

Congestion Control

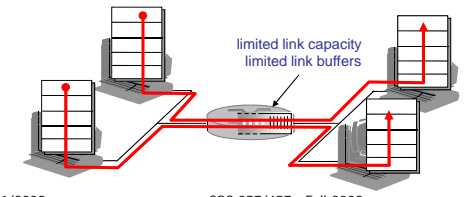
Kai Shen

10/21/2009 CSC 257/457 - Fall 2009 1

Principles of Congestion Control

Congestion:

- informally: "too many sources sending too much data too fast for the network to handle"
- results of congestion:
 - long delays (e.g. queueing in router buffers)
 - lost packets (e.g. buffer overflow at routers)



10/21/2009 CSC 257/457 - Fall 2009 2

Congestion Control

Congestion threatens the growth of the Internet

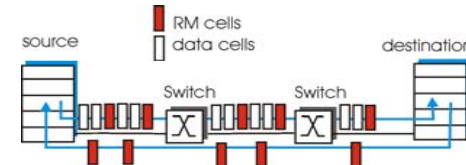
- tragedy of commons

Two broad approaches towards congestion control:

<p>Network-assisted congestion control:</p> <ul style="list-style-type: none"> ■ routers provide feedback to end systems <ul style="list-style-type: none"> ■ mark bit indicating congestion ■ explicit rate sender should send at 	<p>End-end congestion control:</p> <ul style="list-style-type: none"> ■ no explicit feedback from network ■ congestion inferred from end-system observed loss, delay
---	---

10/21/2009 CSC 257/457 - Fall 2009 3

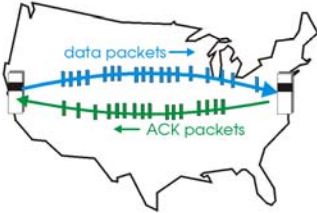
Network-assisted Congestion Control: ATM



- **RM (resource management) cells:**
 - sent by sender, interspersed with data cells
- bits in RM cell set by switches ("network-assisted")
 - **NI bit:** no increase in rate (mild congestion)
 - **CI bit:** congestion indication
 - **ER (explicit rate) field:** congested switch may lower ER value in cell
- RM cells returned to sender by receiver, with bits intact
 - sender control its rate based on information received in RM cells

10/21/2009 CSC 257/457 - Fall 2009 4

TCP Congestion Control



- **Mechanism:** sender controls the sending rate by adjusting the number of packets allowed in flight simultaneously (size of the sliding window)
- **Objective:** figure out the highest rate that does not cause network congestion

10/21/2009 CSC 257/457 - Fall 2009 5

TCP Congestion Control

- **CongWin** is maximum allowed in-flight or un-ACKed bytes
- Roughly:
 - throughput = $\text{CongWin} / \text{RTT}$
- dynamically control **CongWin**, based on perceived network congestion
- **Goal:** maintain a high rate that does not cause congestion

How does sender perceive congestion?

- loss event = timeout *or* triple duplicate ACKs
- timeout before 3 dup ACKs is **more alarming**
 - 3 dup ACKs indicates network capable of delivering some segments

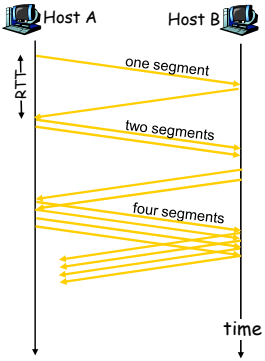
two phases for TCP congestion control:

- start phase (having little idea of where the good rate is)
- steady state phase (being quite close to the good rate)

10/21/2009 CSC 257/457 - Fall 2009 6

Slow Start Phase

- At the beginning, **CongWin** = 1 MSS
- Available bandwidth may be \gg MSS/RTT
 - Example: MSS = 500 bytes & RTT = 200 msec
 \Rightarrow initial rate = 20 kbps
 - desirable to quickly ramp up to respectable rate
- **Exponential increase:** double **CongWin** after every RTT in the absence of loss events
- In this phase, the **CongWin** starts at a slow rate, but it grows fast
- When should we stop exponential growth?



10/21/2009 CSC 257/457 - Fall 2009 7

When does Slow Start become Steady State?

- Stop exponential growth:
 - when there is a loss event
 - or when **CongWin** reaches a **Threshold**.
- **Threshold** divides between two states
 - when you are far away from causing congestion, so you do not need to worry about congestion
 - when you are very close to causing congestion, do need to be very careful
- **Threshold** is maintained based on past history:
 - set to half of **CongWin** at a loss event
 - initially set to a large value so it has no effect;

10/21/2009 CSC 257/457 - Fall 2009 8

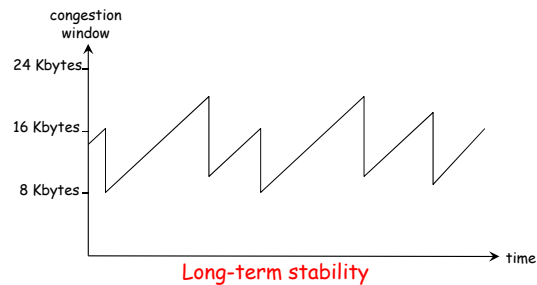
Steady State Phase: AIMD

Additive increase:

increase CongWin by 1 MSS every RTT in the absence of loss events

Multiplicative decrease:

cut CongWin in half after 3 duplicate ACKs



10/21/2009

CSC 257/457 - Fall 2009

9

Handling Loss Events

- After 3 dup ACKs:
 - CongWin is cut in half, it then grows linearly
- But after timeout event:
 - CongWin instead set to 1 MSS, fall back to **slow start** - the very beginning of the connection

Philosophy:

timeout before 3 dup ACKs is **more alarming** because 3 dup ACKs indicates network capable of delivering some segments

10/21/2009

CSC 257/457 - Fall 2009

10

Summary: TCP Congestion Control

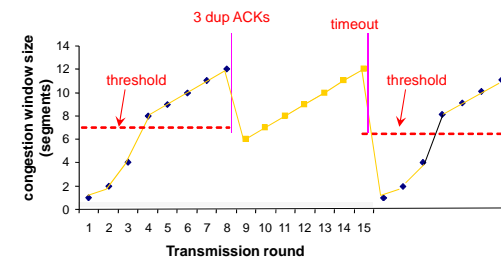
- At **start** phase, window starts at 1 MSS and grows exponentially.
- When CongWin is above **Threshold**, sender is in **steady state** phase, window grows linearly.
- When a **triple duplicate ACK** occurs, CongWin is cut in half.
- When **timeout** occurs, **Threshold** set to CongWin/2 and CongWin is set to 1 MSS (**start** again).

10/21/2009

CSC 257/457 - Fall 2009

11

Congestion Control in Action



10/21/2009

CSC 257/457 - Fall 2009

12

Congestion Control Algorithms

- Reno is what we described.
- Tahoe (a predecessor)
 - does not consider three duplicate ACKs as a loss event.
- New Reno (an enhancement)
 - TCP send window limits the number of unacked segments, many of which are not in-flight;
 - congestion is most related to the number of in-flight segments;
 - send one more segment when receiving a duplicate ACK.
- Vegas

10/21/2009

CSC 257/457 - Fall 2009

13

Outline

- principles of congestion control
- congestion control in TCP
- impact of congestion control on fairness
- impact of congestion control on TCP efficiency

10/21/2009

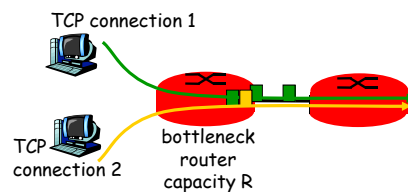
CSC 257/457 - Fall 2009

14

TCP Fairness

Fairness goal:

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



10/21/2009

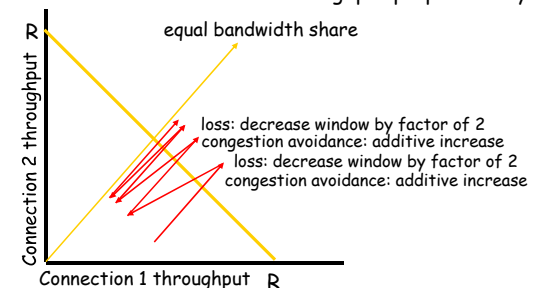
CSC 257/457 - Fall 2009

15

AIMD is Fair

Two competing sessions:

- Additive increase gives slope of 45° , as throughputs of both sessions increase at the same speed
- multiplicative decrease decreases throughput proportionally



10/21/2009

CSC 257/457 - Fall 2009

16

Fairness with UDP

- UDP is not congestion-controlled
 - UDP apps can send data as fast as they want, despite losses (its own losses as well as others)
 - Multimedia applications
- TCP is network friendly while UDP is not
- Regulate UDP traffic within the network

10/21/2009
CSC 257/457 - Fall 2009
17

Fairness with Parallel TCP Connections

- Nothing prevents app from opening parallel connections between 2 hosts.
- Example: link of rate R currently supporting 9 connections;
 - new app uses 1 TCP connection, gets rate R/10
 - new app asks for 9 TCP connections, gets R/2 !

10/21/2009
CSC 257/457 - Fall 2009
18

TCP Efficiency

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

- data transmission delay
- TCP connection establishment
- congestion control (slow start)

Case study:

- assume fixed congestion window: W segments
- more complicated with full TCP congestion control

10/21/2009
CSC 257/457 - Fall 2009
19

Fixed Congestion Window (1)

Notation, assumptions:

- Link rate: R
- Object size: O
- Fixed window size: WS
- No loss, no corruption

Case 1: congestion window is large

- ACK for first segment in window returns before window's worth of data sent

$delay = 2RTT + O/R$

10/21/2009
CSC 257/457 - Fall 2009
20

Fixed Congestion Window (2)

Notation, assumptions:

- Link rate: R
- Object size: O
- Fixed window size: WS
- No loss, no corruption

Case 2: congestion window is not large

- wait for ACK after sending window's worth of data sent

$$\text{delay} = 2RTT + O/R + (K-1)[S/R + RTT - WS/R]$$

where $K = O/WS$

10/21/2009 CSC 257/457 - Fall 2009 21

TCP Efficiency with Slow Start

- small window size in the slow start process \Rightarrow inefficiency
- inefficiency is amortized over long-running connections
 - persistent connection in HTTP/1.1

10/21/2009 CSC 257/457 - Fall 2009 22

Disclaimer

- Parts of the lecture slides contain original work of James Kurose, Larry Peterson, and Keith Ross. The slides are intended for the sole purpose of instruction of computer networks at the University of Rochester. All copyrighted materials belong to their original owner(s).

10/21/2009 CSC 257/457 - Fall 2009 23