

TCP Implementation

TCP Implementation

- **A quick review of TCP**
- TCP implementation options
- Some other implementation considerations
- Optimizing TCP's performance
- TCP over wireless networks

A byte-stream oriented protocol

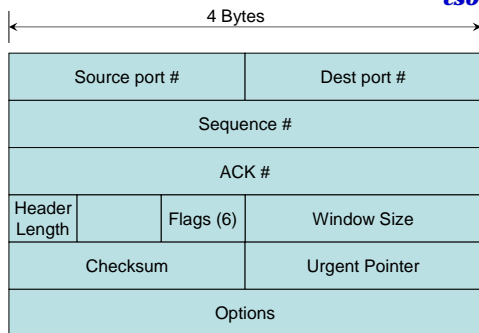
- TCP is designed to treat data as a generic stream of bytes.
- Pays no attention to message boundaries, etc.
- Bytes are sequence-numbered by the sender (32-bit seq number)

TCP transfer protocol

- Extended form of sliding window algorithm
 - Review of SWA:
 - All segs have sequence number
 - Multiple segs can be in flight
 - Sender starts timer when transmitting each seg
 - Receiver sends ACK for seg when all previous segs have been received
 - If send timer times out before seg is ACKed, sender re-sends the seg

Header

cs670



G. W. Cox – Spring 2008

TCP Implementation 5

TCP header flags

cs670

- URG – Indicates seg contains urgent data
- ACK – Indicates this is an ACK seg
- PSH – “Push”. Requests receiver to deliver data to app without buffering
- RST – “Reset”. NACK.
- SYN – Connect request / connect accepted
- FIN – Connection release

G. W. Cox – Spring 2008

TCP Implementation 6

Urgent data

cs670

- A way to embed signaling in the data stream (e.g, to kill the process on the remote machine)
- Process:
 - Sending app gives signal to TCP with “Urgent” flag set
 - Sending TCP sets Urgent Pointer (points to end of urgent data) in seg header and immediately transmits buffer
 - Receiving TCP can interrupt app to transfer urgent data

G. W. Cox – Spring 2008

TCP Implementation 7

TCP flow control

cs670

- Implemented by “Window Size” field in header
- Window size set by receiver to indicate how many bytes (not segs) can be sent before next ACK
- Returned in ACK seg

G. W. Cox – Spring 2008

TCP Implementation 8

Segment size

cs670

- Seg = 20 header bytes
+ ? Option header bytes
+ 0~64KB data
- Max seg size = 64K-20 (max IP payload)
- Actual seg size determined by sending TCP:
 - Local MTU
 - Local sending rules
 - Truncated seg for Urgent Data

G. W. Cox – Spring 2008

TCP Implementation 9

Well-known ports

cs670

- If a remote node wants to obtain a particular service from a server, how does it know which port # to make the request on?
- “Well-known ports” (ports 1-1023) are reserved for standard services:
 - E.g, 21=TCP, 80=HTTP, 110=POP3
 - For the entire list, see <http://www.iana.org/assignments/port-numbers>

G. W. Cox – Spring 2008

TCP Implementation 10

TCP Implementation

cs670

- A quick review of TCP
- **TCP implementation options**
- Some other implementation considerations
- Optimizing TCP's performance
- TCP over wireless networks

G. W. Cox – Spring 2008

TCP Implementation 11

TCP Implementation Options

cs670

- The TCP spec lays out the details of the protocol
- However, some policy options are left to the implementer
- Two implementations that use different options will interoperate, but performance may suffer

G. W. Cox – Spring 2008

TCP Implementation 12

TCP Imp Options: Send Policy

cs670

- The sending TCP accepts bytes from the app and buffers them
- It is free to send them whenever it chooses (except for pushed data, urgent data, or a closed send window)

G. W. Cox – Spring 2008

TCP Implementation 13

TCP Imp Options: Deliver Policy

cs670

- Receiving TCP will receive segments and buffer them, handling errors and in-order considerations.
- It is free to deliver the data to the app wherever it chooses (except for urgent data or pushed data)

G. W. Cox – Spring 2008

TCP Implementation 14

TCP Imp Options: Accept Policy

cs670

- When segs arrive out of order, the receiving TCP can:
 1. Discard any segs that arrive out of order
 2. Accept any seg that has a sequence number within the receive window
- **Note:** Policy 1 is easier to implement and needs less complex buffering, but Policy 2 is better for performance

G. W. Cox – Spring 2008

TCP Implementation 15

TCP Imp Options: Retransmit Policy

cs670

- Governs how the sender will handle a re-transmission
 - First-Only Policy
 - Batch Policy
 - Individual Policy

G. W. Cox – Spring 2008

TCP Implementation 16

Retransmit Policy options: First-Only Policy

cs670

- Keep a send-order queue of unACKed segs.
- Keep a single timeout timer
- When an ACK arrives, remove the ACKed seg(s) from the queue and reset the timer
- When the timer times out, re-send the seg at the head of the queue
- Simple and low-traffic, but can be slow (the timer for the second seg in the queue doesn't start until the first one times out)

G. W. Cox – Spring 2008

TCP Implementation 17

Retransmit Policy Options: Batch Policy

cs670

- Same as First-Only except:
 - When the timer times out, re-send the entire queue
- Basically a go-back-n approach. Simple, but may cause needless additional traffic

G. W. Cox – Spring 2008

TCP Implementation 18

Retransmit Policy Options: Individual Policy

cs670

- Keep a timeout timer for each seg in the queue
- If any timer expires re-transmit just that seg
- Complex implementation. Very traffic efficient.

G. W. Cox – Spring 2008

TCP Implementation 19

Retransmit Policy Options

cs670

- Ideally, the retransmit policy would be selected to be compatible with the receiver's accept policy (e.g., if receiver uses in-order, the best match is a batch retransmit policy).
- But you can't count on that in real networks

G. W. Cox – Spring 2008

TCP Implementation 20

TCP Imp Options: ACK Policy

cs670

- The receiving TCP must ACK segs that are received in order. It can do so:
 - Immediately – Immediately send an empty ACK segment for the received data seg
 - Cumulative – Hold the ACK until it can be piggybacked on outgoing data (recall that ACK applies to the seg received and all before it). Keep a timer to prevent too long a delay.
- Most installations use Cumulative because it yields lower traffic loads, but this is a good bit more complex to implement and manage.

G. W. Cox – Spring 2008

TCP Implementation 21

TCP Implementation

cs670

- A quick review of TCP
- TCP implementation options
- **Some other implementation considerations**
- Optimizing TCP's performance
- TCP over wireless networks

G. W. Cox – Spring 2008

TCP Implementation 22

Setting the segment size

cs670

- Max = 64KB
- Min/default = 556B
- Often restricted to 1460 data bytes:

1460 data bytes
+ 20 TCP header
+ 20 IP header
<hr/>
1500 bytes (One Ethernet payload)
- Actual Max seg size is negotiated by sender and receiver during setup

G. W. Cox – Spring 2008

TCP Implementation 23

A problem with the window size field

cs670

- Window size field is 16 bits => 64KB max window
- Not enough for many purposes:
 - Example (Tanenbaum):
 - On a T3 line (44.7 Mbps), it takes 12 msec to send a 64KB window.
 - If RTT is 50 msec, sender is idle 75% of the time waiting for ACKs

G. W. Cox – Spring 2008

TCP Implementation 24

A window size patch

cs670

- “Window scale” is set during negotiation*
- Sets scale factor used in interpreting the Window Size field:
 - $\text{Act_Win_Size} = \text{Win_Size_Field} \times 2^{\text{Win_Scale}}$
- Allows window sizes up to $\sim 2^{30}$ B

G. W. Cox – Spring 2008

* RFC 1323
TCP Implementation 25

A potential TCP deadlock

cs670

- When the receiver wants the sender to stop sending temporarily, it will advertise a window size of 0. The sender will stop.
- Later, when the receiver can accept data again, it will advertise a larger window size. The sender will re-start.
- A deadlock occurs if the second advertisement is lost.

G. W. Cox – Spring 2008

TCP Implementation 26

Deadlock avoidance

cs670

- When a sender receives a Window Size = 0, it starts a “persistence timer”
- When the persistence timer times out, the sender sends a query to the receiver, and the receiver sends the current advertised window size
- If the size still = 0, the sender re-starts the persistence timer

G. W. Cox – Spring 2008

TCP Implementation 27

TCP Implementation

cs670

- A quick review of TCP
- TCP implementation options
- Some other implementation considerations
- **Optimizing TCP's performance**
- TCP over wireless networks

G. W. Cox – Spring 2008

TCP Implementation 28

Optimizing performance

cs670

- **Deciding when to send a segment**
- The Silly Window Syndrome
- Managing the Congestion Window
- Managing the timeout timer

G. W. Cox – Spring 2008

TCP Implementation 29

Deciding when to send a segment

cs670

- Within the negotiated seg size limits, the sending TCP has to decide when to stop buffering data bytes and send them
- When bytes come in slowly from the app (e.g., a user typing), how does the sending TCP decide when to send the buffered bytes?
 - If you send each one separately, you are using 40 overhead bytes (20 send header and 20 ACK header) to send 1 data byte
 - If you wait to build a large seg, the first bytes typed may be impossibly late

G. W. Cox – Spring 2008

TCP Implementation 30

One way of helping the problem – delayed ACK

cs670

- The receiver can delay ACKing a small seg (typ: 0.5 sec) in the hope of receiving another one that can be ACKed in the same ACK seg
- Fairly common approach, but not practical when fast response needed, and sender is still inefficient

G. W. Cox – Spring 2008

TCP Implementation 31

Another way: Nagle's Algorithm

cs670

- When data comes in a byte at a time from the app,
 - Send the first byte immediately
 - Buffer succeeding bytes until the first byte is ACKed, then send them in one seg
- This is a good strategy when RTT is variable:
 - When network is lightly loaded, the impact of small segs is less. Since ACKs return quickly, more small segs are sent
 - When the network is congested, ACKs return slowly and more data is packed in each seg.
- Note: Nagle's performance may not be good enough for highly interactive applications – sometimes it is disabled

G. W. Cox – Spring 2008

TCP Implementation 32

Optimizing performance

cs670

- Deciding when to send a segment
- **The Silly Window Syndrome**
- Managing the Congestion Window
- Managing the timeout timer

G. W. Cox – Spring 2008

TCP Implementation 33

A related problem: Silly Window Syndrome

cs670

- Occurs when data is sent in large blocks, but receiving app reads a byte at a time:
 - Sending TCP sends until rcv buffer full
 - • Receiving TCP advertises window size = 0
 - Receiving app reads a byte
 - Receiving TCP advertises window size = 1
 - Sending TCP sends a byte

G. W. Cox – Spring 2008

TCP Implementation 34

A fix for the Silly Window Syndrome

cs670

- Receiver is prevented from advertising when only small buffer space is available
- It can only advertise when either:
 - It can handle the negotiated max seg size, or
 - half of the receive buffer is empty

G. W. Cox – Spring 2008

TCP Implementation 35

Optimizing performance

cs670

- Deciding when to send a segment
- The Silly Window Syndrome
- **Managing the Congestion Window**
- Managing the timeout timer

G. W. Cox – Spring 2008

TCP Implementation 36

Congestion control

cs670

- Recall that basic TCP considers two numbers to set the send window size:
 - Advertised window size – set by the receiver to prevent sender from overrunning the receive buffer (flow control)
 - Congestion window size – set by the sender to try to prevent aggravating network congestion (congestion control)
- Send window size is set to the minimum of the two numbers

G. W. Cox – Spring 2008

TCP Implementation 37

Managing congestion window size

cs670

- The basic Slow Start algorithm (“Additive increase, multiplicative decrease”)
 - At start, set $CW = 1$ max seg size
 - If n segs ACKed, $CW = CW + n$ max seg sizes
 - If timeout, $CW = CW / 2$
 - Once CW is reduced, it is not increased again

G. W. Cox – Spring 2008

TCP Implementation 38

A consideration

cs670

- The basic algorithm does not change CW after it is reduced.
- In complex networks, congestion will come and go.
- We’d like to have a scheme that reduces CW when congestion is high, but increases it when the congestion is relieved

G. W. Cox – Spring 2008

TCP Implementation 39

Refined Slow Start

cs670

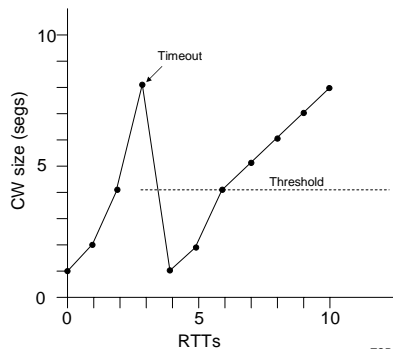
- When timeout occurs:
 - Set threshold = $CW/2$
 - Reset CW to 1 and start Slow Start again
 - After $CW \geq$ threshold, only increase CW by 1 for each ACK received (regardless of how many segs the ACK is for)

G. W. Cox – Spring 2008

TCP Implementation 40

A typical CW profile

cs670



G. W. Cox - Spring 2008

TCP Implementation 41

Optimizing performance

cs670

- Deciding when to send a segment
- The Silly Window Syndrome
- Managing the Congestion Window
- **Managing the timeout timer**

G. W. Cox - Spring 2008

TCP Implementation 42

Managing the timeout timer

cs670

- How do we set the value for the timeout timer?
 - Too long – performance suffers because it takes a long time to discover that a seg has been lost
 - Too short – many unnecessary re-transmissions, increasing network load

G. W. Cox - Spring 2008

TCP Implementation 43

Managing the timeout timer

cs670

- Timeout interval for a connection should be related to the RTT for that connection
- But on the Internet, RTT can vary wildly
- We would like a dynamic measure of RTT

G. W. Cox - Spring 2008

TCP Implementation 44

Dynamic RTT measurement – simple averaging

cs670

- For each connection, TCP maintains a variable *AvgRTT*
- When an ACK arrives, TCP calculates the RTT experienced
- *AvgRTT* is calculated by simple averaging
- Problem: treats long-ago behavior as importantly as recent behavior

G. W. Cox – Spring 2008

TCP Implementation 45

Dynamic RTT measurement – exponential averaging

cs670

$AvgRTT = \alpha \times AvgRTT + (1-\alpha) \times \text{measured RTT}$
where $0 < \alpha < 1$

- This formulation favors recent measurements
- Smaller alpha causes more weight on recent measurements
- Problem: works well for small variance, but does not react well to an abrupt, large change in RTT

G. W. Cox – Spring 2008

TCP Implementation 46

Setting timeout from RTT

cs670

- Generally, timeout is calculated by:
 $\text{timeout} = \beta \times AvgRTT$
- In the early days of the Internet, beta was set to 2
- Problem: when the variance in RTT is wide, a fixed beta may not work
- The fix: Jacobson's Algorithm tracks deviation of the measured RTT and sets beta accordingly.

G. W. Cox – Spring 2008

TCP Implementation 47

Jacobson's Algorithm: The Idea

cs670

- Instead of calculating *AvgRTT* by averaging, take the variability of the RTT samples into account

G. W. Cox – Spring 2008

TCP Implementation 48

Jacobson's Algorithm: Approach

cs670

SampleRTT = ACK_time - SEND_time

Diff = SampleRTT - Current Estimated RTT

EstimatedRTT = EstimatedRTT + $\delta \times \text{Diff}$ (where $0 \leq \delta \leq 1$)

Dev = Dev + $\delta (|\text{Diff}| - \text{Dev})$

THEN:

Timeout = EstimatedRTT + $p \times \text{Dev}$ (typical $p = 4$)

G. W. Cox - Spring 2008

TCP Implementation 49

How Jacobson's Alg acts

cs670

- When variability is small, Timeout tends toward the long-term average of RTT
- When variability is large, Timeout tends toward $p \times$ (magnitude of the variation)

G. W. Cox - Spring 2008

TCP Implementation 50

A problem with timeouts

cs670

- When network is congested, a higher fraction of attempted sends will timeout due to delay and the timed-out senders will re-send
- This increases the traffic that the network is trying to handle, probably increasing the congestion and lengthening delay

G. W. Cox - Spring 2008

TCP Implementation 51

A partial fix: Exponential timeout backoff

cs670

- For each re-transmission, double the timeout
 - For example:
 - Original send - timeout = x
 - First re-send - timeout = $2x$
 - Second re-send - timeout = $4x$

G. W. Cox - Spring 2008

TCP Implementation 52

More problems with setting timeouts

cs670

- Assume you send a seg, timeout, then re-send
- If you get an ACK, is this the ACK for the first send or the second?
- If you guess wrong, you will corrupt the value of *AvgRTT*

G. W. Cox – Spring 2008

TCP Implementation 53

An over-all fix: Karn's Algorithm

cs670

- Use Jacobsons' Alg to dynamically adjust RTT until a re-transmission occurs
- Do not use the RTT of re-transmitted segs in the *AvgRTT* calculation
- Use Exponential Timeout Backoff for all re-transmissions until a non-re-transmitted seg is ACKed
- Then, re-start Jacobson's Alg

G. W. Cox – Spring 2008

TCP Implementation 54

TCP Implementation

cs670

- A quick review of TCP
- **TCP implementation options**
- **Some other implementation considerations**
- **Optimizing TCP's performance**
- **TCP over wireless networks**

G. W. Cox – Spring 2008

TCP Implementation 55

TCP over wireless

cs670

- The earliest versions of TCP assumed that timeouts could be either caused by congestion or by lost packets
- As wired network reliability improved, TCP was optimized to emphasize the congestion case (the retransmit service, which was designed for error handling, has been specialized for congestion control)
- But wireless links are much less reliable and lost packets can occur frequently

G. W. Cox – Spring 2008

TCP Implementation 56

An example of how TCP is not optimized for wireless

cs670

- When segs timeout, TCP re-starts Slow Start and sends more slowly (assuming this will relieve congestion)
- But when data is lost (happens often on a wireless link), the best thing to do for performance is to re-send quickly
- Furthermore, many wireless systems (e.g., 802.11) use L2 re-transmission and slowdown. L2 re-trans causes slowdown which can trigger slow-start.

G. W. Cox – Spring 2008

TCP Implementation 57

TCP and disconnections

cs670

- Mobile users may experience fairly long disconnects
- If a conventional TCP connection was disconnected, TCP using normal settings would timeout and re-transmit repeatedly, doubling interval each time (up to one minute) for a max of 12 times
- If the mobile reconnected during this time, it might have to wait an entire minute before the datastream resumed.

G. W. Cox – Spring 2008

TCP Implementation 58

Some things that TCP needs to tolerate to run on wireless networks

cs670

- High error rates
- Temporary disconnects during handover
- Heterogeneous wired/wireless connections (and, sender may not know the composition)
- NOTE: The fact that conventional TCP is not designed to handle these problem does not mean that it won't work on wireless systems, just that its performance is not optimized.

G. W. Cox – Spring 2008

TCP Implementation 59

TCP mods are possible, but...

cs670

- Some proposals (e.g., RFC 3168) recommend returning the timeout structure to error control and adding an Explicit Congestion Notification (ECN) mechanism
- But this would require changing all TCP installations, most of which run over wired systems and operate fine as-is.

G. W. Cox – Spring 2008

TCP Implementation 60

Some fixes for TCP over wireless

cs670

- Indirect TCP (I-TCP)
- Snooping TCP

G. W. Cox – Spring 2008

TCP Implementation 61

Indirect TCP (I-TCP)

cs670

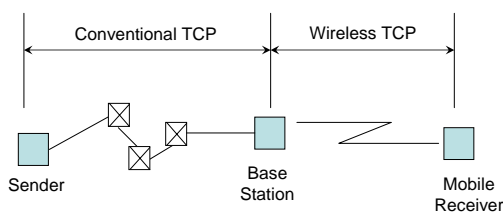
- The idea: Combine conventional TCP on the wired network with a Wireless TCP version tuned to run on the wireless part.

G. W. Cox – Spring 2008

TCP Implementation 62

I-TCP

cs670



G. W. Cox – Spring 2008

TCP Implementation 63

I-TCP: Plusses and Minuses

cs670

- **Plusses:**
 - No changes to existing protocol
 - Wireless errors do not directly affect wired system
- **Minuses:**
 - Breaks the paradigm of TCP operation – an ACK does not mean that the receiver got the seg
 - The base station buffers data for the mobile and ACKs it as it arrives. When a handover occurs, there is no way to get the sender to re-send the buffered data – the base station must forward it.
 - To tolerate disconnections, must have very large buffers – when disconnect is followed by handover, even more data to forward.

G. W. Cox – Spring 2008

TCP Implementation 64

Snooping TCP

cs670

- Uses a single end-to-end TCP connection (so the paradigm is not broken)
- Base station “snoops” on the connection

G. W. Cox – Spring 2008

TCP Implementation 65

“Snooping”

cs670

- The base station buffers and forwards all packets destined to the mobile (without ACKing the sender)
- Monitors the connection from the mobile to detect ACKs and gaps in transmitted data
 - Missing ACKs: Performs local re-transmits based on a short timeout timer
 - Missing data: sends NACK to mobile for re-transmit

G. W. Cox – Spring 2008

TCP Implementation 66

Snooping TCP –Plusses and Minuses

cs670

- **Plusses:**
 - Preserves the end-to-end TCP paradigm
 - No changes in the wired part of the network
 - Handovers to new base stations happen automatically (any buffered data will simply time out at the sender and be re-transmitted)
- **Minuses:**
 - Problems in the wireless segment can affect the wired segment (e.g., lost wireless packets can trigger Slow Start)
 - End-to-end encryption at L3 (e.g., IPSec) foils snooping
 - Does not help in the event of a disconnect (there are no ACKs to snoop)

G. W. Cox – Spring 2008

TCP Implementation 67