

### Reliability Challenges



- Congestion related losses
- Variable packet delays
  - What should the timeout be?
- Reordering of packets
  - How to tell the difference between a delayed packet and a lost one?

Lecture 18: 03-22-2005

### TCP = Go-Back-N Variant



- Sliding window with cumulative acks
  - Receiver can only return a single "ack" sequence number to the sender.
  - · Acknowledges all bytes with a lower sequence number
  - · Starting point for retransmission
  - · Duplicate acks sent when out-of-order packet received
- But: sender only retransmits a single packet.
  - Reason???
    - · Only one that it knows is lost
    - Network is congested → shouldn't overload it
- Error control is based on byte sequences, not packets.
  - Retransmitted packet can be different from the original lost packet – Why?

### Round-trip Time Estimation



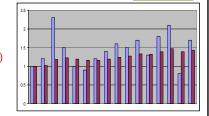
- Wait at least one RTT before retransmitting
- Importance of accurate RTT estimators:
  - · Low RTT estimate
    - · unneeded retransmissions
  - · High RTT estimate
    - poor throughput
- RTT estimator must adapt to change in RTT
  - But not too fast, or too slow!
- Spurious timeouts
  - "Conservation of packets" principle never more than a window worth of packets in flight

Lecture 18: 03-22-2005

Original TCP Round-trip Estimator



- Round trip times exponentially averaged:
  - New RTT =  $\alpha$  (old RTT) + (1 -  $\alpha$ ) (new sample)
  - Recommended value for α: 0.8 0.9
    - . 0.875 for most TCP's



- Retransmit timer set to (b \* RTT), where b = 2
  - · Every time timer expires, RTO exponentially backed-off
- Not good at preventing spurious timeouts
  - Why?

Lecture 18: 03-22-2005

### Jacobson's Retransmission Timeout



- Key observation:
  - At high loads round trip variance is high
- Solution:
  - Base RTO on RTT and standard deviation
    - RTO = RTT + 4 \* rttvar
  - new\_rttvar =  $\beta$  \* dev + (1- $\beta$ ) old\_rttvar
    - Dev = linear deviation
    - Inappropriately named actually smoothed linear deviation

Lecture 18: 03-22-2005

### RTT Sample Ambiguity Sample RTT Sample R

### Timestamp Extension



- Used to improve timeout mechanism by more accurate measurement of RTT
- When sending a packet, insert current timestamp into option
  - 4 bytes for timestamp, 4 bytes for echo
- Receiver echoes timestamp in ACK
  - Actually will echo whatever is in timestamp
- Removes retransmission ambiguity
  - Can get RTT sample on any packet

Lecture 18: 03-22-2005

### **Timer Granularity**



- Many TCP implementations set RTO in multiples of 200,500,1000ms
- Why?
  - Avoid spurious timeouts RTTs can vary quickly due to cross traffic
  - · Make timers interrupts efficient
- What happens for the first couple of packets?
  - Pick a very conservative value (seconds)

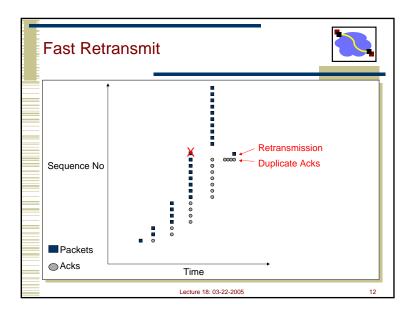
Lecture 18: 03-22-2005

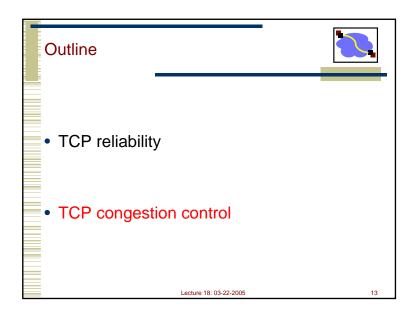
10

### 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





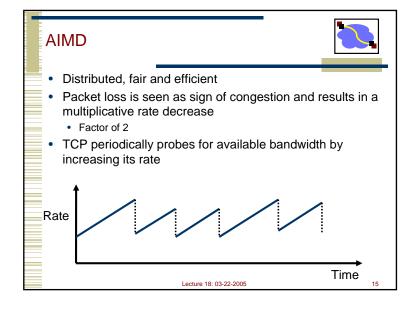
### TCP Congestion Control



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

Lecture 18: 03-22-2005

4.4



### Implementation Issue

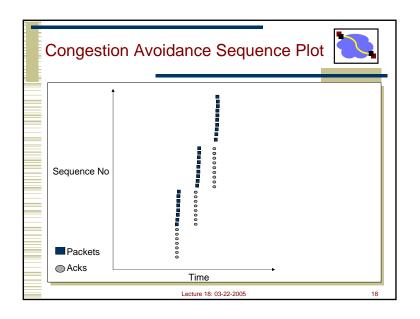


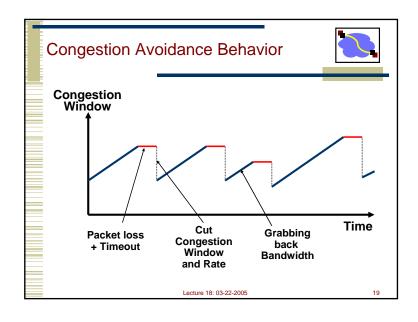
- Operating system timers are very coarse how to pace packets out smoothly?
- Implemented using a congestion window that limits how much data can be in the network.
  - TCP also keeps track of how much data is in transit
- Data can only be sent when the amount of outstanding data is less than the congestion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

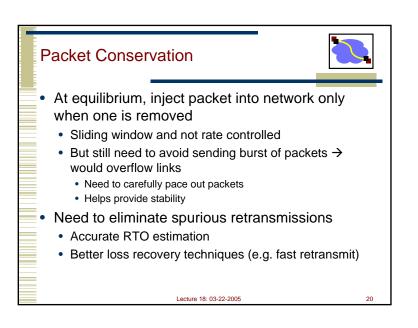
Lecture 18: 03-22-2005

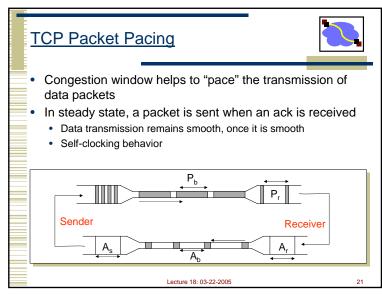
16

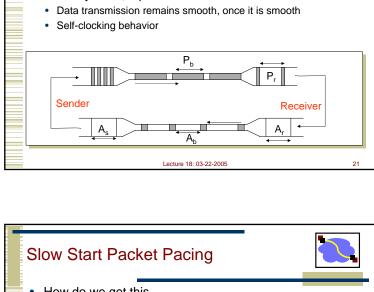
# Congestion Avoidance If loss occurs when cwnd = W Network can handle 0.5W ~ W segments Set cwnd to 0.5W (multiplicative decrease) Upon receiving ACK Increase cwnd by (1 packet)/cwnd What is 1 packet? → 1 MSS worth of bytes After cwnd packets have passed by → approximately increase of 1 MSS Implements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

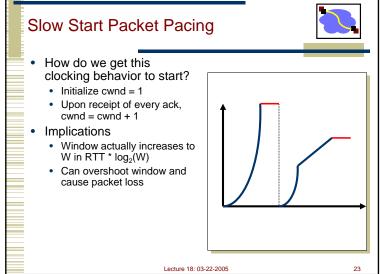


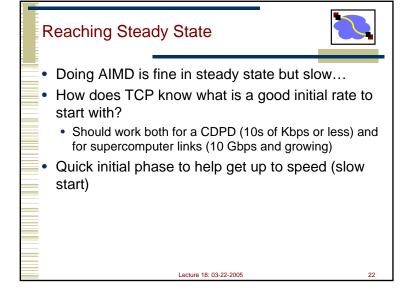


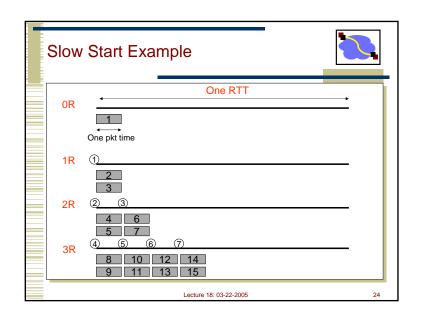


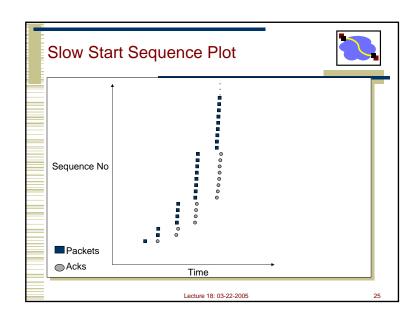


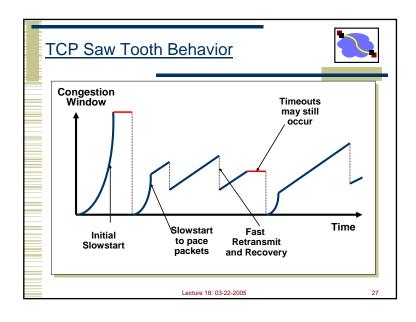












## Peturn to Slow Start If packet is lost we lose our self clocking as well Need to implement slow-start and congestion avoidance together When timeout occurs set ssthresh to 0.5w If cwnd < ssthresh, use slow start Else use congestion avoidance

### Important Lessons



- TCP timeout calculation → how is RTT estimated
- Modern TCP loss recovery
  - Why are timeouts bad?
  - How to avoid them? → e.g. fast retransmit
- How does TCP implement AIMD?
  - Sliding window, slow start & ack clocking
  - How to maintain ack clocking during loss recovery → fast recovery