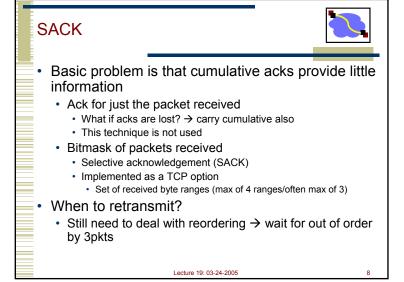
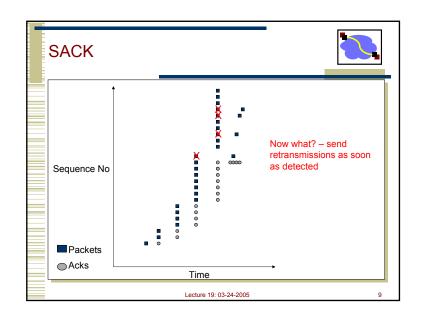
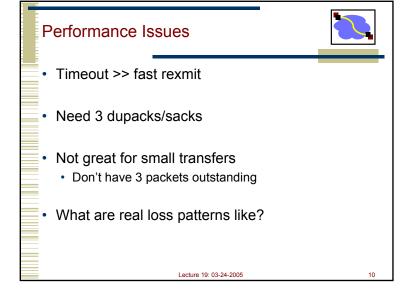
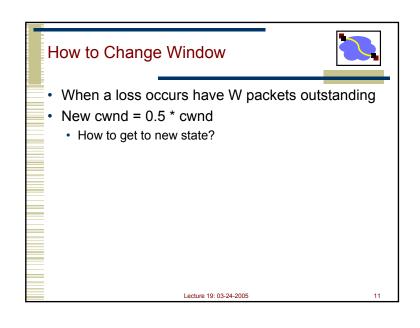


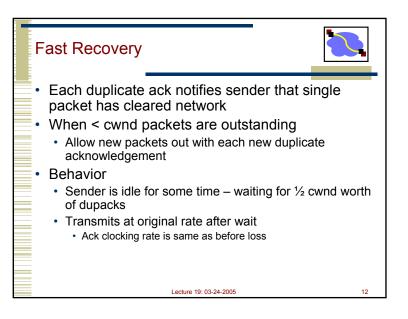
TCP Reno (1990) · All mechanisms in Tahoe Addition of fast-recovery · Opening up congestion window after fast retransmit Delayed acks Header prediction • Implementation designed to improve performance · Has common case code inlined With multiple losses, Reno typically timeouts because it does not see duplicate acknowledgements Lecture 19: 03-24-2005

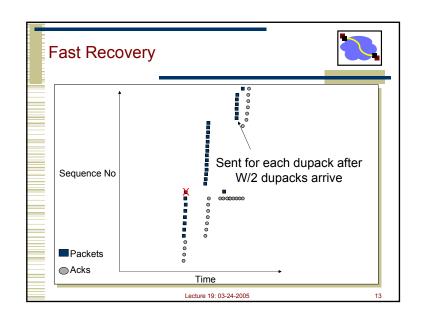


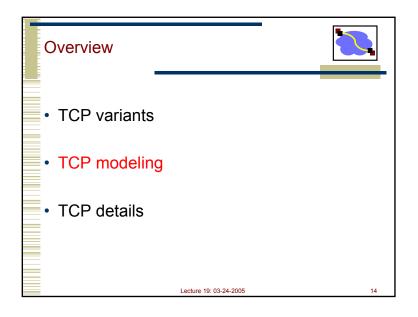


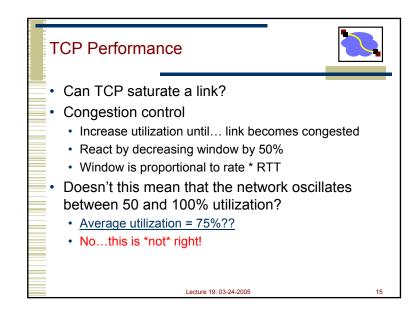


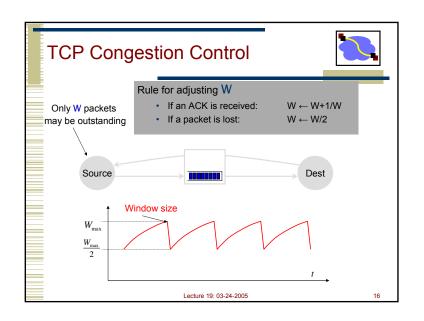


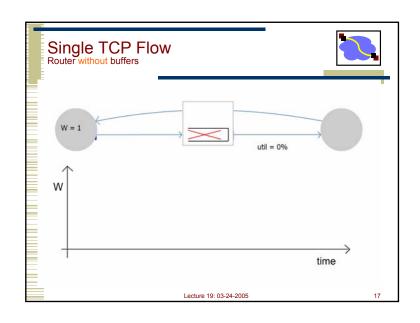


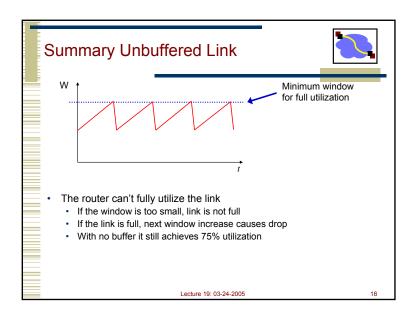




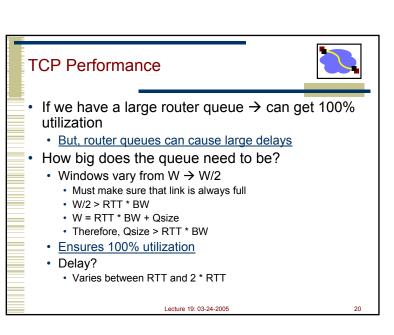


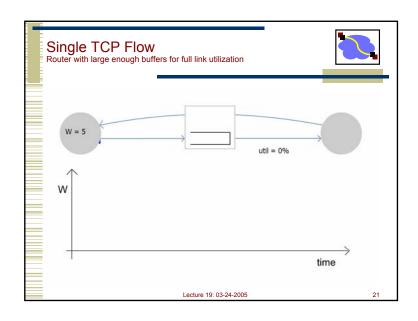


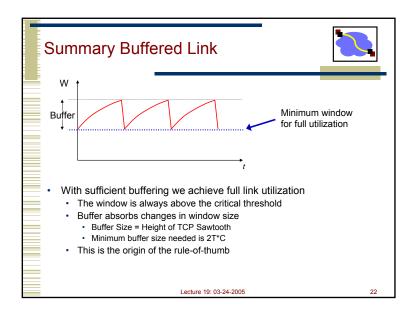


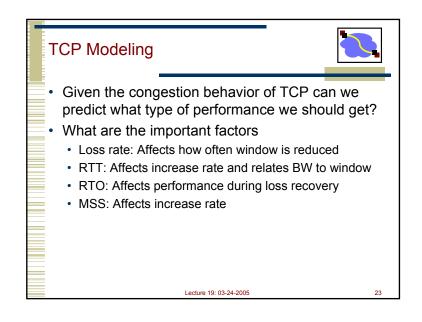


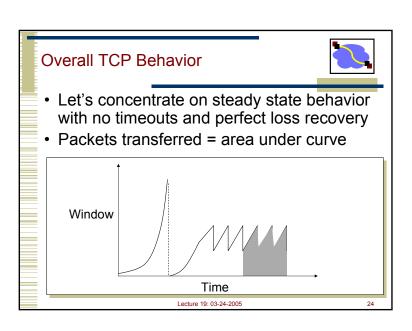
• In the real world, router queues play important role • Window is proportional to rate * RTT • But, RTT changes as well the window • Window to fill links = propagation RTT * bottleneck bandwidth • If window is larger, packets sit in queue on bottleneck link







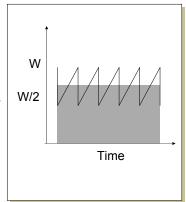




Transmission Rate



- · What is area under curve?
 - W = pkts/RTT, T = RTTs
 - A = avg window * time = ¾
 W * T
- What was bandwidth?
 - BW = A / T = 3/4 W
 - In packets per RTT
 - Need to convert to bytes per second
 - BW = 3/4 W * MSS / RTT
- · What is W?
 - Depends on loss rate



Lecture 19: 03-24-2005

Simple TCP Model



- Some additional assumptions
 - Fixed RTT
 - · No delayed ACKs
- In steady state, TCP losses packet each time window reaches W packets
 - Window drops to W/2 packets
 - Each RTT window increases by 1 packet→W/2 * RTT before next loss

Lecture 19: 03-24-2005

Simple Loss Model



- · What was the loss rate?
 - Packets transferred = (¾ W/RTT) * (W/2 * RTT) = 3W²/8
 - 1 packet lost → loss rate = p = 8/3W²

•
$$W = \sqrt{\frac{8}{3p}}$$

• BW = 3/4 * W * MSS / RTT

$$W = \sqrt{\frac{8}{3p}} = \frac{4}{3} \times \sqrt{\frac{3}{2p}}$$

•
$$BW = \frac{MSS}{RTT \times \sqrt{\frac{2p}{3}}}$$

Lecture 19: 03-24-2005

Fairness



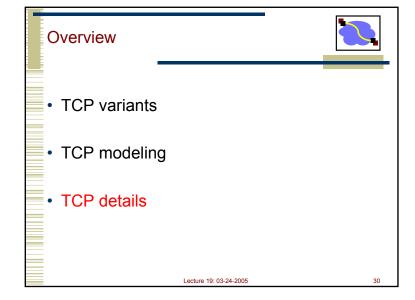
- BW proportional to 1/RTT?
- Do flows sharing a bottleneck get the same bandwidth?
 - NO!
- TCP is RTT fair
 - If flows share a bottleneck and have the same RTTs then they get same bandwidth
 - · Otherwise, in inverse proportion to the RTT

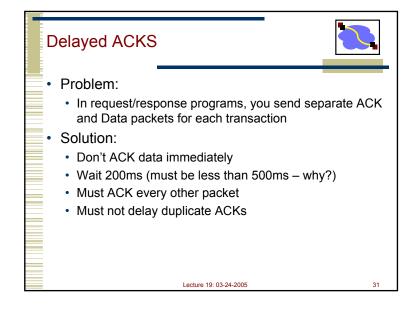
Lecture 19: 03-24-2005

28

What does it mean to be TCP friendly? TCP is not going away Any new congestion control must compete with TCP flows Should not clobber TCP flows and grab bulk of link Should also be able to hold its own, i.e. grab its fair share, or it will never become popular How is this quantified/shown? Has evolved into evaluating loss/throughput behavior If it shows 1/sqrt(p) behavior it is ok But is this really true?

Lecture 19: 03-24-2005





CP ACK Generation [R	FC 1122, RFC 2581]
Event	TCP Receiver action
In-order segment arrival, No gaps, Everything else already ACKed	Delayed ACK. Wait up to 500ms for next segment. If no next segment, send ACK
In-order segment arrival, No gaps, One delayed ACK pending	Immediately send single cumulative ACK
Out-of-order segment arrival Higher-than-expect seq. # Gap detected	Send duplicate ACK, indicating seq. # of next expected byte
Arrival of segment that partially or completely fills gap	Immediate ACK
Lectur	re 19: 03-24-2005

Delayed Ack Impact



- · TCP congestion control triggered by acks
 - If receive half as many acks → window grows half as fast
- Slow start with window = 1
 - · Will trigger delayed ack timer
 - First exchange will take at least 200ms
 - Start with > 1 initial window
 - · Bug in BSD, now a "feature"/standard

Lecture 19: 03-24-2005

33

Nagel's Algorithm



- · Small packet problem:
 - Don't want to send a 41 byte packet for each keystroke
 - How long to wait for more data?
- · Solution:
 - Allow only one outstanding small (not full sized) segment that has not yet been acknowledged
 - Can be disabled for interactive applications

Lecture 19: 03-24-2005

24

Large Windows



- Delay-bandwidth product for 100ms delay
 - 1.5Mbps: 18KB
 - 10Mbps: 122KB
 - 45Mbps: 549KB
 - 100Mbps: 1.2MB
 - 622Mbps: 7.4MB
 - 1.2Gbps: 14.8MB
- · Why is this a problem?
 - 10Mbps > max 16bit window
- Scaling factor on advertised window
- · Specifies how many bits window must be shifted to the left
- · Scaling factor exchanged during connection setup

Lecture 19: 03-24-2005

25

Window Scaling: Example Use of Options "Large window" option (RFC 1323) TCP syn Negotiated by the hosts during connection establishment Option 3 specifies the number of bits by which to shift the value in the 16 bit window field Independently set for the two transmit directions The scaling factor specifies bit shift of the window field in the TCP header Scaling value of 2 translates into a factor of 4 Old TCP implementations will TCP ack simply ignore the option · Definition of an option Lecture 19: 03-24-2005

Maximum Segment Size (MSS)



- Problem: what packet size should a connection use?
- Exchanged at connection setup
 - Uses a TCP option
 - · Typically pick MTU of local link
- · What all does this effect?
 - Efficiency
 - · Congestion control
 - Retransmission
- Path MTU discovery
 - · Why should MTU match MSS?

Lecture 19: 03-24-2005

37

TCP (Summary)



- · General loss recovery
 - Stop and wait
 - Selective repeat
- TCP sliding window flow control
- TCP state machine
- TCP loss recovery
 - · Timeout-based
 - RTT estimation
 - Fast retransmit
 - · Selective acknowledgements

Lecture 19: 03-24-2005

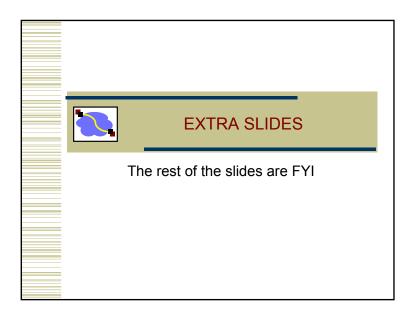
38

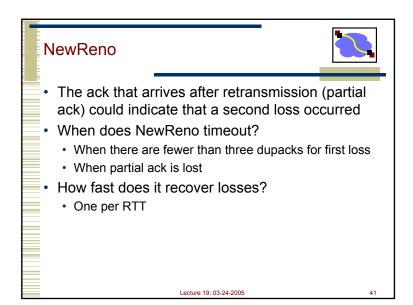
TCP (Summary)

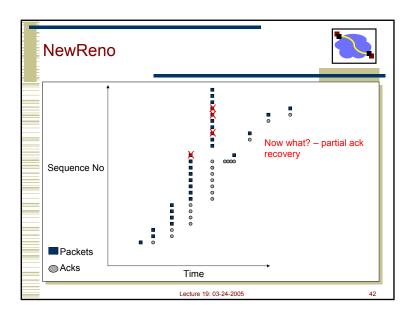


- Congestion collapse
 - Definition & causes
- · Congestion control
 - Why AIMD?
 - Slow start & congestion avoidance modes
 - ACK clocking
 - Packet conservation
- TCP performance modeling
 - · How does TCP fully utilize a link?
 - · Role of router buffers

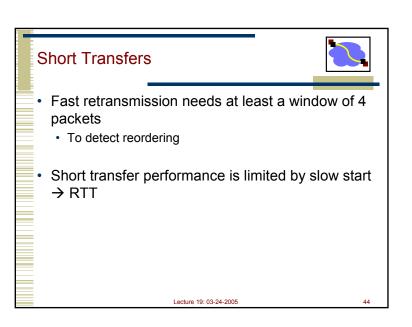
Lecture 19: 03-24-2005







Changing Workloads New applications are changing the way TCP is used 1980's Internet Telnet & FTP → long lived flows Well behaved end hosts Homogenous end host capabilities Simple symmetric routing 2000's Internet Web & more Web → large number of short xfers Wild west – everyone is playing games to get bandwidth Cell phones and toasters on the Internet Policy routing



Short Transfers



- Start with a larger initial window
- · What is a safe value?
 - TCP already burst 3 packets into network during slow start
 - Large initial window = min (4*MSS, max (2*MSS, 4380 bytes)) [rfc2414]
 - Not a standard yet
 - Enables fast retransmission
 - · Only used in initial slow start not in any subsequent slow start

Lecture 19: 03-24-2005

Well Behaved vs. Wild West



- How to ensure hosts/applications do proper congestion control?
- · Who can we trust?
 - Only routers that we control
 - · Can we ask routers to keep track of each flow
 - Per flow information at routers tends to be expensive
 - · Fair-queuing later in the semester

Lecture 19: 03-24-2005

TCP Fairness Issues



- Multiple TCP flows sharing the same bottleneck link do not necessarily get the same bandwidth.
 - · Factors such as roundtrip time, small differences in timeouts, and start time. ... affect how bandwidth is shared
 - The bandwidth ratio typically does stabilize
- Users can grab more bandwidth by using parallel flows.
 - Each flow gets a share of the bandwidth to the user gets more bandwidth than users who use only a single flow

Lecture 19: 03-24-2005

Silly Window Syndrome



- Problem: (Clark, 1982)
 - If receiver advertises small increases in the receive window then the sender may waste time sending lots of small packets
- Solution
 - Receiver must not advertise small window increases
 - Increase window by min(MSS,RecvBuffer/2)

Lecture 19: 03-24-2005

Protection From Wraparound



- · Wraparound time vs. Link speed
 - 1.5Mbps: 6.4 hours
 - 10Mbps: 57 minutes
 - 45Mbps: 13 minutes
 - 100Mbps: 6 minutes
 - 622Mbps: 55 seconds
 - 1.2Gbps: 28 seconds
- Why is this a problem?
 - 55seconds < MSL!
- Use timestamp to distinguish sequence number wraparound

Lecture 19: 03-24-2005

Rule-of-thumb



- Rule-of-thumb makes sense for one flow
- Typical backbone link has > 20,000 flows
- Does the rule-of-thumb still hold?
 - Key assumption → losses are synchronized across all flows
 - All TCP connections halve windows nearly simultaneously
 - · Not necessarily true!

Lecture 19: 03-24-2005

E4

Example



- 10Gb/s linecard
 - · Requires 300Mbytes of buffering.
 - Read and write 40 byte packet every 32ns.
- · Memory technologies
 - DRAM: require 4 devices, but too slow.
 - SRAM: require 80 devices, 1kW, \$2000.
- Problem gets harder at 40Gb/s
 - · Hence RLDRAM, FCRAM, etc.

Lecture 19: 03-24-2005

50