

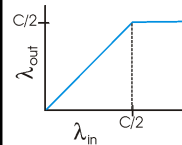
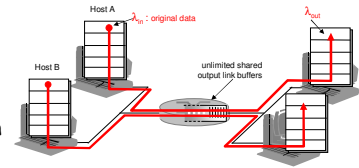
Congestion Control

EECS 122
February 16, 2006

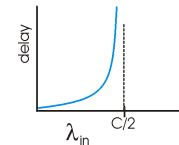
Transport Layer 1

Causes/costs of congestion: scenario 1

- two senders, two receivers
- one router, infinite buffers
- no retransmission



- large delays when congested
- maximum achievable throughput



Transport Layer 4

- Hw 2 due today
- Hw 3 and first phase of project out today

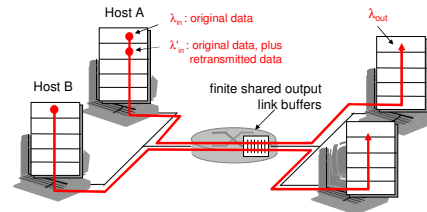
Lecture today:

- Impact of network congestion on end-to-end performance
- Approaches to congestion control
- How TCP does it.

Transport Layer 2

Causes/costs of congestion: scenario 2

- one router, *finite* buffers
- sender retransmission of lost packet



Transport Layer 5

Impact of Network Congestion

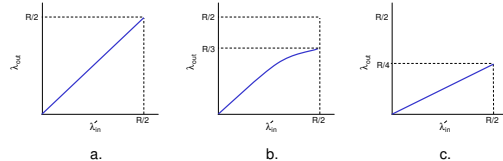
Congestion:

- 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)

Transport Layer 3

Causes/costs of congestion: scenario 2

- always: $\lambda_{in} = \lambda_{out}$ (goodput)
- "perfect" retransmission only when loss: $\lambda'_{in} > \lambda_{out}$, λ'_{in} larger (than perfect case) for same λ_{out}



"costs" of congestion:

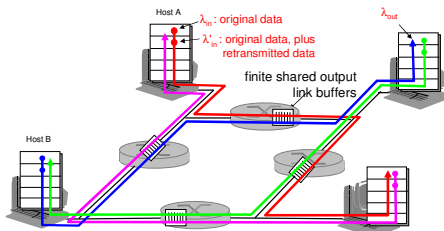
- more work (retrans) for given "goodput"
- unneeded retransmissions: link carries multiple copies of pkt

Transport Layer 6

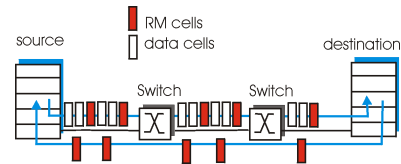
Causes/costs of congestion: scenario 3

- four senders
- multihop paths
- timeout/retransmit

Q: what happens as λ_{in} and λ'_{in} increase?

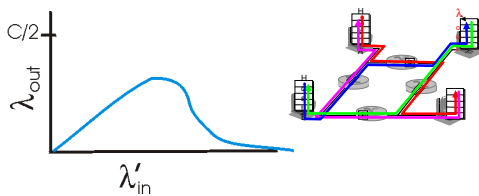


Case study: ATM congestion control



- Virtual circuit architecture
- Switches inside the network cognizant of individual connections.
- Explicit rate notification from each switch fed back to sender.
- Intelligence inside the network vs at the endpoints.

Causes/costs of congestion: scenario 3



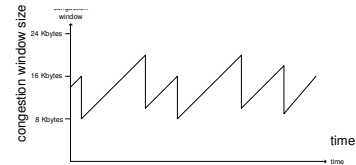
Another "cost" of congestion:

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

TCP congestion control: additive increase, multiplicative decrease

- Approach: increase transmission rate (window size), probing for usable bandwidth, until loss occurs
 - additive increase: increase CongWin by 1 MSS every RTT until loss detected
 - multiplicative decrease: cut CongWin in half after loss

Saw tooth behavior: probing for bandwidth



Approaches towards congestion control

Two broad approaches towards congestion control:

End-end congestion control:

- no explicit feedback from network
- congestion inferred from end-system observed loss, delay
- approach taken by TCP

Network-assisted congestion control:

- routers provide feedback to end systems
 - single bit indicating congestion (SNA, DECbit, TCP/IP ECN, ATM)
 - explicit rate sender should send at

TCP Congestion Control: details

sender limits transmission: $\text{LastByteSent} - \text{LastByteAcked} \leq \text{CongWin}$

Roughly,

$$\text{rate} = \frac{\text{CongWin}}{\text{RTT}} \text{ Bytes/sec}$$

- CongWin is dynamic, function of perceived network congestion

How does sender perceive congestion?

- loss event = timeout or 3 duplicate acks
- TCP sender reduces rate (CongWin) after loss event

three mechanisms:

- AIMD
- slow start
- conservative after timeout events

TCP Slow Start

- When connection begins, $\text{CongWin} = 1 \text{ MSS}$
 - Example: $\text{MSS} = 500$ bytes & $\text{RTT} = 200$ msec
 - initial rate = 20 kbps
- available bandwidth may be $\gg \text{MSS}/\text{RTT}$
 - desirable to quickly ramp up to respectable rate
- When connection begins, increase rate exponentially fast until first loss event

Transport Layer 13

Refinement: inferring loss

- After 3 dup ACKs:
 - CongWin is cut in half
 - window then grows linearly
- **But** after timeout event:
 - CongWin instead set to 1 MSS;
 - window then grows exponentially
 - to a threshold, then grows linearly

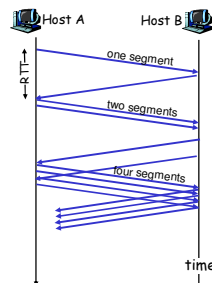
Philosophy:

- 3 dup ACKs indicates network capable of delivering some segments
- timeout indicates a "more alarming" congestion scenario

Transport Layer 16

TCP Slow Start (more)

- When connection begins, increase rate exponentially until first loss event:
 - double CongWin every RTT
 - done by incrementing CongWin for every ACK received
- **Summary:** initial rate is slow but ramps up exponentially fast



Transport Layer 14

Summary: TCP Congestion Control

- When CongWin is below **Threshold**, sender in **slow-start** phase, window grows exponentially.
- When CongWin is above **Threshold**, sender is in **congestion-avoidance** phase, window grows linearly.
- When a **triple duplicate ACK** occurs, **Threshold** set to $\text{CongWin}/2$ and CongWin set to **Threshold**.
- When **timeout** occurs, **Threshold** set to $\text{CongWin}/2$ and CongWin is set to 1 MSS.

Transport Layer 17

Refinement

- Q:** When should the exponential increase switch to linear?
- A:** When CongWin gets to 1/2 of its value before timeout.
- Q:** What happens when there is loss?
- A:** Threshold is set to 1/2 of CongWin just before loss event

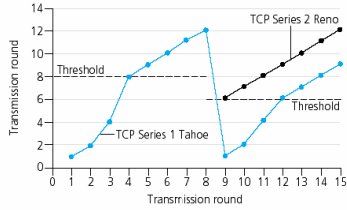
Transport Layer 15

TCP throughput

- What's the average throughput of TCP as a function of window size and RTT?
 - Ignore slow start
- Let W be the window size when loss occurs.
- When window is W , throughput is W/RTT
- Just after loss, window drops to $W/2$, throughput to $W/2\text{RTT}$.
- Average throughput: $.75 W/\text{RTT}$

Transport Layer 18

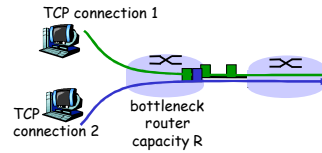
TCP Reno vs Tahoe



Transport Layer 19

TCP Fairness

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



Transport Layer 22

TCP Reno vs Vegas

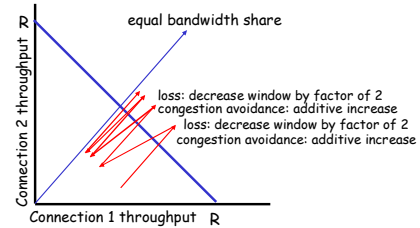
- Quick reaction needed on observing losses => halving the window
- Throughput reduction
- If sender can anticipate losses beforehand, can react more gradually (linear instead of halving).
- Some clues can be obtained by monitoring the RTT's of the segments.

Transport Layer 20

Why is TCP fair?

Two competing sessions:

- Additive increase gives slope of 1, as throughput increases
- multiplicative decrease decreases throughput proportionally



Transport Layer 23

Fast TCP

- Example: 1500 byte segments, 100ms RTT, want 10 Gbps throughput
- Requires window size $W = 83,333$ in-flight segments
- Throughput in terms of loss rate:

$$\frac{1.22 \cdot MSS}{RTT \sqrt{L}}$$

- → $L = 2 \cdot 10^{10}$ *Wow*
- New versions of TCP for high-speed needed!

Transport Layer 21

Fairness (more)

Fairness and UDP

- Multimedia apps often do not use TCP
 - do not want rate throttled by congestion control
- Instead use UDP:
 - pump audio/video at constant rate, tolerate packet loss
- Research area: TCP friendly

Fairness and parallel TCP connections

- nothing prevents app from opening parallel connections between 2 hosts.
- Web browsers do this
- Example: link of rate R supporting 9 cncntions;
 - new app asks for 1 TCP, gets rate R/10
 - new app asks for 11 TCPs, gets R/2 !

Transport Layer 24