

15-744: Computer Networking

L-5 Congestion Control



News



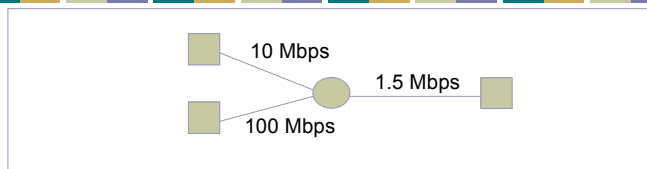
- Note problem set 1 update
 - Part E: Please draw three RTTs after the loss, not just one.
- Assigned reading
 - [JK88] Congestion Avoidance and Control
 - [CJ89] Analysis of the Increase and Decrease Algorithms for Congestion Avoidance in Computer Networks

© Srinivasan Seshan, 2004

L-4; 10-7-04

2

Congestion



- Different sources compete for resources inside network: Link b/w and queue space
- Why is it a problem?
 - Sources are unaware of current state of resource
 - Sources are unaware of each other
 - In many situations will result in < 1.5 Mbps of throughput (congestion collapse)

© Srinivasan Seshan, 2004

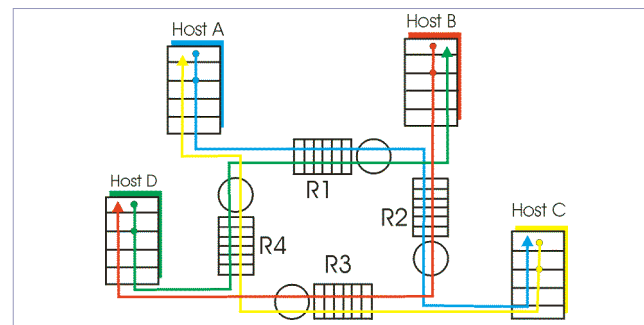
L-4; 10-7-04

3

Causes & Costs of Congestion



- Four senders – multihop paths **Q:** What happens as rate increases?
- Timeout/retransmit

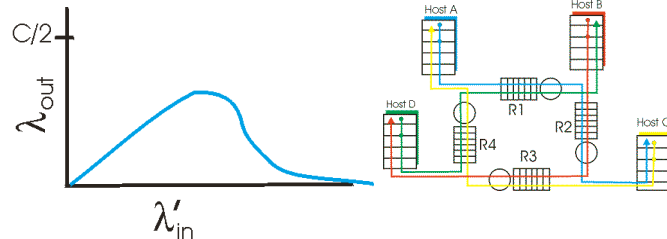


© Srinivasan Seshan, 2004

L-4; 10-7-04

4

Causes & Costs of Congestion



- When packet dropped, any “upstream transmission capacity used for that packet was wasted!

© Srinivasan Seshan, 2004

L -4; 10-7-04

5

Congestion Collapse

- Definition: *Increase in network load results in decrease of useful work done*
- Many possible causes
 - Spurious retransmissions of packets still in flight
 - Classical congestion collapse
 - How can this happen with packet conservation
 - Solution: better timers and TCP congestion control
 - Undelivered packets
 - Packets consume resources and are dropped elsewhere in network
 - Solution: congestion control for ALL traffic

© Srinivasan Seshan, 2004

L -4; 10-7-04

6

Other Congestion Collapse Causes

- Fragments
 - Mismatch of transmission and retransmission units
 - Solutions
 - Make network drop all fragments of a packet (early packet discard in ATM)
 - Do path MTU discovery
- Control traffic
 - Large percentage of traffic is for control
 - Headers, routing messages, DNS, etc.
- Stale or unwanted packets
 - Packets that are delayed on long queues
 - “Push” data that is never used

© Srinivasan Seshan, 2004

L -4; 10-7-04

7

Congestion Control and Avoidance

- Desirable properties:
 - Scalability:
 - # flows, range of capacities, range of delays
 - Do well in entire range!
 - Efficiency: High network utilization
 - Fairness
 - Works-ness: avoids collapse!
- Congestion collapse is not just a theory
 - Has been frequently observed in many networks

© Srinivasan Seshan, 2004

L -4; 10-7-04

8

Fairness



- Jain's fairness index
 - $f = (\sum x_i)^2 / n(\sum x_i^2)$
- All x equal: 1
- k/n get service: k/n
- Max-min fairness
 - No user receives more than their request, π_i
 - No other allocation satisfying (1) has higher min allocation
 - condition 2 holds as we remove the minimal user & reduce total resource accordingly
 - aka: $u_i = \text{MIN}(u \text{ fair}, \pi_i)$
- Goal: Something that works well enough.

© Srinivasan Seshan, 2004

L-4; 10-7-04

9

Objectives



- Simple router behavior
- Distributedness
- Efficiency: $X_{\text{knee}} = \sum x_i(t)$
- Fairness: $(\sum x_i)^2 / n(\sum x_i^2)$
- Power: $(\text{throughput}^\alpha / \text{delay})$
- Convergence: control system must be stable

© Srinivasan Seshan, 2004

L-4; 10-7-04

10

Design questions



- Congestion
 - How is congestion signaled?
 - Either mark or drop packets
 - When is a router congested?
 - Drop tail queues – when queue is full
 - Average queue length – at some threshold
- Control questions:
 - How do senders react to congestion?
 - How do senders determine capacity for flow?

© Srinivasan Seshan, 2004

L-4; 10-7-04

11

Congrol 2



- Upon congestion, flows must reduce rate
- How? Decrease algorithm
- If no congestion, flows might try sending more. Increase algorithm.
- Let's assume window-based flow control
 - sender maintains "cwnd": # of unacknowledged packets in the network at any time
 - Transmission rate: cwnd / rtt
- (Alternate: Rate-based; equation-based)

© Srinivasan Seshan, 2004

L-4; 10-7-04

12

Linear Control

- Many different possibilities for reaction to congestion and probing
 - Examine simple linear controls
 - $\text{Window}(t + 1) = a + b \text{Window}(t)$
 - Different a_i/b_i for increase and a_d/b_d for decrease
- Supports various reaction to signals
 - Increase/decrease additively
 - Increased/decrease multiplicatively
 - Which of the four combinations is optimal?

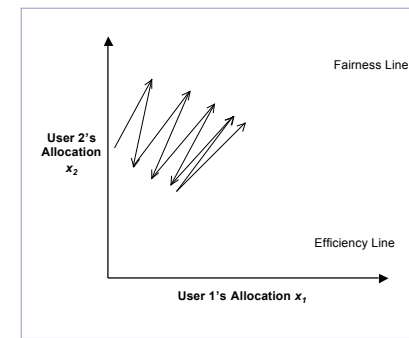
© Srinivasan Seshan, 2004

L-4; 10-7-04

13

Phase plots

- Simple way to visualize behavior of competing connections over time



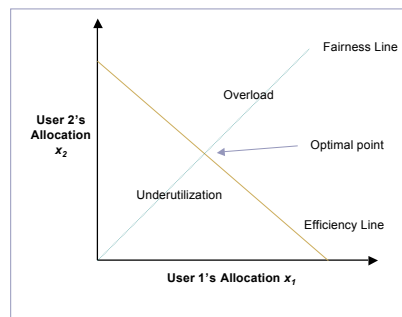
© Srinivasan Seshan, 2004

L-4; 10-7-04

14

Phase plots

- What are desirable properties?
- What if flows are not equal?



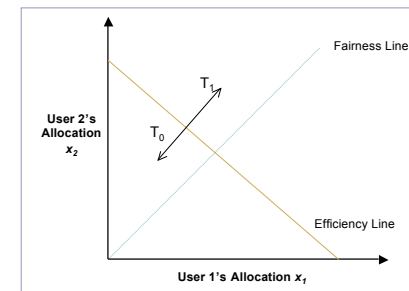
© Srinivasan Seshan, 2004

L-4; 10-7-04

15

Additive Increase/Decrease

- Both x_1 and x_2 increase/decrease by the same amount over time
 - Additive increase improves fairness and additive decrease reduces fairness



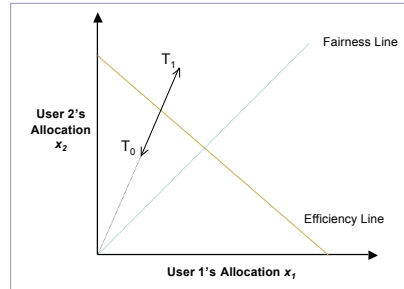
© Srinivasan Seshan, 2004

L-4; 10-7-04

16

Multiplicative Increase/Decrease

- Both x_1 and x_2 increase by the same factor over time
 - Extension from origin – constant fairness



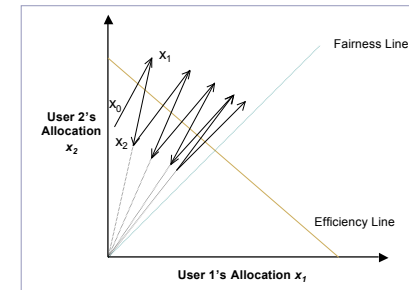
© Srinivasan Seshan, 2004

L-4; 10-7-04

17

What is the Right Choice?

- Constraints limit us to AIMD
 - Can have multiplicative term in increase (MAIMD)
 - AIMD moves towards optimal point



© Srinivasan Seshan, 2004

L-4; 10-7-04

18

TCP and linear controls

- Upon congestion:
 - $w(t+1) = a \cdot w(t)$ $0 < a < 1$
- While probing
 - $w(t+1) = w(t) + b$ $0 < b \ll w_{max}$
- TCP sets $a = 1/2$, $b = 1$ (packet)

© Srinivasan Seshan, 2004

L-4; 10-7-04

19

TCP Congestion Control

- Motivated by ARPANET congestion collapse
- Underlying design principle: packet conservation
 - At equilibrium, inject packet into network only when one is removed
 - Basis for stability of physical systems
- Why was this not working?
 - Connection doesn't reach equilibrium
 - Spurious retransmissions
 - Resource limitations prevent equilibrium

© Srinivasan Seshan, 2004

L-4; 10-7-04

20

TCP Congestion Control - Solutions



- Reaching equilibrium
 - Slow start
- Eliminates spurious retransmissions
 - Accurate RTO estimation
 - Fast retransmit
- Adapting to resource availability
 - Congestion avoidance

© Srinivasan Seshan, 2004

L-4; 10-7-04

21

TCP Congestion Control



- Changes to TCP motivated by ARPANET congestion collapse
- Basic principles
 - AIMD
 - Packet conservation
 - Reaching steady state quickly
 - ACK clocking

© Srinivasan Seshan, 2004

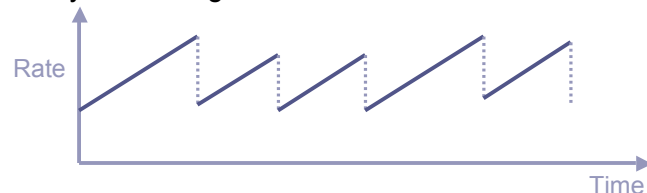
L-4; 10-7-04

22

AIMD: Now you grok the sawtooth



- Distributed, fair and efficient
- Packet loss is seen as sign of congestion and results in a multiplicative rate decrease
 - Factor of 2
- TCP periodically probes for available bandwidth by increasing its rate



© Srinivasan Seshan, 2004

L-4; 10-7-04

23

Congestion Avoidance



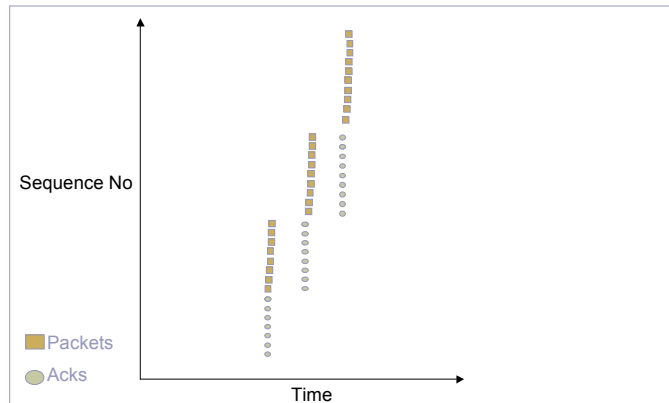
- If loss occurs when $cwnd = W$
 - Network can handle $0.5W \sim W$ segments
 - Set $cwnd$ to $0.5W$ (multiplicative decrease)
- Upon receiving ACK
 - Increase $cwnd$ by $(1 \text{ packet})/cwnd$
 - What is 1 packet? \rightarrow 1 MSS worth of bytes
 - After $cwnd$ packets have passed by \rightarrow approximately increase of 1 MSS
- Implements AIMD

© Srinivasan Seshan, 2004

L-4; 10-7-04

24

Congestion Avoidance Sequence Plot

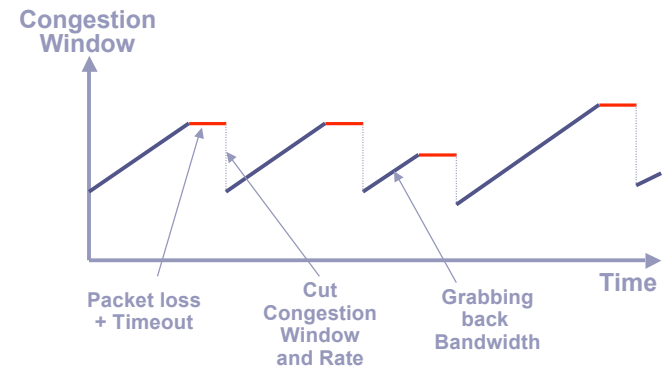


© Srinivasan Seshan, 2004

L-4; 10-7-04

25

Congestion Avoidance Behavior



© Srinivasan Seshan, 2004

L-4; 10-7-04

26

Packet Conservation

- At equilibrium, inject packet into network only when one is removed
 - Sliding window and not rate controlled
 - But still need to avoid sending burst of packets → would overflow links
 - Need to carefully pace out packets (ack clocking)!
 - Helps provide stability
- Need to eliminate spurious retransmissions
 - Accurate RTO estimation
 - Better loss recovery techniques (e.g. fast retransmit)

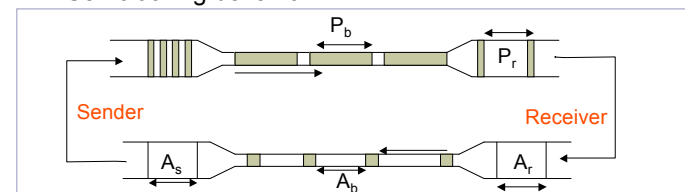
© Srinivasan Seshan, 2004

L-4; 10-7-04

27

TCP Packet Pacing

- Congestion window helps to “pace” the transmission of data packets
- In steady state, a packet is sent when an ack is received
 - Data transmission remains smooth, once it is smooth
 - Self-clocking behavior



© Srinivasan Seshan, 2004

L-4; 10-7-04

28

Reaching Steady State

- Doing AIMD is fine in steady state but slow...
- How does TCP know what is a good initial rate to start with?
 - Should work both for a CDPD (10s of Kbps or less) and for supercomputer links (10 Gbps and growing)
- Quick initial phase to help get up to speed (slow start)

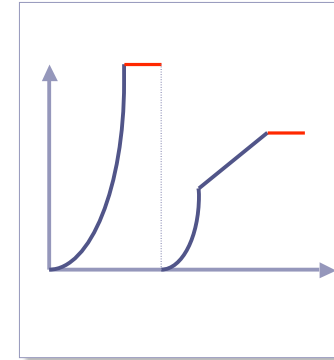
© Srinivasan Seshan, 2004

L-4; 10-7-04

29

Slow Start Packet Pacing

- How do we get this clocking behavior to start?
 - Initialize $\text{cwnd} = 1$
 - Upon receipt of every ack, $\text{cwnd} = \text{cwnd} + 1$
- Implications
 - Window actually increases to W in $\text{RTT} * \log_2(W)$
 - Can overshoot window and cause packet loss

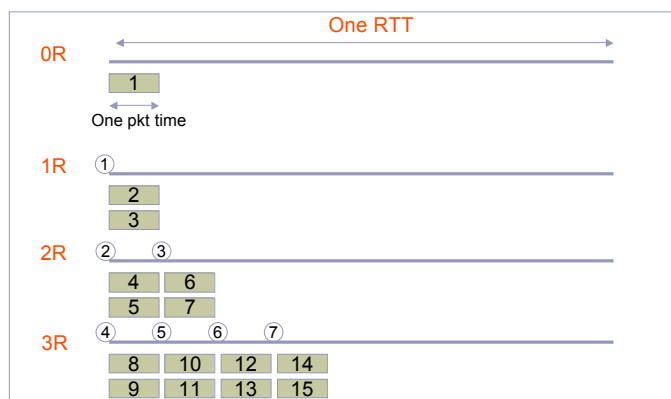


© Srinivasan Seshan, 2004

L-4; 10-7-04

30

Slow Start Example

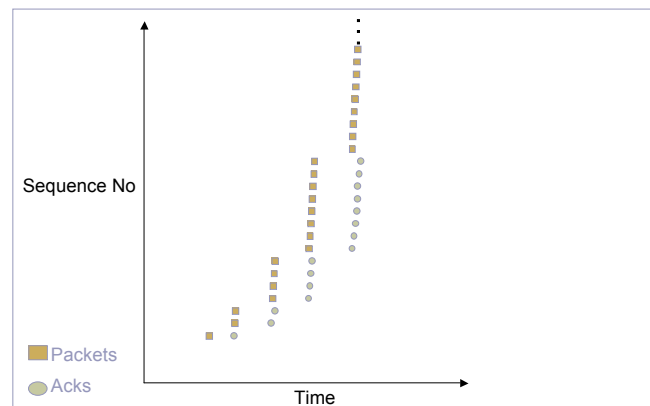


© Srinivasan Seshan, 2004

L-4; 10-7-04

31

Slow Start Sequence Plot



© Srinivasan Seshan, 2004

L-4; 10-7-04

32

Return to Slow Start

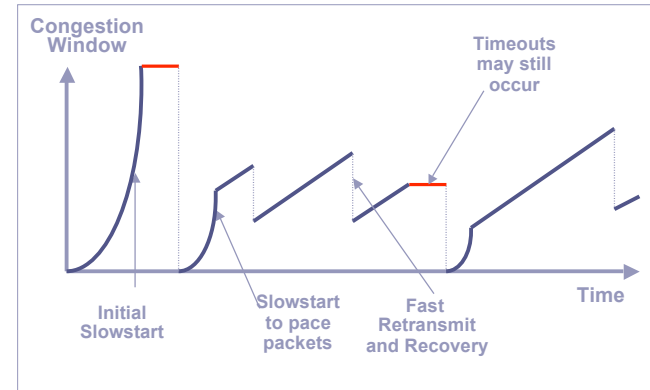
- If packet is lost we lose our self clocking as well
 - Need to implement slow-start and congestion avoidance together
- When timeout occurs set ssthresh to $0.5w$
 - If $cwnd < ssthresh$, use slow start
 - Else use congestion avoidance

© Srinivasan Seshan, 2004

L-4; 10-7-04

33

TCP Saw Tooth Behavior



© Srinivasan Seshan, 2004

L-4; 10-7-04

34

TCP Modeling

- Given the congestion behavior of TCP can we predict what type of performance we should get?
- What are the important factors
 - Loss rate
 - Affects how often window is reduced
 - RTT
 - Affects increase rate and relates BW to window
 - RTO
 - Affects performance during loss recovery
 - MSS
 - Affects increase rate

© Srinivasan Seshan, 2004

L-4; 10-7-04

35

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 $\rightarrow W/2$ * RTT before next loss
 - $BW = MSS * \text{avg window} / RTT = MSS * (W + W/2) / (2 * RTT) = .75 * MSS * W / RTT$

© Srinivasan Seshan, 2004

L-4; 10-7-04

36

Simple Loss Model



- What was the loss rate?
 - Packets transferred = $(.75 W/RTT) * (W/2 * RTT) = 3W^2/8$
 - 1 packet lost \rightarrow loss rate = $p = 8/3W^2$
 - $W = \sqrt{8 / (3 * \text{loss rate})}$
- $BW = .75 * MSS * W / RTT$
 - $BW = MSS / (RTT * \sqrt{2/3p})$

© Srinivasan Seshan, 2004

L-4; 10-7-04

37

TCP Friendliness



- 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?

© Srinivasan Seshan, 2004

L-4; 10-7-04

38

TCP Performance



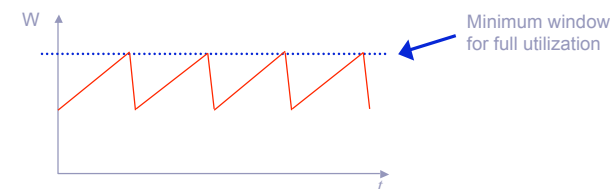
- Can TCP saturate a link?
- Congestion control
 - Increase utilization until... link becomes congested
 - React by decreasing window by 50%
 - Window is proportional to rate * RTT
- Doesn't this mean that the network oscillates between 50 and 100% utilization?
 - Average utilization = 75%??
 - **No...this is **not** right!**

© Srinivasan Seshan, 2004

L-4; 10-7-04

39

Summary Unbuffered Link



- The router can't fully utilize the link
 - If the window is too small, link is not full
 - If the link is full, next window increase causes drop
 - With no buffer it still achieves 75% utilization

© Srinivasan Seshan, 2004

L-4; 10-7-04

40

TCP Performance



- 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

© Srinivasan Seshan, 2004

L-4; 10-7-04

41

TCP Performance



- If we have a large router queue → can get 100% utilization
 - But, router queues can cause large delays
- How big does the queue need to be?
 - Windows vary from $W \rightarrow W/2$
 - Must make sure that link is always full
 - $W/2 > RTT * BW$
 - $W = RTT * BW + Qsize$
 - Therefore, $Qsize > RTT * BW$
 - Ensures 100% utilization
- Delay?
 - Varies between RTT and $2 * RTT$

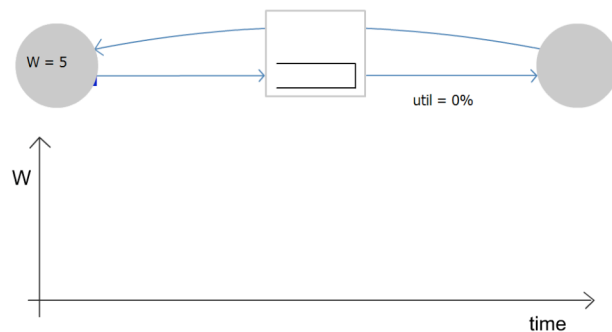
© Srinivasan Seshan, 2004

L-4; 10-7-04

42

Single TCP Flow

Router with large enough buffers for full link utilization



© Srinivasan Seshan, 2004

L-4; 10-7-04

43

Important Lessons



- How does TCP implement AIMD?
 - Sliding window, slow start & ack clocking
 - How to maintain ack clocking during loss recovery → fast recovery
- Modern TCP loss recovery
 - Why are timeouts bad?
 - How to avoid them? → fast retransmit, SACK
- How does TCP fully utilize a link?
 - Role of router buffers

© Srinivasan Seshan, 2004

L-4; 10-7-04

44

Additional Slides for the Curious



© Srinivasan Seshan, 2004

L-4; 10-7-04

45

Example



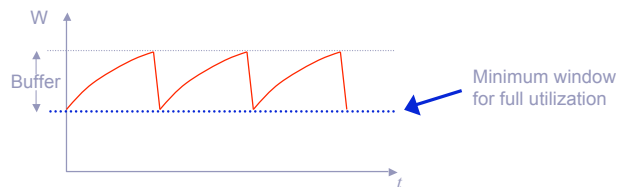
- 10Gb/s linecard
 - Requires 300Mbytes of buffering.
 - Read and write 40 byte packet every 32ns.
- Memory technologies
 - DRAM: require 4 devices, but too slow.
 - SRAM: require 80 devices, 1kW, \$2000.
- Problem gets harder at 40Gb/s
 - Hence RLDRAM, FCRAM, etc.

© Srinivasan Seshan, 2004

L-4; 10-7-04

46

Summary Buffered Link



- With sufficient buffering we achieve full link utilization
 - The window is always above the critical threshold
 - Buffer absorbs changes in window size
 - Buffer Size = Height of TCP Sawtooth
 - Minimum buffer size needed is $2T \cdot C$
 - This is the origin of the rule-of-thumb

© Srinivasan Seshan, 2004

L-4; 10-7-04

47

Rule-of-thumb



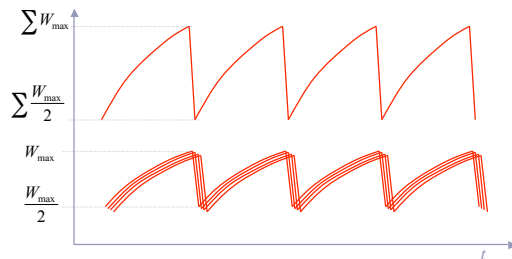
- Rule-of-thumb makes sense for one flow
- Typical backbone link has $> 20,000$ flows
- Does the rule-of-thumb still hold?

© Srinivasan Seshan, 2004

L-4; 10-7-04

48

If flows are synchronized



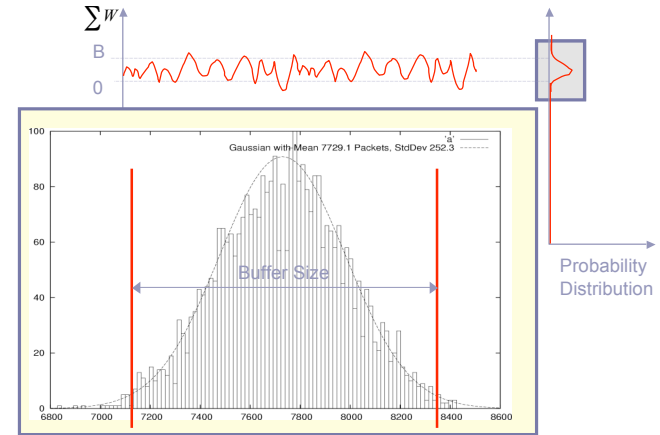
- Aggregate window has same dynamics
- Therefore buffer occupancy has same dynamics
- Rule-of-thumb still holds.

© Srinivasan Seshan, 2004

L-4; 10-7-04

49

If flows are not synchronized



© Srinivasan Seshan, 2004

L-4; 10-7-04

50

Central Limit Theorem

- CLT tells us that the more variables (Congestion Windows of Flows) we have, the narrower the Gaussian (Fluctuation of sum of windows)
- Width of Gaussian decreases with $\frac{1}{\sqrt{n}}$
- Buffer size should also decrease with $\frac{1}{\sqrt{n}}$

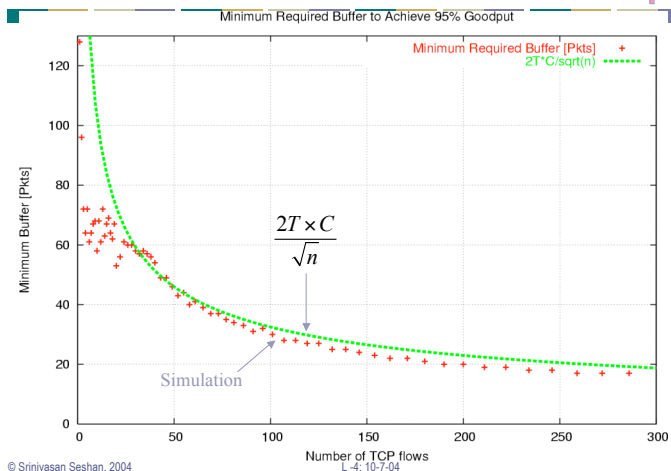
$$B \rightarrow \frac{B_{n=1}}{\sqrt{n}} = \frac{2T \times C}{\sqrt{n}}$$

© Srinivasan Seshan, 2004

L-4; 10-7-04

51

Required buffer size



© Srinivasan Seshan, 2004

L-4; 10-7-04

52

Integrity & Demultiplexing

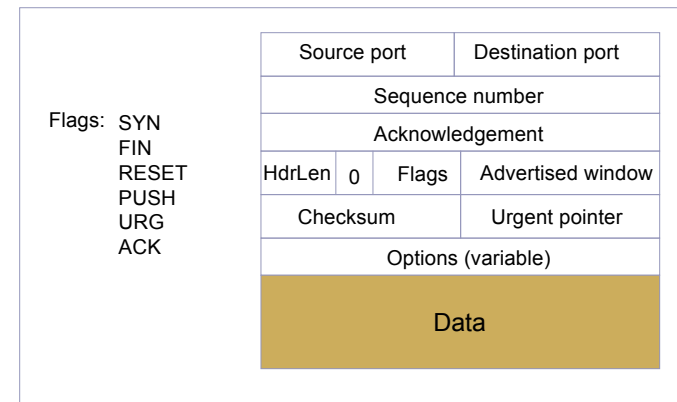
- Port numbers
 - Demultiplex from/to process
 - Servers wait on well known ports (/etc/services)
- Checksum
 - Is it sufficient to just checksum the packet contents?
 - No, need to ensure correct source/destination
 - Pseudoheader – portion of IP hdr that are critical
 - Checksum covers Pseudoheader, transport hdr, and packet body
- UDP provides just integrity and demux

© Srinivasan Seshan, 2004

L-4; 10-7-04

53

TCP Header



© Srinivasan Seshan, 2004

L-4; 10-7-04

54

TCP Flow Control

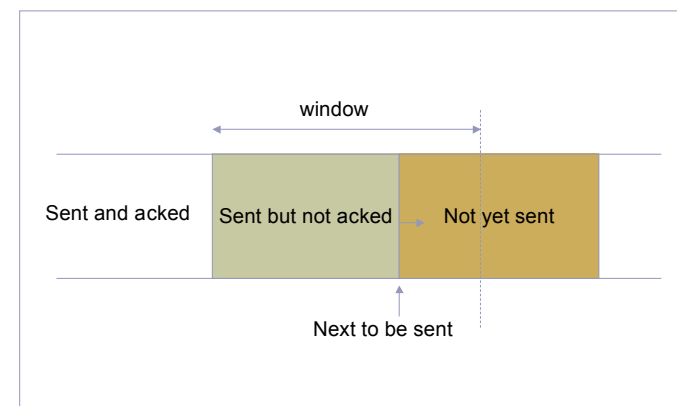
- TCP is a sliding window protocol
 - For window size n , can send up to n bytes without receiving an acknowledgement
 - When the data is acknowledged then the window slides forward
- Each packet advertises a window size
 - Indicates number of bytes the receiver has space for
- Original TCP always sent entire window
 - Congestion control now limits this

© Srinivasan Seshan, 2004

L-4; 10-7-04

55

Window Flow Control: Send Side

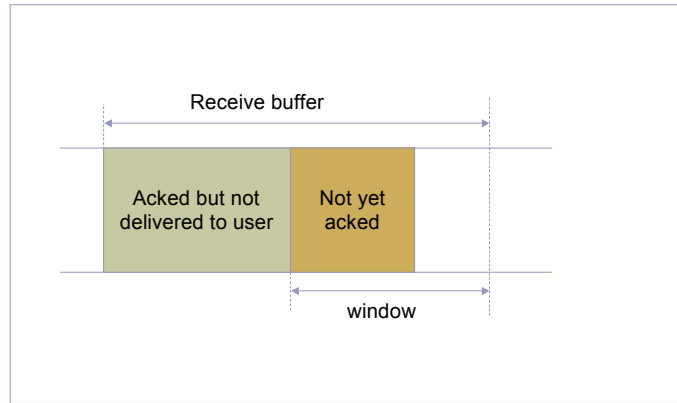


© Srinivasan Seshan, 2004

L-4; 10-7-04

56

Window Flow Control: Receive Side



© Srinivasan Seshan, 2004

L-4; 10-7-04

57

TCP Persist

- What happens if window is 0?
 - Receiver updates window when application reads data
 - What if this update is lost?
- TCP Persist state
 - Sender periodically sends 1 byte packets
 - Receiver responds with ACK even if it can't store the packet

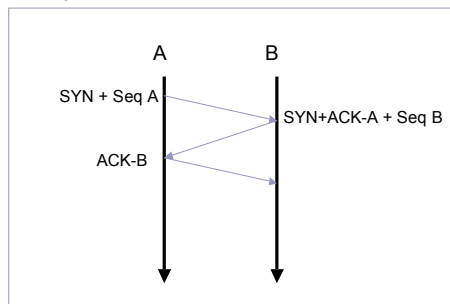
© Srinivasan Seshan, 2004

L-4; 10-7-04

58

Connection Establishment

- A and B must agree on initial sequence number selection
 - Use 3-way handshake



© Srinivasan Seshan, 2004

L-4; 10-7-04

59

Sequence Number Selection

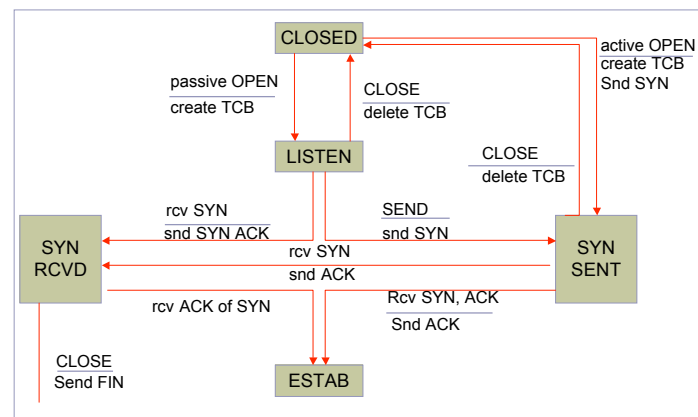
- Why not simply chose 0?
- Must avoid overlap with earlier incarnation

© Srinivasan Seshan, 2004

L-4; 10-7-04

60

Connection Setup



© Srinivasan Seshan, 2004

L-4; 10-7-04

61

Connection Tear-down

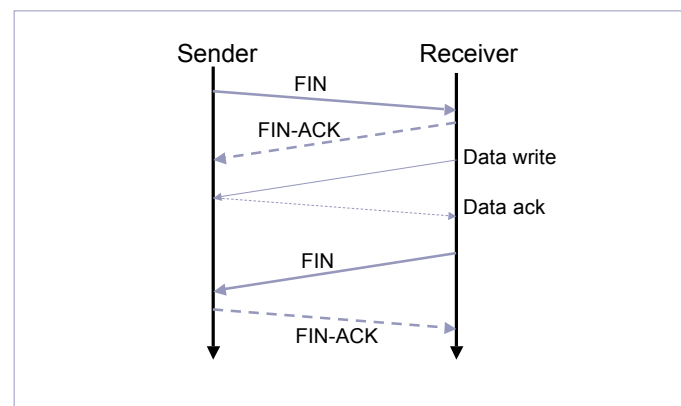
- Normal termination
 - Allow unilateral close
- TCP must continue to receive data even after closing
- Cannot close connection immediately
 - What if a new connection restarts and uses same sequence number?

© Srinivasan Seshan, 2004

L-4; 10-7-04

62

Tear-down Packet Exchange

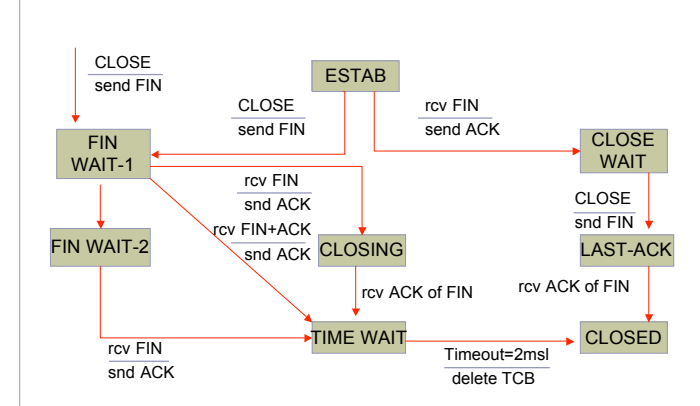


© Srinivasan Seshan, 2004

L-4; 10-7-04

63

Connection Tear-down



© Srinivasan Seshan, 2004

L-4; 10-7-04

64

Detecting Half-open Connections



TCP A

1. (CRASH)
2. CLOSED
3. SYN-SENT → <SEQ=400><CTL=SYN>
4. (!!) ← <SEQ=300><ACK=100><CTL=ACK> ←
5. SYN-SENT → <SEQ=100><CTL=RST>
6. SYN-SENT
7. SYN-SENT → <SEQ=400><CTL=SYN>

TCP B

- (send 300, receive 100)
- ESTABLISHED
- (??)
- ESTABLISHED
- (Abort!!)
- CLOSED
-

© Srinivasan Seshan, 2004

L-4; 10-7-04

65

Observed TCP Problems



- Too many small packets
 - Silly window syndrome
 - Nagel's algorithm
- Initial sequence number selection
- Amount of state maintained

© Srinivasan Seshan, 2004

L-4; 10-7-04

66

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 $\min(\text{MSS}, \text{RecvBuffer}/2)$

© Srinivasan Seshan, 2004

L-4; 10-7-04

67

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

© Srinivasan Seshan, 2004

L-4; 10-7-04

68

Why is Selecting ISN Important?



- Suppose machine X selects ISN based on predictable sequence
- Fred has .rhosts to allow login to X from Y
- Evil Ed attacks
 - Disables host Y – denial of service attack
 - Make a bunch of connections to host X
 - Determine ISN pattern a guess next ISN
 - Fake pkt1: [<src Y><dst X>, guessed ISN]
 - Fake pkt2: desired command

© Srinivasan Seshan, 2004

L -4; 10-7-04

69

Time Wait Issues



- Web servers not clients close connection first
 - Established → Fin-Waits → Time-Wait → Closed
 - Why would this be a problem?
- Time-Wait state lasts for $2 * \text{MSL}$
 - MSL is should be 120 seconds (is often 60s)
 - Servers often have order of magnitude more connections in Time-Wait

© Srinivasan Seshan, 2004

L -4; 10-7-04

70

TCP Extensions



- Implemented using TCP options
 - Timestamp
 - Protection from sequence number wraparound
 - Large windows

© Srinivasan Seshan, 2004

L -4; 10-7-04

71

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 → < MSL!
 - 1.2Gbps: 28 seconds
- Use timestamp to distinguish sequence number wraparound

© Srinivasan Seshan, 2004

L -4; 10-7-04

72

Large Windows



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

Maximum Segment Size (MSS)



- Exchanged at connection setup
 - Typically pick MTU of local link
- What all does this effect?
 - Efficiency
 - Congestion control
 - Retransmission
- Path MTU discovery
 - Why should MTU match MSS?