

# CSE/EE 461

## Congestion Control

---

Or, When to retransmit

---

### This Lecture

---

- Focus
  - How do we decide when to retransmit?
- Topics
  - RTT estimation
  - Karn/Partridge algorithm
  - Jacobson/Karels algorithm

Application
Presentation
Session
<b>Transport</b>
Network
Data Link
Physical

## How Fast Should a Sender Send?

---

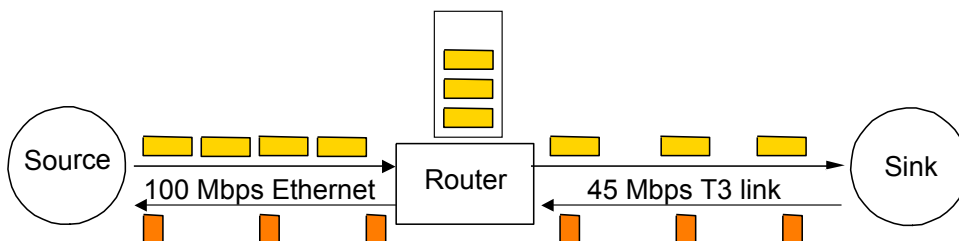
- No greater than the *bottleneck bandwidth*
- Only as fast as packets can be received
  - Network limitation
  - Receiver limitation
- Advertised receiver windows address receiver limitation
- But how to know what the network limitation is?

3

---

## TCP is “Self-Clocking”

---

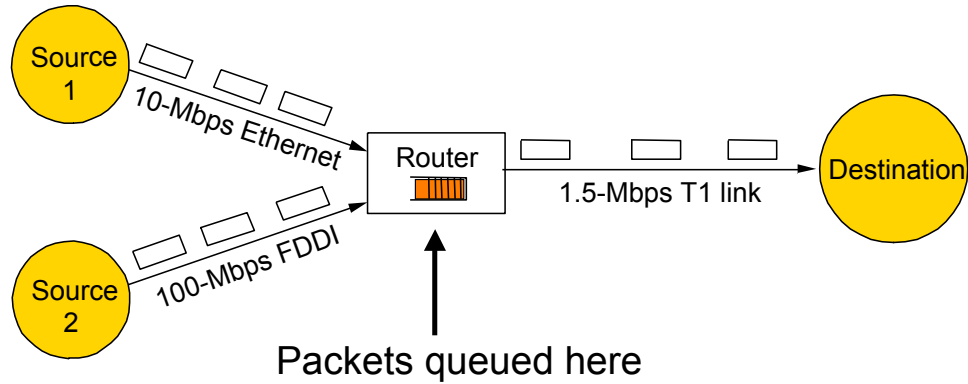


- Observation: acks pace transmissions at approximately the bottleneck rate

4

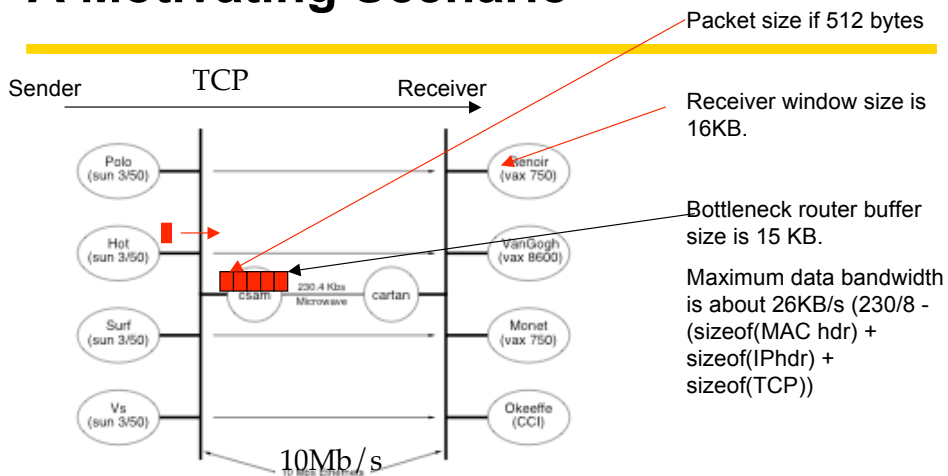
---

# Congestion from Multiple Sources

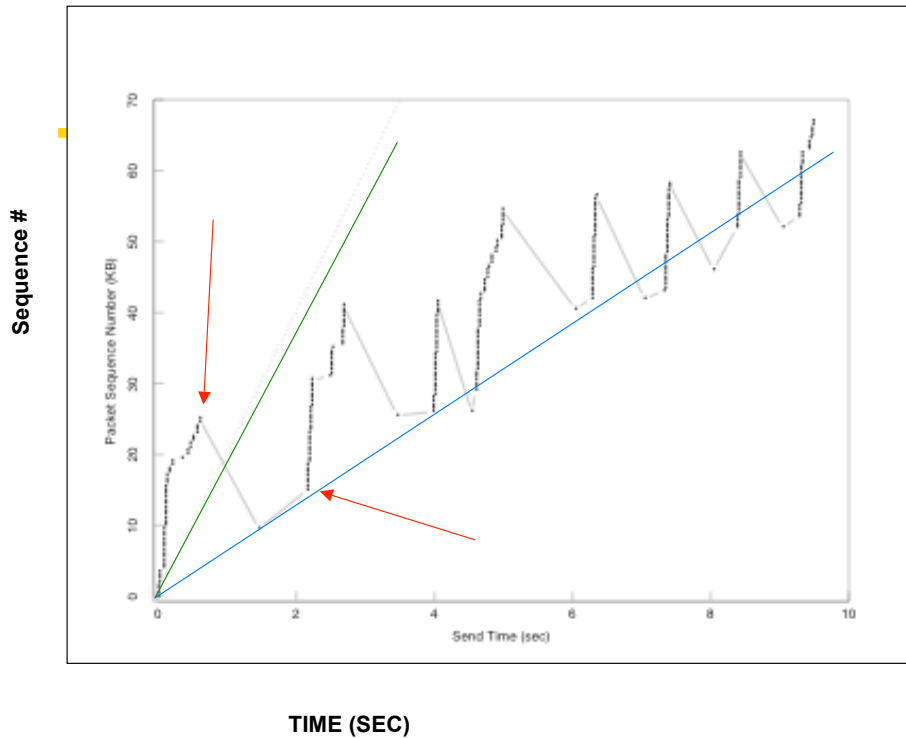


- Buffers at routers used to absorb bursts when input rate > output
- Loss (drops) occur when sending rate is persistently > drain rate

# A Motivating Scenario



Test setup to examine the interaction of multiple, simultaneous TCP conversations sharing a bottleneck link. 1 MByte transfers (2048 512-data-byte packets) were initiated 3 seconds apart from four machines at LBL to four machines at UCB, one conversation per machine pair (the dotted lines above show the pairing). All traffic went via a 230.4 Kbps link connecting IP router **csam** at LBL to IP router **cartan** at UCB. The microwave link queue can hold up to 50 packets. Each connection was given a window of 16 KB (32 512-byte packets). Thus any two connections could overflow the available buffering and the four connections exceeded the queue capacity by 160%.



7

## Congestion Collapse

- In the limit, early retransmissions lead to congestion collapse
  - Sending more packets into the network when it is overloaded exacerbates the problem of congestion
  - Network stays busy but very little useful work is being done
- This happened in real life ~1987
  - Led to Van Jacobson's TCP algorithms, which form the basis of **congestion control** in the Internet today [See "Congestion Avoidance and Control", SIGCOMM'88]
  - Observed 1000x bandwidth reduction between two hosts separated by 400 yards.
  - Led to researchers asking two questions:
    - Was TCP/IP misbehaving?
    - Could TCP/IP be "trained" to work better under 'abysmal network conditions'?

8

# 1988 Observations on Congestion Collapse

---

- Implementation, not the protocol, leads to collapse
- “Obvious” ways of doing things lead to non-obvious and undesirable results
  - “send eff-wind-size # packets, wait rtt, try again”
- Remedial algorithms achieve network stability by forcing the transport connection to obey a ‘packet conservation’ principle.
- For a connection in ‘equilibrium, that is, running stably with a full window of data in transit, the packet flow is “conservative”:
  - a new packet is not put into the network until an old packet leaves.

9

---

## If only...

---

- We knew RTT and Current Router Queue Size,
  - Then we would send  $\text{MIN}(\text{Router Queue Size}, \text{Effective Window Size})$
  - And not resend a packet until it had been sent RTT ago.
- But we don’t know these things, so we have to figure them out.
- And they may change dynamically due to other data sources

10

---

# The Packet Conservation Concept

---

- Connection flow should obey a “conservation of packets” principle
  - For a connection in equilibrium (running stably with a full window of data in transit):
    - New packet is not put into the network until an old packet leaves
- 1. The connection must reach equilibrium.
  - Hurry up and stabilize!
  - When things get wobbly, put on the brakes and reconsider
- 2. A sender must not inject a new packet before an old packet has exited.
  - A packet “exits” when the receiver picks it up or it is lost
  - Ack or packet timeout signals that a packet has “exited.”
  - Acks are easy to detect.
  - Good timeouts are harder.... All about estimating RTT.
- 3. Equilibrium is lost because of resource contention along the way.
  - Competing stream appears/disappears

11

---

## Packet Conservation

---

1. The connection must reach equilibrium.

12

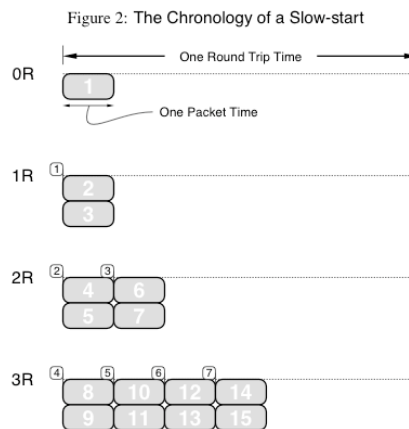
---

# 1. Getting to Equilibrium -- Slow Start

- Goal
  - Quickly determine the appropriate window size
- Strategy
  - Exponential probe
- Tactics
  - Introduce *congestion\_window* (*cwnd*)
  - When starting off, set *cwnd* to 1
  - For each ack received, add 1 to *cwnd*
  - When sending, send the minimum of receiver's advertised window and *cwnd*
  - On timeout, cut *cwnd* in half

13

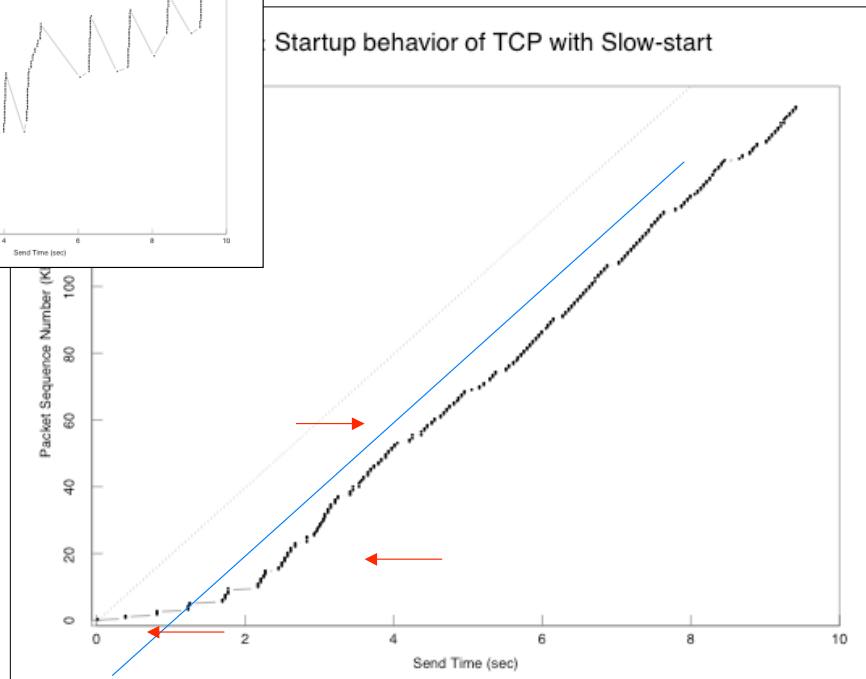
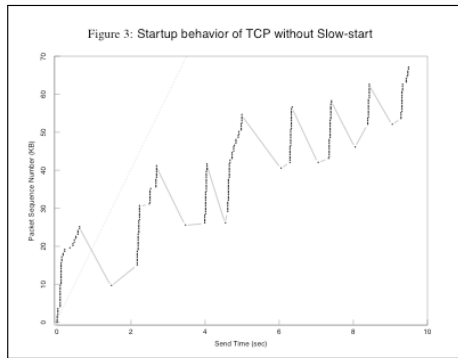
**Cwnd doubles every RTT;**  
**Opening the window of size**  
**W takes time  $(RTT)\log_2 W$ .**



The horizontal direction is time. The continuous time line has been chopped into one-round-trip-time pieces stacked vertically with increasing time going down the page. The grey, numbered boxes are packets. The white numbered boxes are the corresponding acks. As each ack arrives, two packets are generated: one for the ack (the ack says a packet has left the system so a new packet is added to take its place) and one because an ack opens the congestion window by one packet. It may be clear from the figure why an add-one-packet-to-window policy opens the window exponentially in time.

- Will not transmit at more than twice the max bw, and for no more than RTT.
  - (bw delay product)

14



15

- 
2. A sender must not inject a new packet before an old packet has exited.

16



## 2. Packet Injection. Estimating RTTs

---

- Do not inject a new packet until an old packet has left.
  - 1. Ack tells us that an old packet has left.
  - 2. Timeout expires tells us also.
    - *Gotta estimate RTT properly.*
- Strategy 1: Fixed RTT.
  - Simple, but probably wrong. (certainly not adaptive)
- Strategy 2: Estimate based on past behavior.

Tactic 0: Mean

Tactic 1: Mean with exponential decay

Tactic 2: Tactic 1 + *safety margin*

safety margin based on current estimate of error in Tactic 1

17

---

## Simple Estimator (RFC793)

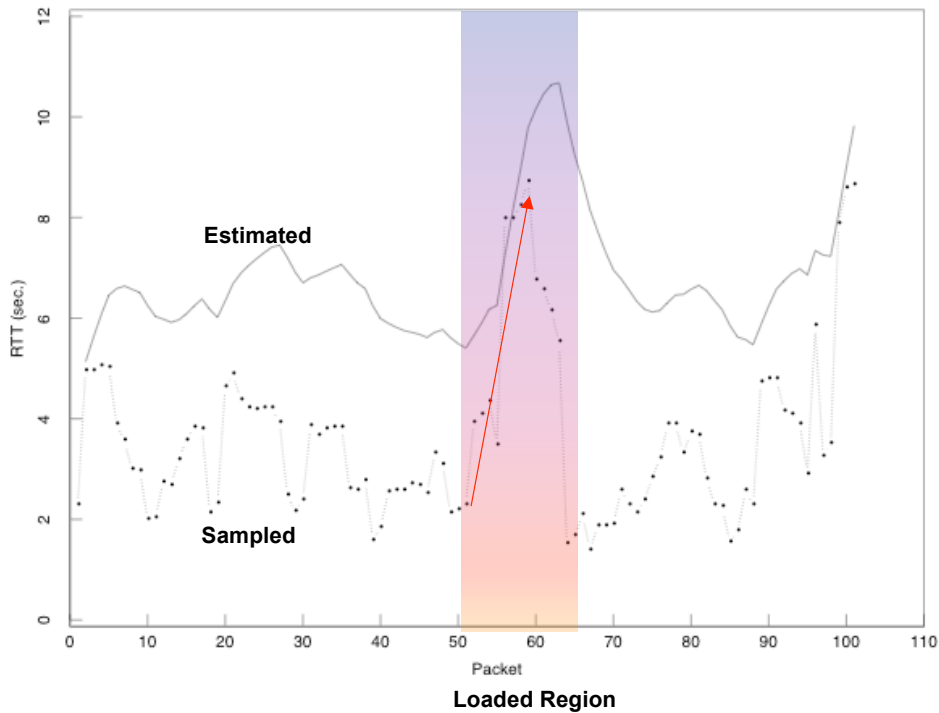
---

- $EstimatedRTT = (1-g)(EstimatedRTT) + g(SampledRTT)$
- Exponentially-weighted moving average ( $0 \leq g \leq 1$ )
- Smooths the samples with a gain of  $g$ 
  - Big  $g$  can be jerky
  - Small  $g$  can be slow to respond
- Stable is better than precise. Typically,  $g = .1$  or  $.2$ 
  - Insight: Respond conservatively to good news, and we already have a way to tell really bad news.
- Conservatively set timeout to small multiple of the estimate in order to account for variance
$$Timeout = 2(EstimatedRTT)$$
- Better to wait “too long” than not long enough. (Why?)

18

---

Figure 5: Performance of an RFC793 retransmit timer



## Bad Estimators and the Bad Things They Do

- Problem:
  - Variance in RTTs gets large as network gets loaded
  - So an average RTT isn't a good predictor when we need it most
    - Time out too soon, unnecessarily drop another packet onto the network.
    - Timing out too soon occurs during load increase
      - if we time out when load increases but packet not yet lost, then we'll inject another packet onto the network which will increase load, which will cause more timeouts, which will increase load, until we actually starting dropping packets!

# Jacobson/Karels Algorithm

---

- EstimatedRTT + “safety margin”
  - large variation in EstimatedRTT --> larger safety margin
- First, estimate how much SampledRTT deviates from EstimatedRTT
  - $DevRTT = (1-b) * DevRTT + b * |SampledRTT - EstimatedRTT|$ 
    - typically,  $b = .25$
- Then, set timeout interval as:
  - $Timeout = EstimatedRTT + k * DevRTT$
  - $k$  is generally set to 4
- Thus,
  - Timeout is close to EstimatedRTT when the Estimate is good,
  - Timeout quickly moves away from EstimatedRTT (4x!) when the Estimate is bad.

21

Figure 5: Performance of an RFC793 retransmit timer

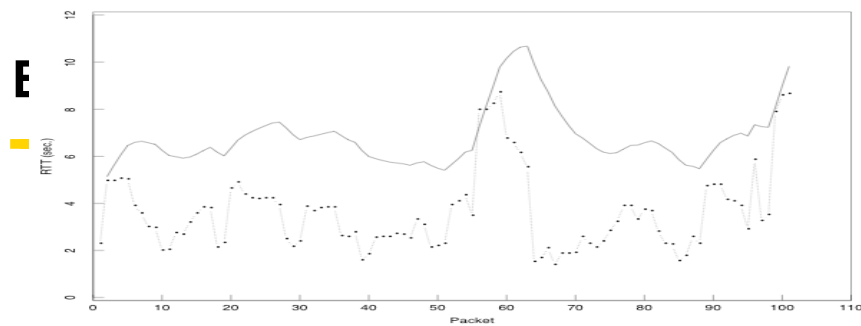
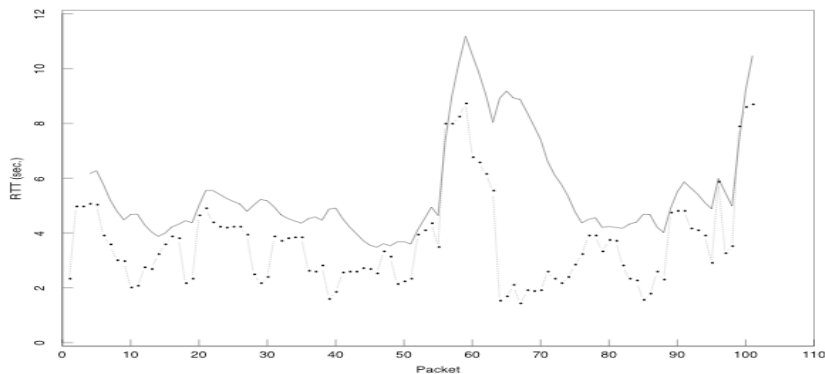


Figure 6: Performance of a Mean+Variance retransmit timer

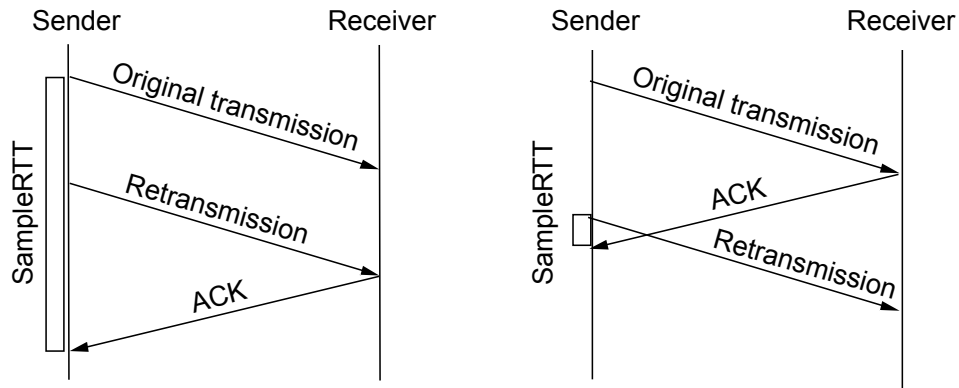


22

## Footnote: Karn/Partridge Algorithm

---

- Problem: RTT for retransmitted packets ambiguous



- Solution: Don't measure RTT for retransmitted packets and do not relax backed off timeout until valid RTT measurements

23

---

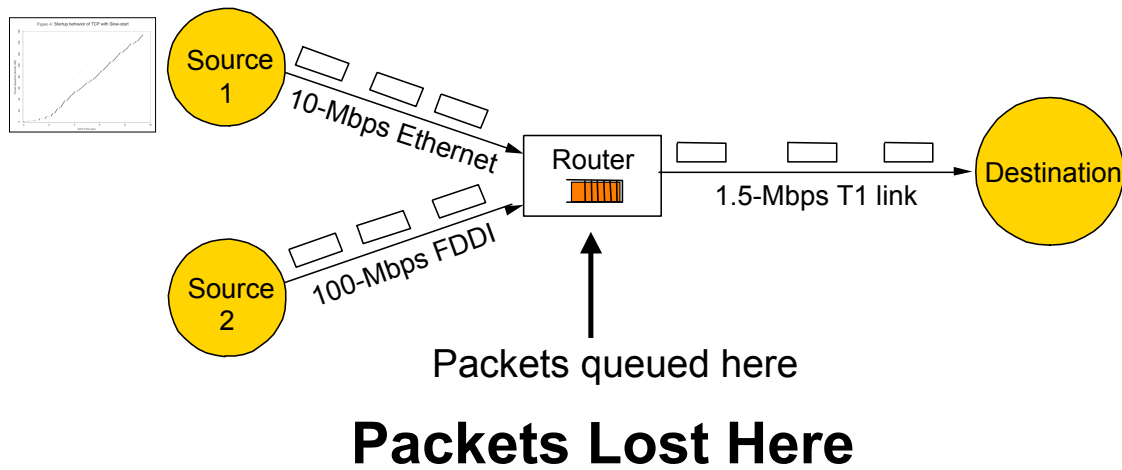
3. Equilibrium is lost because of resource contention along the way.

(room for more? Room for less?)

24

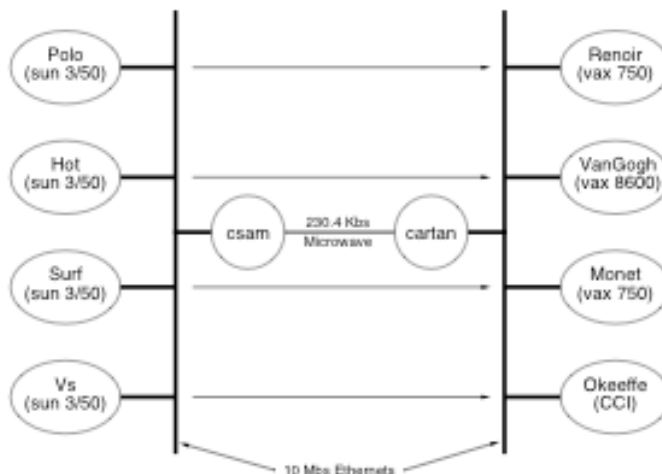
---

# Congestion from Multiple Sources



# In Real Life -- 30KB/s link

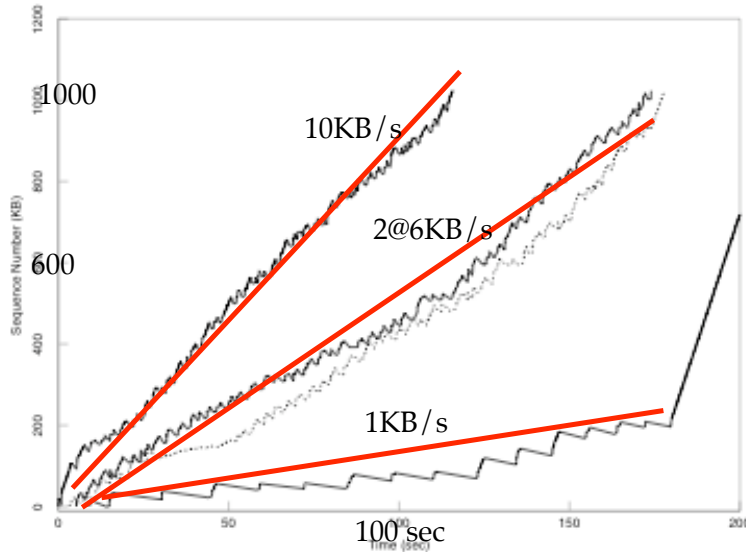
Figure 7: Multiple conversation test setup



## Four Simultaneous Streams

---

Figure 8: Multiple, simultaneous TCPs with no congestion avoidance



27

---

## Implicit Signals from the Network

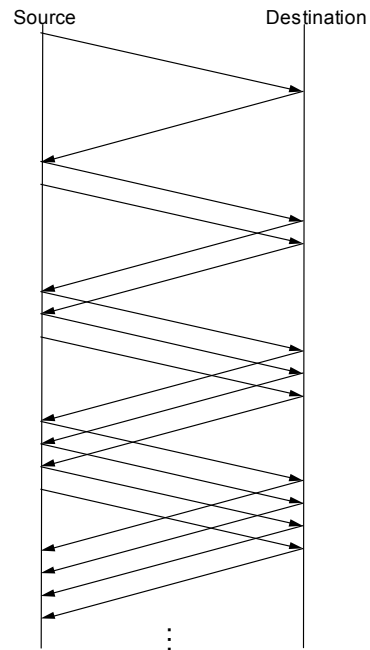
- The network is not saturated: *Send even more*
- The network is saturated: *Send less*
  
- ACK signals that the network is not saturated.
- A Lost packet (no ACK) signals that the network is saturated
  - Assumption here??
- Leads to a simple strategy:
  - On each ack, increase *congestion window (additive increase)*
  - On each lost packet, decrease congestion window (*multiplicative decrease*)
- Why increase slowly and decrease quickly?
  - *Respond to good news conservatively, but bad news aggressively*

28

## AIMD (Additive Increase/Multiplicative Decrease)

---

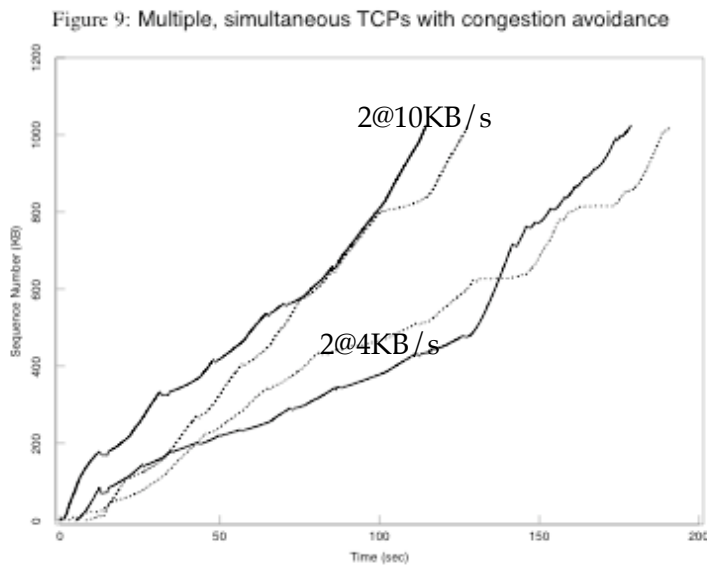
- How to adjust probe rate?
- Increase slowly while we believe there is bandwidth
  - Additive increase per RTT
  - $Cwnd += 1 \text{ packet} / \text{RTT}$
- Decrease quickly when there is loss (went too far!)
  - Multiplicative decrease
  - $Cwnd /= 2$



29

## With Additive Increase/Multiplicative Decrease

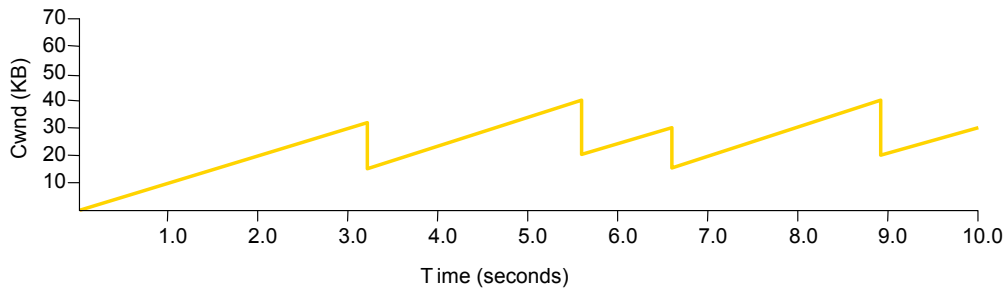
---



30

# TCP Sawtooth Pattern

---

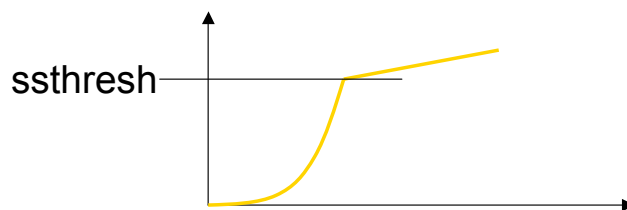


31

---

# Combining Slow Start and AIMD

---



- Slow start is used whenever the connection is not running with packets outstanding
  - initially, and after timeouts indicating that there's no data on the wire
- But we don't want to overshoot our ideal cwnd on next slow start, so remember the last cwnd that worked with no loss
  - $Ssthresh = cwnd \text{ after } cwnd / = 2 \text{ on loss}$
  - Switch to AIMD once cwnd passes ssthresh

32

---



## Key Concepts

---

- Packet conservation is a fundamental concept in TCP's congestion management
  - Get to equilibrium
    - *Slow Start*
  - Do nothing to get out of equilibrium
    - *Good RTT Estimate*
  - Adapt when equilibrium has been lost due to other's attempts to get to stay in equilibrium
    - *Additive Increase/Multiplicative Decrease*
- The Network Reveals Its Own Behavior

33

---

## Key Concepts (next level down)

---

- TCP probes the network for bandwidth, assuming that loss signals congestion
- The congestion window is managed to be additive increase / multiplicative decrease
  - It took fast retransmit and fast recovery to get there
- Slow start is used to avoid lengthy initial delays
  - Ramp up to near target rate and then switch to AIMD

34

---

# A Fast Algorithm for RTT Mean and Variation

---

- Let  $a$  = estimated round trip time,  $v$  = estimated error,  $g$  = gain ( $0 < g < 1$ ),  $m$  = new sampled round trip time
- $a = (1-g)a + gm$  // compute new estimate using gain
- $a = a + g(m-a)$  // rearrange terms:
  - $a$  is a prediction of next measurement, and  $(m-a)$  is the "error" in that prediction.
  - so, the new prediction is the old prediction plus some fraction of the prediction error.
  - The prediction error is the sum of two components:
    - $E_r$  = noise (random unpredictable effects like fluctuations in competing traffic)
    - $E_e$  = bad choice of  $a$
    - $a = a + g E_r + g E_e$ 
      - The term  $g E_e$  kicks  $a$  in the right direction towards the real estimate
      - The term  $g E_r$  kicks it off in the random direction
      - Over many samples, the random errors cancel each other so we get closer and closer to the real estimate
      - **But,  $g$  represents a compromise.**
        - » **Big 'g' means that we get a lot of value out of a prediction error, but it also means that the random errors introduce a lot of noise.**
      - **Since  $g E_e$  moves  $a$  in the right direction regardless of  $g$ , we're better off using a small  $g$  and waiting a bit longer to get a better estimate than to very quickly get a lousy estimate**
- Or,
  - $Err = (m - a)$  // Sampled Error
  - $a = a + g (Err)$  // Estimate of round trip time
  - $v = v + g(|Err| - v)$  // Estimate of error
- Not necessary to use same gain; in general want to force timer to go up