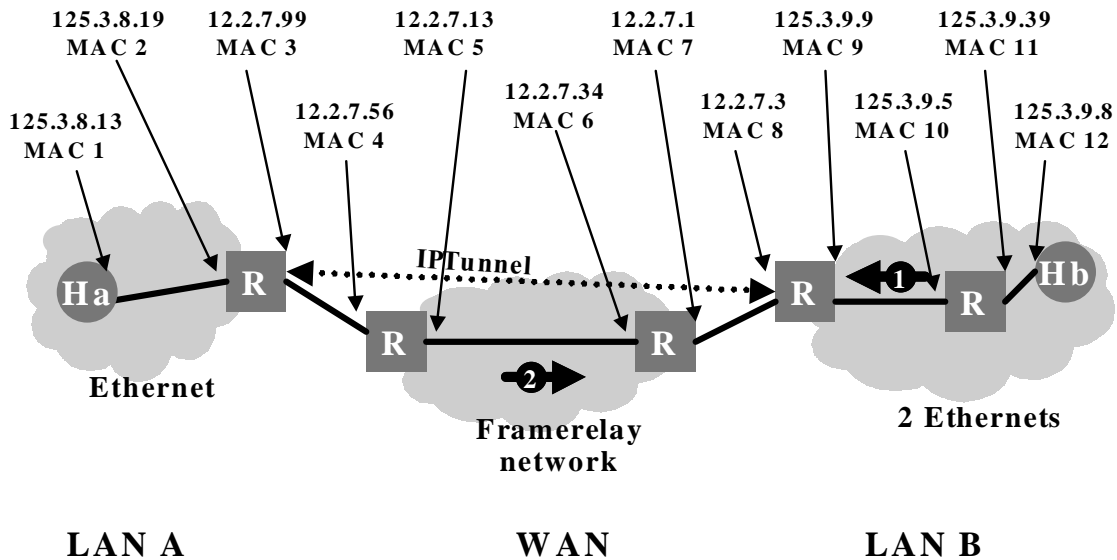


15-441:Networking

Homework 3 Solutions

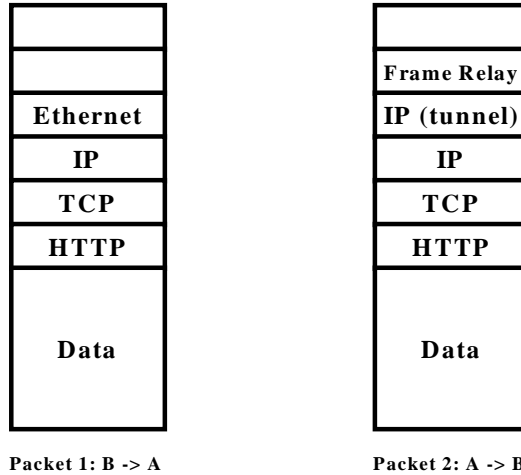
Problem 1: Packet Headers



A corporation has two local area networks connected over a wide area network using an IP tunnel, as is shown in the picture above. The two edge routers that connect the LANs to the WAN forward any packets destined for the other LAN through the tunnel. Tunneling is used so that the corporation can more easily use non-standard protocol features in its corporate network. The first LAN is a simple Ethernet while the second LAN consists of two Ethernets connected by an internal router. Part of the WAN runs through a Frame Relay network. In the figure, R stands for Router and H for Host. The Frame Relay switches are not shown in the picture, and Ethernet hubs and switches are also not shown.

We want to look in more detail at the headers of packets being exchanged by a web browser executing on the client host H_a in LAN A and the web server H_b executing in LAN B.

a) Please show the type and the order of the headers of the packets being exchanged by the web browser and server if they are captured at locations “1” and “2” as indicated on the network picture. When showing the “type” of header, be specific (e.g., Don’t just write ‘transport’, but rather indicate ‘TCP’ or ‘UDP’. Don’t just write ‘datalink’, but rather indicate ‘Frame Relay’, or ‘Ethernet’ or something else.) Include all headers of the following protocol layers: network layer, application layer, datalink layer, and transport layer. Use the packet outlines shown below, i.e., fill in the header types in the right box. Note that there may be more boxes than you need. The top of each packet outline corresponds to the first byte in the packet while the bottom corresponds to the last byte.



b) For the packet traveling in LAN B from LAN B to LAN A, identify the primary fields of the protocol headers, meaning transport-layer source and destination port numbers, network-layer source and destination IP addresses, and link-layer addresses or connection identifiers. There is no need to fill in anything for application headers. Use the IP and datalink address information shown in the figure. The URL that the browser is trying to retrieve is "http://www.internal.supercom/index.html". If you do not have enough information to specify a field, fill in an "X" and briefly explain. If a field does not exist, fill in "-" and explain briefly.

Hdr Type	Src	Dest
Ethernet	MAC 10	MAC 9
IP	125.3.9.8	125.3.8.13
TCP	80	X
HTTP		
Data		

The source port in the TCP header is 80, since that is the default port used by HTTP (the application). The destination port is unknown. It was assigned by the operating system of host Ha when it open the connection to the server.

c) The same question for the packet traveling from LAN A to LAN B through the Frame Relay network.

Hdr Type	Src	Dest
Frame Relay	-	X
IP (tunnel)	12.2.7.99	12.2.7.3
IP	125.3.8.13	125.3.9.8
TCP	X	80
HTTP		
Data		

Because the two packets belong to the same connection, the IP and TCP headers for packets 1 and 2 are similar. The difference is that the source and destination fields are flipped.

The second IP header is for the tunnel. The source and destination IP addresses are for the entry and exit point of the tunnel.

Communication in frame relay is based on virtual connections. This means that the destination field will hold a connection identifier (unknown) and that there is no source field.

Problem 2: DNS

The Andrew Linux machines provide a program `dig` that allows you to query Domain Name Service (DNS) servers around the Internet (some documentation is available if you type `man dig`). When running `dig` for the purposes of this question, you should use the following format: `dig +norecurse @name.of.dns.server record-type domain-name` where

- *name.of.dns.server* is the hostname of the DNS server you wish to query, such as A.ROOT-SERVERS.NET.
- *record-type* is the type of DNS record you wish to retrieve, such as ANY, MX, HINFO, A, or SOA.
- *domain-name* is the name of the host or domain you seek information on.

The DNS is a distributed architecture that uses hierarchical delegation. At the top of the system are the “root” name servers, who know which DNS server is responsible for each second-level domain (such as CMU.EDU). If you send a root server a query for a particular machine, you will receive a reply listing the servers that have been delegated authority for that machine’s second-level domain. It is common for a large domain such as CMU.EDU to further delegate to “departmental” or workgroup DNS servers, which you can discover by querying the second-level servers.

(a) In order to discover the chain of delegation in use at CMU, run a series of NS queries for UX15.SP.CS.CMU.EDU. You may start with any of the root servers, and you should continue your sequence of queries until you stop getting new delegations (in some domains, this is indicated by a DNS server returning you a delegation pointing to itself, and in other domains this is indicated by a DNS server returning you a SOA record instead).

Delegation chain for: AOL.COM

Server queried

NS delegations to

A.ROOT-SERVERS.NET
K.GTLD-SERVERS.NET
DNS-01.NS.AOL.COM

A.GTLD-SERVERS.NET, K.GTLD-SERVERS.NET
DNS-01.NS.AOL.COM, DNS-02.NS.AOL.COM
DNS-01.NS.AOL.COM, DNS-02.NS.AOL.COM

This was produced by running the following commands:

```
% dig +norecurse @a.root-servers.net NS aol.com
% dig +norecurse @k.gtld-servers.net NS aol.com
% dig +norecurse @dns-01.ns.aol.com NS aol.com
```

Generate the delegation chain for UX15.SP.CS.CMU.EDU. Present your results in the table form shown above. Each NS query will typically return two or more answers; choose among them at random. If you query a server and get a timeout, choose an alternate server.

Server queried	NS delegations to
-----	-----
A.ROOT-SERVERS.NET	L3.NSTLD.COM, H3.NSTLD.COM
H3.NSTLD.COM	CABBAGE.SRV.CS.CMU.EDU, T-NS1.NET.CMU.EDU
T-NS1.NET.CMU.EDU	SPINACH.SRV.CS.CMU.EDU, LETTUCE.SRV.CS.CMU.EDU
LETTUCE.SRV.CS.CMU.EDU	returns SOA record

Note: the second column only shows the first and last server returned, not the full list. The

(b) The DNS is also used to translate IP addresses into hostnames. Again, the database is distributed in a hierarchical fashion, with a wrinkle. The most-specific part of a domain name is on the left (i.e., UX15 in UX15.SP.CS.CMU.EDU), but the reverse is true of IP addresses (i.e., in 128.2.203.134, 128 is “top-level”, 128.2 is CMU.EDU, and 128.2.203 belongs to CS.CMU.EDU). Thus, address-to-name mapping is handled by reversing the bytes of the IP address and making queries in a special domain. To turn 128.2.203.134 into a hostname, various servers are sent queries seeking PTR records for 134.203.2.128.in-addr.arpa. The first query would be:

```
% dig @a.root-servers.net PTR 134.203.2.128.in-addr.arpa
```

You will know you’re done when your query gives you back a PTR record in addition to (or instead of) NS records.

Fill in a table like the one above showing a query chain for the IP address 64.91.109.37.

Server queried	NS delegations (or PTR record)
-----	-----
A.ROOT-SERVERS.NET	chia.ARIN.NET, figwort.ARIN.NET
figwort.ARIN.NET	ns2.centurytel.net, ns1.centurytel.net
ns1.centurytel.net	returns PTR 9netics.com

Problem 3: Mars Exploration

You are part of an exploratory mission to Mars. When, after a long trip, the mission arrives on Mars it discovers the remains of an ancient civilization. Shockingly, it appears that the previous inhabitants of Mars used 4Mb/s token ring for their computer infrastructure. Sadly, however, while your team easily locates computer interface cards, cabling, and hubs, it appears that all the transmit tokens are missing.

When you try to report your findings (including beautiful color pictures of the writing on the Martian token-ring hubs) to Mission Control on Earth, you discover that transmission is extremely slow. It takes almost three hours (10,000 seconds, to be exact) to send just one picture

(of 1000x1000 pixels, 8 bits per pixel). The radio link transmits at 800 KBit/second, but the average bit rate you are achieving to Mission Control is only a fraction of this speed.

Your mission colleagues ask you to apply your 15-441 expertise to diagnose the problem. After some digging, you discover that the communication link uses a very simple point-to-point data-link protocol: it uses flow control and assumes the transmission channel is reliable. The packet size is 1000 bytes. You also discover that the radio transmission system was originally built for the Apollo 11 mission to the moon. You immediately suspect that the engineers who designed your Mars spacecraft did not take 15-441!

	Moon	Mars
Distance to Earth (km)	390,000	390,000,000
Diameter (km)	3475	6000
Distance to Sun (km)	150 M	230 M

A: Do your best to explain the problem. Please be precise (i.e., give numbers). We have provided some information (which may or may not be relevant) above about distances at the time of the mission. Also, the speed of radio waves in a vacuum is 300,000 km/second.

Answer: Given the picture size (8,000,000 bits) and the transmission rate, you would have expected a transmission time of 10 seconds per picture (obviously there is also a long propagation delay). The only possibility is that the datalink protocol is slowing things down. The maximum transmission rate of a flow-controlled protocol is the buffer size at the receiver, divided by the roundtrip time. Given that the distance to Mars is 1000 times the distance to the moon, it seems like the datalink protocol was engineered to support full transmission rates to the Moon (you can calculate that the system has 260 packet buffers: transmission rate x roundtrip time), but it only achieves 1/1000th of its throughput to Mars.

B: How would you fix the problem? Again, be precise. How practical do you think this is?

Answer: You need to increase the buffer size by a factor of 1000 to make up for the longer distance, i.e. 260,000 packet buffers of 1000 Bytes each.

Note that since you are sending data to earth, you need the larger buffers on earth. That is a lucky break since they can presumably go to the store and buy it. It is unlikely that you packed that much extra memory for your trip to Mars so it is less likely that you can speed up transmission in the other direction.

C: Once you have fixed the throughput problem, mission control is reporting corrupted pictures. It seems that a radio transmission system designed to work from the moon does not quite cut it from Mars. What changes would you make to the datalink protocol to recover from the transmission errors?

Answer:

First, you have the receiver (earth) send acknowledgments for each correct packet you receive. Next, you add timeouts on the sender (Mars) so you can retransmit lost packet - note that you need long timeouts! Finally, you add sequence numbers so the sender and receiver can agree on what packets are missing.

Problem 4: TCP Handshake

After leaving CMU you are hired by a new client/server middleware startup called Super-Tricky.com. After you've been there a few days, your manager drops by your office to ask

if you to evaluate a cool new idea the company has been considering.

The problem is that applications based on their middleware frequently need results from multiple servers. For an application to contact one server and get a response TCP requires five packets and two round-trip times. Your manager suggests that this process could be streamlined if a client would begin each TCP conversation by sending a packet with the SYN and FIN flags set, plus application data in the segment body. Then the server can return one packet containing the response and the TCP connection can consist of exactly two packets.

(a) If this scheme is to be implemented, which flags should be set in the server's response to the client? Why?

Answer: The server would also need to set the SYN and FIN flags. It would also set the ACK flag to acknowledge the SYN and FIN flags and sequence numbers sent by the sender.

(b) Explain one way in which the Berkeley Socket API makes it inconvenient to pursue this approach.

Answer: The Berkeley Socket API uses separate calls to establish a connection (`connect()`, `accept()`) and to send and receive data (`read()`, `send()`, `write()`, `recv()`). This is particularly inconvenient on the server side—there would need to be a way to do the equivalent of `accept()` and `read()` without any packets being sent back to the client, and a way to do `write()` and `close()` as a single operation.

(c) As you are writing up a report for your manager, one of your co-workers drops by. Leslie claims that this approach has a fatal flaw: neither client queries nor server responses can be more than 1500 bytes long. Is this true? Explain.

Answer: Leslie has a point—traditional Ethernets have a maximum payload size of 1500 bytes. But that's not the whole story. Gigabit Ethernet can be configured to allow 9000 bytes of payload, and a popular IP-over-ATM arrangement allowed 9180, so 1500 isn't a magic number appearing at the heart of all IP protocols. Furthermore, IPv4 supports fragmentation, which in theory should allow requests and responses to fill a maximal-size IP datagram, meaning a payload size of just under 64 KByte. However, most existing TCP implementations are tuned to *avoid* fragmentation, so there would be messy code changes.

By the way, your manager's proposal is a stripped-down version of "Transactional TCP" (T/TCP), RFC 1644. Though the ideas were worked out in 1994, T/TCP has not taken over the world—so far.

Problem 5: Sized for Speed

You are hired to design a reliable byte-stream protocol that uses a sliding window (like TCP). This protocol will run over a 100 Mb/s network. The round-trip time will be bounded by 100 ms and you may assume the maximum segment lifetime is 60 seconds.

(a) How many bits wide should the AdvertisedWindow header field be? Show your work.

Answer: There should be enough bits represent a full bandwidth-delay product, assuming maximum a maximum data rate and latency. This corresponds to 10 Mbits (i.e. 10^7). Thus must be rounded up to a power of two, i.e. 16Mbits (i.e. 2^{24}). Assuming the window is expressed in bytes, this will require 21 bits to represent.

(b) How many bits wide should the SequenceNumber field be? Show your work.

Answer: The sequence number field must be large enough so that "old" packets cannot be confused with "new" packets after the sequence number wraps around. This means that the sequence number space must be larger than 100 Mb/sec * 60 sec, or $6 * 10^9$ bits. This means that you need at least 30 bits of sequence number space.

(c) Your manager walks into your office and asks you whether your design can be made "future-proof." For each of the three design parameters, briefly explain whether or not you

think it makes sense to design for a factor of 100 increase (rate = 10 Gb/s, RTT = 10s, MSL = 6000s).

Answer: Gigabit Ethernet interfaces are readily available for many hardware platforms (e.g., some PCs ship with Gigabit Ethernet by default), and it is easy to come by SONET hardware faster than that. So 100 Mb/s does not seem like a good “ceiling” to build into a new protocol. Ten seconds is more than enough to get to the Moon and back but nowhere near enough time for a round trip to Mars. It is possible to argue that a 10-second RTT is either plausible or esoteric, as you choose. An MSL of 6000 seconds is harder to justify, though—that’s 1.6 *hours*. Except for serious deep-space networking, it’s hard to imagine a good justification for a packet taking hours to arrive. In summary, protocol designers should be wary of setting short-sighted throughput limits, but the nature of the universe can provide us with useful guidance on round-trip-time and segment-lifetime limits.