
CNP3-Practice Documentation

Release 0.0

Olivier Bonaventure Mickael Hoerd, Laurent Vanbever and Virg

October 28, 2010

CONTENTS

1	The Alternating Bit Protocol	3
2	Go-back-n and selective repeat	7
3	Programming project	11
3.1	Deliverables	12
3.2	Demonstration	12
4	The Transmission Control Protocol	15
4.1	Packet trace analysis tools	15
4.2	Emulating a network with Netkit	16
4.3	Questions	19
5	TCP congestion control	23
5.1	Trace analysis	24
6	Indices and tables	27
	Index	29

This document contains the questions for the practical part of the INGI2141 course during the 2010-2011 academic year. The html and pdf versions will be updated every week. These exercises have been written by [Olivier Bonaventure](#) with the help of [Mickael Hoerd](#), [Damien Saucez](#) and [Laurent Vanbever](#) and [Virginie van den Schriek](#)

THE ALTERNATING BIT PROTOCOL

The objective of this set of exercises is to better understand the basic mechanisms of the alternating bit protocol and the utilisation of the socket interface with the connectionless transport service.

1. Consider the Alternating Bit Protocol as described in the book.
 - How does the protocol recover from the loss of a data segment ?
 - How does the protocol recovers from the loss of an acknowledgement ?
2. A student proposed to optimise the Alternating Bit Protocol by adding a negative acknowledgment, i.e. the receiver sends a *NAK* control segment when it receives a corrupted data segment. What kind of information should be placed in this control segment and how should the sender react when receiving such a *NAK* ?
3. Transport protocols rely on different types of checksums to verify whether segments have been affected by transmission errors. The most frequently used checksums are :
 - the Internet checksum used by UDP, TCP and other Internet protocols which is defined in **RFC 1071** and implemented in various modules, e.g. <http://ilab.cs.byu.edu/cs460/code/ftp/ichchecksum.py> for a python implementation
 - the 16 bits or the 32 bits Cyclical Redundancy Checks (CRC) that are often used on disks, in zip archives and in datalink layer protocols. See <http://docs.python.org/library/binascii.html> for a python module that contains the 32 bits CRC
 - the Alder checksum defined in **RFC 2920** for the SCTP protocol but replaced by a CRC later
 - the Fletcher checksum

By using your knowledge of the Internet checksum, can you find a transmission error that will not be detected by the Internet checksum ?
4. The CRCs are efficient error detection codes that are able to detect :
 - all errors that affect an odd number of bits
 - all errors that affect a sequence of bits which is shorter than the length of the CRC

Carry experiments with one implementation of CRC-32 to verify that this is indeed the case.
5. Checksums and CRCs should not be confused with secure hash functions such as MD5 defined in **RFC 1321** or SHA-1 described in **RFC 4634**. Secure hash functions are used to ensure that files or sometimes packets/segments have not been modified. Secure hash functions aim at detecting malicious changes while checksums and CRCs only detect random transmission errors. Perform some experiments with hash functions such as those defined in the <http://docs.python.org/library/hashlib.html> python hashlib module to verify that this is indeed the case.
6. A version of the Alternating Bit Protocol supporting variable length segments uses a header that contains the following fields :
 - a number (0 or 1)
 - a length field that indicates the length of the data

- a CRC

To speedup the transmission of the segments, a student proposes to compute the CRC over the data part of the segment but not over the header. What do you think of this optimisation ?

7. On Unix hosts, the `ping(8)` command can be used to measure the round-trip-time to send and receive packets from a remote host. Use `ping(8)` to measure the round-trip to a remote host. Chose a remote destination which is far from your current location, e.g. a small web server in a distant country. There are implementations of ping in various languages, see e.g. <http://pypi.python.org/pypi/ping/0.2> for a python implementation of ping
8. How would you set the retransmission timer if you were implementing the Alternating Bit Protocol to exchange files with a server such as the one that you measured above ?
9. What are the factors that affect the performance of the Alternating Bit Protocol ?
10. Links are often considered as symmetrical, i.e. they offer the same bandwidth in both directions. Symmetrical links are widely used in Local Area Networks and in the core of the Internet, but there are many asymmetrical link technologies. The most common example are the various types of ADSL and CATV technologies. Consider an implementation of the Alternating Bit Protocol that is used between two hosts that are directly connected by using an asymmetric link. Assume that a host is sending segments containing 10 bytes of control information and 90 bytes of data and that the acknowledgements are 10 bytes long. If the round-trip-time is negligible, what is the minimum bandwidth required on the return link to ensure that the transmission of acknowledgements is not a bottleneck ?
11. Derive a mathematical expression that provides the *goodput* achieved by the Alternating Bit Protocol assuming that :
 - Each segment contains D bytes of data and c bytes of control information
 - Each acknowledgement contains c bytes of control information
 - The bandwidth of the two directions of the link is set to B bits per second
 - The delay between the two hosts is s seconds in both directions

The goodput is defined as the amount of SDUs (measured in bytes) that is successfully transferred during a period of time

12. The socket interface allows you to use the UDP protocol on a Unix host. UDP provides a connectionless unreliable service that in theory allows you to send SDUs of up to 64 KBytes.
 - Implement a small UDP client and a small UDP server (in python, you can start from the example provided in <http://docs.python.org/library/socket.html> but you can also use C or java)
 - run the client and the servers on different workstations to determine experimentally the largest SDU that is supported by your language and OS. If possible, use different languages and Operating Systems in each group.
13. By using the socket interface, implement on top of the connectionless unreliable service provided by UDP a simple client that sends the following message shown in the figure below.

In this message, the bit flags should be set to `01010011b`, the value of the 16 bits field must be the square root of the value contained in the 32 bits field, the character string must be an ASCII representation (without any trailing `0`) of the number contained in the 32 bits character field. The last 16 bits of the message contain an Internet checksum that has been computed over the entire message.

Upon reception of a message, the server verifies that :

- the flag has the correct value
- the 32 bits integer is the square of the 16 bits integer
- the character string is an ASCII representation of the 32 bits integer
- the Internet checksum is correct

If the verification succeeds, the server returns a SDU containing `11111111b`. Otherwise it returns `01010101b`

Inside each group, implement two different clients and two different servers (both using different languages). The clients and the servers must run on both the Linux workstations and the Sun server (*sirius*). Verify the interoperability of the clients and the servers inside the group. You can use C, Java or python to write these implementations.

```

0           1           2           3
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|  Bit flags   |          16 bits field          |          Zero   |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|          32 bits field          |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
| Type=1       | Len (8 bits)   |  Character string ...   |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|          Character string (cont.)          |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|  16 bits Internet checksum  |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

14. Consider an Alternating Bit Protocol that is used over a link that suffers from deterministic errors. When the error ratio is set to $\frac{1}{p}$, this means that $p - 1$ bits are transmitted correctly and the p^{th} bit is corrupted. Discuss the factors that affect the performance of the Alternating Bit Protocol over such a link.

GO-BACK-N AND SELECTIVE REPEAT

Go-back-n and selective repeat are the basic mechanisms used in reliable window-based transport-layer protocols. These questions cover these two mechanisms in details. You do not need to upload the time sequence diagrams on the svn repository, bring them with you on paper.

1. Amazon provides the [S3 storage service](#) where companies and researchers can store lots of information and perform computations on the stored information. Amazon allows users to send files through the Internet, but also by sending hard-disks. Assume that a 1 Terabyte hard-disk can be delivered within 24 hours to Amazon by courier service. What is the minimum bandwidth required to match the bandwidth of this courier service ?
2. Several large datacenters operators (e.g. [Microsoft](#) and [google](#)) have announced that they install servers as containers with each container hosting up to 2000 servers. Assuming a container with 2000 servers and each storing 500 GBytes of data, what is the time required to move all the data stored in one container over one 10 Gbps link ? What is the bandwidth of a truck that needs 10 hours to move one container from one datacenter to another.
3. What are the techniques used by a go-back-n sender to recover from :
 - transmission errors
 - losses of data segments
 - losses of acknowledgements
4. Consider a b bits per second link between two hosts that has a propagation delay of t seconds. Derive a formula that computes the time elapsed between the transmission of the first bit of a d bytes segment from a sending host and the reception of the last bit of this segment on the receiving host.
5. Consider a go-back-n sender and a go-back receiver that are directly connected with a 10 Mbps link that has a propagation delay of 100 milliseconds. Assume that the retransmission timer is set to three seconds. If the window has a length of 4 segments, draw a time-sequence diagram showing the transmission of 10 segments (each segment contains 10000 bits):
 - when there are no losses
 - when the third and seventh segments are lost
 - when the second, fourth, sixth, eighth, ... acknowledgements are lost
 - when the third and fourth data segments are reordered (i.e. the fourth arrives before the third)
6. Same question when using selective repeat instead of go-back-n. Note that the answer is not necessarily the same.
7. Consider two high-end servers connected back-to-back by using a 10 Gbps interface. If the delay between the two servers is one millisecond, what is the throughput that can be achieved by a transport protocol that is using 10,000 bits segments and a window of
 - one segment

- ten segments
 - hundred segments
8. Consider two servers are directly connected by using a b bits per second link with a round-trip-time of r seconds. The two servers are using a transport protocol that sends segments containing s bytes and acknowledgements composed of a bytes. Can you derive a formula that computes the smallest window (measured in segments) that is required to ensure that the servers will be able to completely utilise the link ?
 9. Same question as above if the two servers are connected through an asymmetrical link that transmits bu bits per second in the direction used to send data segments and bd bits per second in the direction used to send acknowledgements.
 10. The Trivial File Transfer Protocol is a very simple file transfer protocol that is often used by diskless hosts when booting from a server. Read the TFTP specification in [RFC 1350](#) and explain how TFTP recovers from transmission errors and losses.
 11. Is it possible for a go-back- n receiver to interoperate with a selective-repeat sender ? Justify your answer.
 12. Is it possible for a selective-repeat receiver to interoperate with a go-back- n sender ? Justify your answer.
 13. The go-back- n and selective repeat mechanisms that are described in the book exclusively rely on cumulative acknowledgements. This implies that a receiver always returns to the sender information about the last segment that was received in-sequence. If there are frequent losses or reordering, a selective repeat receiver could return several times the same cumulative acknowledgment. Can you think of other types of acknowledgements that could be used by a selective repeat receiver to provide additional information about the out-of-sequence segments that it has received. Design such acknowledgements and explain how the sender should react upon reception of this information.
 14. The *goodput* achieved by a transport protocol is usually defined as the number of application layer bytes that are exchanged per unit of time. What are the factors that can influence the *goodput* achieved by a given transport protocol ?
 15. The Transmission Control Protocol (TCP) attaches a 40 bytes header to each segment sent. Assuming an infinite window and no losses nor transmission errors, derive a formula that computes the maximum TCP goodput in function of the size of the segments that are sent.
 16. A go-back- n sender uses a window size encoded in a n bits field. How many segments can it send without receiving an acknowledgement ?
 17. Consider the following situation. A go-back- n receiver has sent a full window of data segments. All the segments have been received correctly and in-order by the receiver, but all the returned acknowledgements have been lost. Show by using a time sequence diagram (e.g. by considering a window of four segments) what happens in this case. Can you fix the problem on the go-back- n sender ?
 18. Same question as above, but assume now that both the sender and the receiver implement selective repeat. Note the the answer will be different from the above question.
 19. Consider a transport that supports window of one hundred 1250 Bytes segments. What is the maximum bandwidth that this protocol can achieve if the round-trip-time is set to one second ? What happens if, instead of advertising a window of one hundred segments, the receiver decides to advertise a window of 10 segments ?
 20. Explain under which circumstances a transport entity could advertise a window of 0 segments ?
 21. The socket library is also used to develop applications above the reliable bytestream service provided by TCP. We have installed on the *sirius.info.ucl.ac.be* server a simple server that provides a simple client-server service. The service operates as follows :
 - the server listens on port *62141* for a TCP connection
 - upon the establishment of a TCP connection, the server sends an integer by using the following TLV format :
 - the first two bits indicate the type of information (01 for ASCII, 10 for boolean)
 - the next six bits indicate the length of the information (in bytes)

- An ASCII TLV has a variable length and the next bytes contain one ASCII character per byte. A boolean TLV has a length of one byte. The byte is set to *00000000b* for *true* and *00000001b* for *false*.
- the client replies by sending the received integer encoded as a 32 bits integer in *network byte order*
- the server returns a TLV containing *true* if the integer was correct and a TLV containing *false* otherwise and closes the TCP connection

Each group of two students must implement a client to interact with this server in C, Java or python.

PROGRAMMING PROJECT

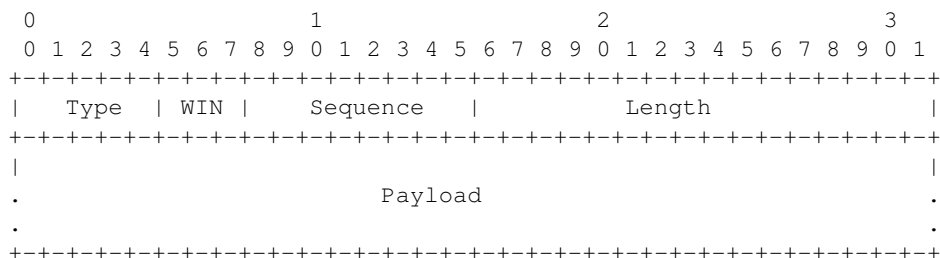
Your objective in this project is to implement a simple reliable transport protocol by groups of 2 students. These groups must be subgroups of the main groups that are registered on icampus. If the number of students in an icampus group is odd, there can be one group of three students.

The protocol uses a sliding window to transmit more than one segment without being forced to wait for an acknowledgment. Your implementation must support variable size sliding window as the other end of the flow can send its maximum window size. The window size is encoded as a three bits unsigned integer.

The protocol identifies the DATA segments by using sequence numbers. The sequence number of the first segment must be 0. It is incremented by one for each new segment. The receiver must acknowledge the delivered segments by sending an ACK segment. The sequence number field in the ACK segment always contains the sequence number of the next expected in-sequence segment at the receiver. The flow of data is unidirectional, meaning that the sender only sends DATA segments and the receiver only sends ACK segments.

To deal with segments losses, the protocol must implement a recovery technique such as go-back-n or selective repeat and use retransmission timers. The project will partially be evaluated on the quality of the recovery technique. Groups of three must implement the selective repeat technique while groups of two can implement a simpler recovery scheme such as go-back-n.

Segment format



- *Type*: segment type
 - 0x1 DATA segment.
 - 0x2 ACK segment
- *WIN*: the size of the current window (an integer encoded as a 3 bits field). In DATA segments, this field indicates the size of the sending window of the sender. In ACK segments, this field indicates the current value of the receiving window.
- *Sequence*: Sequence number (8 bits unsigned integer), starts at 0. The sequence number is incremented by 1 for each new DATA segment sent by the sender. Inside an ACK segment, the sequence field carries the sequence number of the next in-sequence segment that is expected by the receiver.
- *Length*: length of the payload in multiple of one byte. All DATA segments contain a payload with 512 bytes of data, except the last DATA segment of a transfert that can be shorter. The reception of a DATA segment whose length is different than 512 indicates the end of the data transfert.
- *Payload*: the data to send

3.1 Deliverables

Before October 22nd, 2010 at 23:59 each sub-group must submit its commented source code (with a Makefile) on the SVN and a short report (up to four pages in pdf format) describing the chosen recovery technique, the architecture of the client and server and the tests that have been carried out. Each group must implement both a receiver and a sender. The implementation language can be chosen among C, Java and Python.

The client and the server exchange UDP datagrams that contain the DATA and ACK segments. They must be command-line tools that allow to transmit one binary file and support the following parameters :

```
sender <destination_DNS_name> <destination_port_number> <window_size> <input_file>
```

```
receiver <listening_port_number> <window_size> <output_file>
```

A demo session will be organised on Tuesday October 26th. During the demo session, you will be invited to demonstrate that your implementation is operational and is interoperable with another. You also need to perform tests to show that your implementations works well in case of segment losses. For these tests, you can use a random number generator to probabilistically drop received segments and introduce random delays upon the arrival of a segment.

3.2 Demonstration

The demonstration of your project will take place in the intel room at :

- group 1, 2, 3: October 26th 10:45 am
- group 4, 5, 6: October 26th 11:45 am

We will test your implementations. For that, we will link you with each sub-group of the other groups with the same subgroup number. You personally assign the sub-group number inside your group. The receivers will be tested on sirius, the senders from the Intel room. Group 1 (resp. 4) sends to group 2 (resp. 5); Group 1 (resp. 4) receives from 3 (resp. 6); group 2 (resp. 5) sends to group 3 (resp. 6). The sender window size will be set to 6 at startup. The receiver window size will be set to 3 at startup.

We will provide you 6 files to each sender. The receiver must receive all of them and prove the transfer correctness by giving the md5 hashes (digest -v -a md5 <file>.).

The receiver port number is computed is defined as 6SGsg, where

- S: sender sub-group number
- G: sender group number
- s: receiver sub-group number
- g: receiver group number

e.g., port 61213 means that sub-group 1 of group 2 sends traffic to sub-group 1 of group 3.

Evaluation

Code: 50%

- no compilation -> 0
- increment of sequence number/ack (1)
- ack processing should be decoupled from data processing (1)
- window management (2)
- timer management (2)
- network byte order (1)
- architecture (2)

- code readability/documentation/synopsis (1)

Report: 25%

- architecture description (2)
- pitfall highlighting (2)
- validation (5)
- further work (1)

Demo: 25%

- correct transfers (6)
- support random loss/delay (2)
- ability to explain the code/events (2)

THE TRANSMISSION CONTROL PROTOCOL

The Transmission Control Protocol plays a key role in the TCP/IP protocol suite by providing a reliable byte stream service on top of the unreliable connectionless service provided by IP. During this exercise session, you will learn how to establish correctly a TCP connection and analyse packet traces that contain TCP segments. Note that some of exercises involve the creation of non-standard TCP segments. These exercises cannot be performed outside the [netkit](#) environment that is described below.

4.1 Packet trace analysis tools

When debugging networking problems or to analyse performance problems, it is sometimes useful to capture the segments that are exchanged between two hosts and to analyse them.

Several packet trace analysis tools are available, either as commercial or open-source tools. These tools are able to capture all the packets exchanged on a link. Of course, capturing packets require administrator privileges. They can also analyse the content of the captured packets and display information about them. The captured packets can be stored in a file for offline analysis.

[tcpdump](#) is probably one of the most well known packet capture software. It is able to both capture packets and display their content. [tcpdump](#) is a text-based tool that can display the value of the most important fields of the captured packets. Additional information about [tcpdump](#) may be found in *tcpdump(1)*. The text below is an example of the output of [tcpdump](#) for the first TCP segments exchanged on an scp transfer between two hosts

```
21:05:56.230737 IP 192.168.1.101.54150 > 130.104.78.8.22: S 1385328972:1385328972(0) win 65535 <m
21:05:56.251468 IP 130.104.78.8.22 > 192.168.1.101.54150: S 3627767479:3627767479(0) ack 13853289
21:05:56.251560 IP 192.168.1.101.54150 > 130.104.78.8.22: . ack 1 win 65535 <nop,nop,timestamp 27
21:05:56.279137 IP 130.104.78.8.22 > 192.168.1.101.54150: P 1:21(20) ack 1 win 49248 <nop,nop,timestamp 2
21:05:56.279241 IP 192.168.1.101.54150 > 130.104.78.8.22: . ack 21 win 65535 <nop,nop,timestamp 2
21:05:56.279534 IP 192.168.1.101.54150 > 130.104.78.8.22: P 1:22(21) ack 21 win 65535 <nop,nop,t
21:05:56.303527 IP 130.104.78.8.22 > 192.168.1.101.54150: . ack 22 win 49248 <nop,nop,timestamp 1
21:05:56.303623 IP 192.168.1.101.54150 > 130.104.78.8.22: P 22:814(792) ack 21 win 65535 <nop,nop
```

You can easily recognise in the output above the *SYN* segment containing the *MSS*, *window scale*, *timestamp* and *sackOK* options, the *SYN+ACK* segment whose *wscale* option indicates that the server uses window scaling for this connection and then the first few segments exchanged on the connection.

[wireshark](#) is more recent than [tcpdump](#). It evolved from the [ethereal](#) packet trace analysis software. It can be used as a text tool like [tcpdump](#). For a TCP connection, [wireshark](#) would provide almost the same output as [tcpdump](#). The main advantage of [wireshark](#) is that it also includes a graphical user interface that allows to perform various types of analysis on a packet trace.

The [wireshark](#) window is divided in three parts. The top part of the window is a summary of the first packets from the trace. By clicking on one of the lines, you can show the detailed content of this packet in the middle part of the window. The middle of the window allows you to inspect all the fields of the captured packet. The bottom part

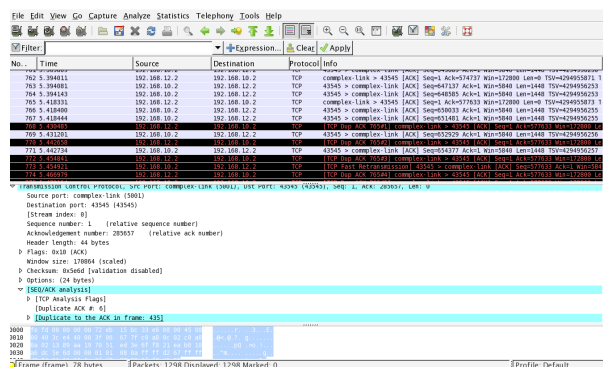


Figure 4.1: Wireshark : default window

of the window is the hexadecimal representation of the packet, with the field selected in the middle window being highlighted.

wireshark is very good at displaying packets, but it also contains several analysis tools that can be very useful. The first tool is *Follow TCP stream*. It is part of the *Analyze* menu and allows you to reassemble and display all the payload exchanged during a TCP connection. This tool can be useful if you need to analyse for example the commands exchanged during a SMTP session.

The second tool is the flow graph that is part of the *Statistics* menu. It provides a time sequence diagram of the packets exchanged with some comments about the packet contents. See blow for an example.

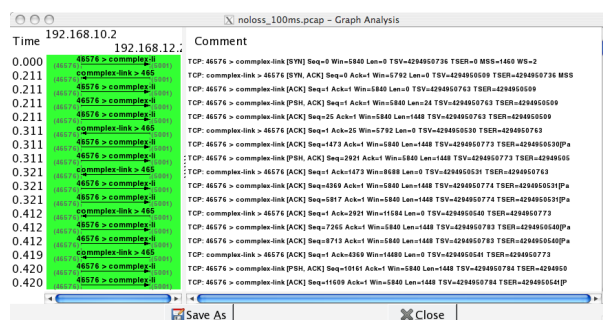


Figure 4.2: Wireshark : flow graph

The third set of tools are the *TCP stream graph* tools that are part of the *Statistics menu*. These tools allow you to plot various types of information extracted from the segments exchanged during a TCP connection. A first interesting graph is the *sequence number graph* that shows the evolution of the sequence number field of the captured segments with time. This graph can be used to detect graphically retransmissions.

A second interesting graph is the *round-trip-time* graph that shows the evolution of the round-trip-time in function of time. This graph can be used to check whether the round-trip-time remains stable or not. Note that from a packet trace, **wireshark** can plot two *round-trip-time* graphs, One for the flow from the client to the server and the other one. **wireshark** will plot the *round-trip-time* graph that corresponds to the selected packet in the top **wireshark** window.

4.2 Emulating a network with Netkit

Netkit is network emulator based on User Mode Linux. It allows to easily set up an emulated network of Linux machines, that can act as end-host or routers.

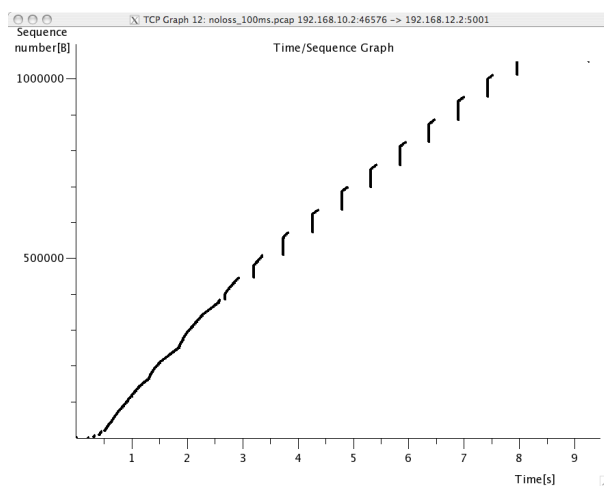


Figure 4.3: Wireshark : sequence number graph

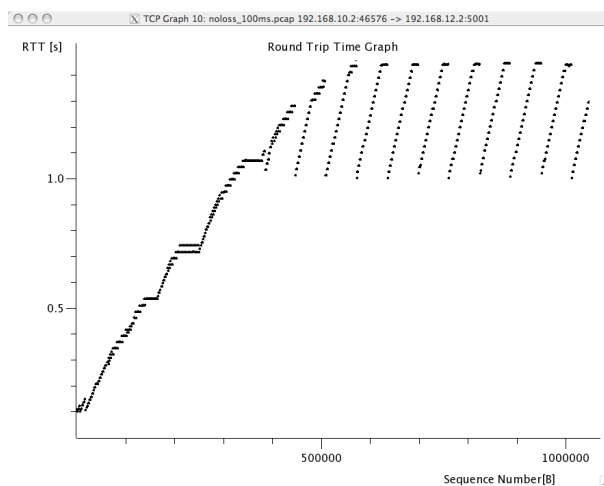


Figure 4.4: Wireshark : round-trip-time graph

4.2.1 Where can I find Netkit?

Netkit is available online. Files can be downloaded from http://wiki.netkit.org/index.php/Download_Official, and instructions for the installations are available here : <http://wiki.netkit.org/download/netkit/INSTALL> .

Netkit has already been installed in the student labs, in */etinfo/applications/netkit* . All you have to do in order to use it is to set the following environment variables :

```
export NETKIT_HOME=/etinfo/applications/netkit
export MANPATH=$NETKIT_HOME/man
export PATH=$NETKIT_HOME/bin:$PATH
```

It is usually convenient to put those lines in your shell initialization file.

4.2.2 How do I use Netkit?

There are two ways to use Netkit : The manual way, and by using pre-configured labs. In the first case, you boot and control each machine individually, using the commands starting with a “v” (for virtual machine). In the second case, you can start a whole network in a single operation. The commands for controlling the lab start with a “l”. The man pages of those commands is available from <http://wiki.netkit.org/man/man7/netkit.7.html>

You must be careful not to forgot to stop your virtual machines and labs, using either *vhalt* or *lhalt*.

4.2.3 Example of using a lab

A lab is simply a directory containing at least a configuration file called *lab.conf*, and one directory for each virtual machine. In the case the lab available on iCampus, the network is composed of two pc, pc1 and pc2, both of them being connected to a router r1. The lab.conf file contains the following lines :

```
pc1[0]=A
pc2[0]=B
r1[0]=A
r1[1]=B
```

This means that pc1 and r1 are connected to a “virtual LAN” named A via their interface eth0, while pc2 and r1 are connected to the “virtual LAN” B via respectively their interfaces eth0 and eth1.

The directory of each device is initially empty, but will be used by Netkit to store their filesystem.

The lab directory can contain optional files. In the lab provided to you, the “pc1.startup” file contains the shell instructions to be executed on startup of the virtual machine. In this specific case, the script configures the interface eth0 to allow traffic exchanges between pc1 and r1, as well as the routing table entry to join pc2.

Starting a lab consists thus simply in unpacking the provided archive, going into the lab directory and typing *lstart* to start the network.

4.2.4 File sharing between virtual machines and host

Virtual machines can access to the directory of the lab they belong to. This repertory is mounted in their filesystem at the path */hostlab*.

4.2.5 Tools available on Netkit

As the virtual machines run Linux, standard networking tools such as hping, tcpdump, netstats etc. are available as usual.

Note that capturing network traces can be facilitated by using the *uml_dump* extension available at <http://kartoch.msi.unilim.fr/blog/?p=19> . This extension is already installed in the Netkit installation on the student lab. In order to capture the traffic exchanged on a given ‘virtual LAN’, you simply need to issue the command

`vdump <LAN name>` on the host. If you want to pipe the trace to Wireshark, you can use `vdump A | wireshark -i - -k`

In the lab provided in iCampus, you can find a simple [Python](#) client/server application that establishes TCP connections. Feel free to re-use this code to perform your analysis.

4.3 Questions

1. A TCP/IP stack receives a SYN segment with the sequence number set to 1234. What will be the value of the acknowledgement number in the returned SYN+ACK segment ?
2. Is it possible for a TCP/IP stack to return a SYN+ACK segment with the acknowledgement number set to 0 ? If no, explain why. If yes, what was the content of the received SYN segment.
3. Open the `tcpdump` packet trace `traces/trace.5connections_opening_closing.pcap` and identify the number of different TCP connections that are established and closed. For each connection, explain by which mechanism they are closed. Analyse the initial sequence numbers that are used in the SYN and SYN+ACK segments. How do these initial sequence numbers evolve ? Are they increased every 4 microseconds ?
4. The `tcpdump` packet trace `traces/trace.5connections.pcap` contains several connection attempts. Can you explain what is happening with these connection attempts ?
5. The `tcpdump` packet trace `traces/trace.ipv6.google.com.pcap` was collected from a popular website that is accessible by using IPv6. Explain the TCP options that are supported by the client and the server.
6. The `tcpdump` packet trace `traces/trace.sirius.info.ucl.ac.be.pcap` Was collected on the departmental server. What are the TCP options supported by this server ?
7. A TCP implementation maintains a Transmission Control Block (TCB) for each TCP connection. This TCB is a data structure that contains the complete “state” of each TCP connection. The TCB is described in [RFC 793](#). It contains first the identification of the TCP connection :
 - `localip` : the IP address of the local host
 - `remoteip` : the IP address of the remote host
 - `remoteport` : the TCP port used for this connection on the remote host
 - `localport` : the TCP port used for this connection on the local host. Note that when a client opens a TCP connection, the local port will often be chosen in the ephemeral port range ($49152 \leq \text{localport} \leq 65535$).
 - `sndnxt` : the sequence number of the next byte in the byte stream (the first byte of a new data segment that you send will use this sequence number)
 - `snduna` : the earliest sequence number that has been sent but has not yet been acknowledged
 - `rcvnxt` : the sequence number of the next byte that your implementation expects to receive from the remote host. For this exercise, you do not need to maintain a receive buffer and your implementation can discard the out-of-sequence segments that it receives
 - `sndwnd` : the current sending window
 - `rcvwnd` : the current window advertised by the receiver

Using the `traces/trace.sirius.info.ucl.ac.be.pcap` packet trace, what is the TCB of the connection on host `130.104.78.8` when it sends the third segment of the trace ?

1. The `tcpdump` packet trace `traces/trace.maps.google.com` was collected by containing a popular web site that provides mapping information. How many TCP connections were used to retrieve the information from this server ?

2. Some network monitoring tools such as `ntop` collect all the TCP segments sent and received by a host or a group of hosts and provide interesting statistics such as the number of TCP connections, the number of bytes exchanged over each TCP connection, ... Assuming that you can capture all the TCP segments sent by a host, propose the pseudocode of an application that would list all the TCP connections established and accepted by this host and the number of bytes exchanged over each connection. Do you need to count the number of bytes contained inside each segment to report the number of bytes exchanged over each TCP connection ?
3. There are two types of firewalls ¹ : special devices that are placed at the border of campus or enterprise networks and software that runs on endhosts. Software firewalls typically analyse all the packets that are received by a host and decide based on the packet's header and contents whether it can be processed by the host's network stack or must be discarded. System administrators often configure firewalls on laptop or student machines to prevent students from installing servers on their machines. How would you design a simple firewall that blocks all incoming TCP connections but still allows the host to establish TCP connections to any remote server ?
4. Using the `netkit` lab explained above, perform some tests by using `hping3(8)`. `hping3(8)` is a command line tool that allows anyone (having system administrator privileges) to send special IP packets and TCP segments. `hping3(8)` can be used to verify the configuration of firewalls ¹ or diagnose problems. We will use it to test the operation of the Linux TCP stack running inside `netkit`.
 1. On the server host, launch `tcpdump(1)` with `-vv` as parameter to collect all packets received from the client and display them. Using `hping3(8)` on the client host, send a valid SYN segment to one unused port on the server host (e.g. `12345`). What are the contents of the segment returned by the server ?
 2. Perform the same experiment, but now send a SYN segment towards port 7. This port is the default port for the discard service (see `services(5)`) launched by `xinetd(8)`). What segment does the server sends in reply ? What happens upon reception of this segment ? Explain your answer.
1. The Linux TCP/IP stack can be easily configured by using `sysctl(8)` to change kernel configuration variables. See <http://fasterdata.es.net/TCP-tuning/ip-sysctl-2.6.txt> for a recent list of the `sysctl` variables on the Linux TCP/IP stack. Try to disable the selective acknowledgements and the RFC1323 timestamp and large window options and open a TCP connection on port 7 on the server by using `:manpage:telnet(1)`. Check by using `tcpdump(1)` the effect of these kernel variables on the segments sent by the Linux stack in `netkit`.
2. Network administrators sometimes need to verify which networking daemons are active on a server. When logged on the server, several tools can be used to verify this. A first solution is to use the `netstat(8)` command. This command allows you to extract various statistics from the networking stack on the Linux kernel. For TCP, `netstat` can list all the active TCP connections with the state of their FSM. `netstat` supports the following options that could be useful during this exercises :
 - `-t` requests information about the TCP connections
 - `-n` requests numeric output (by default, `netstat` sends DNS queries to resolve IP addresses in hosts and uses `/etc/services` to convert port number in service names, `-n` is recommended on `netkit` machines)
 - `-e` provides more information about the state of the TCP connections
 - `-o` provides information about the timers
 - `-a` provides information about all TCP connections, not only those in the Established state

On the `netkit` lab, launch a daemon and start a TCP connection by using `telnet(1)` and use `netstat(8)` to verify the state of these connections.

A second solution to determine which network daemons are running on a server is to use a tool like `nmap(1)`. `nmap(1)` can be run remotely and thus can provide information about a host on which the system administrator cannot login. Use `tcpdump(1)` to collect the segments sent by `nmap(1)` running on the client and explain how `nmap(1)` operates.

¹ A firewall is a software or hardware device that analyses TCP/IP packets and decides, based on a set of rules, to accept or discard the packets received or sent. The rules used by a firewall usually depend on the value of some fields of the packets (e.g. type of transport protocols, ports, ...). We will discuss in more details the operation of firewalls in the network layer chapter.

1. Long lived TCP connections are susceptible to the so-called *RST attacks*. Try to find additional information about this attack and explain how a TCP stack could mitigate such attacks.

TCP CONGESTION CONTROL

The TCP congestion control mechanisms, defined in [RFC 5681](#) plays a key role in today's Internet. Without this mechanism that was first defined and implemented in the late 1980s, the Internet would not have been able to continue to work until now. The objective of this exercise is to allow you to have a better understanding of the operation of TCP's congestion control mechanism by analysing all the segments exchanged over a TCP connection.

1. To understand the operation of the TCP congestion control mechanism, it is useful to draw some time sequence diagrams. Let us consider a simple scenario of a web client connected to the Internet that wishes to retrieve a simple web page from a remote web server. For simplicity, we will assume that the delay between the client and the server is 0.5 seconds and that the packet transmission times on the client and the servers are negligible (e.g. they are both connected to a 1 Gbps network). We will also assume that the client and the server use 1 KBytes segments.
1. Compute the time required to open a TCP connection, send an HTTP request and retrieve a 16 KBytes web page. This page size is typical of the results returned by search engines like [google](#) or [bing](#). An important factor in this delay is the initial size of the TCP congestion window on the server. Assume first that the initial window is set to 1 segment as defined in [RFC 2001](#), 4 KBytes (i.e. 4 segments in this case) as proposed in [RFC 3390](#) or 16 KBytes as proposed in a recent [paper](#).
2. Perform the same analysis with an initial window of one segment is the third segment sent by the server is lost and the retransmission timeout is fixed and set to 2 seconds.
3. Same question as above but assume now that the 6th segment is lost.
4. Same question as above, but consider now the loss of the second and seventh acknowledgements sent by the client.
5. Does the analysis above changes if the initial window is set to 16 KBytes instead of one segment ?
2. Several MBytes have been sent on a TCP connection and it becomes idle for several minutes. Discuss which values should be used for the congestion window, slow start threshold and retransmission timers.
3. To operate reliably, a transport protocol that uses Go-back-n (resp. selective repeat) cannot use a window that is larger than $2^n - 1$ (resp. 2^{n-1}) segments. Does this limitation affects TCP ? Explain your answer.
4. Consider the simple network shown in the figure below. In this network, the router between the client and the server can only store on each outgoing interface one packet in addition to the packet that it is currently transmitting. It discards all the packets that arrive while its buffer is full. Assuming that you can neglect the transmission time of acknowledgements and that the server uses an initial window of one segment and has a retransmission timer set to 500 milliseconds, what is the time required to transmit 10 segments from the client to the server. Does the performance increase if the server uses an initial window of 16 segments instead ?

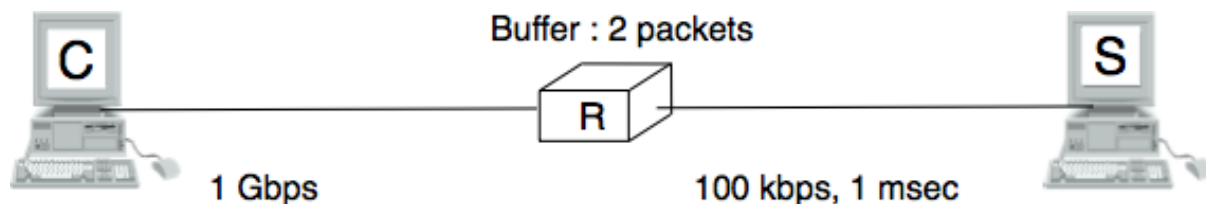


Figure 5.1: Simple network

5.1 Trace analysis

1. For the exercises below, we have performed measurements in an emulated ¹ network similar to the one shown below.

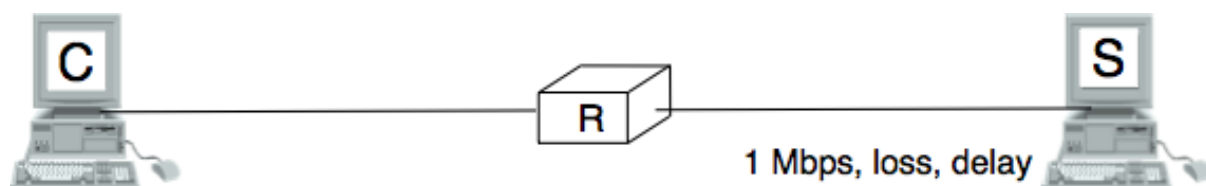


Figure 5.2: Emulated network

The emulated network is composed of three UML machines ²: a client, a server and a router. The client and the server are connected via the router. The client sends data to the server. The link between the router and the client is controlled by using the `netem` Linux kernel module. This module allows us to insert additional delays, reduce the link bandwidth and insert random packet losses.

We used `netem` to collect several traces :

- `traces/trace0.pcap`
- `traces/trace1.pcap`
- `traces/trace2.pcap`
- `traces/trace3.pcap`

Each team of two students will analyse these traces by using `wireshark` or `tcpdump`. For each trace, you should

1. Identify the TCP options that have been used on the TCP connection
2. Try to find explanations for the evolution of the round-trip-time on each of these TCP connections. For this, you can use the *round-trip-time* graph of `wireshark`, but be careful with their estimation as some versions of `wireshark` are buggy
3. Verify whether the TCP implementation used implemented *delayed acknowledgements*
4. Inside each packet trace, find :
 1. one segment that has been retransmitted by using *fast retransmit*. Explain this retransmission in details.
 2. one segment that has been retransmitted thanks to the expiration of TCP's retransmission timeout. Explain why this segment could not have been retransmitted by using *fast retransmit*.

¹ With an emulated network, it is more difficult to obtain quantitative results than with a real network since all the emulated machines need to share the same CPU and memory. This creates interactions between the different emulated machines that do not happen in the real world. However, since the objective of this exercise is only to allow the students to understand the behaviour of the TCP congestion control mechanism, this is not a severe problem.

² For more information about the TCP congestion control schemes implemented in the Linux kernel, see <http://linuxgazette.net/135/pfeiffer.html> and <http://www.cs.helsinki.fi/research/iwtcp/papers/linuxtcp.pdf> or the source code of a recent Linux. A description of some of the sysctl variables that allow to tune the TCP implementation in the Linux kernel may be found in <http://fasterdata.es.net/TCP-tuning/linux.html>. For this exercise, we have configured the Linux kernel to use the NewReno scheme [RFC 3782](#) that is very close to the official standard defined in [RFC 5681](#)

5. [wireshark](#) contain several two useful graphs : the *round-trip-time* graph and the *time sequence* graph. Explain how you would compute the same graph from such a trace .
 6. When displaying TCP segments, recent versions of [wireshark](#) contain *expert analysis* heuristics that indicate whether the segment has been retransmitted, whether it is a duplicate ack or whether the retransmission timeout has expired. Explain how you would implement the same heuristics as [wireshark](#).
 7. Can you find which file has been exchanged during the transfer ?
2. You have been hired as an networking expert by a company. In this company, users of a networked application complain that the network is very slow. The developers of the application argue that any delays are caused by packet losses and a buggy network. The network administrator argues that the network works perfectly and that the delays perceived by the users are caused by the applications or the servers where the application is running. To resolve the case and determine whether the problem is due to the network or the server on which the application is running. The network administrator has collected a representative packet trace that you can download from `traces/trace9.pcap`. By looking at the trace, can you resolve this case and indicate whether the network or the application is the culprit ?

INDICES AND TABLES

- *genindex*
- *search*

These notes are licensed under the creative commons attribution share-alike license 3.0. You can obtain detailed information about this license at <http://creativecommons.org/licenses/by-sa/3.0/>

INDEX

R

RFC

- RFC 1071, 3
- RFC 1321, 3
- RFC 1350, 8
- RFC 2001, 23
- RFC 2920, 3
- RFC 3390, 23
- RFC 3782, 24
- RFC 4634, 3
- RFC 5681, 23, 24
- RFC 793, 19