
Chapter 21:

TCP Timeout and Retransmission

Introduction

- ❑ **Two examples of timeout and retransmission had already seen:**
 - ❖ 1. In the ICMP port unreachable (TFTP, Section 6.5)
 - ❖ 2. In the ARP example to a nonexistent host (Section 4.5)
- ❑ **TCP manage four different timers for each connection:**
 - ❖ 1. A retransmission timer (this chapter)
 - ❖ 2. A persist timer (Chapter 22)
 - ❖ 3. A keepalive timer (Chapter 23)
 - ❖ 4. A 2MSL timer (Section 18.6)
- ❑ **Simple Timeout and Retransmission Example:**
 - ❖ Line 4 is the transmission of “hello,world” and line 5 is its ACK.
 - ❖ Line 6 shows “and hi”. Line 7-18 are 12 retransmissions
 - ❖ Line 19 is the TCP finally gives up and sends a reset

Round-Trip Time Measurement (Cont.)

❑ Karn's Algorithm:

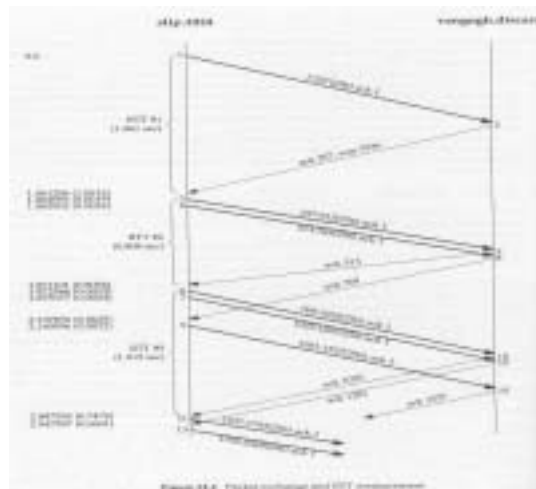


- ❖ A packet is transmitted, a timeout occurs, the packet is retransmitted with the longer RTO, and an acknowledgment is received
- ❖ Is the ACK for the first transmission or the second?
- ❖ This is called the *retransmission ambiguity problem*
- ❖ Karn and Partridge: don't count retransmitted packets into the RTT estimator.

❑ An RTT Example:

- ❖ sent 32768 bytes of data from *slip* to *vangogh.cs.berkeley.edu*
 - `slip% sock -D -i -n32 vangogh.cs.berkeley.edu. discard`
- ❖ slip is connected to the 140.252.1 Ethernet by two SLIP links
- ❖ MTU between slip and bsd is 296
- ❖ 32 1024-byte => 128 segment with 256 bytes of user data

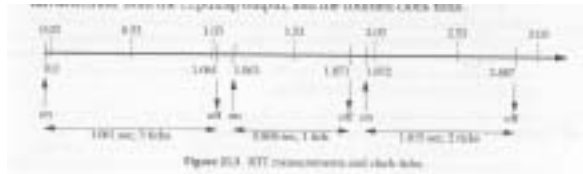
An RTT Example



An RTT Example

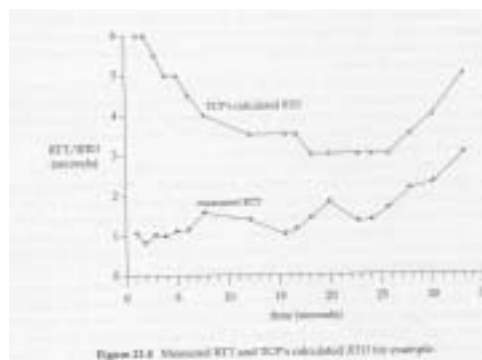
□ Round-Trip Time Measurements:

- ❖ One connection, one timer
- ❖ RTT#1 is 1.061 seconds \Rightarrow 3 clock ticks
- ❖ RTT#2 is 0.808 seconds \Rightarrow 1 clock tick
- ❖ RTT#3 is 1.015 seconds \Rightarrow 2 clock ticks
- ❖ Segment 4,7,9 cannot be timed, since the timer is already being used by segment 3 and 6



An RTT Example

- ❖ In this complete example 18 RTT samples were collected



An RTT Example

RTT Estimator Calculations:

- ❖ The initial $RTO = A + 2D \Rightarrow 0 + 2 \times 3 = 6$ seconds
- ❖ After 5.802 seconds $RTO = A + 4D \Rightarrow 0 + 4 \times 3 = 12$ seconds
- ❖ The ACK arrives 467 ms after the retransmission. The A and D are not updated because of retransmission ambiguity
- ❖ The ACK on line 4 is not timed since it is only an ACK



An RTT Example

RTO calculations

- ❖ first segment arrives $\Rightarrow RTO = 6$ seconds
- ❖ second segment arrives $\Rightarrow RTO = 6.3125$ seconds
- ❖ Fixed-point calculations that are actually used $\Rightarrow RTO$ is 6 seconds (not 6.3125)

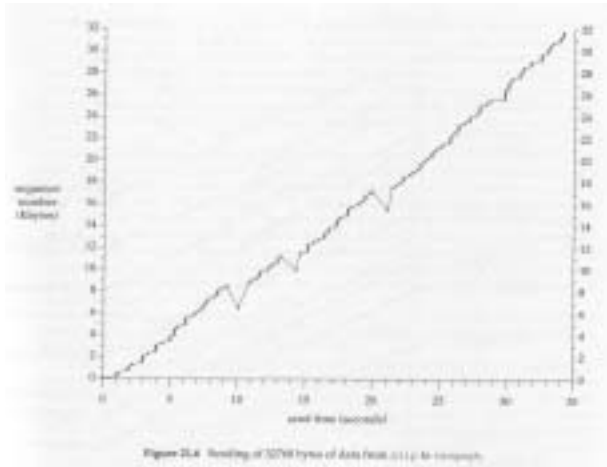
Slow Start

- ❖ See the slow start algorithm in Section 20.6

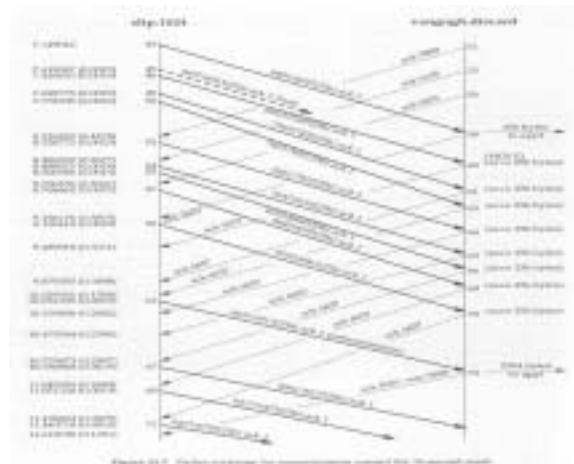
Congestion Example:

- ❖ Normally the data points should move up and to the right, with the slope of the points being the transfer rate.
- ❖ Retransmissions will appear as motion down and to the right.

Congetion Example



Congetion Example



Congestion Example

❑ The Jacobson's fast retransmit algorithm:



- ❖ It is followed by his fast recovery algorithm. The third of the duplicate ACKs was received that forces to retransmit
- ❖ Berkeley-derived implementation when the third one is received, assume that a segment has been lost and retransmit only one segment

❑ When the missing data arrives (segment 63):

- ❖ The receiving TCP now has data bytes 6657-8960 in its buffer, and passes these 2304 bytes to the user process.
- ❖ All 2304 bytes are acknowledged in segment 72

Congestion Avoidance Algorithm

❑ What's congestion avoidance?



- ❖ It is a way to deal with lost packets
- ❖ Assumption: packet loss caused by damage is very small ($< 1\%$)
- ❖ The loss of a packet signals congestion somewhere in the path between the source and the destination.

❑ Two indications of packet loss:



- ❖ A timeout occurring
- ❖ The receipt of duplicate ACKs

❑ Congestion avoidance and slow start are independent algorithms with different objectives.

❑ Two variables be maintained for each connection: *cwnd* and *ssthresh*



Congestion Avoidance Algorithm (Cont.)

□ Congestion avoidance algorithm operates:

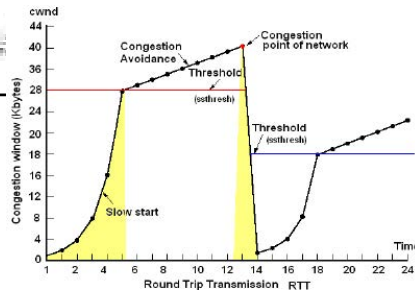
- ❖ 1. Initialization \Rightarrow $cwnd$ is one segment, $ssthresh$ is 65535 bytes
- ❖ 2. TCP output never sends more than the minimum of $cwnd$ and the receiver's advertised window
- ❖ 3. When congestion occurs, one-half of the current window size is saved in $ssthresh$. If timeout, $cwnd$ is set to one segment
- ❖ 4. When new data is acknowledged by the other end, increase $cwnd$

□ If $cwnd$ is less than or equal to $ssthresh$, doing slow start. Otherwise, we're doing congestion avoidance.



□ Congestion avoidance dictates that $cwnd$ be increased by $1/cwnd$, plus a small fraction of segment size, each time an ACK is received: **an additive increase.**

□ Increase $cwnd$ by at most **one segment each round-trip time** (regardless how many ACKs are received in that RTT)



Network & System Lab, NSYSU

15



Congestion Avoidance Algorithm

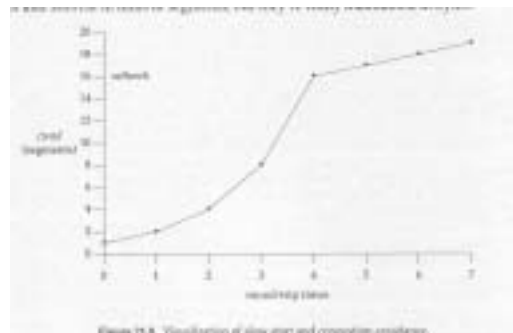


Figure 23.8 Visualization of slow start and congestion avoidance.



Network & System Lab, NSYSU

16

Fast Retransmit and Fast Recovery Algorithms

- ❑ TCP is required to generate an immediate ACK when an out-of-order segment is received.
- ❑ **Duplicate ACK:** a segment is received out of order.
 - ❖ Two possible situations:
 - Packet loss – if two or more duplicate ACKs are received – a strong indication of segment loss.
 - First out last in
- ❑ The receipt of the duplicate ACKs tells us more than just a packet has been lost:
 - ❖ A segment leaves the network into receiver's buffer.
 - ❖ There still is data flowing between the two ends – DON'T reduce the flow abruptly by going into slow start.

Fast Retransmit and Fast Recovery Algorithms (Cont.)

- ❑ **Fast retransmit algorithm:**
 - ❖ If three or more duplicate ACKs are received in a row, indicate a segment has been lost, then retransmission the missing segment
- ❑ **Fast recovery algorithm:**
 - ❖ Next, congestion avoidance, but not slow start is performed
- ❑ **Congestion Example**
 - ❖ In congestion avoidance:
 - 考 ➢ $cwnd = cwnd + (segsz \times segsz) / cwnd + segsz / 8$
 - ❖ By fast retransmit and fast recovery, we can send a new data segment when $cwnd > unacknowledged\ bytes$
 - ❖ $cwnd$ is allowed to keep increasing while the duplicate ACKs are received, since each duplicate ACK means the a segment has left the network.

Congestion Example (Continued)

Sequence (Figure 11.2)	Action			Totals	
	Send	Receive	Comment	count	window
	99%		initial	20	1000
	99%		transmit	20	101
	ACK	59% ACK			
1	1.037286				
2		ACK 307	clear start	101	101
3	22.724286				
4	10.576286				
5		ACK 311	clear start	101	101
6	100.027286				
7	905.136286				
8		ACK 369	cong. avoid	985	101
9	100.040286				
10		ACK 320	cong. avoid	991	101
11	917.776286				
12		ACK 330	cong. avoid	1000	101

Figure 11.2 Example of congestion avoidance.

Sequence (Figure 11.3)	Action		Totals	
	Send	Receive	Count	Window
36		ACK 407	ACK at new data	1126
39	870.896286			
40		ACK 407	duplicate ACK #1	1430
41		ACK 407	duplicate ACK #1	1430
42		ACK 407	duplicate ACK #1	1742
43	947.991286			
44		ACK 407	duplicate ACK #1	1946
45		ACK 407	duplicate ACK #1	2104
46		ACK 407	duplicate ACK #1	2304
47	990.1017286			
48		ACK 407	duplicate ACK #1	2616
49	1027.947286			
50		ACK 407	duplicate ACK #1	2872
51	1470.070286			
52		ACK 898	ACK at new data	3280

Figure 11.3 Example of congestion avoidance (continued).

Congestion Example (Continued)

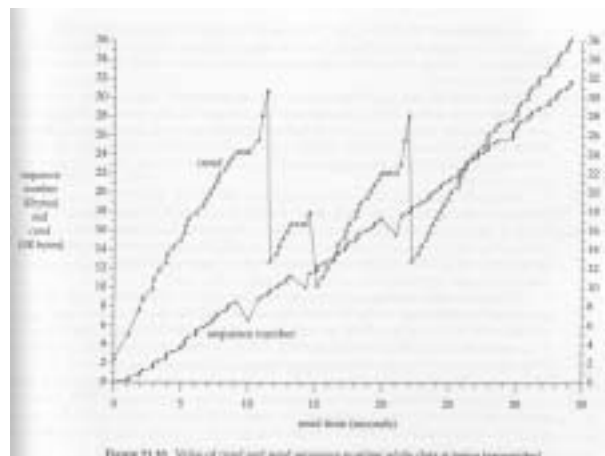


Figure 11.10 Value of send and send sequence number while data is being transmitted.

ICMP Errors

❑ Berkeley-based implementations handle ICMP errors as follows:

- ❖ A received source quench causes the *cwnd* set to one segment to initiate slow start, but the *ssthresh* is not changed
- ❖ A received host unreachable or network unreachable is effectively ignored, since these two errors are considered transient.

❑ An Example

- slip% sock aix echo
- test line
- test line
-
- another line
-
-
- down

SLIP link is brought down here

type this line and retransmissions

SLIP link is reestablished here

after the last line, SLIP link is brought

➤ read error: No route to host

TCP finally gives up



ICMP Errors

The PCAP shows a sequence of events where a SLIP link is brought down and then reestablished. The capture shows multiple retransmissions of data segments and the eventual timeout of the connection. The events are as follows:

- Initial connection setup and data transfer.
- SLIP link is brought down (indicated by a 'read error: No route to host' message).
- Multiple retransmissions of data segments are shown, with the sender increasing its congestion window.
- The link is eventually reestablished, and data transfer resumes.
- The connection eventually times out due to the link being brought down again.



