# CS268: Beyond TCP **Congestion Control**

Kevin Lai February 4, 2003

#### **TCP** Problems

- When TCP congestion control was originally designed in 1988:
  - Maximum link bandwidth: 10Mb/s
  - Users were mostly from academic and government organizations (i.e., well-behaved)
- Almost all links were wired (i.e., negligible error rate)

2

4

- Thus, current problems with TCP: - High bandwidth-delay product paths

  - Selfish users
  - Wireless (or any high error links)

### **High Bandwidth-Delay Product Paths**

Motivation

- 10Gb/s links now common in Internet core as a result of Wave Division Multiplexing (WDM), link bandwidth 2x/9 months
- Some users have access to and need for 10Gb/s end-to-end · e.g., very large scientific, financial databases
- Satellite/Interplanetary links have a high delay
- Problems
  - slow start
- Additive increase, multiplicative decrease (AIMD)
- Congestion Control for High Bandwidth-Delay Product Networks. Dina Katabi, Mark Handley, and Charlie Rohrs. Proceedings on ACM Sigcomm 2002.

3

# **Slow Start**

- TCP throughput controlled by congestion window (cwnd) size
- In slow start, window increases exponentially, but may not be enough
- example: 10Gb/s, 200ms RTT, 1460B payload, assume no loss
  - Time to fill pipe: 18 round trips = 3.6 seconds
  - Data transferred until then: 382MB
  - Throughput at that time: 382MB / 3.6s = 850Mb/s - 8.5% utilization  $\rightarrow$  not very good
- Loose only one packet  $\rightarrow$  drop out of slow start into AIMD (even worse)



5

7

#### Proposed Solution:

Decouple Congestion Control from Fairness

Coupled because a *single* mechanism controls both <u>Example:</u> In TCP, Additive-Increase Multiplicative-Decrease (AIMD) controls both

How does decoupling solve the problem?

- 1. To control congestion: use  $\mathsf{MIMD}$  which shows fast response
- 2. <u>To control fairness:</u> use AIMD which converges to fairness

### **Characteristics of Solution**

- 1. Improved Congestion Control (in high bandwidth-delay & conventional environments):
  - Small queues
  - Almost no drops
- 2. Improved Fairness
- 3. Scalable (no per-flow state)
- 4. Flexible bandwidth allocation: min-max fairness, proportional fairness, differential bandwidth allocation,...

8





















### **XCP Summary**

- XCP
  - Outperforms TCP
  - Efficient for any bandwidth
  - Efficient for any delay
  - Scalable (no per flow state)
- Benefits of Decoupling
  - Use MIMD for congestion control which can grab/release large bandwidth quickly
  - Use AIMD for fairness which converges to fair bandwidth allocation

19

# Selfish Users

- Motivation
- Many users would sacrifice overall system efficiency for more performance
- Even more users would sacrifice fairness for more performance
- Users can modify their TCP stacks so that they can receive data from a normal server at an un-congestion controlled rate.
- Problem
  - How to prevent users from doing this?
  - General problem: How to design protocols that deal with lack of trust?
- OT ITUST? Congestion Control with a Misbehaving Receiver. Stefan Savage, Neal Cardwell, David Wetherall and Tom Anderson. ACM Computer Communications Review, pp. 71-78, v 29, no 5, October, 1999. Robust Congestion Signaling, David Vetherall, David Ely, Neil Spring, Stefan Savage and Tom Anderson. IEEE International Conference on Network Protocols, November 2001

20













#### Solutions

- Modify transport protocol
- Modify link layer protocol

Hybrid

### **Modify Transport Protocol**

#### Explicit Loss Signal

- Distinguish non-congestion losses
- Explicit Loss Notification (ELN) [BK98]
- If packet lost due to interference, set header bit - Only needs to be deployed at wireless router

30

32

- Need to modify end hosts - How to determine loss cause?
- What if ELN gets lost?

29

31

### **Modify Transport Protocol**

- TCP SACK
  - TCP sends cumulative ack only $\rightarrow$ cannot distinguish multiple losses in a window
  - Selective acknowledgement: indicate exactly which
  - packets have not been received - Allows filling multiple "holes" in window in one RTT
  - Quick recovery from a burst of wireless losses
  - Still causes TCP to reduce window

# Modify Link Layer

- How does IP convey reliability requirements to link layer? - not all protocols are willing to pay for reliability
  - Read IP TOS header bits(8)?

  - must modify hosts
    TCP = 100% reliability, UDP = doesn't matter? what about other degrees?
  - consequence of lowest common denominator IP architecture
- Link layer retransmissions
  - Wireless link adds seq. numbers and acks below the IP layer
  - If packet lost, retransmit it May cause reordering
  - Causes at least one additional link RTT delay

  - Some applications need low delay more than reliability e.g. IP telephony
  - easy to deploy

# **Modify Link Layer**

- Forward Error Correction (FEC) codes
  - k data blocks, use code to generate n>k coded blocks
  - can recover original k blocks from any k of the n blocks
  - n-k blocks of overhead
  - trade bandwidth for loss
  - can recover from loss in time independent of link RTT
  - useful for links that have long RTT (e.g. satellite)
  - pay n-k overhead whether loss or not • need to adapt n, k depending on current channel conditions

33

# Hybrid

- Indirect TCP [BB95]
- Split TCP connection into two parts - regular TCP from fixed host (FH) to base station
- modified TCP from base station to mobile host (MH)
- base station fails?
- wired path faster than wireless path?
- TCP Snoop [BSK95]
  - Base station snoops TCP packets, infers flow
  - cache data packets going to wireless side
  - If dup acks from wireless side, suppress ack and retransmit from cache
     soft state

  - what about non-TCP protocols?
  - what if wireless not last hop?

34

# Conclusion

- Transport protocol modifications not deployed - SACK was deployed because of general utility
- Cellular, 802.11b
  - link level retransmissions
  - 802.11b: acks necessary anyway in MAC for collision avoidance
  - additional delay is only a few link RTTs (<5ms)
- Satellite
  - FEC because of long RTT issues
- Link layer solutions give adequate, predictable performance, easily deployable

35