to linear?
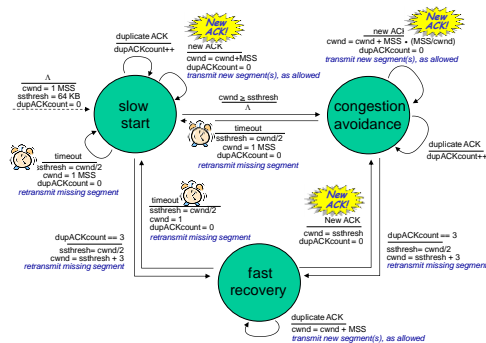
A: when **cwnd** gets to 1/2 of its value before timeout.



### Implementation:
- variable **ssthresh**
- on loss event, **ssthresh** is set to 1/2 of **cwnd** just before loss event

*For metrics such as cwnd and ssthresh, check out the structures in /usr/include/netinet/tcp.h*
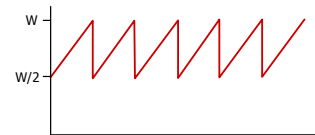
## Summary: TCP Congestion Control

## TCP throughput

- avg. TCP thruput 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. thruput is 3/4W per RTT

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

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

- example: 1500 byte segments, 100ms RTT, want 10 Gbps throughput (1500 bytes = 12,000 bits seg, 100 ms can carry 83,333 segments at 10Gbps)
- requires W = 83,333 in-flight segments
- throughput in terms of segment loss probability, L [Mathis 1997]:

$$\text{TCP throughput} = \frac{1.22 \cdot MSS}{RTT\sqrt{L}}$$
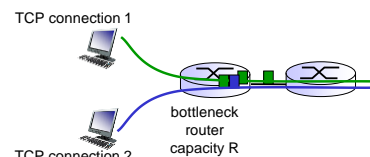
  → to achieve 10 Gbps throughput, need a loss rate of L = $2 \cdot 10^{-10}$ – *a very small loss rate!*

- these observations led to new versions of TCP for high-speed [Jin 2004; RFC 3649; Kelly 2003; Ha 2008].

## TCP Fairness

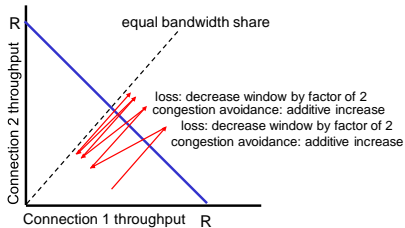*fairness goal:* if K TCP sessions share same bottleneck link of bandwidth R, each should have average rate of R/K



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



Connection 2 throughput

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 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 9 TCPs, gets R/2

## Examine some source code

- Linux 2.6 implementation of TCP congestion control:
  - http://lxr.free-electrons.com/source/net/ipv4/tcp_cong.c
- Look for
  - snd_cwnd
  - tcp_slow_start
  - tcp_cong_avoid_ai
  - tcp_reno_cong_avoid

## Chapter 3: 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"