

15-744: Computer Networking

L-4 TCP



This Lecture: Congestion Control



- Congestion Control
- Assigned Reading
 - [Chiu & Jain] Analysis of Increase and Decrease Algorithms for Congestion Avoidance in Computer Networks
 - [Jacobson and Karels] Congestion Avoidance and Control

2

Introduction to TCP



- Communication abstraction:
 - Reliable
 - Ordered
 - Point-to-point
 - Byte-stream
 - Full duplex
 - Flow and congestion controlled
- Protocol implemented entirely at the ends
 - Fate sharing
- Sliding window with cumulative acks
 - Ack field contains last in-order packet received
 - Duplicate acks sent when out-of-order packet received

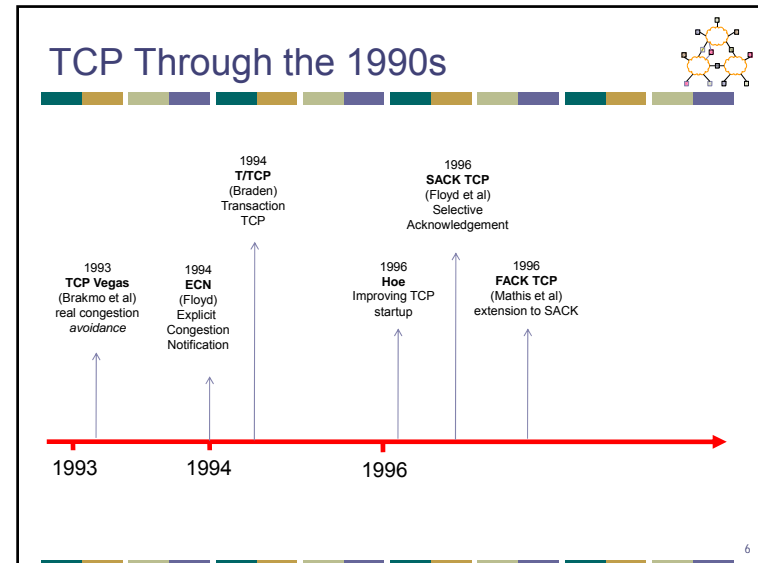
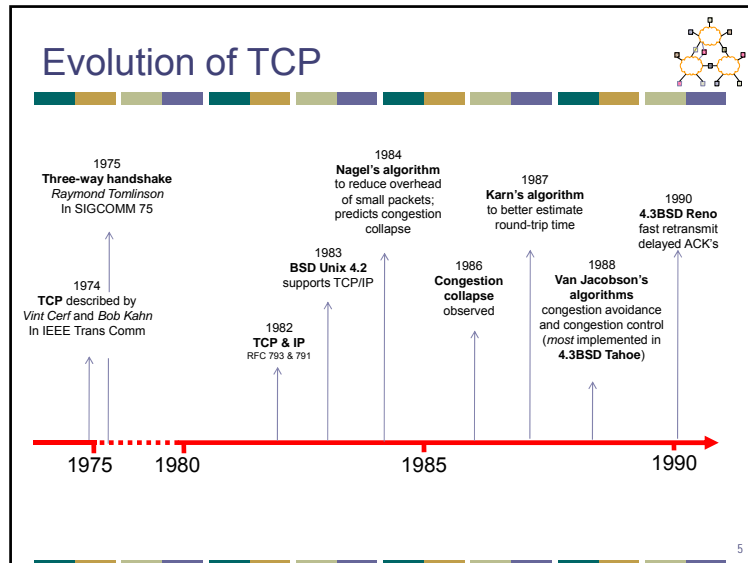
3

Key Things You Should Know Already

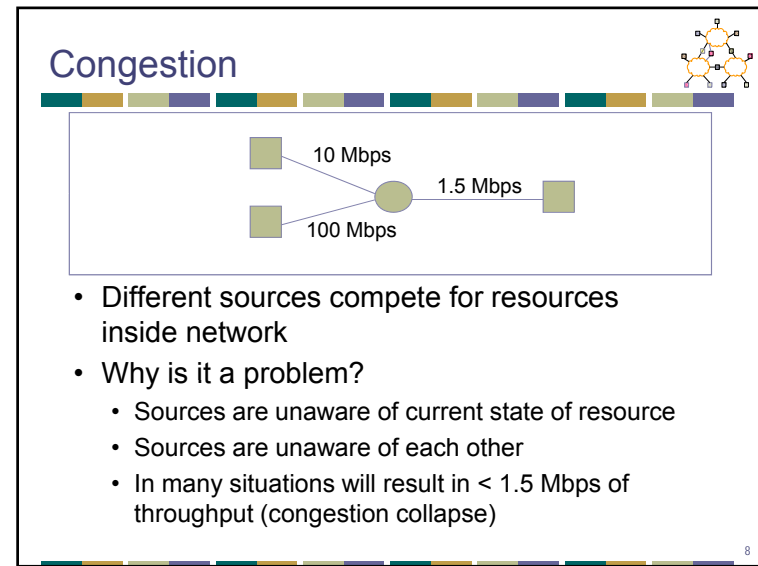


- Port numbers
- TCP/UDP checksum
- Sliding window flow control
 - Sequence numbers
- TCP connection setup
- TCP reliability
 - Timeout
 - Data-driven

4



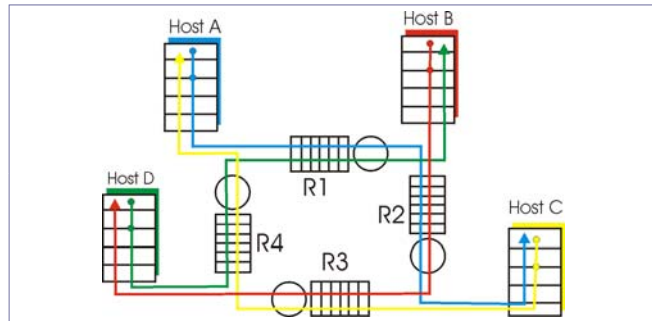
- ## Overview
- Congestion sources and collapse
 - Congestion control basics
 - TCP congestion control
 - TCP modeling



Causes & Costs of Congestion

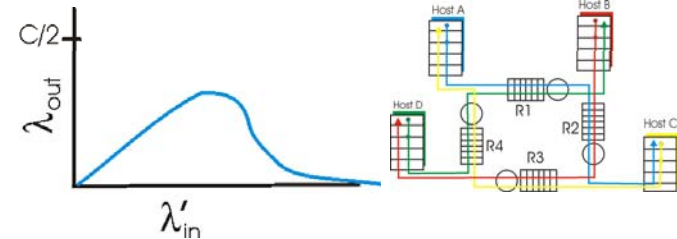


- Four senders – multihop paths **Q:** What happens as rate increases?
- Timeout/retransmit



9

Causes & Costs of Congestion



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

10

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

11

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

12

Where to Prevent Collapse?



- Can end hosts prevent problem?
 - Yes, but must trust end hosts to do right thing
 - E.g., sending host must adjust amount of data it puts in the network based on detected congestion
- Can routers prevent collapse?
 - No, not all forms of collapse
 - Doesn't mean they can't help
 - Sending accurate congestion signals
 - Isolating well-behaved from ill-behaved sources

13

Congestion Control and Avoidance



- A mechanism which:
 - Uses network resources efficiently
 - Preserves fair network resource allocation
 - Prevents or avoids collapse
- Congestion collapse is not just a theory
 - Has been frequently observed in many networks

14

Overview



- Congestion sources and collapse
- Congestion control basics
- TCP congestion control
- TCP modeling

15

Objectives



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

16

Basic Control Model



- Let's assume window-based control
- Reduce window when congestion is perceived
 - 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
- Increase window otherwise
 - Probe for available bandwidth – how?

17

Linear Control



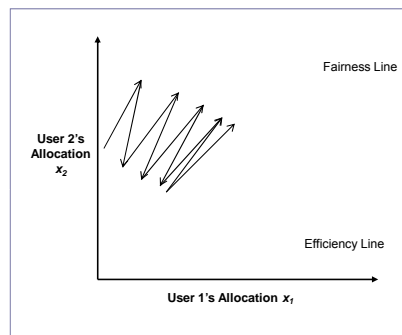
- Many different possibilities for reaction to congestion and probing
 - Examine simple linear controls
 - $Window(t + 1) = a + b \cdot Window(t)$
 - Different a_i/b_i for increase and a_d/b_d for decrease
- Supports various reaction to signals
 - Increase/decrease additively
 - Increased/decrease multiplicatively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

18

Phase plots



- Simple way to visualize behavior of competing connections over time

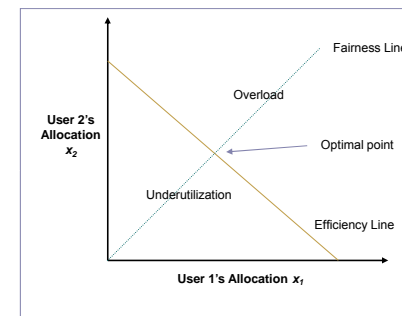


19

Phase plots



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

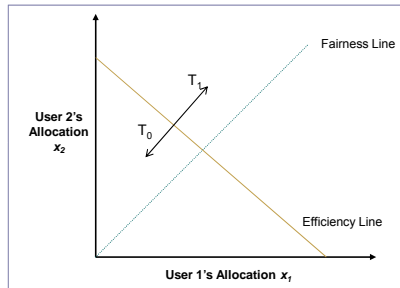


20

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

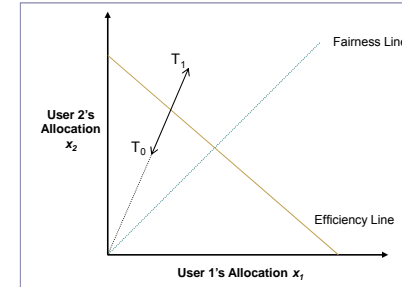


21

Multiplicative Increase/Decrease

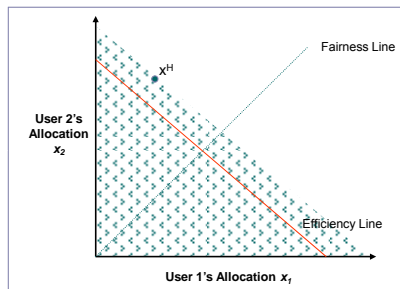


- Both X_1 and X_2 increase by the same factor over time
 - Extension from origin – constant fairness



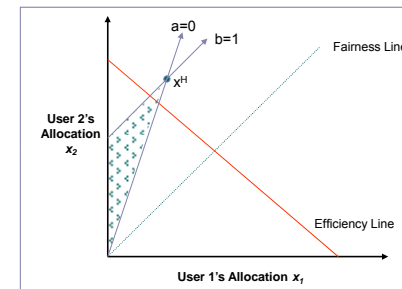
22

Convergence to Efficiency



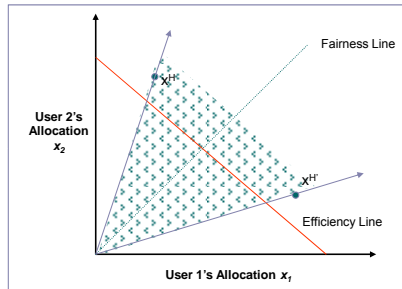
23

Distributed Convergence to Efficiency



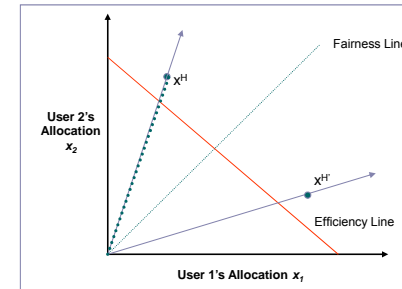
24

Convergence to Fairness



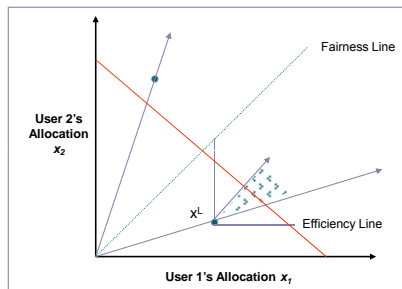
25

Convergence to Efficiency & Fairness



26

Increase



27

Constraints



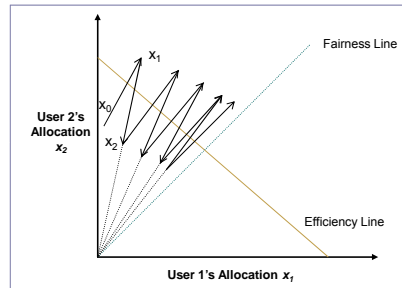
- Distributed efficiency
 - I.e., $\sum \text{Window}(t+1) > \sum \text{Window}(t)$ during increase
 - $a_i > 0$ & $b_i \geq 1$
 - Similarly, $a_d < 0$ & $b_d \leq 1$
- Must never decrease fairness
 - a & b 's must be ≥ 0
 - $a_i/b_i > 0$ and $a_d/b_d \geq 0$
- Full constraints
 - $a_d = 0$, $0 \leq b_d < 1$, $a_i > 0$ and $b_i \geq 1$

28

What is the Right Choice?



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



29

Overview



- Congestion sources and collapse
- Congestion control basics
- **TCP congestion control**
- TCP modeling

30

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

31

TCP Congestion Control - Solutions



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

32

TCP Congestion Control



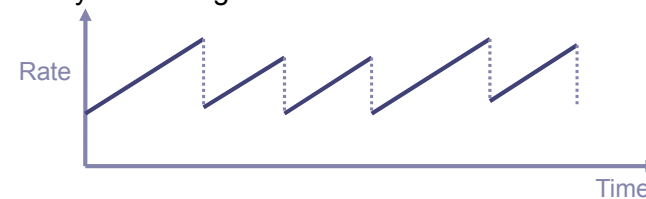
- Changes to TCP motivated by ARPANET congestion collapse
- Basic principles
 - AIMD
 - Packet conservation
 - Reaching steady state quickly
 - ACK clocking

33

AIMD



- 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



34

Implementation Issue



- Operating system timers are very coarse – how to pace packets out smoothly?
- Implemented using a congestion window that limits how much data can be in the network.
 - TCP also keeps track of how much data is in transit
- Data can only be sent when the amount of outstanding data is less than the congestion window.
 - The amount of outstanding data is increased on a “send” and decreased on “ack”
 - $(\text{last sent} - \text{last acked}) < \text{congestion window}$
- Window limited by both congestion and buffering
 - $\text{Sender's maximum window} = \text{Min}(\text{advertised window}, \text{cwnd})$

35

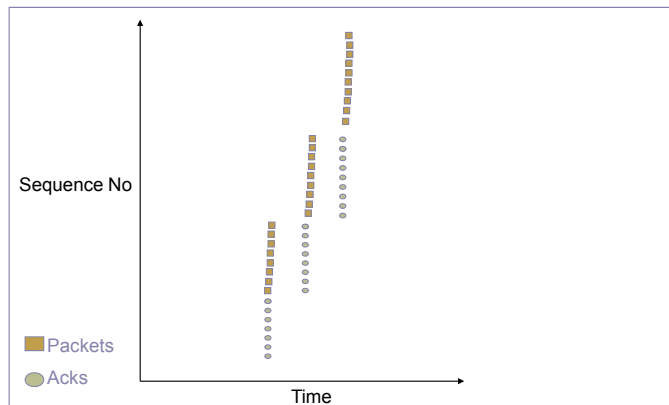
Congestion Avoidance



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

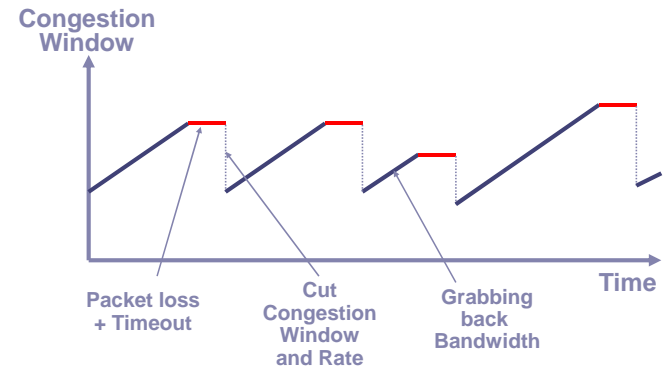
36

Congestion Avoidance Sequence Plot



37

Congestion Avoidance Behavior



38

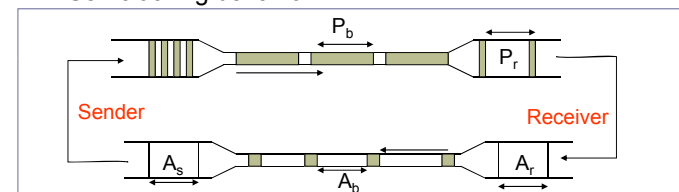
Packet Conservation

- 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
 - Helps provide stability
- Need to eliminate spurious retransmissions
 - Accurate RTO estimation
 - Better loss recovery techniques (e.g. fast retransmit)

39

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



40

Reaching Steady State



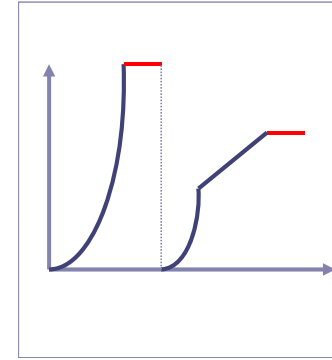
- Doing AIMD is fine in steady state but slow...
- How does TCP know what is a good initial rate to start with?
 - Should work both for a CDPD (10s of Kbps or less) and for supercomputer links (10 Gbps and growing)
- Quick initial phase to help get up to speed (slow start)

41

Slow Start Packet Pacing

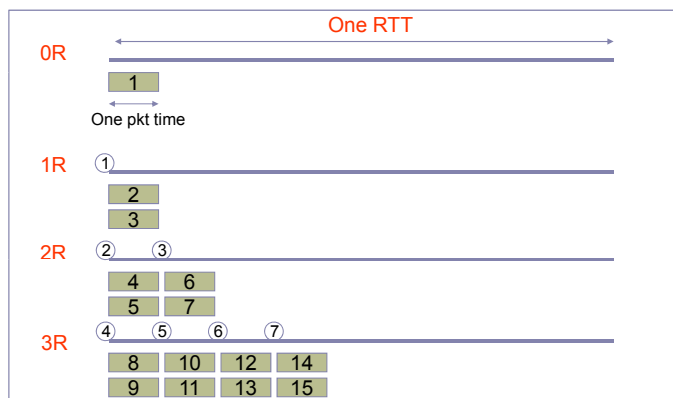


- How do we get this clocking behavior to start?
 - Initialize $cwnd = 1$
 - Upon receipt of every ack, $cwnd = cwnd + 1$
- Implications
 - Window actually increases to W in $RTT * \log_2(W)$
 - Can overshoot window and cause packet loss



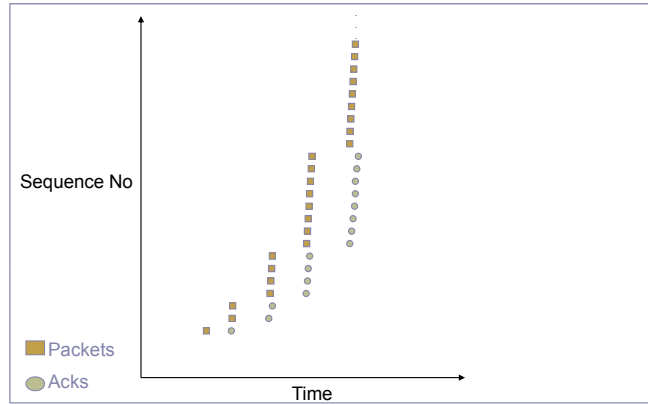
42

Slow Start Example



43

Slow Start Sequence Plot



44

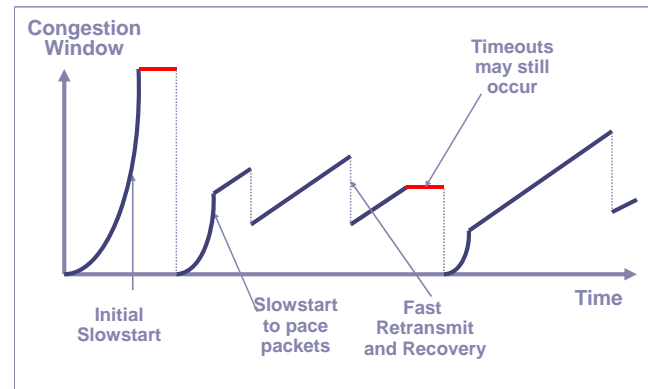
Return to Slow Start



- 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

45

TCP Saw Tooth Behavior



46

How to Change Window



- When a loss occurs have W packets outstanding
- New $cwnd = 0.5 * cwnd$
 - How to get to new state?

47

Fast Recovery



- Each duplicate ack notifies sender that single packet has cleared network
- When $< cwnd$ packets are outstanding
 - Allow new packets out with each new duplicate acknowledgement
- Behavior
 - Sender is idle for some time – waiting for $\frac{1}{2}$ $cwnd$ worth of dupacks
 - Transmits at original rate after wait
 - Ack clocking rate is same as before loss

48

Simple TCP Model



- Some additional assumptions
 - Fixed RTT
 - No delayed ACKs
- In steady state, TCP loses 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
 - $BW = MSS * \text{avg window} / RTT =$
 - $MSS * (W + W/2) / (2 * RTT)$
 - $.75 * MSS * W / RTT$

53

Simple Loss Model



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

54

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 $1/\text{sqrt}(p)$ behavior it is ok
 - But is this really true?

55

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!**

56

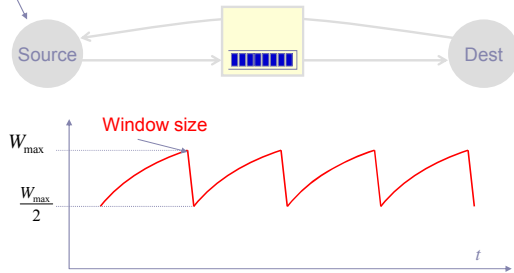
TCP Congestion Control



Only W packets may be outstanding

Rule for adjusting W

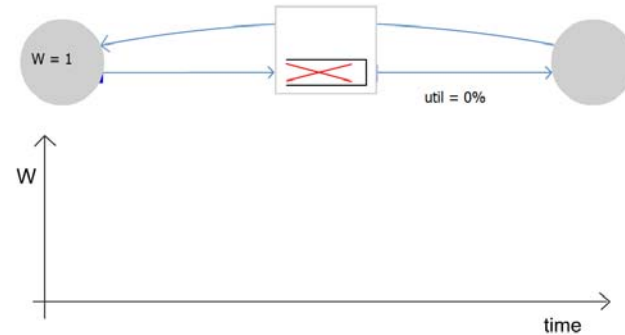
- If an ACK is received: $W \leftarrow W+1/W$
- If a packet is lost: $W \leftarrow W/2$



57

Single TCP Flow

Router **without** buffers



58

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

59

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

60

TCP Performance

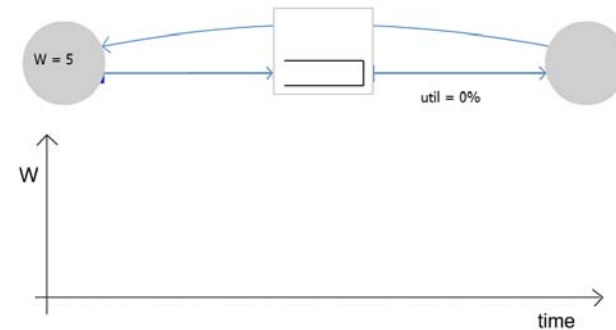


- If we have a large router queue \rightarrow can get 100% utilization
 - But, router queues can cause large delays
- How big does the queue need to be?
 - Windows vary from $W \rightarrow W/2$
 - Must make sure that link is always full
 - $W/2 > RTT * BW$
 - $W = RTT * BW + Qsize$
 - Therefore, $Qsize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

61

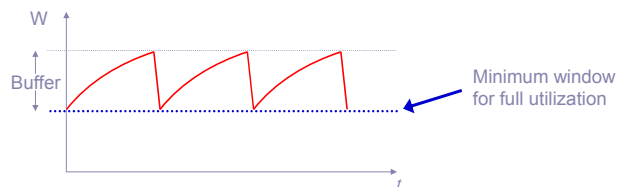
Single TCP Flow

Router with large enough buffers for full link utilization



62

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 $2T * C$
 - This is the origin of the rule-of-thumb

63

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.

64

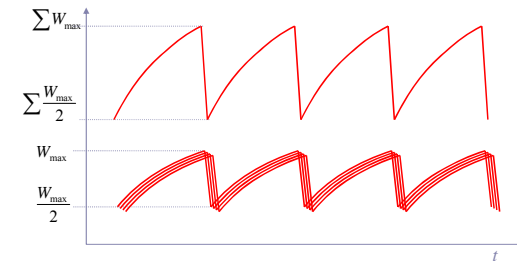
Rule-of-thumb



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

65

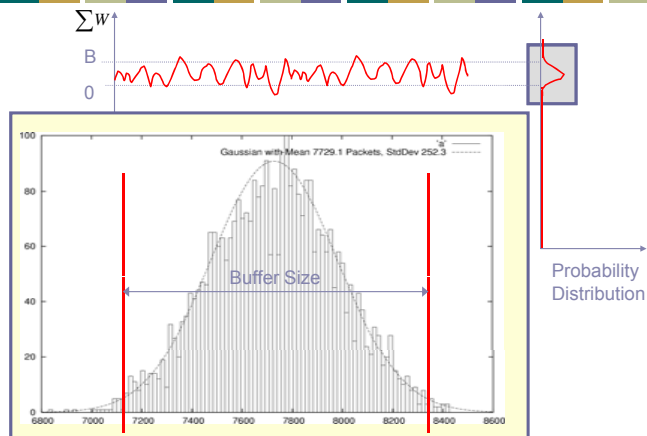
If flows are synchronized



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

66

If flows are not synchronized



67

Central Limit Theorem

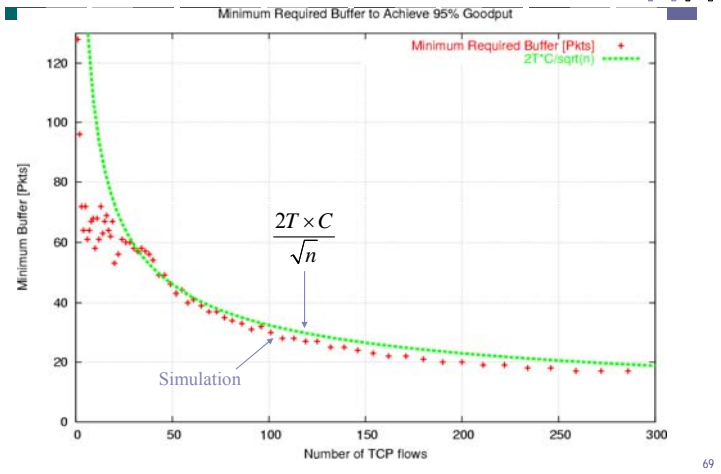


- CLT tells us that the more variables (Congestion Windows of Flows) we have, the narrower the Gaussian (Fluctuation of sum of windows)
- Width of Gaussian decreases with $\frac{1}{\sqrt{n}}$
- Buffer size should also decrease with $\frac{1}{\sqrt{n}}$

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

68

Required buffer size



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

Next Lecture

- Fair-queueing
- Assigned reading
 - [Demers, Keshav, Shenker] Analysis and Simulation of a Fair Queueing Algorithm
 - [Stoica, Shenker, Zhang] Core-Stateless Fair Queueing: Achieving Approximately Fair Bandwidth Allocations in High Speed Networks*