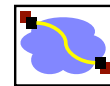# 15-441 Computer Networking
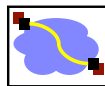
Lecture 19 – TCP Performance

---

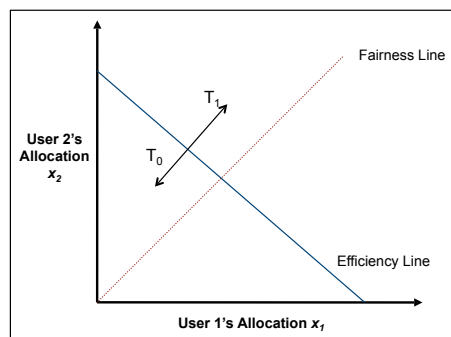# Outline

- TCP congestion avoidance
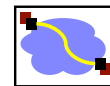
- TCP slow start
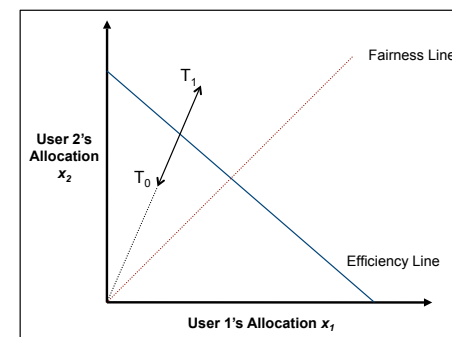
- TCP modeling

---

# Additive Increase/Decrease

- Both $X_1$ and $X_2$ increase/ decrease by the same amount over time
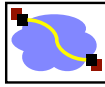  - Additive increase improves fairness and additive decrease reduces fairness



Fairness Line

$T_1$

$T_0$

User 2's Allocation $x_2$

Efficiency Line

User 1's Allocation $x_1$

---

# Muliplicative Increase/Decrease

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



Fairness Line

$T_1$

$T_0$

User 2's Allocation $x_2$
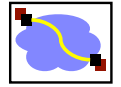
Efficiency Line

User 1's Allocation $x_1$

## What is the Right Choice?

- Constraints limit us to AIMD
  - Improves or keeps fairness constant at each step
  - AIMD moves towards optimal point



Fairness Line

User 2's Allocation $x_2$

$x_1$

$x_0$

$x_2$

Efficiency Line
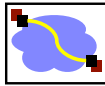
User 1's Allocation $x_1$

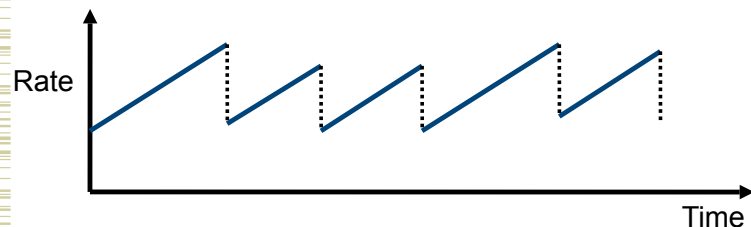---

## TCP Congestion Control

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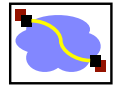  - Reaching steady state quickly
  - ACK clocking

---

## 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
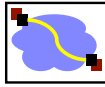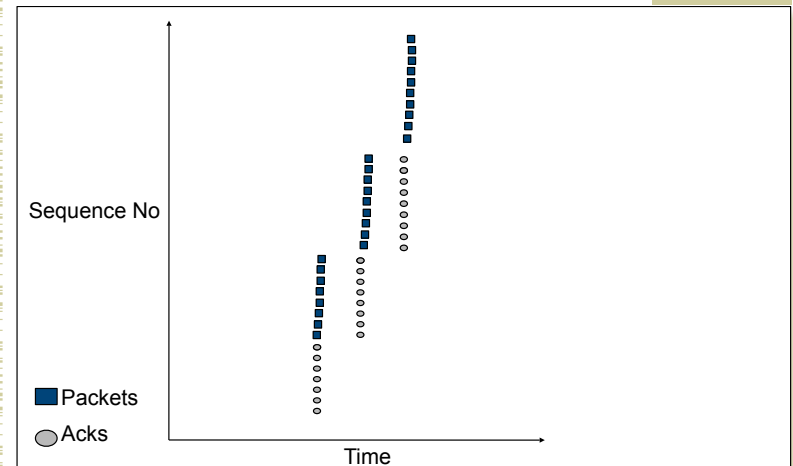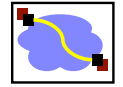


Rate

Time

---

## 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"
  - (last sent – last acked) < congestion window
- Window limited by both congestion and buffering
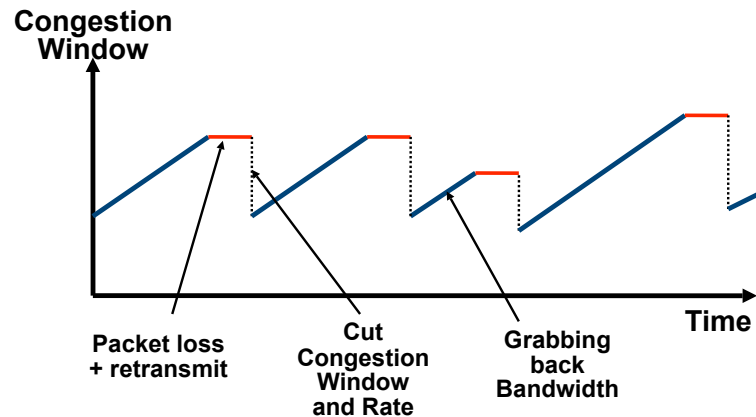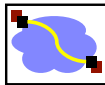  - Sender's maximum window = Min (advertised window, cwnd)

## Congestion Avoidance

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

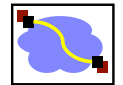## Congestion Avoidance Sequence Plot



Sequence No

- Packets
- Acks

Time

## Congestion Avoidance Behavior



Congestion Window

Time

Packet loss + retransmit

Cut Congestion Window and Rate
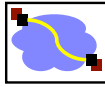
Grabbing back Bandwidth
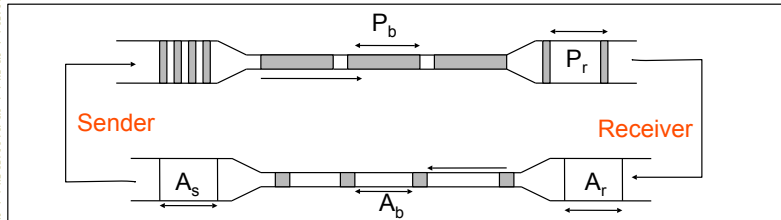
## Packet Conservation

- 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
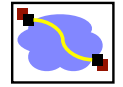  - Better loss recovery techniques (e.g. fast retransmit)

## 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
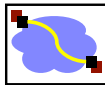  - Self-clocking behavior



$P_b$

$P_r$

Sender

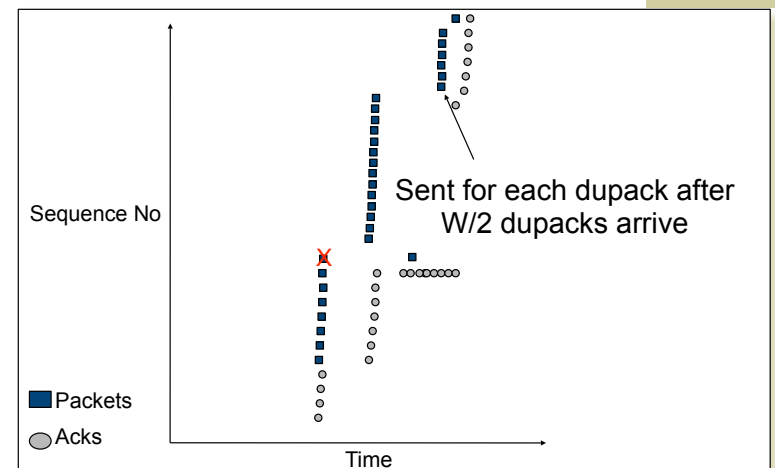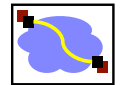Receiver

$A_s$

$A_b$

$A_r$

## How to Change Window

- When a loss occurs have W packets outstanding
- New cwnd = 0.5 * cwnd
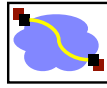  - How to get to new state without losing ack clocking?

## 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 ½ cwnd worth of dupacks
  - Transmits at original rate after wait
    - Ack clocking rate is same as before loss

## Fast Recovery



Sequence No

Sent for each dupack after W/2 dupacks arrive
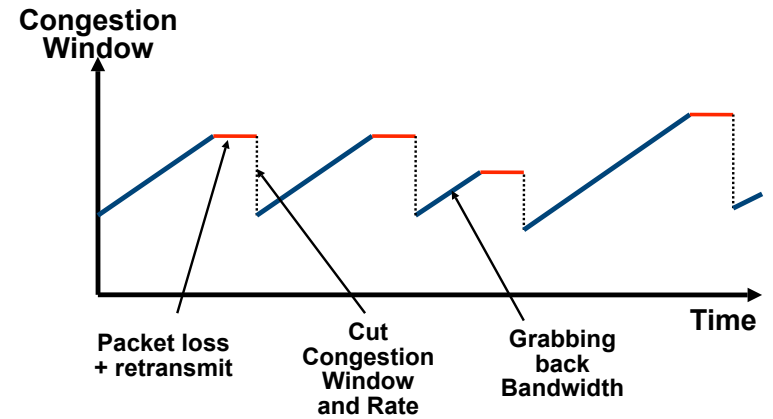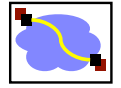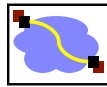
■ Packets

● Acks

Time

## Outline

- TCP congestion avoidance

- TCP slow start

- TCP modeling

## Congestion Avoidance Behavior

**Congestion Window**

**Time**

Packet loss + retransmit

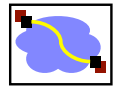Cut Congestion Window and Rate
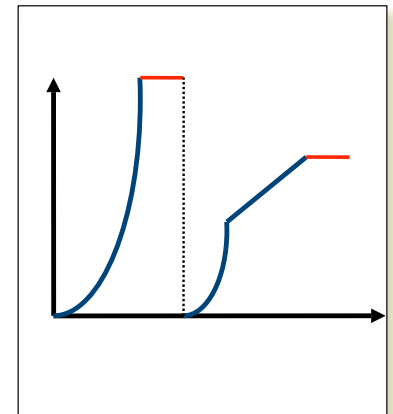
Grabbing back Bandwidth

## 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)
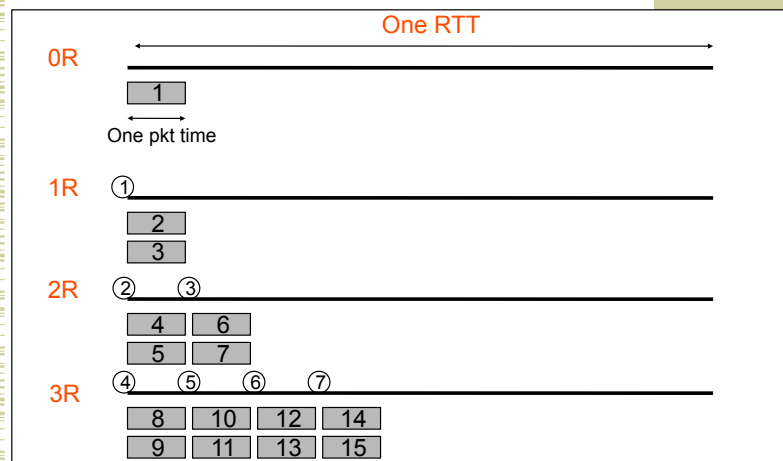- Quick initial phase to help get up to speed (slow start)
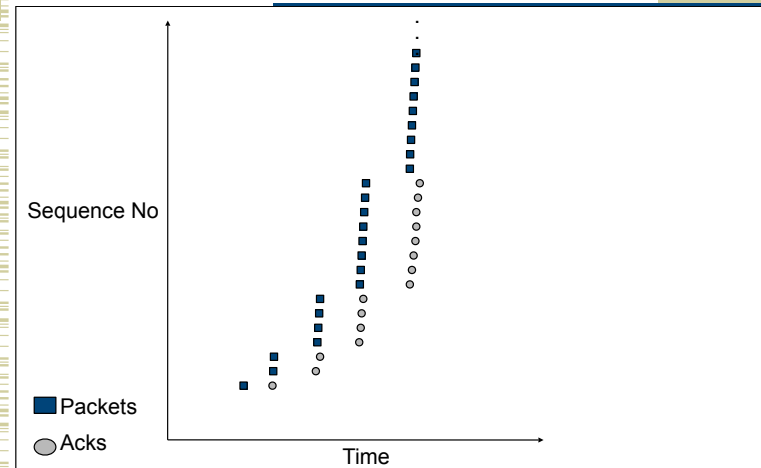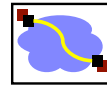
## Slow Start Packet Pacing

- 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

## Slow Start Example

One RTT

0R ————————————————————————
  [1]
  One pkt time

1R ①————————————————————————
  [2]
  [3]

2R ② ③————————————————————————
  [4] [6]
  [5] [7]

3R ④ ⑤ ⑥ ⑦————————————————————————
  [8] [10] [12] [14]
  [9] [11] [13] [15]

## Slow Start Sequence Plot



Sequence No
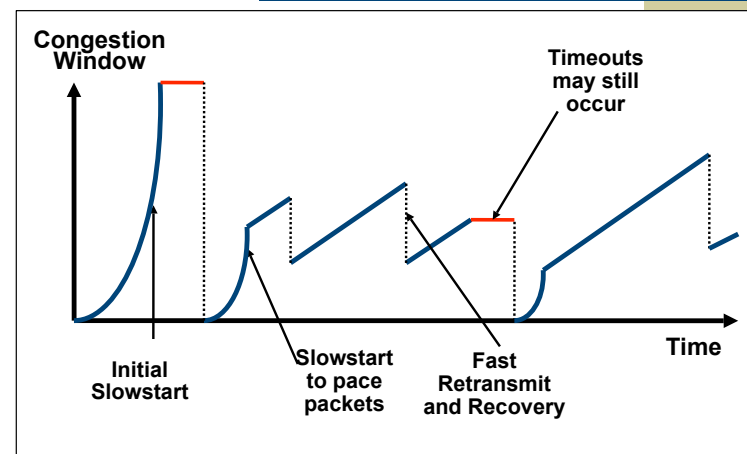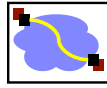
■ Packets
◯ Acks

Time

## 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 retransmission occurs set ssthresh to 0.5w
  - If cwnd < ssthresh, use slow start
  - Else use congestion avoidance

## TCP Saw Tooth Behavior



Congestion Window

Timeouts may still occur

Time

Initial Slowstart

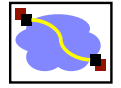Slowstart to pace packets
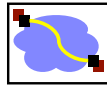
Fast Retransmit and Recovery

## Outline

- TCP congestion avoidance

- TCP slow start

- TCP modeling

## 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?
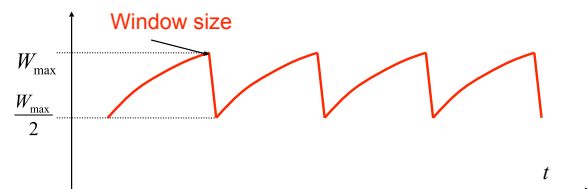  - Average utilization = 75%??
  - No…this is *not* right!

## TCP Congestion Control
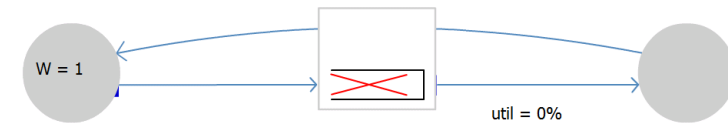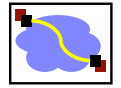
Only $W$ packets may be outstanding

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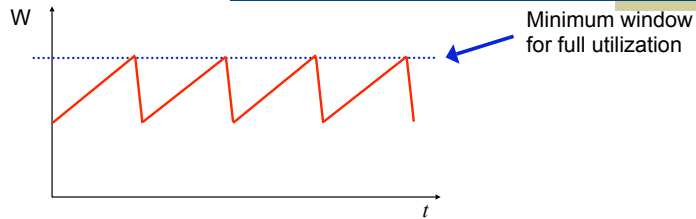- If a packet is lost: $W \leftarrow W/2$

Source

Dest

Window size

$W_{max}$

$\dfrac{W_{max}}{2}$

$t$

## Single TCP Flow
Router without buffers

$W = 1$

util = 0%

$W$
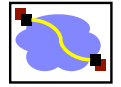
time

## Summary Unbuffered Link
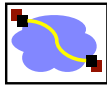
W

Minimum window for full utilization

t

- 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

## 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
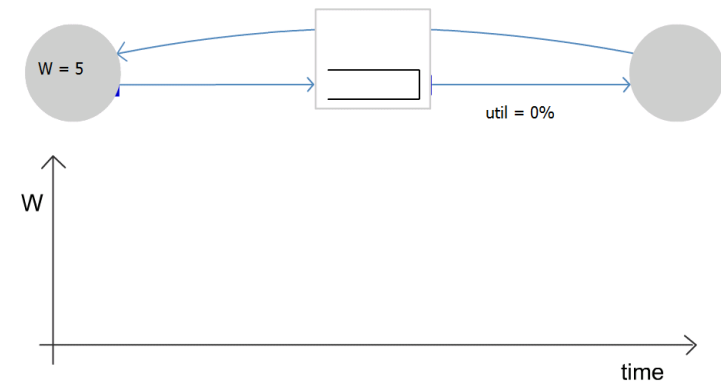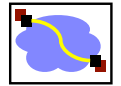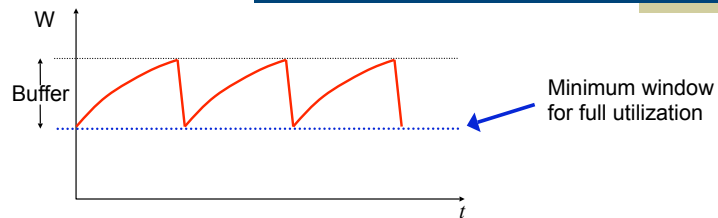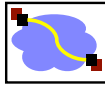    - If window is larger, packets sit in queue on bottleneck link

## 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 → 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

## Single TCP Flow
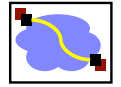Router with large enough buffers for full link utilization

W = 5

util = 0%

W

time

## Summary Buffered Link



W

Buffer

Minimum window
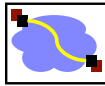for full utilization

t

- 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

## 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

## 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