

# User Datagram Protocol (UDP) Transmission Control Protocol (TCP)

Matti Siekkinen

22.09.2009

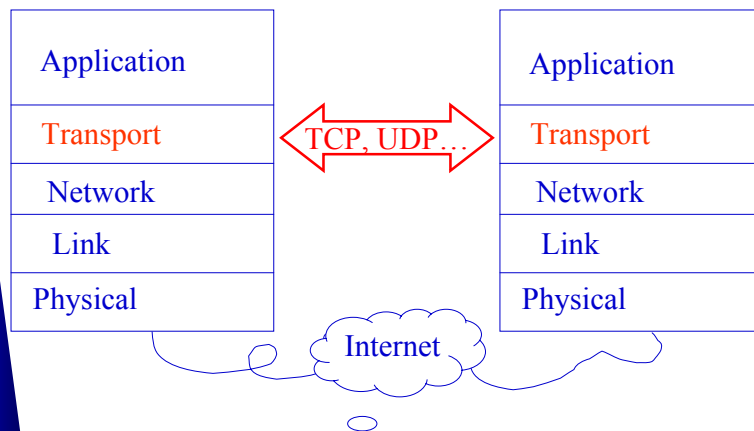
## Outline

- Background
- UDP
  - Role and Functioning
- TCP
  - Basics
  - Error control
  - Flow control
  - Congestion control

22 September 2009



## Transport layer



22 September 2009



## Transport layer (cont.)

- Offers end-to-end transport of data for applications
- Different characteristics
  - Reliable vs. unreliable
  - Forward error correction (FEC) vs. Automatic Repeat-request (ARQ)
  - TCP friendly or not
  - Structured vs. unstructured stream
  - ...

22 September 2009



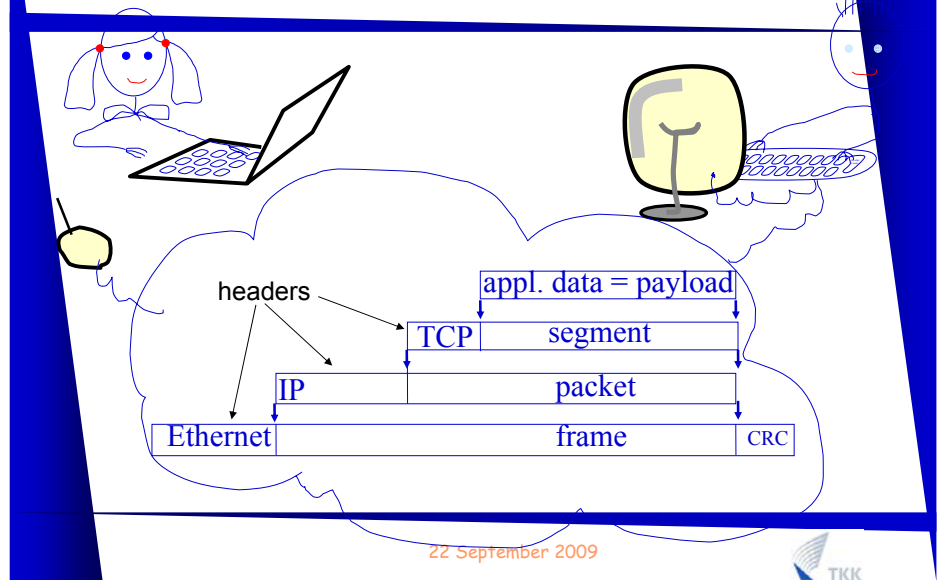
## Reliable vs. best effort

- ❑ TCP - reliable transport
  - Guarantees ordered delivery of packets
  - Important for e.g.
    - Signaling messages
    - File transfer
- ❑ UDP - best effort transport
  - No guarantees of packet delivery
  - Non-critical data delivery, e.g. VoIP

22 September 2009



## Encapsulation

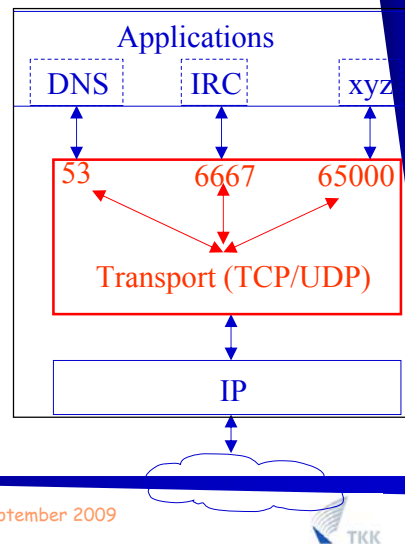


22 September 2009



## Role of ports

- ❑ Well-known port numbers
  - RFC 2780 (&4443)
  - 0-1023
- ❑ Registered port numbers
  - 1024-49151
- ❑ Other port numbers
  - 49152-65535



22 September 2009



## Checksums

- ❑ For detecting damaged packets
  - Compute at sender, check at receiver
- ❑ Computed from pseudo-header and transport segment
  - Pseudo header includes
    - source and destination IP addresses
    - protocol number
    - TCP/UDP length
    - Slightly different method for IPv4 (RFC 768/793) and IPv6 (RFC 2460)
    - Included for protection against misrouted segments
  - Divide into 16-bit words and compute one's complement of the one's complement sum of all the words

22 September 2009



## Part 2: UDP - User Datagram Protocol

## User Datagram Protocol (UDP)

- ❑ Lightweight protocol
  - Just add port numbering and integrity checking (checksums) to IP
  - No segmentation
- ❑ Unreliable connectionless transport service
  - No acknowledgments and no retransmissions
  - Checksum optional in IPv4 and mandatory in IPv6

22 September 2009



## UDP datagram

0	16	31
UDP SOURCE PORT	UDP DESTINATION PORT	
UDP MSG LENGTH	UDP CHECKSUM	
DATA ...		

- ❑ Source port and checksum are optional
  - Checksum mandatory with IPv6
- ❑ Length: header and data in bytes
- ❑ Ports provide application multiplexing within a host (single IP)

22 September 2009



## Part 3: TCP - Transmission Control Protocol

## Outline

- TCP general overview
- TCP-header
- Connection management
- Error control
- Flow control
- Congestion control

22 September 2009



## TCP properties

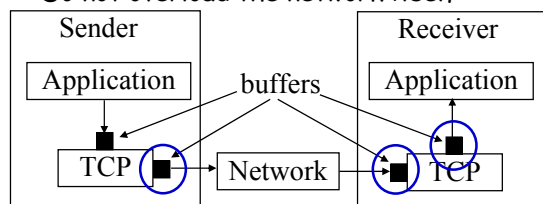
- End-to-end
- Connection oriented
  - State maintained at both ends
  - Identified by a four-tuple
    - Formed by the two end point's IP address and TCP port number
- Reliable
  - Try to guarantee in order delivery of each packet
  - Buffered transfer
- Full Duplex
  - Data transfer simultaneously in both directions

22 September 2009



## TCP properties

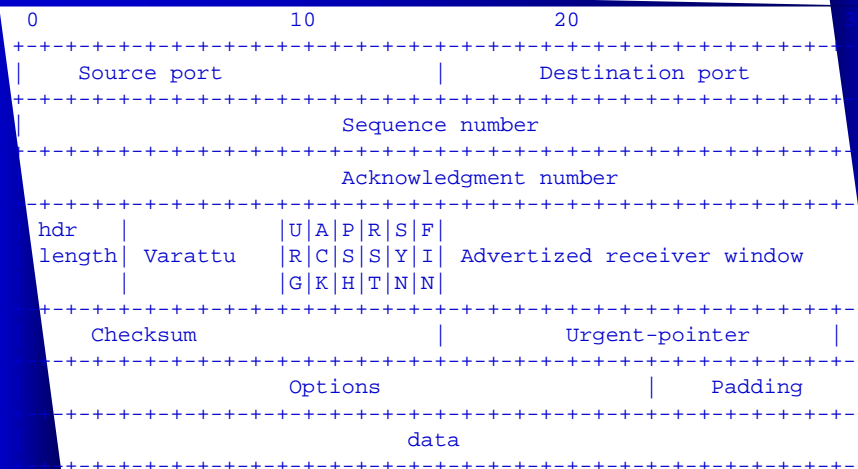
- Three main functionalities for active connection
  1. Error control
    - Deal with the best effort unreliable network
  2. Flow control
    - Do not overload the receiving application
  3. Congestion control
    - Do not overload the network itself



22 September 2009



## TCP-header (RFC 793)



22 September 2009



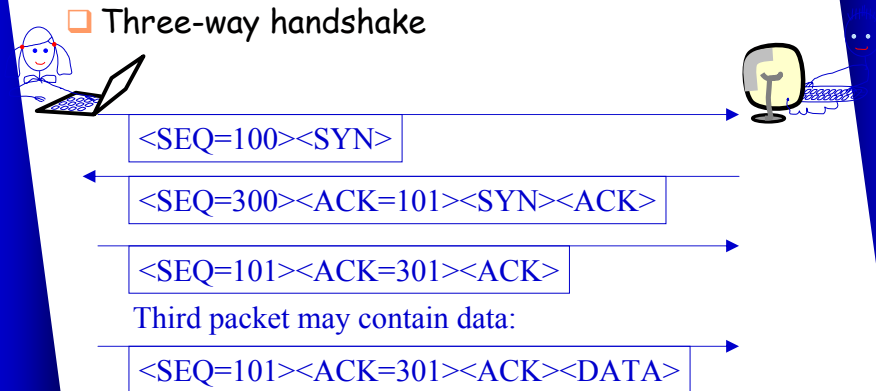
## TCP options

- ❑ 3 = window scaling
- ❑ 8,10 = Timestamp and echo of previous timestamp
  - Improve accuracy of RTT computation
  - Protect against wrapped sequence numbers
- ❑ 2 = Maximum Segment Size (MSS)
  - Negotiated while establishing connection
  - Try to avoid fragmentation
- ❑ 1 = No-operation
  - Sometimes between options, align option fields
- ❑ 0 = End of options

22 September 2009



## Connection establishment



22 September 2009



## Terminating connection

- ❑ Modified three-way handshake
- ❑ If other end has no more data to send, can be terminated one way:
  - Send a packet with FIN flag set
  - Recipient acks the FIN packet
- ❑ After done with the data transfer to the other direction
  - FIN packet and ack to the inverse direction

22 September 2009



## Outline

- ❑ TCP general overview
- ❑ TCP-header
- ❑ Connection management
- ➔ ❑ Error control
- ❑ Flow control
- ❑ Congestion control

22 September 2009



## Error control

- ❑ Mechanisms to detect and recover from lost packets
  - Positive acknowledgments (ARQ)
  - Retransmissions
  - Error detection
    - Timers
    - Checksums

22 September 2009



## Sequence numbers

- ❑ Used in acknowledgments
  - Identify the packets that are acknowledged
- ❑ Marks the number of octets from the start value
  - Start value initiated at connection establishment
- ❑ Distinct values to both directions

22 September 2009



## Cumulative Acknowledgments

- ❑ Acknowledge only the next expected packet in sequence
  - E.g. received 1,2,3,4,6 → ACK 5
- ❑ Advantages
  - Single ACK for multiple packets
    - Delayed ACKs scheme = one ACK for  $2 \times \text{MSS}$  data
  - Lost ACK does not necessarily trigger retransmission
- ❑ Drawback
  - Cannot tell if lost only first or all of a train of packets

22 September 2009



## Retransmission timeout (RTO)

- ❑ RTO associated to each transmitted packet
- ❑ Retransmit packet if no ACK is received before RTO has elapsed
- ❑ RTO is adjusted based on observed delays between sent packets and received ACKs
  - Round trip time (RTT)

22 September 2009



## Computing RTO

- ❑ Original algorithm:
  - $RTT = (\alpha * oldRTT) + ((1 - \alpha) * newRTTsample)$  (recommended  $\alpha = 0.9$ )
  - $RTO: \beta * RTT, \beta > 1$  (recommended  $\beta = 2$ )
- ❑ Does not take into account large variation in RTT

22 September 2009



## Modified algorithm

- ❑ Initialize:  $RTO = 3$
- ❑ Two variables: SRTT (smoothed round-trip time) and RTTVAR (round-trip time variation)
  - First measurement R:
    - $SRTT = R$
    - $RTTVAR = R/2$
  - For subsequent measurement R:
    - $RTTVAR = (1 - \beta) * RTTVAR + \beta * |SRTT - R|$
    - $SRTT = (1 - \alpha) * SRTT + \alpha * R$
    - Use  $\alpha = 1/8, \beta = 1/4$
- ❑  $RTO = SRTT + 4 * RTTVAR$
- ❑ If computed  $RTO < 1s \rightarrow$  round it up to 1s

22 September 2009



## Karn's algorithm

- ❑ Receiving ACK for retransmitted packet
  - Is the ACK for original packet or retransmission??
  - No way to know...
  - ➔ Do not update RTO for retransmitted packets
- ❑ Timer backoff also needed
  - At timeout:  $new\_timeout = 2 * timeout$  (exponential backoff)
- ❑ TCP timestamps can also help disambiguate ACKs

22 September 2009



## Fast Retransmit

- ❑ Introduced by Van Jacobson 1988
- ❑ TCP always ACKs the next expected missing packet
- ❑ Duplicate ACKs are a sign of lost packet(s)
- ❑ Do not wait for timeout but retransmit after 3 duplicate ACKs
  - Wait for reordered packets, don't do go-back-n

22 September 2009



# Selective Acknowledgments (SACK)

- ❑ RFC 2018
- ❑ Helps recovery when multiple packets are lost
- ❑ Receiver reports which segments were lost using TCP SACK (Selective Acknowledgment) options
- ❑ Sender can retransmit several packets per RTT

22 September 2009



# Outline

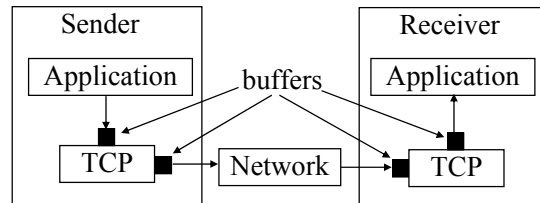
- ❑ TCP general overview
- ❑ TCP-header
- ❑ Connection management
- ❑ Error control
- ➔ ❑ Flow control
- ❑ Congestion control

22 September 2009



# Flow control

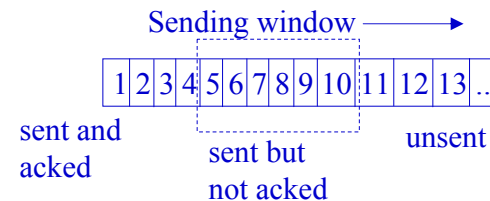
- ❑ Goal: do not overflow the receiving application
- ❑ Window based mechanism to limit transmission rate
- ❑ Receiver advertised window



22 September 2009



# Sliding Window



- ❑ Multiple packets simultaneously "in flight", i.e. outstanding
  - Improve efficiency
- ❑ Buffer sent unacked packets

22 September 2009





## Receiver advertised window

- ❑ Receiver advertises the maximum window size the sender is allowed to use
- ❑ Enables receiver TCP to signal sending TCP to backoff
  - Receiving application not consuming received data fast enough
- ❑ Value is included in each ACK
  - Can change dynamically

22 September 2009



## Silly Window Syndrome (SWS)

- ❑ Problem in worst case:
  - Receiver buffer between TCP and application fills up
  - Receiving application read a single byte -> TCP advertises a receiver window of size one
  - Sender transmits a single byte
- ❑ Lot of overhead due to packet headers

22 September 2009



## Avoiding SWS

- ❑ Dave Clark's solution
- ❑ Window update only with significant size
  - At least MSS worth of data or
  - Half of its buffer
- ❑ Analogy at sender side
  - Application gives small chunks of data to TCP -> send small packets
  - Nagle's algorithm: Delay sending data until have MSS worth of it
    - Does not work for all applications, e.g. delay sensitive apps
    - Need also mechanism to tell TCP to transmit immediately -> Push flag

22 September 2009



## Outline

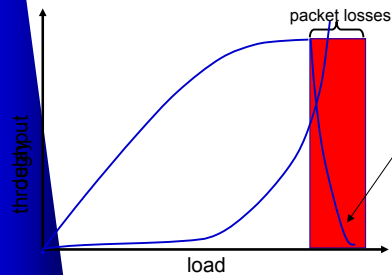
- ❑ TCP general overview
- ❑ TCP-header
- ❑ Connection management
- ❑ Error control
- ❑ Flow control
- ➔ ❑ Congestion control
  - Background and motivation
  - Evolution of congestion control in TCP
  - Most important mechanisms
  - Recent developments
  - Conclusions

22 September 2009



# Why we need congestion control

- Flow control in TCP prevents overwhelming the receiving application
- Problem: Multiple TCP senders sharing a link can still overwhelm it



Congestion collapse due to:

- Retransmitting lost packets
  - Further increases the load
- Spurious retransmissions of packets still in flight
  - Unnecessary retransmissions lead to even more load!
  - Like pouring gasoline on a fire

22 September 2009



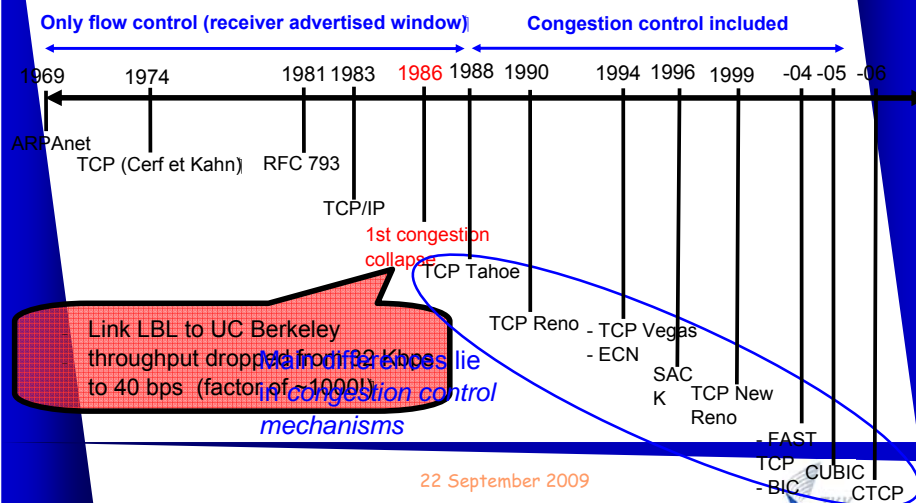
# Congestion control (cont.)

- Principle:
  - Continuously throttle TCP sender's transmission rate
  - Probe the network by increasing the rate when all is fine
  - Decrease rate when signs of congestion (e.g. packet loss)
- How?
  - Introduce *congestion window* (cwnd):
    - $\#outstanding\ bytes = \min(cwnd, rwnd)$  ← flow control
  - Adjust cwnd size to control the transmission rate
    - Adjustment strategy depends on TCP version

22 September 2009



# Glimpse into the past



22 September 2009



# TCP Tahoe

- 1988 Van Jacobson
- The basis for TCP congestion control
- Lost packets are sign of congestion
  - Detected with *timeouts*: no ACK received in time
- Two modes:
  - Slow Start
  - Congestion Avoidance
- New retransmission timeout (RTO) calculation
  - Incorporates variance in RTT samples
  - Timeout really means a lost packet (=congestion)
- Fast Retransmit

22 September 2009



## Slow Start (SS)

- In two cases:
  - Beginning of connection
  - After a timeout
- On each ACK for new data, increase cwnd by one packet
  - Exponential increase in the size of cwnd
  - Ramp up a new TCP connection fast (not slow!)
  - Kind of a misnomer...

22 September 2009



## Congestion Avoidance (CA)

- Approach the rate limit of the network more conservatively
  - Easy to drive the net into saturation but hard for the net to recover
- AIMD (additive increase / multiplicative decrease)
  - On each ACK for new data: increase cwnd by 1/cwnd
  - On timeout: set cwnd to half the current window size

22 September 2009



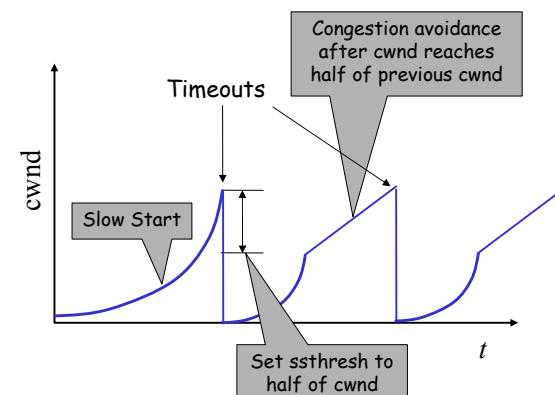
## Combining SS and CA

- Introduce Slow start threshold (ssthresh)
- On timeout:
  - $ssthresh = 0.5 \times cwnd$
  - $cwnd = 1 \text{ packet}$
- On new ACK:
  - If  $cwnd < ssthresh$ : do Slow Start
  - Else: do Congestion Avoidance

22 September 2009



## TCP Tahoe: adjusting cwnd



22 September 2009



## TCP Reno

□ Van Jacobson 1990

□ Fast retransmit with Fast recovery

- Duplicate ACKs tell sender that packets still go through
- Do less aggressive back-off:
  - $ssthresh = 0.5 \times cwnd$
  - $cwnd = ssthresh + 3$  packets
  - Increment cwnd by one for each additional duplicate ACK
  - When the next new ACK arrives:  $cwnd = ssthresh$

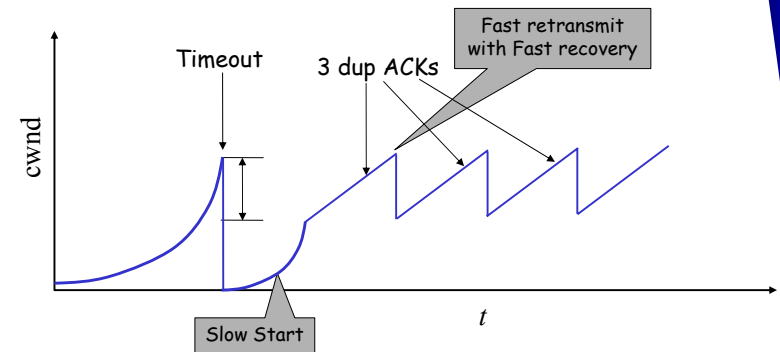
Fast recovery

Nb of packets that were delivered

22 September 2009



## TCP Reno: adjusting cwnd



22 September 2009



## TCP New Reno

□ 1999 by Sally Floyd

□ Modification to Reno's Fast Recovery phase

□ Problem with Reno:

- Multiple packets lost in a window
- Sender out of Fast Recovery after retransmission of only one packet
  - cwnd closed up
  - no room in cwnd to generate duplicate ACKs for additional Fast Retransmits
  - eventual timeout

□ New Reno continues Fast Recovery until all lost packets from that window are recovered

22 September 2009



## More recent developments

□ Delay-based congestion control

- TCP Vegas

□ Wireless networks

- Take into account random packet loss due to bit errors (not congestion!)
- E.g. TCP Veno

□ Paths with high *bandwidth\*delay*

- These "long fat pipes" require large cwnd to be saturated
- SS and CA provide too slow response
- TCP CUBIC
- Compound TCP (CTCP)

22 September 2009



## TCP Vegas

- 1994 by Brakmo et Peterson
- Issue: Tahoe and Reno RTO clock is very coarse grained
  - "ticks" each 500ms
- Increasing delay is a sign of congestion
  - Packets start to fill up queues
  - Expected throughput =  $cwnd / BaseRTT$  ← minimum of all measured round trip times
  - Compare expected to actual throughput
  - Adjust rate accordingly before packets are lost
- Also some modifications to Slow start and Fast Retransmit
- Potentially up to 70% better throughput than Reno
- Fairness with Reno?
  - Reno grabs larger share due to late congestion detection

22 September 2009



## BIC and CUBIC

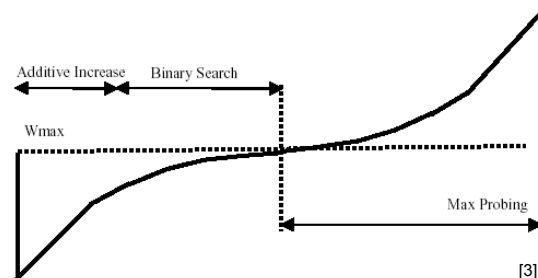
- 2004, 2005 by Xu and Rhee
- Both for paths with high (*bandwidth x delay*)
  - These "long fat pipes" lead to large cwnd
  - SS and CA provide too slow response
  - Scale up to tens of Gb/s
- BIC TCP
  - No AIMD
  - Window growth function is combination of *binary search* and *linear increase*
  - Aim for TCP friendliness and RTT fairness

22 September 2009



## BIC and CUBIC

- BIC window growth function

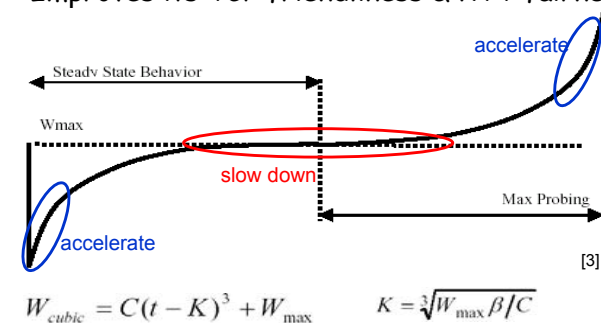


22 September 2009



## BIC and CUBIC (cont.)

- CUBIC TCP
  - Enhanced version of BIC
  - Simplifies BIC window control using a *cubic function*
  - Improves its TCP friendliness & RTT fairness



$$W_{cubic} = C(t - K)^3 + W_{max} \quad K = \sqrt[3]{W_{max} \beta / C}$$

22 September 2009



## Compound TCP (CTCP)

- ❑ From Microsoft research, 2006
- ❑ Tackles same problems as BIC and CUBIC
  - High speed and long distance networks
  - RTT fairness, TCP friendliness
- ❑ Loss-based vs. delay-based approaches
  - Loss-based (e.g. HSTCP, BIC...) too aggressive
  - Delay-based (e.g. Vegas) too timid
- ❑ Compound approach
  - Use delay metric to sense the network congestion
  - Adaptively adjust aggressiveness based on network congestion level
  - Loss-based component: *cwnd* (standard TCP Reno)
  - Scalable delay-based component: *dwnd*
  - TCP sending window is  $Win = cwnd + dwnd$

22 September 2009



## Explicit Congestion Notification (ECN)

- ❑ Routers flag packets upon congestion
  - Active queue management
- ❑ TCP sender consequently reduces *cwnd*

22 September 2009



## Deployment

- ❑ Windows
  - Server 2008 uses Compound TCP (CTCP) by default
  - Vista, XP support CTCP, New Reno by default
- ❑ Linux
  - TCP BIC default in kernels 2.6.8 through 2.6.18
  - TCP CUBIC since 2.6.19

22 September 2009



## Conclusions

- ❑ Transport layer
  - End-to-end transport of data for applications
  - Application multiplexing through port numbers
  - Reliable (TCP) vs. unreliable (UDP)
- ❑ UDP
  - Unreliable, no state
  - Optionally integrity checking
- ❑ TCP
  - Connection management
  - Error control: deal with unreliable network path
  - Flow control: Prevent overwhelming receiving application
  - Congestion control: Prevent overwhelming the network
  - Many TCP versions exist today
    - Main differences in congestion control
    - Loss-based and delay-based congestion detection
    - More and less aggressive rate control
    - Suitable for different network types

22 September 2009



# References

- [1] IETF's RFC page: <http://www.ietf.org/rfc.html>
- [2] V. Jacobson: **Congestion Avoidance and Control**. *In proceedings of SIGCOMM '88*.
- [3] L. Brakmo et al.: **TCP Vegas: New techniques for congestion detection and avoidance**. *In Proceedings of SIGCOMM '94*.
- [4] **RFC2582/RFC3782 - The NewReno Modification to TCP's Fast Recovery Algorithm**.
- [5] L. Hu et al.: **Binary Increase Congestion Control for Fast, Long Distance Networks**, *IEEE Infocom, 2004*.
- [6] S. Ha et al.: **CUBIC: A New TCP-Friendly High-Speed TCP Variant**, *ACM SIGOPS, 2008*.
- [7] K. Tan et al.: **Compound TCP: A Scalable and TCP-friendly Congestion Control for High-speed Networks**, *In IEEE Infocom, 2006*.
- [8] W. John et al.: **Trends and Differences in Connection Behavior within Classes of Internet Backbone Traffic**, *In PAM 2008*.
- [9] A. Medina et al.: **Measuring the evolution of transport protocols in the internet**, *SIGCOMM CCR, 2005*.