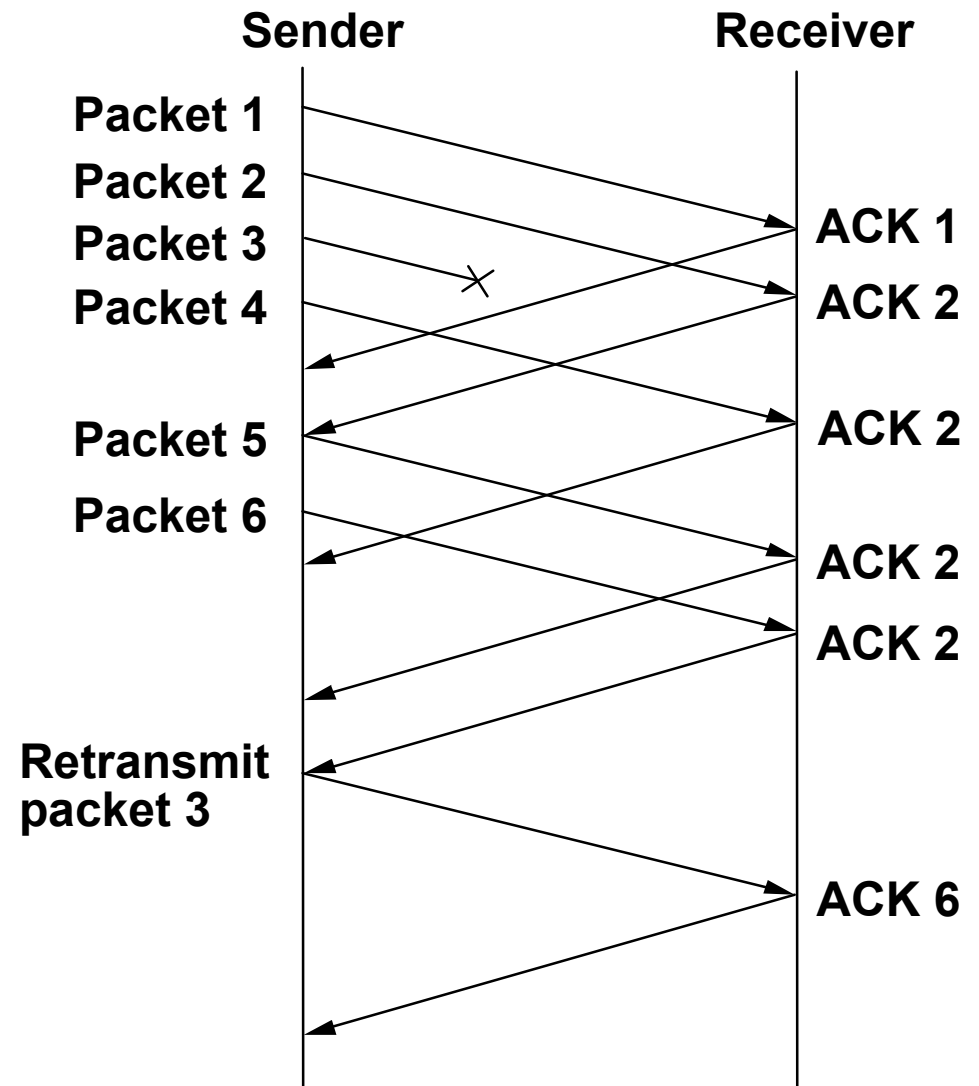


# Fast Retransmit

- ▶ Problem: coarse-grain TCP timeouts lead to idle periods
- ▶ Fast retransmit: use duplicate ACKs to trigger retransmission



# Fast Retransmit

- ▶ If we get 3 duplicate acks for segment N
  - Retransmit segment N
  - Set ssthresh to  $0.5 * \text{cwnd}$
  - Set cwnd to ssthresh + 3
- ▶ For every subsequent duplicate ack
  - Increase cwnd by 1 segment
- ▶ When new ack received
  - Reset cwnd to ssthresh (resume congestion avoidance)



# Congestion Avoidance

- ▶ TCP needs to create congestion to find the point where congestion occurs
- ▶ Coarse grained timeout as loss indicator
- ▶ If loss occurs when  $cwnd = W$ 
  - Network can absorb  $0.5W \sim W$  segments
  - Set  $cwnd$  to  $0.5W$  (multiplicative decrease)
  - Needed to avoid exponential queue buildup
- ▶ Upon receiving ACK
  - Increase  $cwnd$  by  $1/cwnd$  (additive increase)
  - Multiplicative increase  $\rightarrow$  non-convergence



# Slow Start and Congestion Avoidance

- ▶ 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

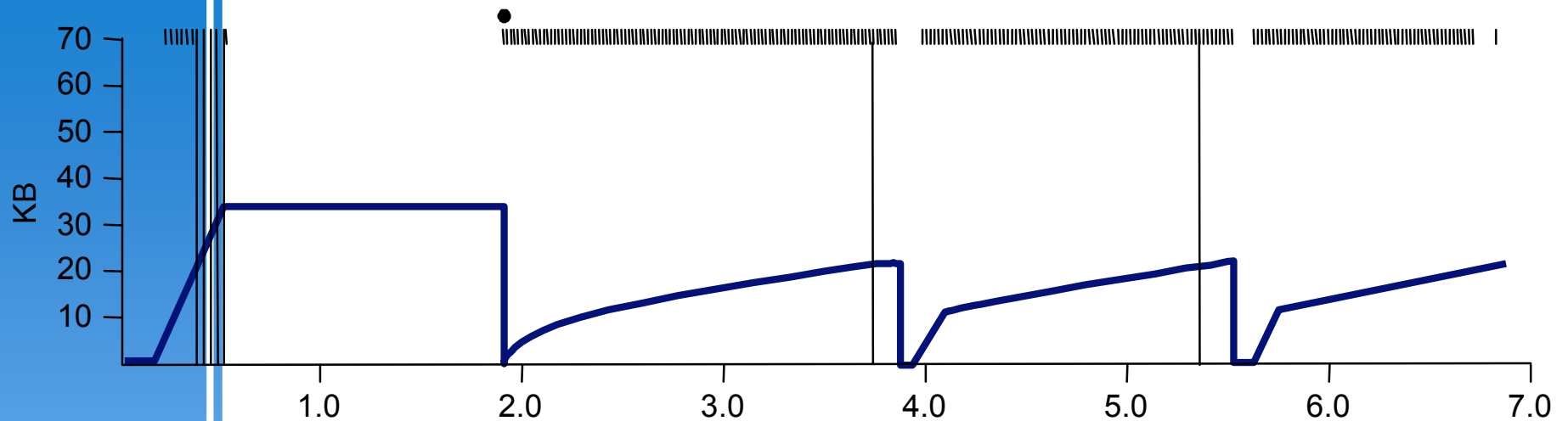


# Fast Recovery

- ▶ In congestion avoidance mode, if duplicate acks are received, reduce cwnd to half
- ▶ If  $n$  successive duplicate acks are received, we know that receiver got  $n$  segments after lost segment:
  - Advance cwnd by that number



# Results



## ► Fast recovery

- skip the slow start phase
- go directly to half the last successful `CongestionWindow` (`sssthresh`)



# Impact of Timeouts

- ▶ Timeouts can cause sender to
  - Slow start
  - Retransmit a possibly large portion of the window
- ▶ Bad for lossy high bandwidth-delay paths
- ▶ Can leverage duplicate acks to:
  - Retransmit fewer segments (fast retransmit)
  - Advance cwnd more aggressively (fast recovery)



# Summary

- ▶ Three way handshake to initiate TCP. Either side can tear down connection using FIN sequence.
- ▶ Initial sequence number chosen to avoid packets from earlier incarnations
- ▶ AIMD to slowly probe network. Slow start to speed up the probing process (at the expense of hard congestion for the last step). Use ssthresh to decide whether to restart slow start or AIMD after loss
- ▶ Fast retransmit by inferring duplicate ACKs = lost packet, Fast recovery by inferring duplicate ACKs = increasing CWND





# TCP Extensions

- ▶ Implemented using TCP options
  - Timestamp
  - Protection from sequence number wraparound
  - Large windows



# Timestamp Extension

- ▶ Used to improve timeout mechanism by more accurate measurement of RTT
- ▶ When sending a packet, insert current timestamp into option
- ▶ Receiver echoes timestamp in ACK



# Protection Against Wrap Around

- ▶ 32-bit **SequenceNum**

Bandwidth	Time Until Wrap Around
T1 (1.5 Mbps)	6.4 hours
Ethernet (10 Mbps)	57 minutes
T3 (45 Mbps)	13 minutes
FDDI (100 Mbps)	6 minutes
STS-3 (155 Mbps)	4 minutes
STS-12 (622 Mbps)	55 seconds
STS-24 (1.2 Gbps)	28 seconds

- ▶ Use timestamp to distinguish sequence number wraparound



# Keeping the Pipe Full

## ▶ 16-bit Advertised Window

Bandwidth	Delay x Bandwidth Product
T1 (1.5 Mbps)	18KB
Ethernet (10 Mbps)	122KB
T3 (45 Mbps)	549KB
FDDI (100 Mbps)	1.2MB
STS-3 (155 Mbps)	1.8MB
STS-12 (622 Mbps)	7.4MB
STS-24 (1.2 Gbps)	14.8MB



# Large Windows

- ▶ Apply scaling factor to advertised window
  - Specifies how many bits window must be shifted to the left
- ▶ Scaling factor exchanged during connection setup



# TCP Flavors

- ▶ Tahoe, Reno, Vegas
- ▶ TCP Tahoe (distributed with 4.3BSD Unix)
  - Original implementation of van Jacobson's mechanisms (VJ paper)
  - Includes:
    - Slow start (exponential increase of initial window)
    - Congestion avoidance (additive increase of window)
    - Fast retransmit (3 duplicate acks)



# TCP Reno

▶ 1990: includes:

- All mechanisms in Tahoe
- Addition of fast-recovery (opening up window after fast retransmit)
- Delayed acks (to avoid silly window syndrome)
  - Silly window syndrome is when one byte is ack'd causing advertisedwindow=1
- Header prediction (to improve performance)



# SACK TCP

(RFC 2018)





# What's Wrong with Current TCP?

- ▶ TCP uses a cumulative acknowledgment scheme, in which the receiver identifies the last byte of data successfully received.
- ▶ Received segments that are not at the left window edge are not acknowledged.
- ▶ This scheme forces the sender to either wait a roundtrip time to find out a segment was lost, or unnecessarily retransmit segments which have been correctly received.
- ▶ Results in significantly reduced overall throughput.



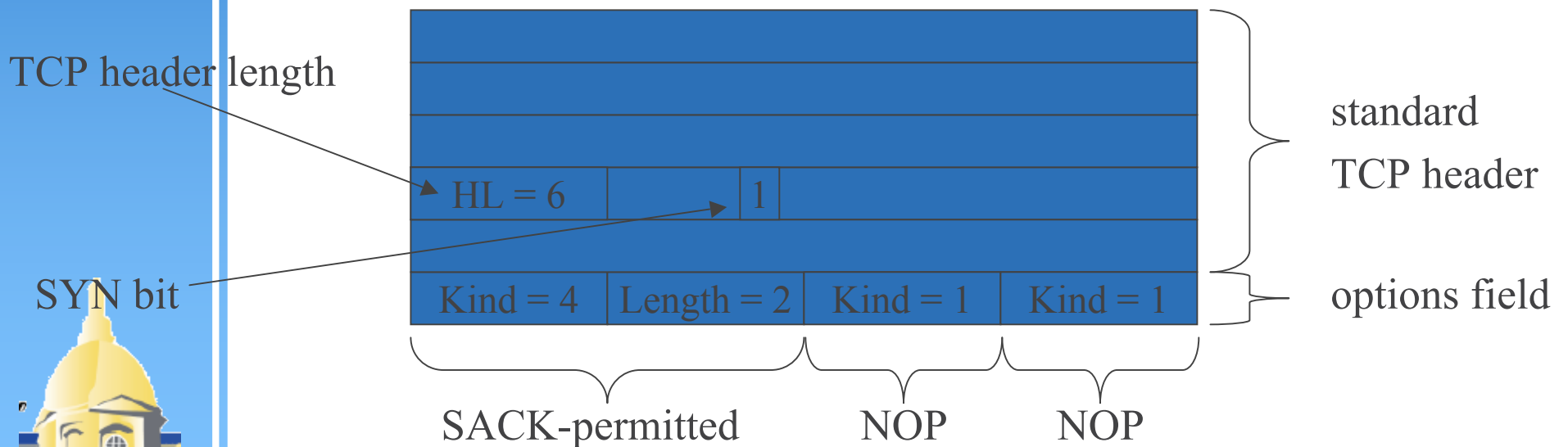
# Selective Acknowledgment TCP

- ▶ Selective Acknowledgment (SACK) allows the receiver to inform the sender about all segments that have been successfully received.
- ▶ Allows the sender to retransmit only those segments that have been lost.
- ▶ SACK is implemented using two different TCP options.



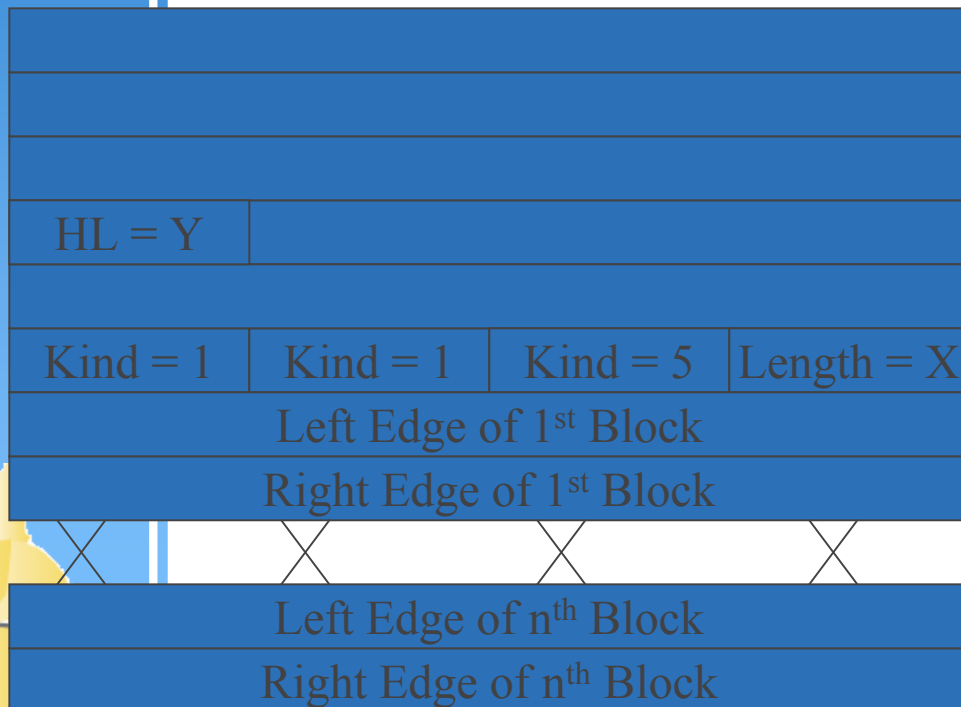
# The SACK-Permitted Option

- ▶ The first TCP option is the enabling option, “SACK-permitted,” allowed only in a SYN segment.
- ▶ This indicates that the sender can handle SACK data and the receiver should send it, if possible. (Both sides can enable SACK, but each direction of the TCP connection is treated independently.)



# The SACK Option

- ▶ If the SACK-permitted option is received, the receiver may send the SACK option.



What is a simple formula for the SACK option length field (based on n, the number of blocks in the option)?

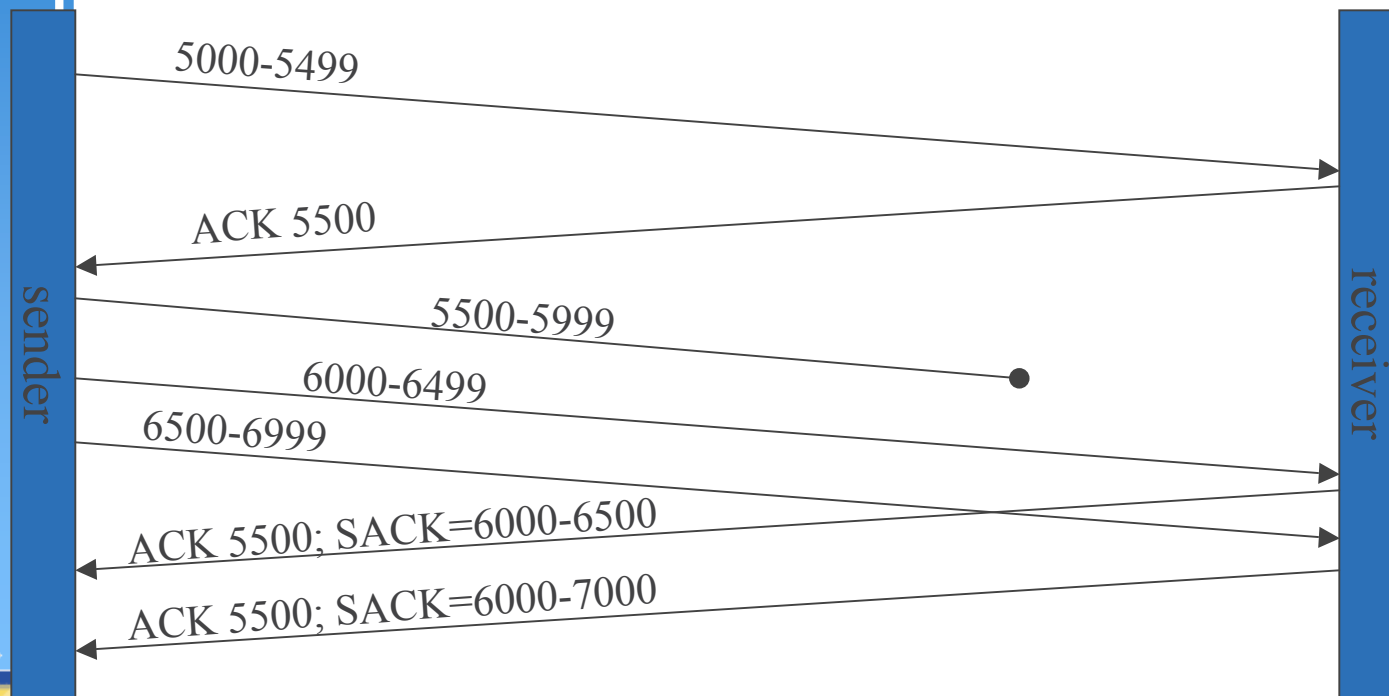
**$(2 + 8 * n)$  bytes**

What is the maximum number of SACK blocks possible? Why?

The maximum size of the options field is 40 bytes, giving a maximum of **4** SACK blocks (barring no other TCP options).

# The SACK Option

- ▶ Each block in a SACK represents bytes successfully received that are contiguous and isolated (the bytes immediately to the left and the right have not yet been received).



# 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.

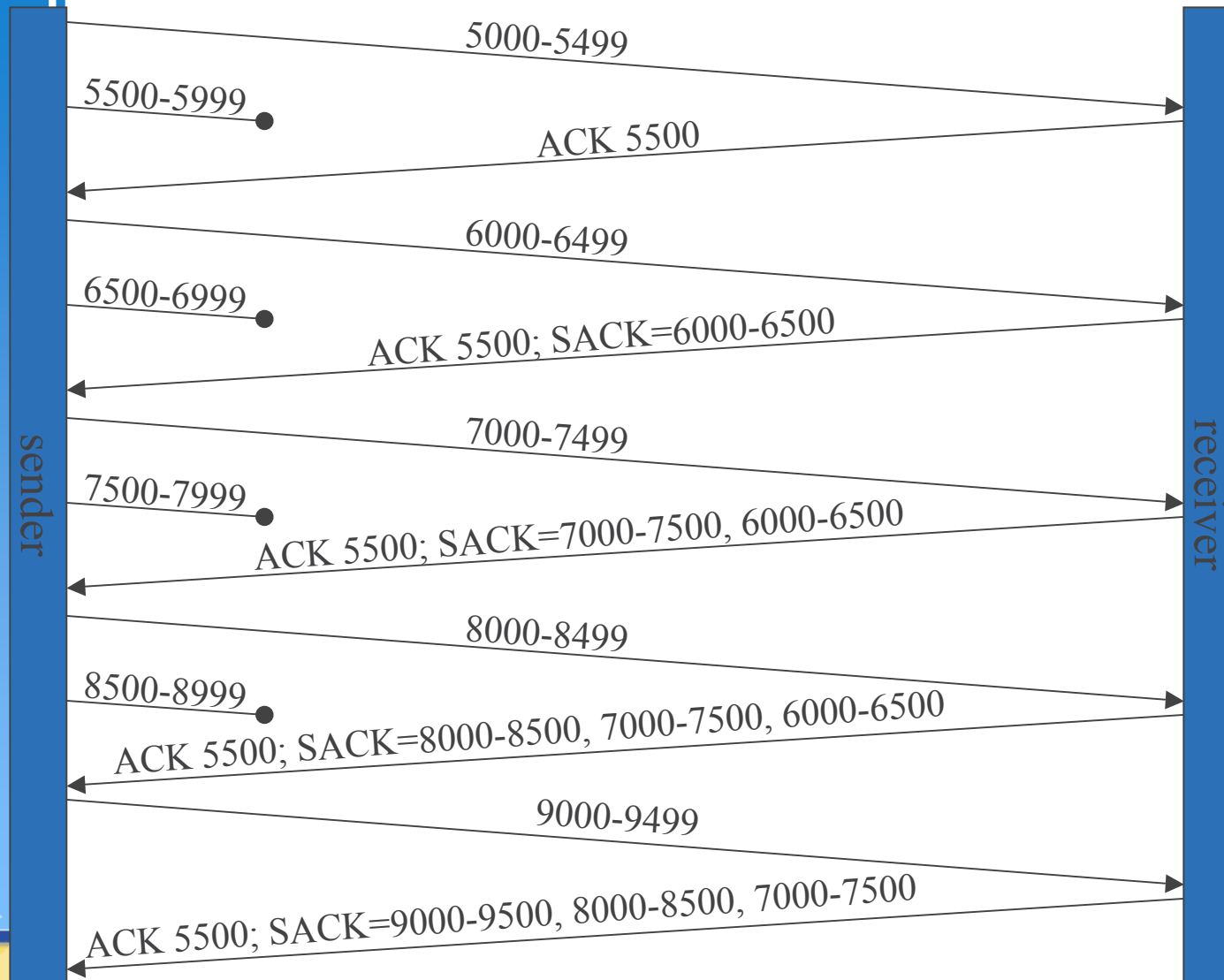


# SACK TCP Rules

- ▶ The first SACK block must contain the most recently received segment that is to be SACKed.
- ▶ The second block must contain the second most recently received segment that is to be SACKed, and so forth.
- ▶ Notice this can result in some data in the receiver's buffers which should be SACKed but is not (if there are more segments to SACK than available space in the TCP header).



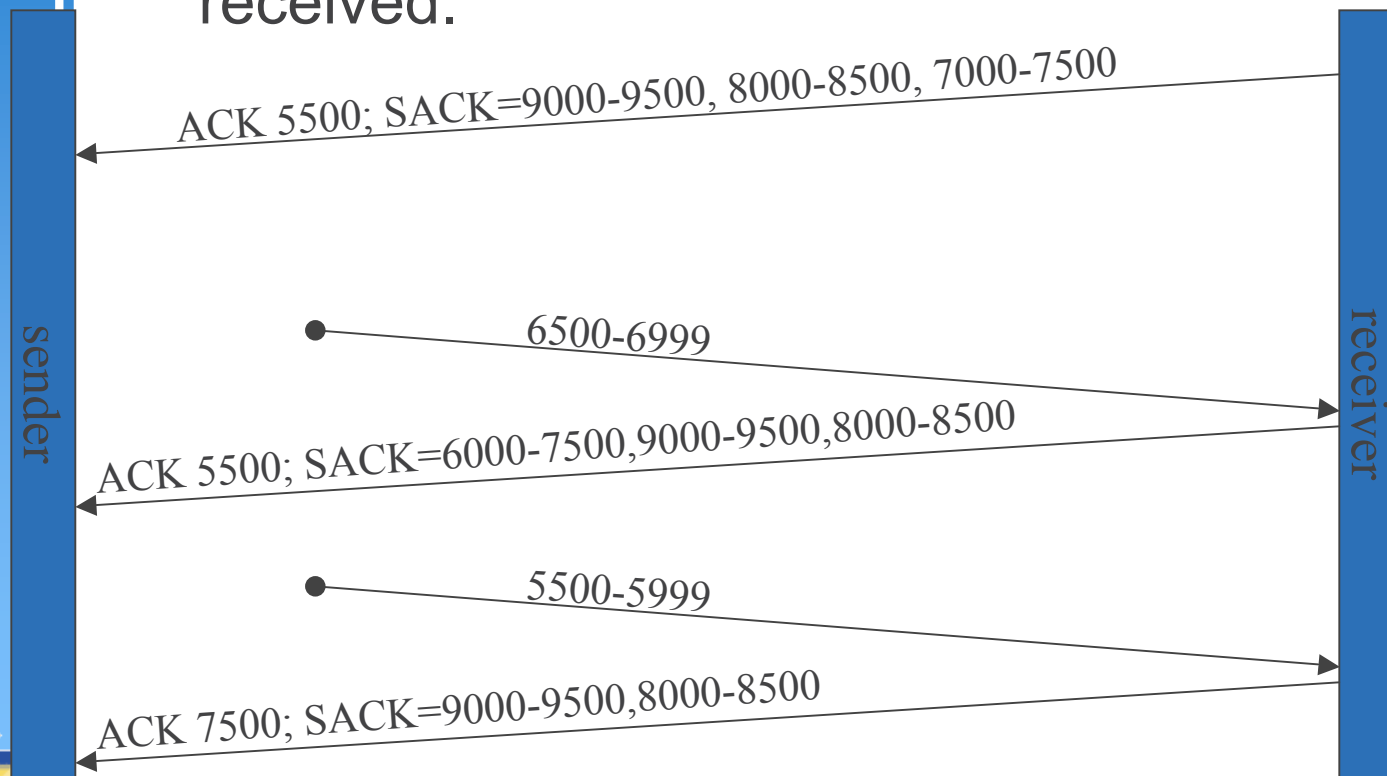
# SACK TCP Example (assuming a maximum of 3 blocks)





# SACK TCP Example (continued)

- ▶ At this point, the 4th segment (6500-6999) is received. After the receiver acknowledges this reception, the 2nd segment (5500-5999) is received.



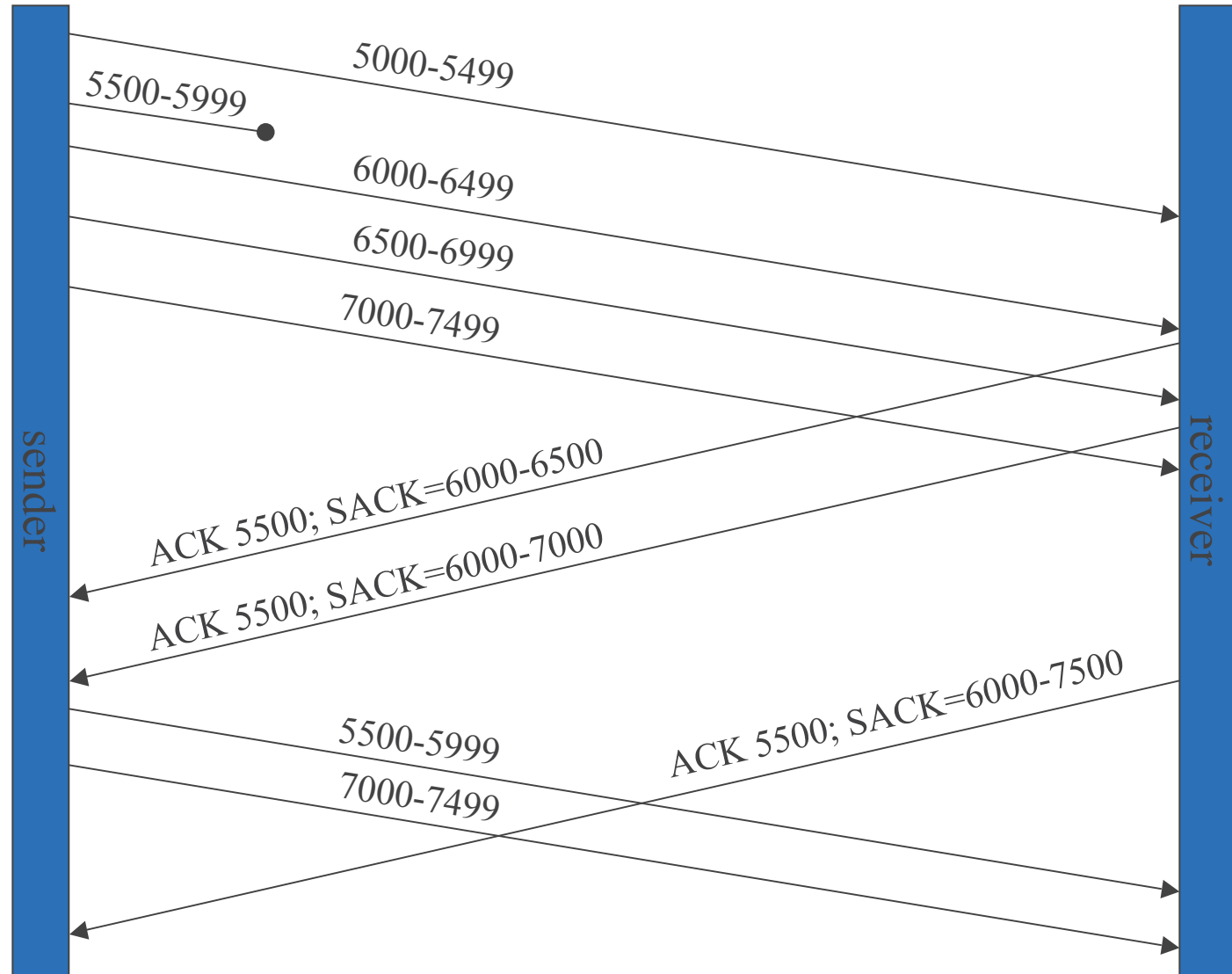
# What Should the Sender do?

- ▶ The sender must keep a buffer of unacknowledged data. When it receives a SACK option, it should turn on a SACK-flag bit for all segments in the transmit buffer that are wholly contained within one of the SACK blocks.
- ▶ After this SACK flag bit has been turned on, the sender should skip that segment during any later retransmission.

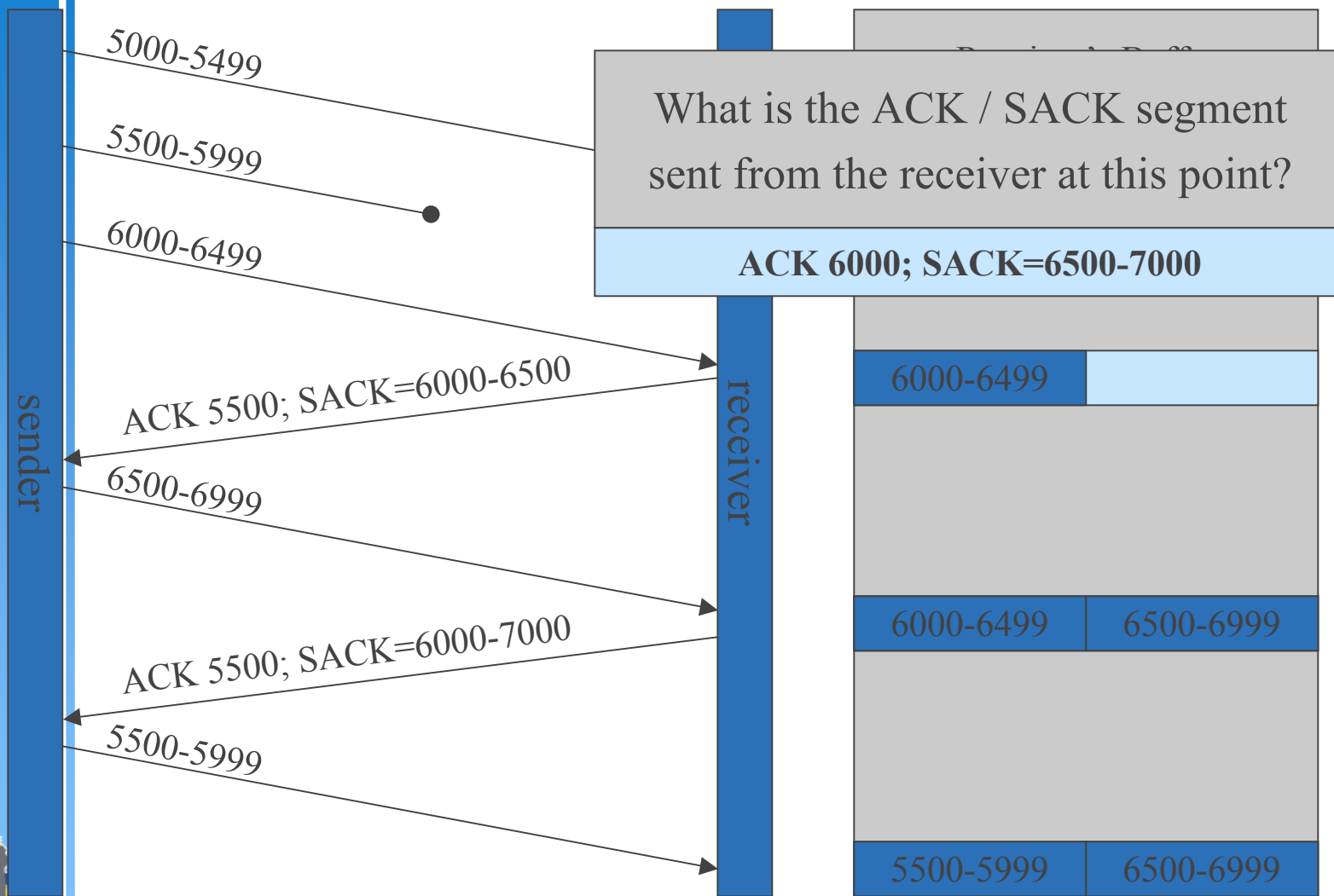


# SACK TCP at the Sender Example

SENDER  
TIMEOUT



# Receiver Has A Two-Segment Buffer (A Problem?)



# Reneging in SACK TCP

- ▶ It is possible for the receiver to SACK some data and then later discard it. This is referred to as reneging. This is discouraged, but permitted if the receiver runs out of buffer space.
- ▶ If this occurs,
  - The first SACK block must still reflect the newest segment, i.e. contain the left and right edges of the newest segment, even if that segment is going to be discarded.
  - Except for the newest segment, all SACK blocks must not report any old data that has been discarded.



# 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 Observations

- ▶ 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

- ▶ 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).





# 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.



# D-SACK TCP

(RFC 2883)

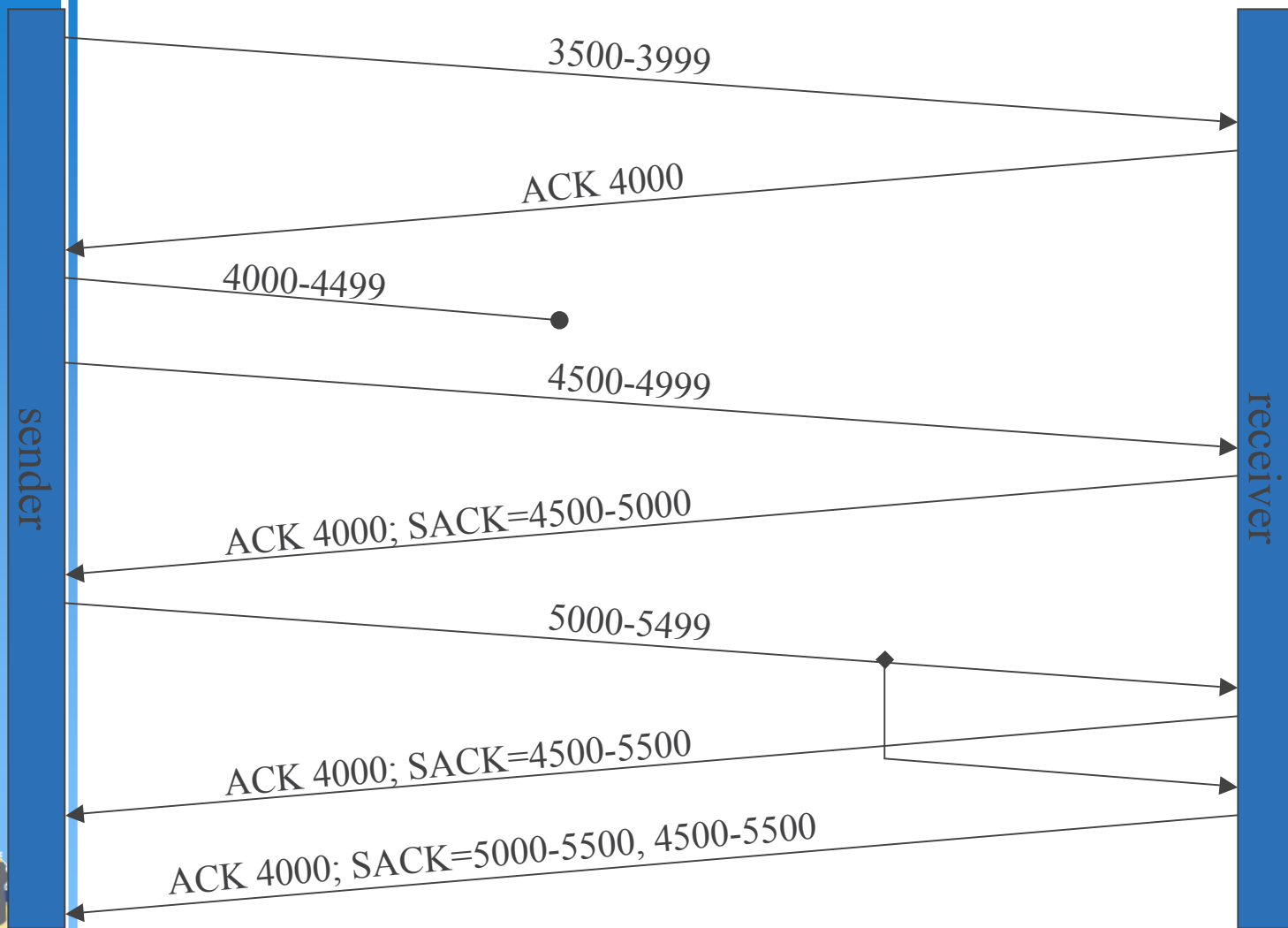


# One Step Further: D-SACK TCP

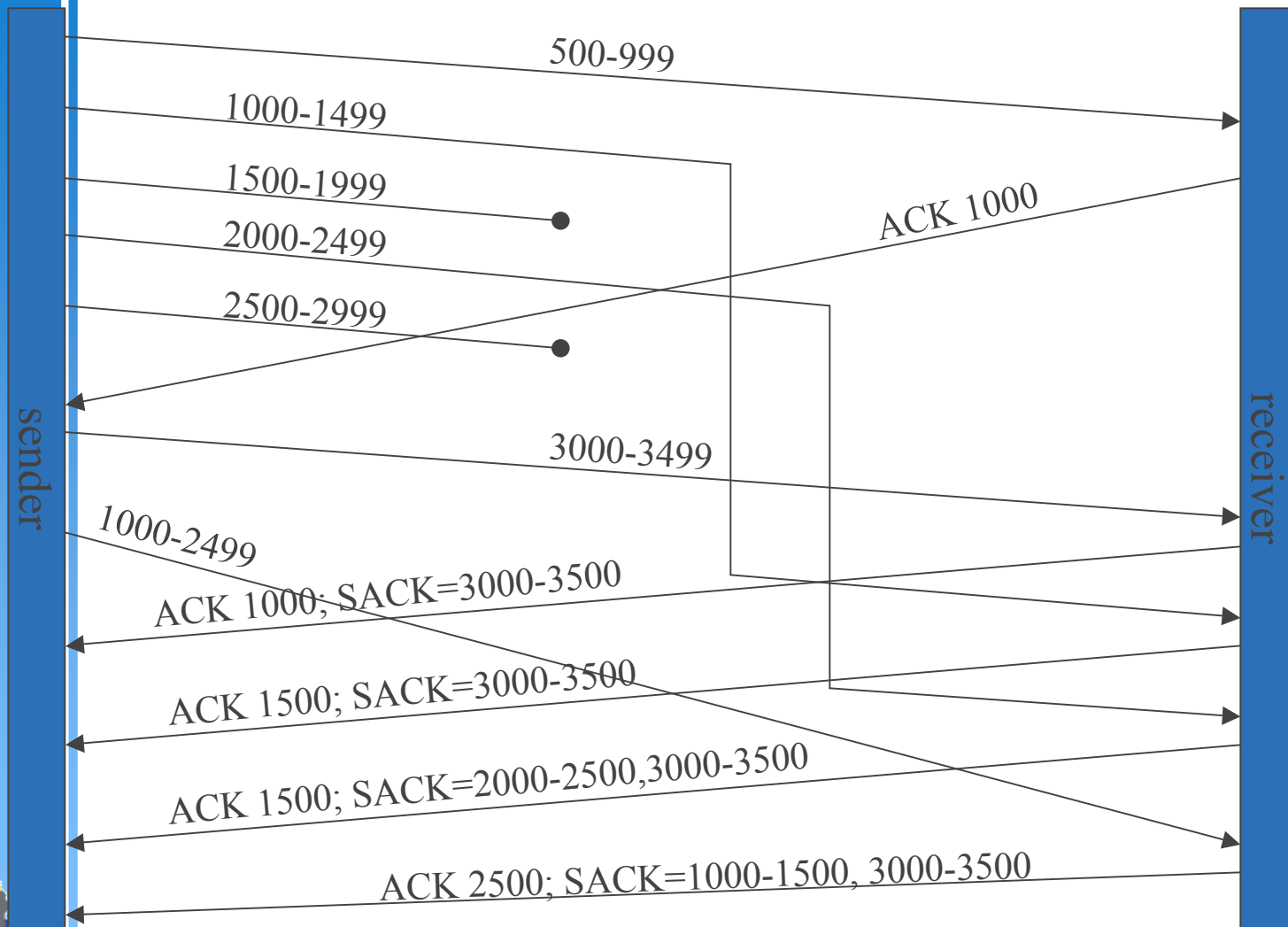
- ▶ Duplicate-SACK, or D-SACK is an extension to SACK TCP which uses the first block of a SACK option is used to report duplicate segments that have been received.
- ▶ A D-SACK block is only used to report a duplicate contiguous sequence of data received by the receiver in the most recent segment.
- ▶ Each duplicate is reported at most once.
- ▶ This allows the sender TCP to determine when a retransmission was not necessary. It may not have been necessary due to the retransmit timer expiring prematurely or due to a false Fast Retransmit (3 duplicate ACKs received due to network reordering).



# D-SACK Example (packet replicated by the network)



# D-SACK Example (losses, and the sender changes the segment size)



# 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.



# D-SACK TCP and Retransmissions

- ▶ D-SACK allows TCP to determine when a retransmission was not necessary (it receives a D-SACK after it retransmitted a segment). When this determination is made, the sender can “undo” the halving of the congestion window, as it will do when a segment is retransmitted (as it assumes net congestion).
- ▶ D-SACK also allows TCP to determine if the network is duplicating packets (it will receive a D-SACK for a segment it only sent once).
- ▶ D-SACK’s weakness is that it does not allow a sender to determine if both the original and retransmitted segment are received, or the original is lost and the retransmitted segment is duplicated by the network.



# 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 inherent 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).





# Increasing the Maximum TCP Initial Window Size

(RFC 2414)



# Increasing the Initial Window

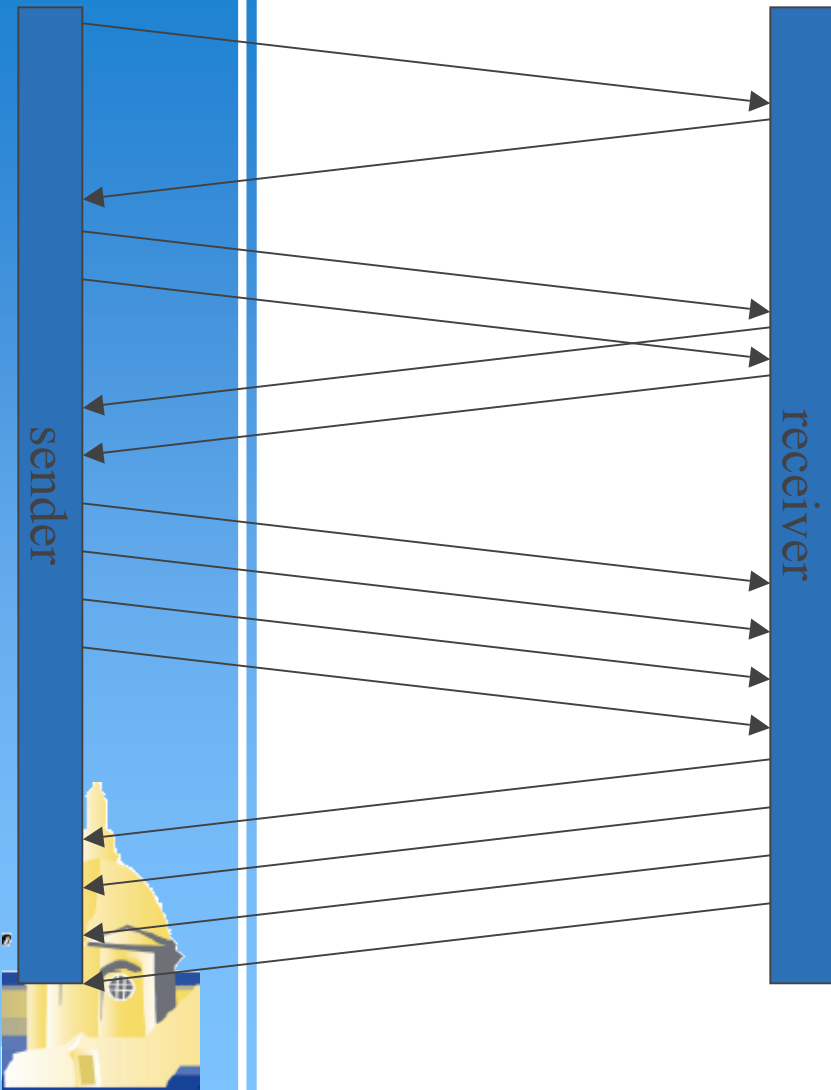
- ▶ RFC 2414 specifies an experimental change to TCP, the increasing of the maximum initial window size, from one segment to a larger value.
- ▶ This new larger value is given as:  
$$\min ( 4 * \text{MSS}, \max ( 2 * \text{MSS}, 4380 \text{ bytes} ) )$$
- ▶ This translates to:

Maximum Segment Size (MSS)	Maximum Initial Window Size
$\leq 1095$ bytes	$\leq 4 * \text{MSS}$
$1095 \text{ bytes} < \text{MSS} < 2190$ bytes	$\leq 4380$ bytes
$\geq 2190$ bytes	$\leq 2 * \text{MSS}$

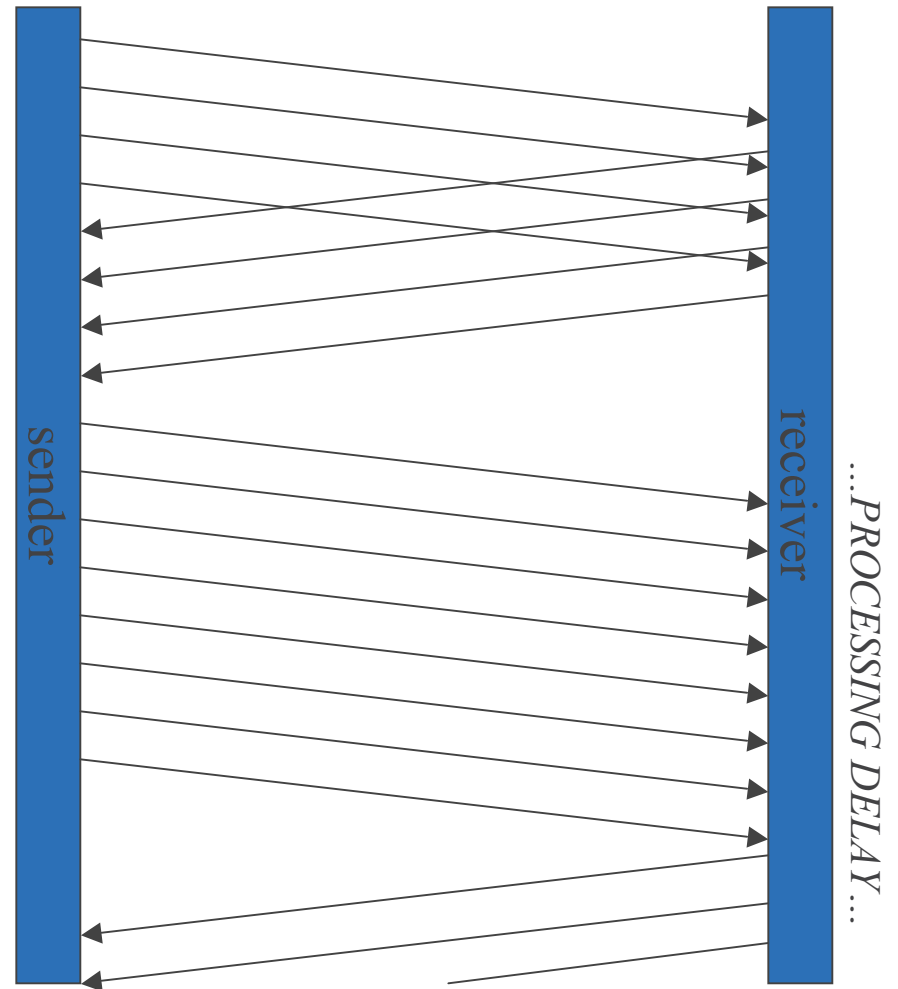


# Increasing the Initial Window

## Slow-Start TCP



## RFC 2414 TCP



# Advantages of an Increased Initial Window Size

- ▶ This change is in contrast to the slow start mechanism, which initializes the initial window size to one segment. This mechanism is in place to implement sender-based congestion control (see RFC 2001 for a complete discussion).
- ▶ This new larger window offers three distinct advantages:
  - With slow start, a receiver which uses delayed ACKs is forced to wait for a timeout before generating an ACK. With an initial window of at least two segments, the receiver will generate an ACK after the second segment arrives, causing a speedup in data acknowledgement.



# Advantages of an Increased Initial Window Size

- For TCP connections transferring a small amount of data (such as SMTP and HTTP requests), the larger initial window will reduce the transmission time, as more data can be outstanding at once.
- For TCP connections transferring a large amount of data with high propagation delays (long haul pipes; such as backbone connects and satellite links), this change eliminates up to three round-trip times (RTTs) and a delayed ACK timeout during the initial slow start.



# 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).



# 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 sender-side congestion control window (even if the experimental higher initial window is considered).



# Summary

- ▶ SACK TCP provides additional information to the sender, allowing the reduction of needless retransmissions. There is no danger in providing this information, it simply serves to make a “smarter” TCP sender.
- ▶ D-SACK TCP allows the sender to determine when it has needlessly resent segments. This will allow the sender to continuously refine its retransmission strategy and undo unnecessary and incorrect congestion control mechanisms.
- ▶ Increasing the initial TCP window is a slight change that has advantages for both small and large data transfers, without significantly affecting the congestion control a smaller window provides.

