







# SACK TCP Rules

- A SACK cannot be sent unless the SACK-permitted option has been received (in the SYN).
- If a receiver has chosen to send SACKs, it must send them whenever it has data to SACK at the time of an ACK.
- The receiver should send an ACK for every valid segment it receives containing new data (standard TCP behavior), and each of these ACKs should contain a SACK, assuming there is data to SACK.

🛕 Mar-6-03

4/598N: Computer Networks 7















# Reneging in SACK TCP Therefore, the sender must maintain normal TCP timeouts. A segment cannot be considered received until an ACK is received for it. The sender must retransmit the segment at the left window edge after a retransmit timeout, even if the SACK bit is on for that segment. A segment cannot be removed from the transmit buffer until the left window edge is advanced over it, via the receiving of an ACK.

SACK TCP follows standard TCP congestion control; it should not damage the network.
 SACK TCP has an advantage over other implementations (Reno, Tahoe, Vegas, and NewReno) as it has added information due to the SACK data.
 This information allows the sender to better decide what it needs to retransmit and what it does not. This can only serve to help the sender, and should not adversely affect other TCPs.

SACK TCP Observations

# SACK TCP Observations

4/598N: Computer Networks

15

17

Mar-6-03

Mar-6-03

- While it is still possible for a SACK TCP to needlessly retransmit segments, the number of these retransmissions has been shown to be quite low in simulations, relative to Reno and Tahoe TCP.
- In any case, the number of needless retransmissions must be strictly less than Reno/Tahoe TCP. As the sender has additional information from which to devise its retransmission scheme, worse performance is not possible (barring a flawed implementation).

4/598N: Computer Networks

SACK TCP Implementation Progress Current SACK TCP implementations: - Windows 2000 - Windows 98 / Windows ME - Solaris 7 and later - Linux kernel 2.1.90 and later - FreeBSD and NetBSD have optional modules · ACIRI has measured the behavior of 2278 random web servers that claim to be SACK-enabled. Out of these, 2133 (93.6%) appeared to ignore SACK data and only 145 (6.4%) appeared to actually use the SACK data. Mar-6-03 4/598N: Computer Networks 18









# **D-SACK TCP Rules**

- If the D-SACK block reports a duplicate sequence from a (possibly larger) block of data in the receiver buffer above the cumulative acknowledgement, the second SACK block (the first non D-SACK block) should specify this block.
- As only the first SACK block is considered to be a D-SACK block, if multiple sequences are duplicated, only the first is contained in the D-SACK block.

4/598N: Computer Networks

23

À

Mar-6-03



# SACK and D-SACK Interaction

- There is no difference between SACK and D-SACK, except that the first SACK block is used to report a duplicate segment in D-SACK.
- There is no separate negotiation/options for D-SACK.
- There are no inherit problems with having the receiver use D-SACK and having the sender use traditional SACK. As the duplicate that is being reported is still being SACKed (for the second or greater time), there is no problem with a SACK TCP using this extension with a D-SACK TCP (although the D-SACK specific data is not used).

Mar-6-03

4/598N: Computer Networks

25



Increasing the Initial Window							
<ul> <li>RFC 2 increating segment</li> <li>This n</li> </ul>	414 specifies a sing of the max ent to a larger v ew larger value	in experir imum init alue. is given	nental change to TCP, the ial window size, from one as:				
<ul> <li>This tr</li> </ul>	min ( 4*MS anslates to:	SS, max ( 1	2*MSS, 4380 bytes) )				
	Maximum Segment	Size (MSS)	Maximum Initial Window Size				
	<= 1095 bytes		<= 4 * MSS				
	1095 bytes < MSS < 2190 bytes		<= 4380 bytes				
	>= 2190 bytes		<= 2 * MSS				



Advantages of an						
<ul> <li>This change mechanism to one seggimplement RFC 2001</li> <li>This new la advantage:         <ul> <li>With slow forced to With an irreceiver warrives, carrives, carrives</li> </ul> </li> </ul>	e is in contrast to the slow Size e is in contrast to the slow star n, which initializes the initial wir ment. This mechanism is in pla sender-based congestion contr for a complete discussion). arger window offers three distin s: start, a receiver which uses delayed wait for a timeout before generating itial window of at least two segment <i>i</i> ll generate an ACK after the secon- ausing a speedup in data acknowled	t ndow size ace to rol (see ct d ACKs is an ACK. s, the d segment lgement.				
A Mar-6-03	4/598N: Computer Networks	29				



## Disadvantages of an Increased Initial Window Size

- This approach also has disadvantages:
- This approach could cause increased congestion, as multiple segments are transmitted at once, at the beginning of the connection. As modern routers tend to not handle bursty traffic well (Drop Tail queue management), this could increase the drop rate.
- ACIRI research on this topic concludes that there is no more danger from increasing the initial TCP window size to a maximum of 4KB than the presence of UDP communications (that do not have end-to-end congestion control).

Mar-6-03

Mar-6-03

4/598N: Computer Networks

31

## Increased Initial Window Size Implementation Progress

- Looking at ACIRI observations, current web servers use a wide range of initial TCP window sizes, ranging from one segment (slow start) to seventeen segments.
- This is a clear violation of RFC 2414, not to mention RFC 2001 (the currently approved IETF/ISOC standard).
- Such large initial window sizes seem to indicate a greedy TCP, not conforming to the required senderside congestion control window (even if the experimental higher initial window is considered).

4/598N: Computer Networks

32

À

Mar-6-03



4/598N: Computer Networks















CHAN Details					
Lost message (request, reply, or ACK)					
<ul> <li>set RETRANSMIT timer</li> </ul>					
<ul> <li>use message id (MID) field to distinguish</li> </ul>					
<ul> <li>Slow (long running) server</li> </ul>					
<ul> <li>client periodically sends "are you alive" probe, or</li> </ul>					
<ul> <li>server periodically sends "I'm alive" notice</li> </ul>					
Want to support multiple outstanding calls					
<ul> <li>use channel id (CID) field to distinguish</li> </ul>					
<ul> <li>Machines crash and reboot</li> </ul>					
<ul> <li>use boot id (BID) field to distinguish</li> </ul>					
Mar-6-03 4/598N: Computer Networks	42				

CHAN Header Format							
<pre>typedef stru u_short int int int int } ChanHdr;</pre>	CID; /* CID; /* MID; /* BID; /* Length; /* ProtNum; /*	REQ, REP, ACK, PROBE */ unique channel id */ unique message id */ unique boot id */ length of message */ high-level protocol */					
<pre>typedef stru u_char u_char int int XkReturn Msg Semaphore int int } ChanState;</pre>	<pre>ttype; status; retries; timeout; ret_val; *request; *reply; a reply_sem; mid; bid;</pre>	<pre>/* CLIENT or SERVER */ /* BUSY or IDLE */ /* number of retries */ /* timeout value */ /* return value */ /* request message */ /* reply message */ /* client semaphore */ /* message id */ /* boot id */</pre>					
Mar-6-03	4	/598N: Computer Networks	43				



















