

Flow Control

- Fast sender can overrun receiver:
 - Packet loss, unnecessary retransmissions
- Possible solutions:
 - Sender transmits at pre-negotiated rate
 - Sender limited to a window's worth of unacknowledged data
- Flow control different from congestion control



Mar-23-04

4/598N: Computer Networks

Flow Control

- Send buffer size: MaxSendBuffer
- Receive buffer size: MaxRcvBuffer
- Receiving side
 - $\text{LastByteRcvd} - \text{LastByteRead} \leq \text{MaxRcvBuffer}$
 - $\text{AdvertisedWindow} = \text{MaxRcvBuffer} - (\text{NextByteExpected} - \text{NextByteRead})$
- Sending side
 - $\text{LastByteSent} - \text{LastByteAcked} \leq \text{AdvertisedWindow}$
 - $\text{EffectiveWindow} = \text{AdvertisedWindow} - (\text{LastByteSent} - \text{LastByteAcked})$
 - $\text{LastByteWritten} - \text{LastByteAcked} \leq \text{MaxSendBuffer}$
 - block sender if $(\text{LastByteWritten} - \text{LastByteAcked}) + y > \text{MaxSenderBuffer}$
- Always send ACK in response to arriving data segment
- Persist when $\text{AdvertisedWindow} = 0$



Mar-23-04

4/598N: Computer Networks

Round-trip Time Estimation

- Wait at least one RTT before retransmitting
- Importance of accurate RTT estimators:
 - Low RTT -> unneeded retransmissions
 - High RTT -> poor throughput
- RTT estimator must adapt to change in RTT
 - But not too fast, or too slow!



Mar-23-04

4/598N: Computer Networks

Initial Round-trip Estimator

Round trip times exponentially averaged:

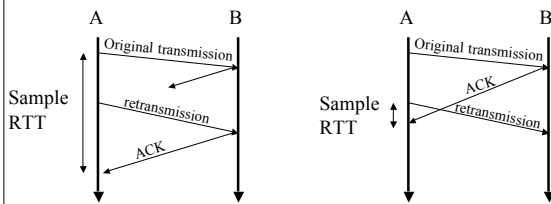
- New RTT = α (old RTT) + $(1 - \alpha)$ (new sample)
- Recommended value for α : 0.8 - 0.9
- Retransmit timer set to β RTT, where $\beta = 2$
- Every time timer expires, RTO exponentially backed-off



Mar-23-04

4/598N: Computer Networks

Retransmission Ambiguity



Mar-23-04

4/598N: Computer Networks

Karn's Retransmission Timeout Estimator

- Accounts for retransmission ambiguity
- If a segment has been retransmitted:
 - Don't count RTT sample on ACKs for this segment
 - Keep backed off time-out for next packet
 - Reuse RTT estimate only after one successful transmission



Mar-23-04

4/598N: Computer Networks

Karn/Partridge Algorithm

- Do not sample RTT when retransmitting
- Double timeout after each retransmission

The diagram shows two scenarios between a Sender and a Receiver. In the first, an 'Original transmission' is sent, followed by an 'ACK' received. A vertical line marks the start of a 'Sample RTT' period. In the second scenario, an 'Original transmission' is sent but no 'ACK' is received. A 'Retransmission' is sent. A second 'ACK' is received, but the 'Sample RTT' period is shown as a longer interval, indicating a doubled timeout.

Mar-23-04 4/598N: Computer Networks

Jacobson's Retransmission Timeout Estimator

- Key observation:
 - Using β RTT for timeout doesn't work
 - At high loads round trip variance is high
- Solution:
 - If D denotes mean variation
 - Timeout = RTT + 4D

Mar-23-04 4/598N: Computer Networks

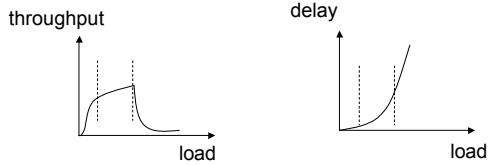
Congestion

The diagram shows a network topology where two input links with capacities of 10 Mbps and 100 Mbps converge into a single output link with a capacity of 1.5 Mbps. This configuration is prone to congestion.

- If both sources send full windows, we may get congestion collapse
- Other forms of congestion collapse:
 - Retransmissions of large packets after loss of a single fragment
 - Non-feedback controlled sources

Mar-23-04 4/598N: Computer Networks

Congestion Response



Avoidance keeps the system performing at the *knee*
Control kicks in once the system has reached a congested state



Mar-23-04

4/598N: Computer Networks

Separation of Functionality

- Sending host must adjust amount of data it puts in the network based on detected congestion
- Routers can help by:
 - Sending accurate congestion signals
 - Isolating well-behaved from ill-behaved sources



Mar-23-04

4/598N: Computer Networks

6.3 TCP Congestion Control

- Idea
 - assumes best-effort network (FIFO or FQ routers) each source determines network capacity for itself
 - uses implicit feedback
 - ACKs pace transmission (*self-clocking*)
- Challenge
 - determining the available capacity in the first place
 - adjusting to changes in the available capacity




Mar-23-04

4/598N: Computer Networks


TCP Congestion Control

- A collection of interrelated mechanisms:
 - Slow start
 - Congestion avoidance
 - Accurate retransmission timeout estimation
 - Fast retransmit
 - Fast recovery

 Mar-23-04 4/598N: Computer Networks


Congestion Control

- Underlying design principle: packet conservation
 - At equilibrium, inject packet into network only when one is removed
 - Basis for stability of physical systems
- A mechanism which:
 - Uses network resources efficiently
 - Preserves fair network resource allocation
 - Prevents or avoids collapse
- Congestion collapse is not just a theory
 - Has been frequently observed in many networks

 Mar-23-04 4/598N: Computer Networks

Congestion Under Infinite Buffering

- Nagle (RFC 970) showed that congestion will not go away even with infinite buffers
- Basic argument
 - A datagram network must have TTL
 - With infinite buffering queuing delays increase
 - Even if buffers are not dropped for lack of buffering, they will be dropped because TTL expires

 Mar-23-04 4/598N: Computer Networks

Additive Increase/Multiplicative Decrease

- Objective: adjust to changes in the available capacity
- New state variable per connection: `CongestionWindow`
 - limits how much data source has in transit

$$\text{MaxWin} = \text{MIN}(\text{CongestionWindow}, \text{AdvertisedWindow})$$

$$\text{EffWin} = \text{MaxWin} - (\text{LastByteSent} - \text{LastByteAcked})$$

- Idea:
 - increase `CongestionWindow` when congestion goes down
 - decrease `CongestionWindow` when congestion goes up



Mar-23-04

4/598N: Computer Networks

AIMD (cont)

- Question: how does the source determine whether or not the network is congested?
- Answer: a timeout occurs
 - timeout signals that a packet was lost
 - packets are seldom lost due to transmission error
 - lost packet implies congestion

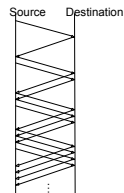


Mar-23-04

4/598N: Computer Networks

AIMD (cont)

- Algorithm
 - increment `CongestionWindow` by one packet per RTT (*linear increase*)
 - divide `CongestionWindow` by two whenever a timeout occurs (*multiplicative decrease*)



- In practice: increment a little for each ACK
$$\text{Increment} = (\text{MSS} * \text{MSS}) / \text{CongestionWindow}$$
$$\text{CongestionWindow} += \text{Increment}$$

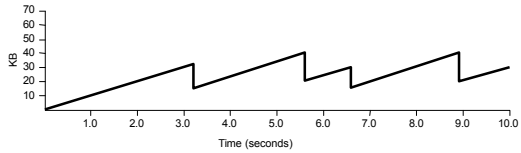


Mar-23-04

4/598N: Computer Networks

AIMD (cont)

- Trace: sawtooth behavior



Mar-23-04

4/598N: Computer Networks

Self-clocking

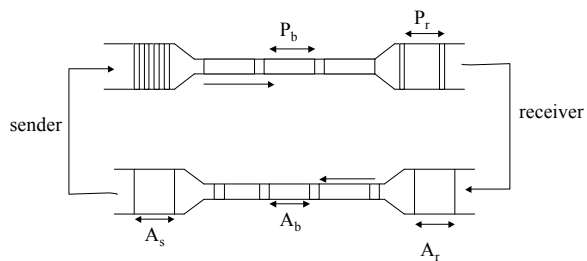
- If we have large actual window, should we send data in one shot?
 - No, use acks to clock sending new data



Mar-23-04

4/598N: Computer Networks

..Self-clocking

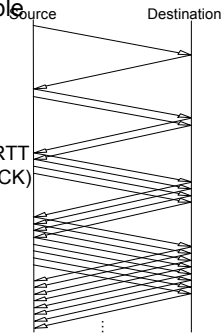


Mar-23-04

4/598N: Computer Networks

Slow Start

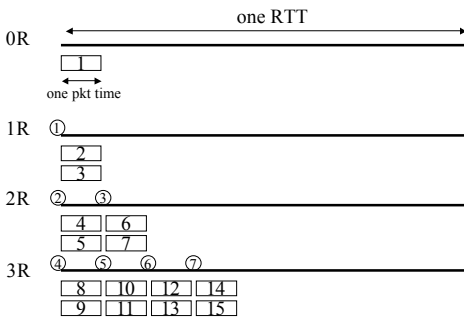
- Objective: determine the available capacity in the first
- Idea:
 - begin with CongestionWindow = 1 packet
 - double CongestionWindow each RTT (increment by 1 packet for each ACK)



Mar-23-04

4/598N: Computer Networks

Slow Start Example

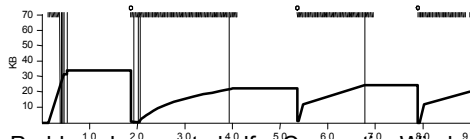


Mar-23-04

4/598N: Computer Networks

Slow Start (cont)

- Exponential growth, but slower than all at once
- Used...
 - when first starting connection
 - when connection goes dead waiting for timeout
- Trace



- Problem: lose up to half a CongestionWindow's worth of data



Mar-23-04

4/598N: Computer Networks

Congestion Avoidance

- 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



Mar-23-04

4/598N: Computer Networks

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



Mar-23-04

4/598N: Computer Networks

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)



Mar-23-04

4/598N: Computer Networks

Fast Retransmit and Fast Recovery

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

Sender

Packet 1
Packet 2
Packet 3
Packet 4

Packet 5
Packet 6

Retransmit packet 3

Receiver

ACK 1
ACK 2
ACK 2
ACK 2
ACK 2

ACK 6

The diagram illustrates a sequence of events between a Sender and a Receiver. The Sender transmits Packet 1, Packet 2, Packet 3, and Packet 4. Packet 3 is marked with an 'X', indicating it is lost. The Receiver receives ACK 1, then ACK 2, followed by three more duplicate ACK 2s. This triggers the Sender to retransmit Packet 3. After the retransmission, the Receiver receives ACK 6, indicating that Packets 1 through 6 have been successfully received.

Mar-23-04
4/598N: Computer Networks

Fast Retransmit and Recovery

- If we get 3 duplicate acks for segment N
 - Retransmit segment N
 - Set ssthresh to $0.5 * cwnd$
 - 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)

Mar-23-04
4/598N: Computer Networks

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

Mar-23-04
4/598N: Computer Networks

Results

- Fast recovery
 - skip the slow start phase
 - go directly to half the last successful `congestionWindow` (`ssthresh`)

Mar-23-04
4/598N: Computer Networks

TCP Extensions

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

Mar-23-04
4/598N: Computer Networks

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

Mar-23-04
4/598N: Computer Networks

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



Mar-23-04

4/598N: Computer Networks

Keeping the Pipe Full

- 16-bit AdvertisedWindow

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



Mar-23-04

4/598N: Computer Networks

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



Mar-23-04

4/598N: Computer Networks

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)



Mar-23-04

4/598N: Computer Networks

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)
 - Header prediction (to improve performance)



Mar-23-04

4/598N: Computer Networks

SACK TCP

(RFC 2018)



Mar-23-04

4/598N: Computer Networks

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.



Mar-23-04

4/598N: Computer Networks

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.

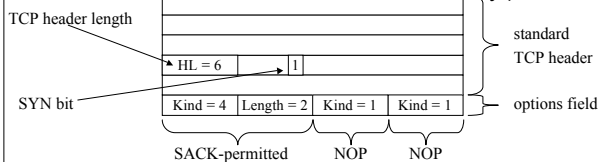


Mar-23-04

4/598N: Computer Networks

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.)



Mar-23-04

4/598N: Computer Networks

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).

Mar-23-04
4/598N: Computer Networks

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).

sender

receiver

Mar-23-04
4/598N: Computer Networks

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-23-04
4/598N: Computer Networks

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).

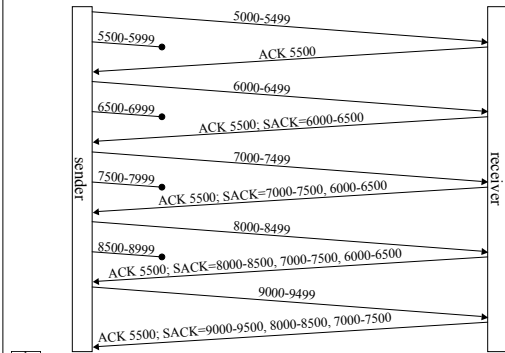


Mar-23-04

4/598N: Computer Networks

SACK TCP Example

(assuming a maximum of 3 blocks)

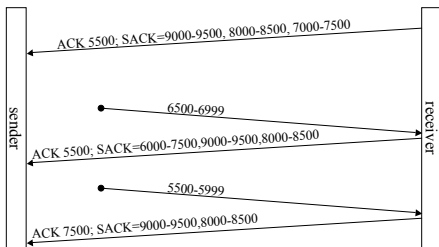


Mar-23-04

4/598N: Computer Networks

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.



Mar-23-04

4/598N: Computer Networks

What Should the Sender do?

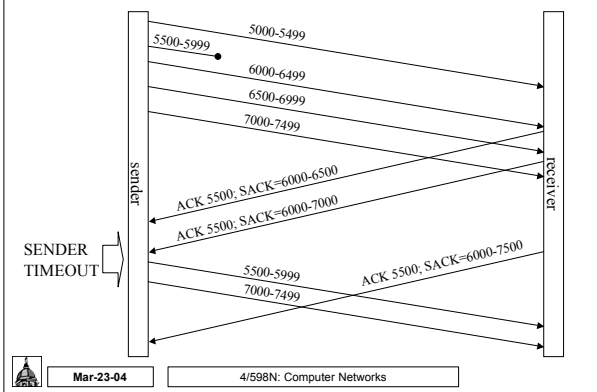
- The sender must keep a buffer of unacknowledged data. When it receives a SACK option, it should turn on a SACK-flag bit for all segments in the transmit buffer that are wholly contained within one of the SACK blocks.
- After this SACK flag bit has been turned on, the sender should skip that segment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



Mar-23-04

4/598N: Computer Networks

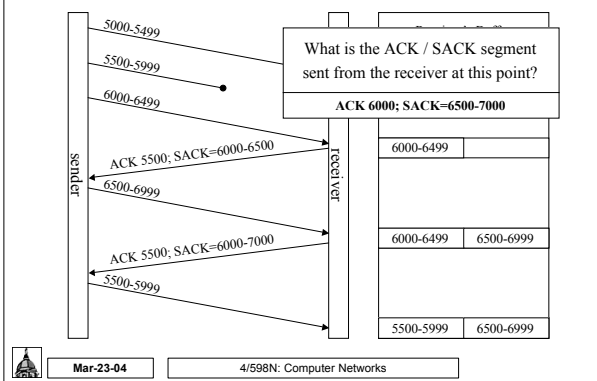
SACK TCP at the Sender Example



Mar-23-04

4/598N: Computer Networks

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



Mar-23-04

4/598N: Computer Networks

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.



Mar-23-04

4/598N: Computer Networks

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.



Mar-23-04

4/598N: Computer Networks

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.




Mar-23-04

4/598N: Computer Networks


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).

 Mar-23-04 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-23-04 4/598N: Computer Networks

D-SACK TCP

(RFC 2883)

 Mar-23-04 4/598N: Computer Networks

One Step Further: D-SACK TCP

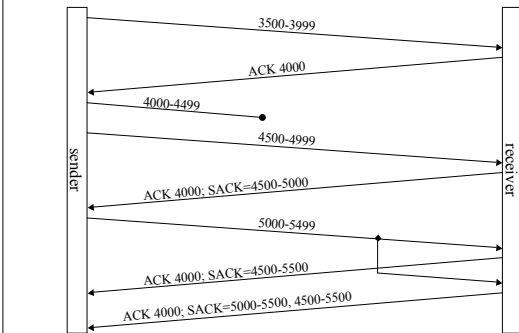
- Duplicate-SACK, or D-SACK is an extension to SACK TCP which uses the first block of a SACK option is used to report duplicate segments that have been received.
- A D-SACK block is only used to report a duplicate contiguous sequence of data received by the receiver in the most recent segment.
- Each 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).



Mar-23-04

4/598N: Computer Networks

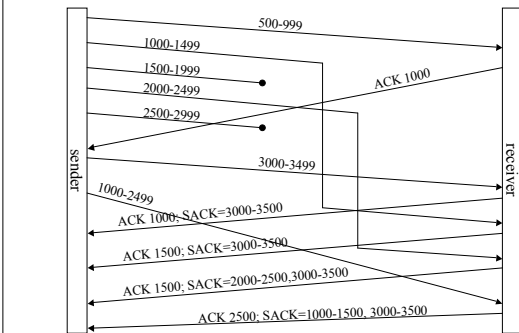
D-SACK Example (packet replicated by the network)



Mar-23-04

4/598N: Computer Networks

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



Mar-23-04

4/598N: Computer Networks

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.



Mar-23-04

4/598N: Computer Networks

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.



Mar-23-04

4/598N: Computer Networks

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).



Mar-23-04

4/598N: Computer Networks

Increasing the Maximum TCP Initial Window Size

(RFC 2414)

Mar-23-04
4/598N: Computer Networks

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:
 $\min(4 * MSS, \max(2 * MSS, 4380 \text{ bytes}))$
- This translates to:

Maximum Segment Size (MSS)	Maximum Initial Window Size
<= 1095 bytes	<= 4 * MSS
1095 bytes < MSS < 2190 bytes	<= 4380 bytes
>= 2190 bytes	<= 2 * MSS

Mar-23-04
4/598N: Computer Networks

Increasing the Initial Window

Slow-Start TCP

RFC 2414 TCP

Mar-23-04
4/598N: Computer Networks

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.



Mar-23-04

4/598N: Computer Networks

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.



Mar-23-04

4/598N: Computer Networks

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-23-04

4/598N: Computer Networks

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).



Mar-23-04

4/598N: Computer Networks

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.



Mar-23-04

4/598N: Computer Networks

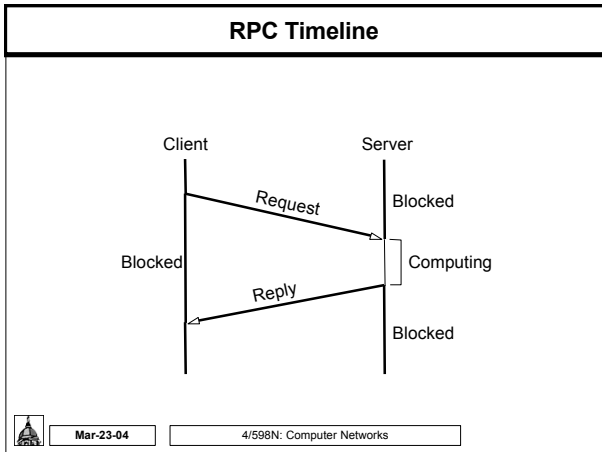
Remote Procedure Call

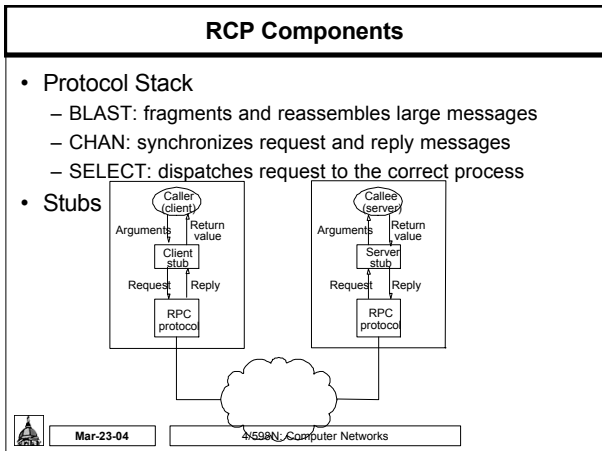
- Outline
 - Protocol Stack
 - Presentation Formatting

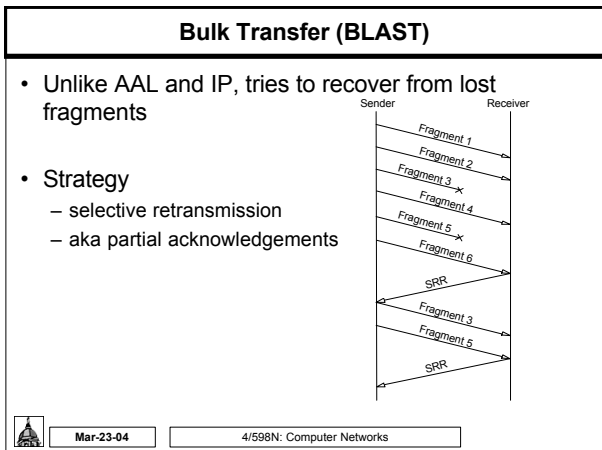


Mar-23-04

4/598N: Computer Networks







BLAST Details

- Sender:
 - after sending all fragments, set timer DONE
 - if receive SRR, send missing fragments and reset DONE
 - if timer DONE expires, free fragments



Mar-23-04

4/598N: Computer Networks

BLAST Details (cont)

- Receiver:
 - when first fragments arrives, set timer LAST_FRAG
 - when all fragments present, reassemble and pass up
 - four exceptional conditions:
 - if last fragment arrives but message not complete
 - send SRR and set timer RETRY
 - if timer LAST_FRAG expires
 - send SRR and set timer RETRY
 - if timer RETRY expires for first or second time
 - send SRR and set timer RETRY
 - if timer RETRY expires a third time
 - give up and free partial message

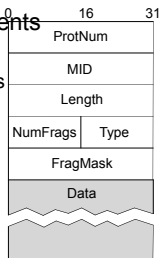


Mar-23-04

4/598N: Computer Networks

BLAST Header Format

- MID must protect against wrap around
- TYPE = DATA or SRR
- NumFrag indicates number of fragments
- FragMask distinguishes among fragments
 - if Type=DATA, identifies this fragment
 - if Type=SRR, identifies missing fragments



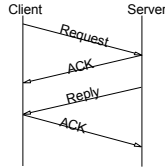
Mar-23-04

4/598N: Computer Networks

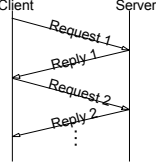
Request/Reply (CHAN)

- Guarantees message delivery
- Synchronizes client with server
- Supports at-most-once semantics

• Simple case



Implicit Acks



Mar-23-04

4/598N: Computer Networks

CHAN Details

- Lost message (request, reply, or ACK)
 - set RETRANSMIT timer
 - use message id (MID) field to distinguish
- Slow (long running) server
 - client periodically sends “are you alive” probe, or
 - server periodically sends “I’m alive” notice
- Want to support multiple outstanding calls
 - use channel id (CID) field to distinguish
- Machines crash and reboot
 - use boot id (BID) field to distinguish



Mar-23-04

4/598N: Computer Networks

CHAN Header Format

```
typedef struct {
    u_short Type; /* REQ, REP, ACK, PROBE */
    u_short CID; /* unique channel id */
    int MID; /* unique message id */
    int BID; /* unique boot id */
    int Length; /* length of message */
    int ProtNum; /* high-level protocol */
} ChanHdr;

typedef struct {
    u_char type; /* CLIENT or SERVER */
    u_char status; /* BUSY or IDLE */
    int retries; /* number of retries */
    int timeout; /* timeout value */
    XkReturn ret_val; /* return value */
    Msg *request; /* request message */
    Msg *reply; /* reply message */
    Semaphore reply_sem; /* client semaphore */
    int mid; /* message id */
    int bid; /* boot id */
} ChanState;
```



Mar-23-04

4/598N: Computer Networks

Synchronous vs Asynchronous Protocols

- Asynchronous interface

```
xPush(Sessn s, Msg *msg)
xPop(Sessn s, Msg *msg, void *hdr)
xDemux(Protl hlp, Sessn s, Msg *msg)
```

- Synchronous interface

```
xCall(Sessn s, Msg *req, Msg *rep)
xCallPop(Sessn s, Msg *req, Msg *rep, void *hdr)
xCallDemux(Protl hlp, Sessn s, Msg *req, Msg *rep)
```

- CHAN is a hybrid protocol

- synchronous from above: **xCall**
- asynchronous from below: **xPop/xDemux**



Mar-23-04

4/598N: Computer Networks

```
chanCall(Sessn self, Msg *msg, Msg *rmsg){
    ChanState *state = (ChanState *)self->state;
    ChanHdr    *hdr;
    char       *buf;

    /* ensure only one transaction per channel */
    if ((state->status != IDLE))
        return XK_FAILURE;
    state->status = BUSY;

    /* save copy of req msg and ptr to rep msg*/
    msgConstructCopy(&state->request, msg);
    state->reply = rmsg;
    /* fill out header fields */
    hdr = state->hdr_template;
    hdr->Length = msgLen(msg);
    if (state->mid == MAX_MID)
        state->mid = 0;
    hdr->MID = ++state->mid;
}
```



Mar-23-04

4/598N: Computer Networks

```
/* attach header to msg and send it */
buf = msgPush(msg, HDR_LEN);
chan_hdr_store(hdr, buf, HDR_LEN);
xPush(xGetDown(self, 0), msg);

/* schedule first timeout event */
state->retries = 1;
state->event = evSchedule(retransmit, self, state->timeout);

/* wait for the reply msg */
semWait(&state->reply_sem);

/* clean up state and return */
flush_msg(state->request);
state->status = IDLE;
return state->ret_val;
}
```



Mar-23-04

4/598N: Computer Networks

```

retransmit(Event ev, int *arg){
    Sessn      s = (Sessn)arg;
    ChanState  *state = (ChanState *)s->state;
    Msg        tmp;

    /* see if event was cancelled */
    if ( evIsCancelled(ev) ) return;

    /* unblock client if we've retried 4 times */
    if (++state->retries > 4) {
        state->ret_val = XK_FAILURE;
        semSignal(state->rep_sem);
        return;
    }

    /* retransmit request message */
    msgConstructCopy(&tmp, &state->request);
    xPush(xGetDown(s, 0), &tmp);

    /* reschedule event with exponential backoff */
    evDetach(state->event);
    state->timeout = 2*state->timeout;
    state->event = evSchedule(retransmit, s,
        state->timeout);
}

```

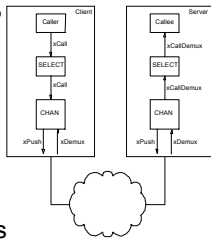


Mar-23-04

4/598N: Computer Networks

Dispatcher (SELECT)

- Dispatch to appropriate procedure
- Synchronous counterpart to UDP



- Address Space for Procedures
 - flat: unique id for each possible procedure
 - hierarchical: program + procedure number



Mar-23-04

4/598N: Computer Networks

Example Code

```

Client side
static XkReturn
selectCall(Sessn self, Msg *req, Msg *rep)
{
    SelectState *state=(SelectState *)self->state;
    char        *buf;

    buf = msgPush(req, HLEN);
    select_hdr_store(state->hdr, buf, HLEN);
    return xCall(xGetDown(self, 0), req, rep);
}

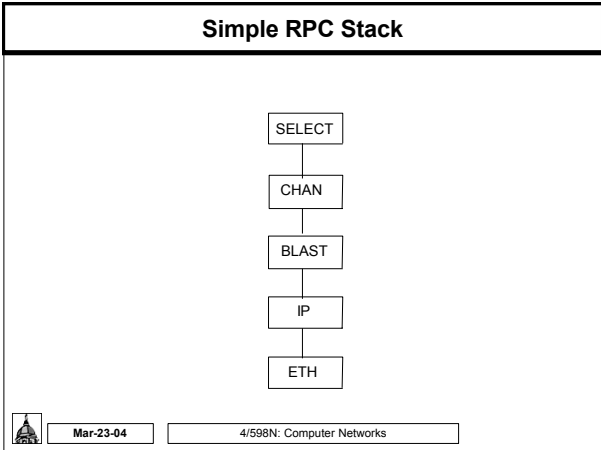
Server side
static XkReturn
selectCallPop(Sessn s, Sessn lls, Msg *req, Msg *rep, void *inHdr)
{
    return xCallDemux(xGetUp(s), s, req, rep);
}

```



Mar-23-04

4/598N: Computer Networks



VCHAN: A Virtual Protocol

```

static XkReturn
vchanCall(Sessn s, Msg *req, Msg *rep)
{
    Sessn      chan;
    XkReturn   result;
    VchanState *state=(VchanState *)s->state;

    /* wait for an idle channel */
    semWait(&state->available);
    chan = state->stack[--state->tos];

    /* use the channel */
    result = xCall(chan, req, rep);

    /* free the channel */
    state->stack[state->tos++] = chan;
    semSignal(&state->available);
    return result;
}
  
```

Mar-23-04 4/598N: Computer Networks

SunRPC

- IP implements BLAST-equivalent
 - except no selective retransmit
- SunRPC implements CHAN-equivalent
 - except not at-most-once
- UDP + SunRPC implement SELECT-equivalent
 - UDP dispatches to program (ports bound to programs)
 - SunRPC dispatches to procedure within program

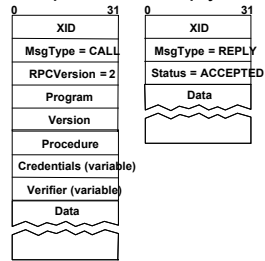
```

graph TD
    SunRPC[SunRPC] --- UDP[UDP]
    UDP --- IP[IP]
    IP --- ETH[ETH]
  
```

Mar-23-04 4/598N: Computer Networks

SunRPC Header Format

- XID (transaction id) is similar to CHAN's MID
- Server does not remember last XID it serviced
- Problem if client retransmits request while reply is in transit



Mar-23-04

4/598N: Computer Networks
