



## 15-441 Computer Networking

### Lecture 18 – More TCP & Congestion Control

## Overview



- TCP congestion control
- TCP modern loss recovery
- TCP modeling

Lecture 18: 04-23-2004

2

## TCP Congestion Control



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

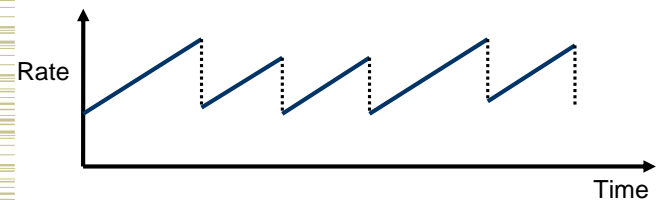
Lecture 18: 04-23-2004

3

## 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



Lecture 18: 04-23-2004

4

## 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
  - Sender's maximum window =  $\text{Min}(\text{advertised window}, \text{cwnd})$

Lecture 18: 04-23-2004

5

## Congestion Avoidance

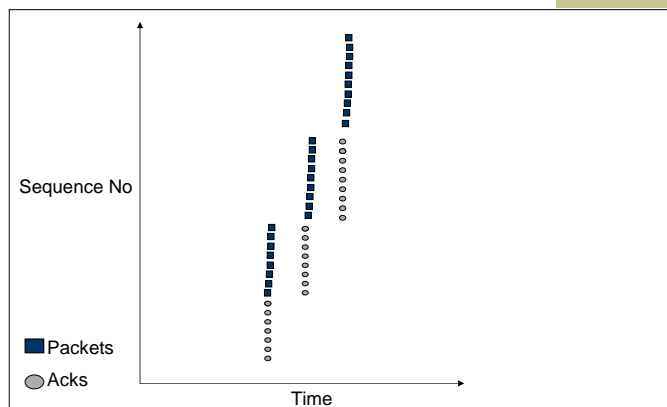


- If loss occurs when  $\text{cwnd} = W$ 
  - Network can handle  $0.5W \sim W$  segments
  - Set  $\text{cwnd}$  to  $0.5W$  (multiplicative decrease)
- Upon receiving ACK
  - Increase  $\text{cwnd}$  by  $(1 \text{ packet})/\text{cwnd}$ 
    - What is 1 packet?  $\rightarrow 1 \text{ MSS}$  worth of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

Lecture 18: 04-23-2004

6

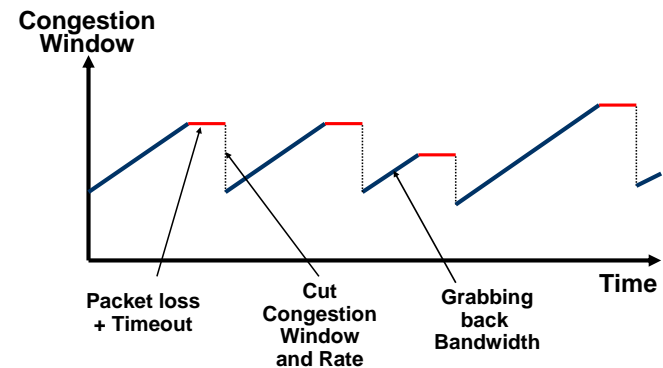
## Congestion Avoidance Sequence Plot



Lecture 18: 04-23-2004

7

## Congestion Avoidance Behavior



Lecture 18: 04-23-2004

8

## 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)

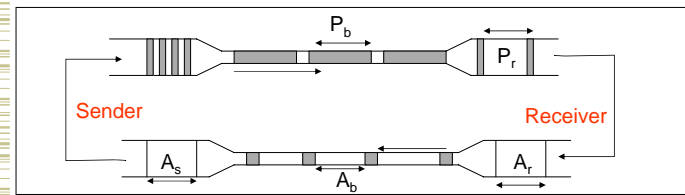
Lecture 18: 04-23-2004

9

## 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



Lecture 18: 04-23-2004

10

## 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)

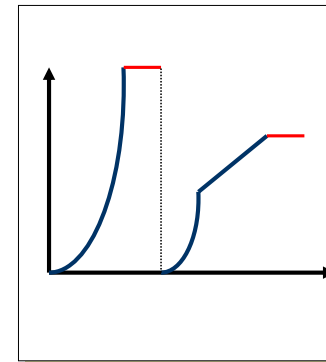
Lecture 18: 04-23-2004

11

## Slow Start Packet Pacing



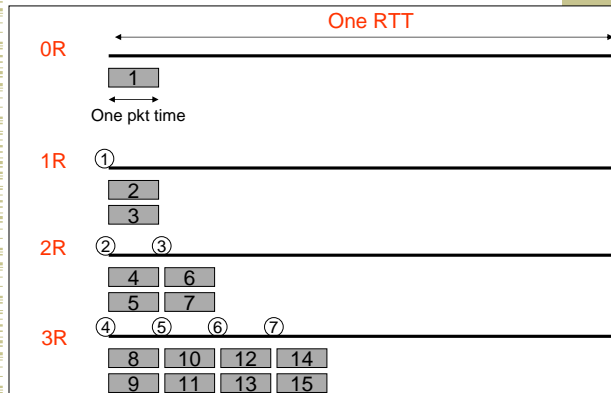
- How do we get this clocking behavior to start?
  - Initialize  $cwnd = 1$
  - Upon receipt of every ack,  $cwnd = cwnd + 1$
- Implications
  - Window actually increases to  $W$  in  $RTT * \log_2(W)$
  - Can overshoot window and cause packet loss



Lecture 18: 04-23-2004

12

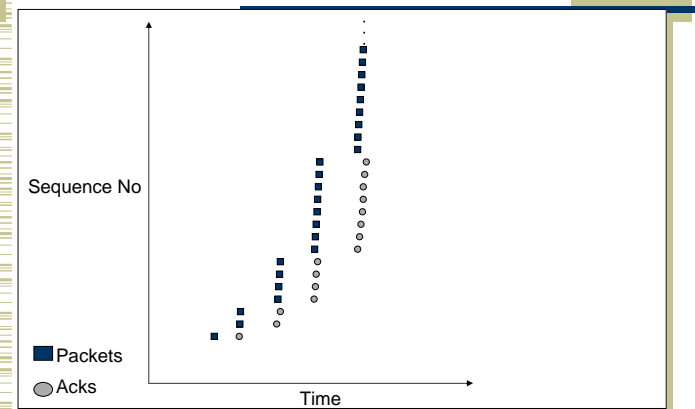
## Slow Start Example



Lecture 18: 04-23-2004

13

## Slow Start Sequence Plot



Lecture 18: 04-23-2004

14

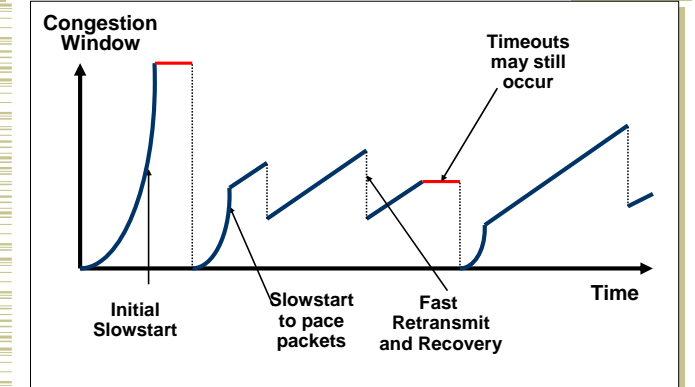
## Return to Slow Start

- 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

Lecture 18: 04-23-2004

15

## TCP Saw Tooth Behavior



Lecture 18: 04-23-2004

16

## Overview

- TCP congestion control
- TCP modern loss recovery
- TCP modeling

Lecture 18: 04-23-2004

17

## TCP Flavors

- Tahoe, Reno, Vegas
- TCP Tahoe (distributed with 4.3BSD Unix)
  - Original implementation of Van Jacobson's mechanisms (VJ paper)
  - Includes:
    - Slow start
    - Congestion avoidance
    - Fast retransmit

Lecture 18: 04-23-2004

18

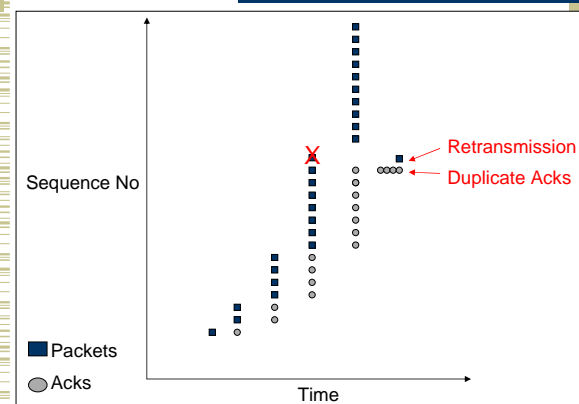
## Fast Retransmit

- What are duplicate acks (dupacks)?
  - Repeated acks for the same sequence
- When can duplicate acks occur?
  - Loss
  - Packet re-ordering
  - Window update – advertisement of new flow control window
- Assume re-ordering is infrequent and not of large magnitude
  - Use receipt of 3 or more duplicate acks as indication of loss
  - Don't wait for timeout to retransmit packet

Lecture 18: 04-23-2004

19

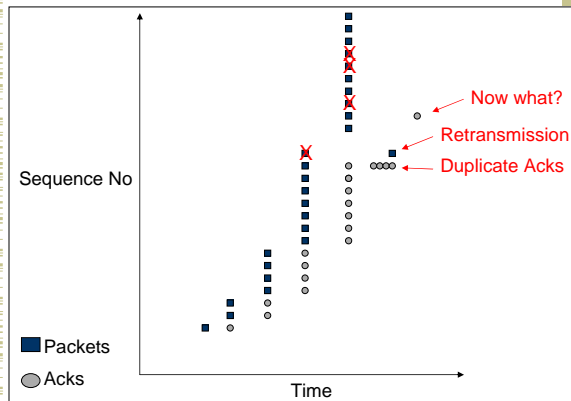
## Fast Retransmit



Lecture 18: 04-23-2004

20

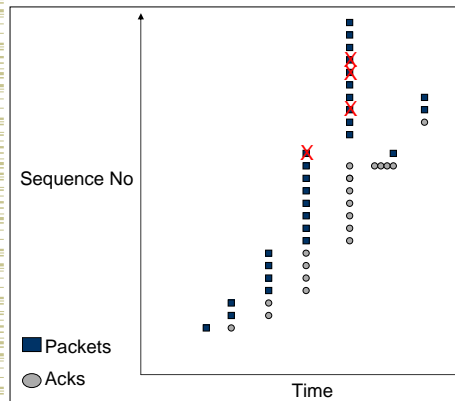
## Multiple Losses



Lecture 18: 04-23-2004

21

## Tahoe



Lecture 18: 04-23-2004

22

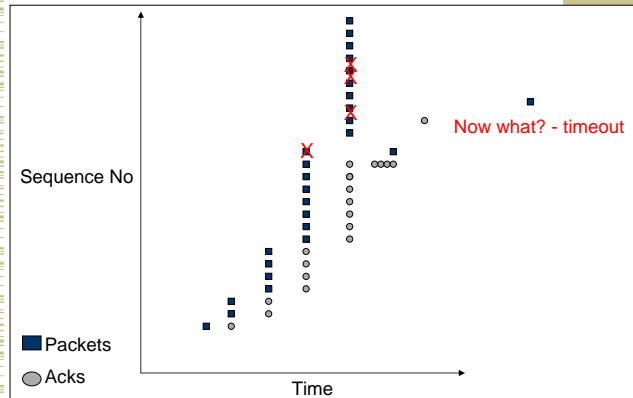
## TCP Reno (1990)

- 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 18: 04-23-2004

23

## Reno



Lecture 18: 04-23-2004

24

## NewReno

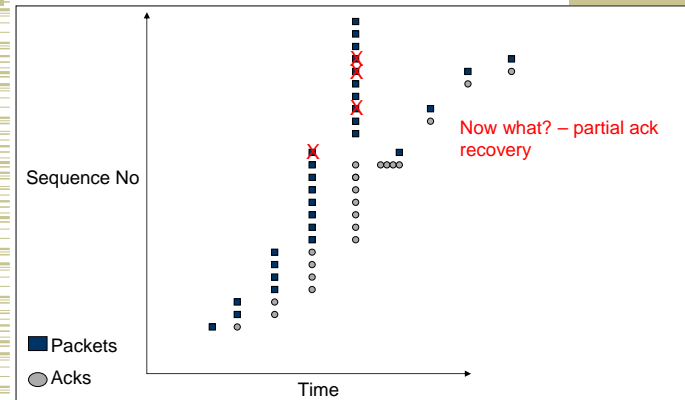


- The ack that arrives after retransmission (partial ack) could indicate that a second loss occurred
- When does NewReno timeout?
  - When there are fewer than three dupacks for first loss
  - When partial ack is lost
- How fast does it recover losses?
  - One per RTT

Lecture 18: 04-23-2004

25

## NewReno



Lecture 18: 04-23-2004

26

## SACK

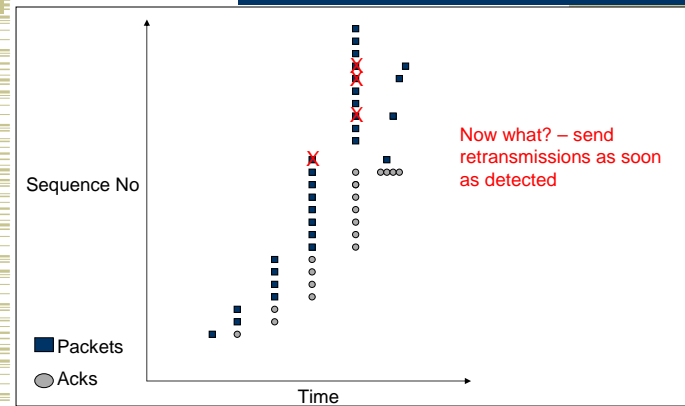


- Basic problem is that cumulative acks provide little information
  - Ack for just the packet received
    - What if acks are lost? → carry cumulative also
    - This technique is not used
  - Bitmask of packets received
    - Selective acknowledgement (SACK)
- When to retransmit?
  - Still need to deal with reordering → wait for out of order by 3pkts

Lecture 18: 04-23-2004

27

## SACK



Lecture 18: 04-23-2004

28

## Performance Issues

- Timeout  $\gg$  fast retransmit
  - Need 3 dupacks/sacks
  - Not great for small transfers
    - Don't have 3 packets outstanding
  - What are real loss patterns like?
- How to deal with reordering?

Lecture 18: 04-23-2004

29

## How to Change Window

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

Lecture 18: 04-23-2004

30

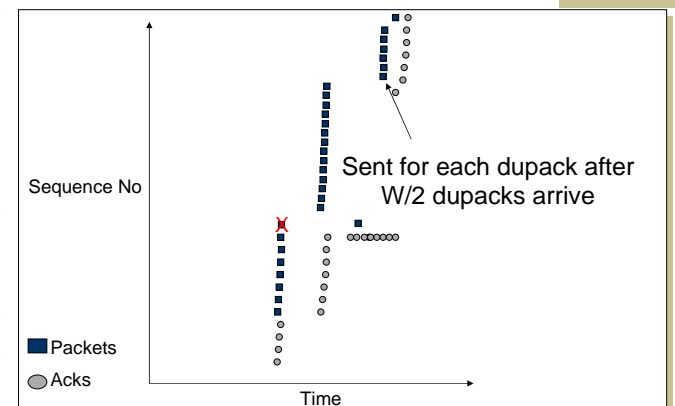
## Fast Recovery

- 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

Lecture 18: 04-23-2004

31

## Fast Recovery



Lecture 18: 04-23-2004

32



## Overview



- TCP congestion control
- TCP modern loss recovery
- TCP modeling

Lecture 18: 04-23-2004

33

## 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!

Lecture 18: 04-23-2004

34

## 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

Lecture 18: 04-23-2004

35

## 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  $2 * RTT$

Lecture 18: 04-23-2004

36

## 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