

Computer Networking and IT Security (INHN0012)

Tutorial 11

Problem 1 Flow and congestion control with TCP

The most widely used transport protocol on the Internet is TCP. This implements mechanisms for flow and congestion control.

a)* Discuss the differences between flow and congestion control. What are the objectives of each mechanism?

- **Flow control:**
Prevention of overload situations at the receiver
- **Congestion control:**
Adaption to overload situations in the network

b) Assign each of the following terms to TCP flow or congestion control:

- Slow-Start
- Receive window
- Congestion-Avoidance
- Multiplicative-Decrease

Only the receive window is part of the flow control, since the receiver can use it to tell the sender the maximum amount of data it can send at once. The remaining terms all belong to congestion control, with slow-start and congestion-avoidance being the two congestion control phases of a TCP connection. Multiplicative decrease, on the other hand, is the halving of the congestion control window when a segment is lost.

To analyze the data rate that can be achieved with TCP, we consider the course of a contiguous data transmission in which the slow-start phase has already been completed. TCP is therefore in the congestion-avoidance phase. We designate the individual windows as follows:

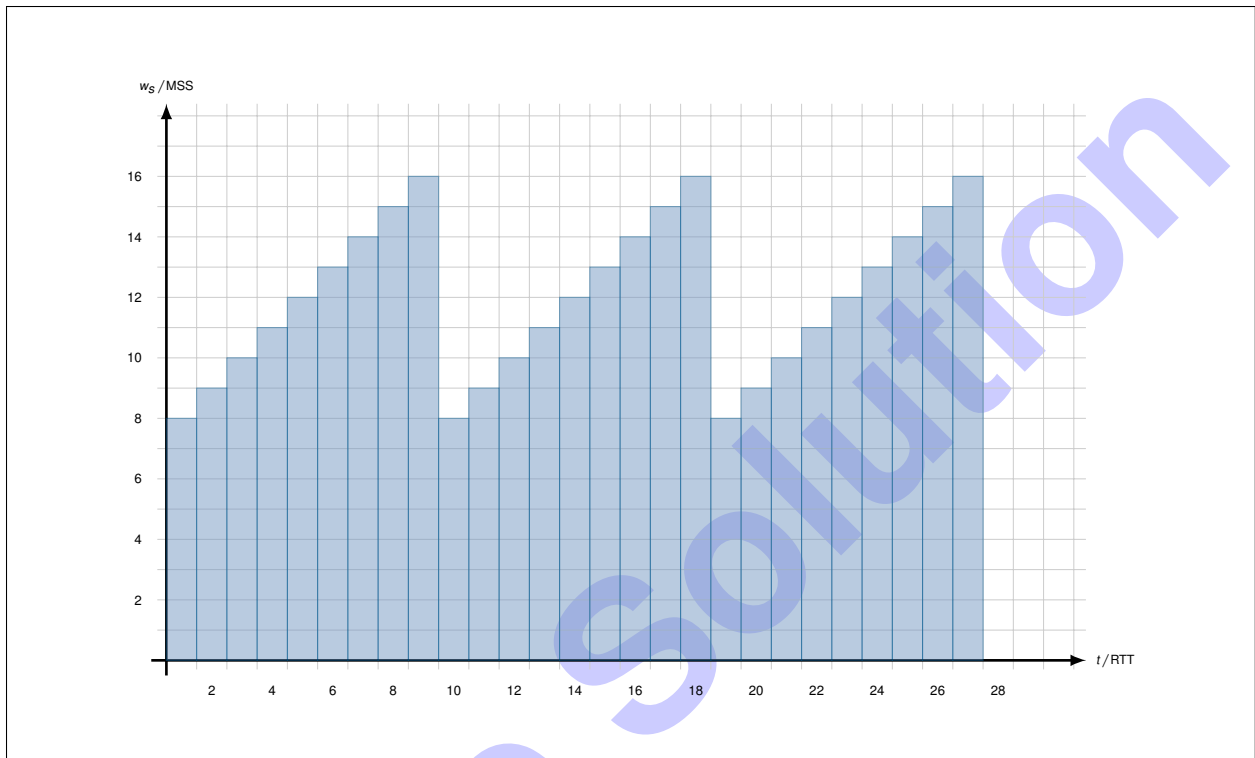
- Send window W_s , $|W_s| = w_s$
- Receive window W_r , $|W_r| = w_r$
- Traffic-control window W_c , $|W_c| = w_c$

We assume that the receive window is arbitrarily large, so that the transmit window is determined solely by the congestion control window, i. e., h. $W_s = W_c$. No losses occur as long as the transmit window is smaller than a maximum value x , i. e., $w_s < x$.

If a full send window is confirmed, the currently used window will increase by exactly 1 MSS. If the send window has reached the value x , exactly one of the sent TCP segments is lost. The sender detects the loss by receiving the same ACK number multiple times. The station then halves the congestion control window, but still remains in the congestion-avoidance phase, i. e., no new slow-start takes place. This approach corresponds to a simplified variant of TCP-Reno (cf. lecture).

As concrete numerical values, we assume that the maximum TCP segment size (MSS) is 1460 B and the RTT is 200 ms. Let the serialization time of segments be negligible compared to the propagation delay. Segment loss occurs as soon as the transmit window reaches a size of $w_s \geq x = 16 \text{ MSS}$.

c)* Create a graph plotting the current size of the transmit window w_s measured in MSS over the time axis t measured in RTT. In your diagram, at the time $t_0 = 0$ s, you want the transmit window size to have just been halved, so $w_s = x/2$ applies. Draw the diagram in the time interval $t = \{0, \dots, 27\}$.



d)* How much time elapses before the congestion control window is reduced again after a segment loss as a result of another segment loss?

After a segment loss, w_c is reduced to $x/2$ and then increased again by 1 MSS per fully acknowledged window. Thus, since the serialization time is negligible, a total of w_c segments can be sent at time t_0 , which are confirmed at time $t_0 + RTT$. Consequently, for the time until the maximum value is reached again, we get

$$T = \left(\frac{x}{2} + 1\right) \cdot RTT = 9 \cdot 200 \text{ ms} = 1.8 \text{ s.}$$

e)* In general, determine the average loss rate L . Note: Since the behavior of TCP is periodic in this idealized model, it is sufficient to consider only one period. Set the total number of transmitted segments in relation to the number of lost segments (specification as a truncated fraction is sufficient).

First, we determine the number n of segments that are transmitted during each „sawtooth“

$$\begin{aligned}
 n &= \sum_{i=x/2}^x i = \sum_{i=1}^x i - \sum_{i=1}^{x/2-1} i = \frac{x \cdot (x+1)}{2} - \frac{\left(\frac{x}{2}-1\right) \cdot \frac{x}{2}}{2} \\
 &= \frac{x^2+x}{2} - \frac{x^2}{8} + \frac{x}{4} \\
 &= \frac{3}{8}x^2 + \frac{3}{4}x \\
 &\stackrel{x=16}{=} 108
 \end{aligned}$$

Exactly one segment is lost per „sawtooth“. We therefore obtain for the loss rate

$$L = \frac{1}{\frac{3}{8}x^2 + \frac{3}{4}x} = \frac{1}{108} \approx 9.26 \cdot 10^{-3}$$

f) Using the results from subtasks (c) and (e), determine the average achievable transmission rate in kB/s during the TCP transmission phase under consideration. **Note:** Use the exact value (fraction) from subtask e).

For the data rate we get

$$\begin{aligned}
 r_{TCP} &= \frac{n \cdot MSS}{T} \cdot (1 - \theta) \\
 &= \frac{108 \cdot 1460 \text{ B}}{1.8 \text{ s}} \cdot \frac{107}{108} \\
 &= \frac{107 \cdot 1460 \text{ B}}{1.8 \text{ s}} \\
 &= \frac{1562200}{18} \text{ B/s} \\
 &\approx \frac{1562}{18} \text{ kB/s} \approx 86.79 \text{ kB/s.}
 \end{aligned}$$

g)* What is the maximum transmission rate you could send over the channel with UDP without creating a congestion? Take into account that the UDP header 12 B is smaller than the TCP header without options.

Apparently, 15 MSS can be reliably transferred. In addition, a segment carries 12 B more payload data in UDP than in TCP. So we get

$$\begin{aligned}
 r_{UDP} &= \frac{15 \cdot (MSS + 12 \text{ B})}{RTT} \\
 &= \frac{15 \cdot (1460 \text{ B} + 12 \text{ B})}{0.2 \text{ s}} \\
 &= \frac{15 \cdot 1472 \text{ B}}{0.2 \text{ s}} \\
 &= 110.40 \text{ kB/s.}
 \end{aligned}$$

Problem 2 Network Address Translation

This task will look at the forwarding of IP packets (IPv4) when using a NAT-enabled router. For the mapping between public and private IP addresses, a NAT-enabled router has a mapping table that stores the relationship between the local and the global port. Many NAT-enabled devices also store additional information such as the remote IP address or the router's own global IP address (e.g., if the router has more than one global IP). We will refrain from doing this here.

Figure 2.1 shows the network topology. Router R1 has NAT enabled, using a private IP address on IF1 and a public IP address on IF2. Router R2 does not use NAT. PC2 had already communicated with server 2, which created the entry in R1's NAT table (see the 2.1 figure). Where you have the freedom, choose sensible values for the IP addresses and port numbers.

a)* Give PC1 and interface 1 of R1 a suitable IP address. The subnet is 10.0.0.0/24.

Possible are e.g.

- PC1: 10.0.0.1
- R1 IF1: 10.0.0.254

b)* PC1 now establishes an HTTP connection to server 2. Specify the source IP, destination IP, source port, destination port, and TTL fields of the IP or TCP header for the packets in the three highlighted locations in Figure 2.1. Also, specify newly created entries in the NAT table of R1.

See figure 2.1.

- **Between PC1 and R1:** TTL = 64
Important for the source port is that it is larger than 1023 (values smaller than 1024 are *well-known-ports* and can not be used as a source port). Furthermore it should not be larger than 65535, because port numbers are 16 bit long. The destination port is given with TCP 80 (HTTP).
- **R1 and R2** TTL = 63
R1 replaces the private source IP with its own public IP address. The source port will usually (if not already used otherwise) stay the same. Otherwise it will also be changed, e.g. incremented. The choice of port numbers depends on the respective NAT type. If possible, we keep the same port number. At this point a new entry in the NAT table is created: [10.0.0.1, 3627, 3627].
- **Between R2 and Server 2** TTL = 62
No change because a normal router does not change IP addresses or port numbers. The TTL field is obviously decremented.

c) Server 2 now answers PC1. In Figure 2.2, specify the header fields at the three named locations, as well as newly created entries in the NAT table of R1, analogous to the previous subtask.

We assume that the server send out packets with a TTL = 64

- **Between Server 2 and R2** TTL = 64
The server first addresses the response to R1 (where else?).
- **Between R2 and R1** TTL = 63
R2 does not change anything (except TTL).
- **Between R1 and PC1:** TTL = 62
R1 uses the entry in its own NAT table to find the IP address of the real receiver. Afterwards the destination IP address and port number are replaced (if needed) and the packet is forwarded.

Sample Solution

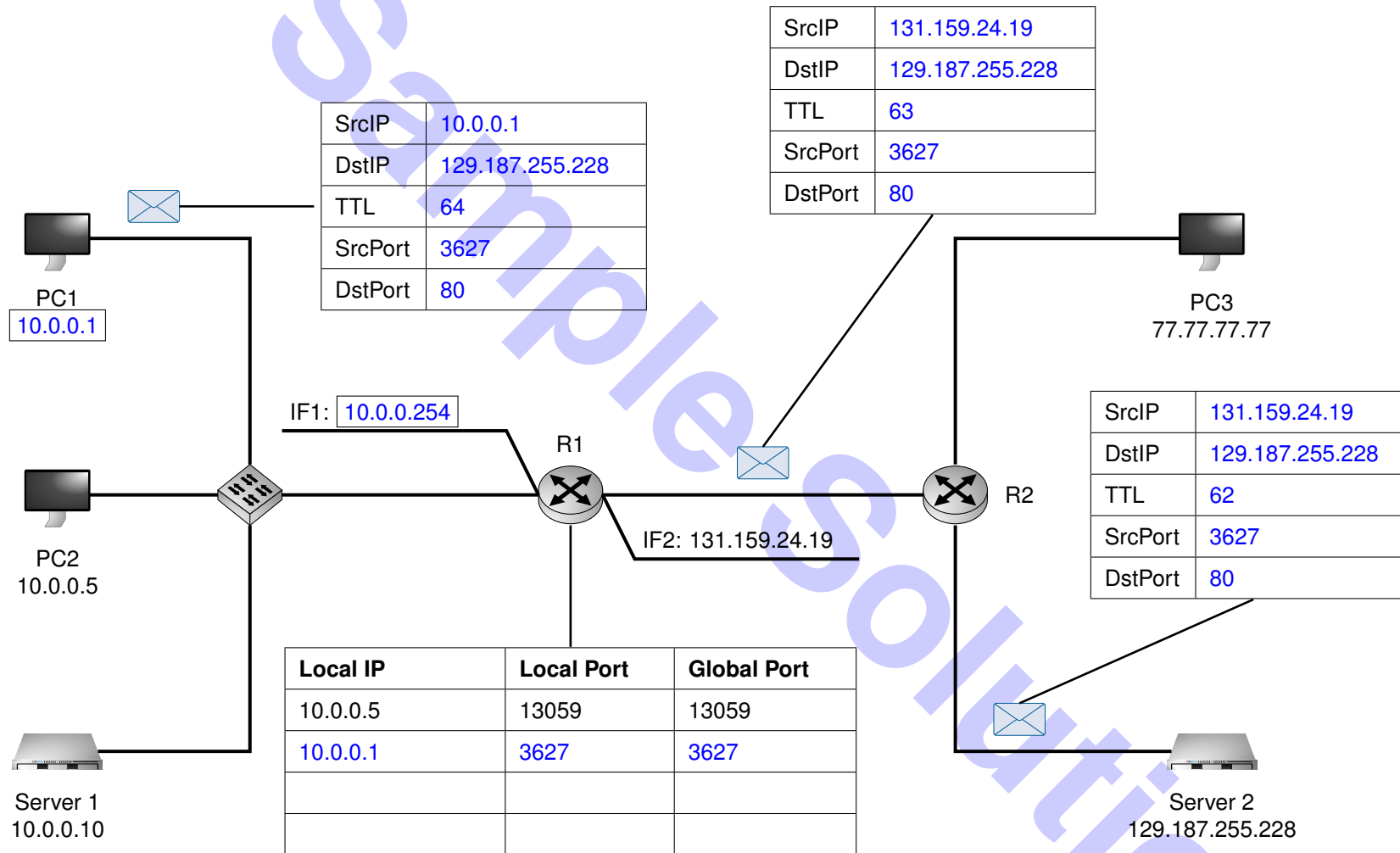


Figure 2.1: Lösungsblatt für Aufgabe 2a)/b)

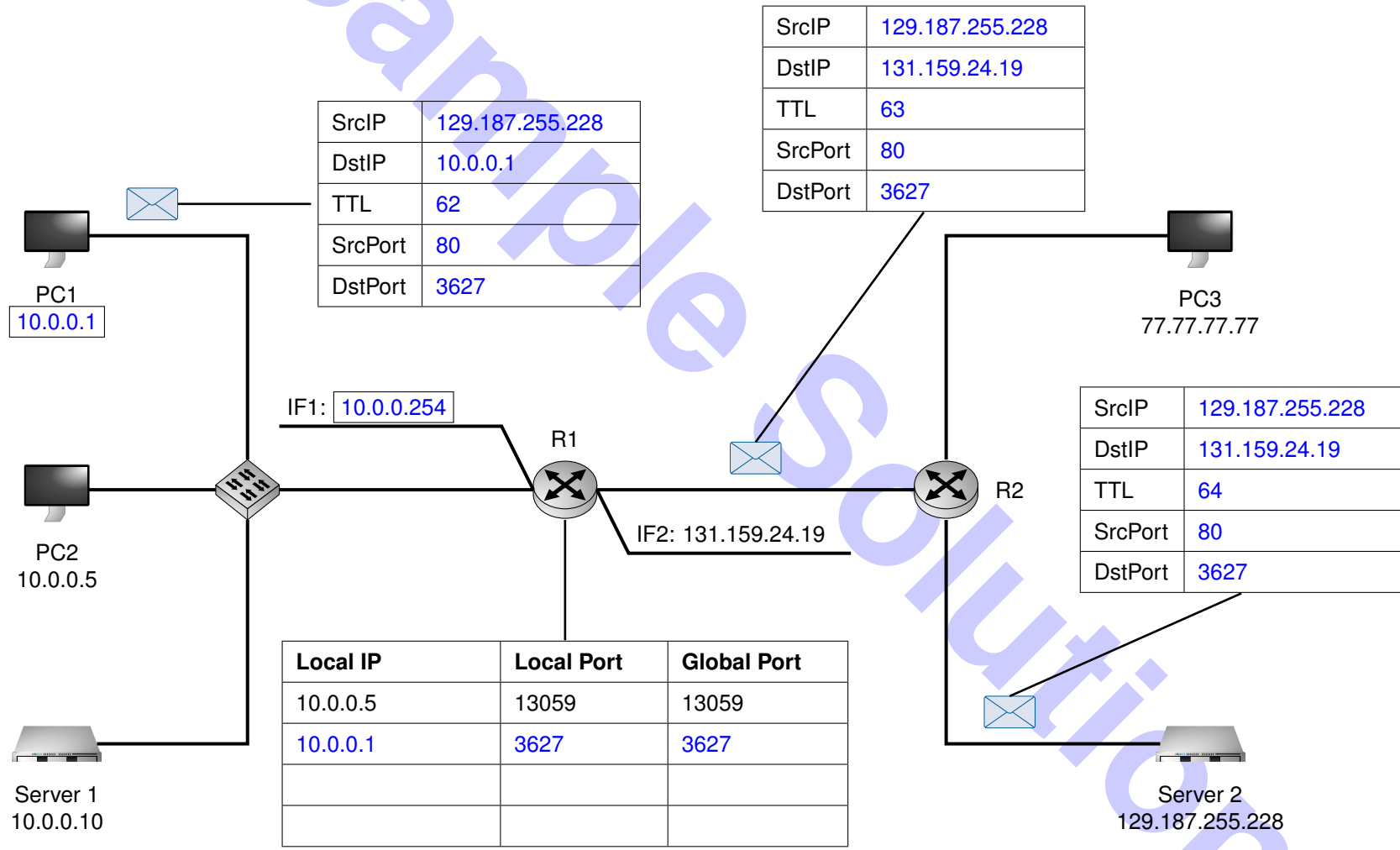


Figure 2.2: Solution sheet for problem 2c)

d)* Server 1 now also establishes a TCP connection to server 2 on port 80. In doing so, it randomly chooses the sender port 13059. Describe the problem that occurs on the NAT and how it is solved.

There is a collision with the first entry in the NAT table: The NAT router can no longer distinguish whether responses from server 2 are destined for PC1 or server 2, since the only distinguishing feature is the global port number.

The solution is that the NAT router checks whether the respective port is already in use before creating new entries. If it is, the NAT server chooses a random port number from the ephemeral port range (or increments the port number) and stores both the local and the new global port number. For incoming packets, the port number is translated back in the L4 PDUs.

e)* R1 receives a packet from PC3 addressed to 131.159.24.19:13059. What will R1 do with this packet? What problems may arise from it?

R1 will translate the destination address of the packet according to the NAT table and forward it to PC2, even though the original entry was created for Server2. PC2 receives an „unexpected“ packet and must be able to handle it. The erroneously often assumed firewall function of NAT cannot be enabled here.

f) Does a problem arise for PC2 when it receives a „random“ packet with TCP payload on a port with an existing connection?

The packet probably has a different sender IP and source port and thus is not associated with the existing connection. If the sender IP and source port „coincidentally“ match, the packet's sequence number (with high probability) does not fall within the currently valid receive window and is thus discarded.

g)* What other distinguishing criteria could be used by a NAT router?

global IP (if the router has multiple interfaces/IP addresses configured), remote IP, remote port and the protocol number (TCP or UDP).

h)* What problem arises when PC1 sends an echo request to server 2?

ICMP does not use port numbers therefore the NAT router can't create an entry. The response is thus discarded.

i) Describe a possible solution for the problem which arose in the previous subproblem.

The NAT router could in case of ICMP packet additionally to the protocol number use the ICMP identifier as a replacement for the missing port numbers. In this case, however, the NAT router has to also differentiate between the different L3-SDU protocols (TCP, UDP, ICMP, etc.).

j) What problem arises if a NAT router receives ICMP TTL-exceeded messages and wants to forward it to the intended receiver (sender who caused the problem)? How can this problem be worked around?

TTL Exceeded messages are own ICMP messages whose identifier was not entered in the NAT (messages are not generated in the own network, but from computers outside). An assignment to the recipient is therefore not possible. ICMP TTL Exceeded contain besides the ICMP header also the IP header and the first 8 payload bytes of the triggering packet¹. This allows the NAT to identify the triggering connection. For TCP and UDP the port numbers can be found here, for ICMP messages the original identifier.

k)* Now PC3 wants to establish a HTTP connection with Server 1. Can this happen under the given circumstances? (Explain!)

PC3 can't address the packet directly to the address 10.0.0.10 because it's a private IP address which is not routed in the public internet. If PC3 knows the public IP of R1, behind which Server 1 is located, it can address the packet to the IP address of R1 and TCP80; but R1 does (as far as the problem statement says) not have a suitable entry in its NAT table and therefore can't identify the correct receiver of the packet.

l) How can this problem be avoided while maintaining a NAT?

In the NAT a static forwarding (so called portforwarding) can be configured.
Example: 10.0.0.10 80 80 With this Server 1 can be reached with the ip address of R1 via the router R1 on Port 80 from the outside.

Problem 3 TCP and Long Fat Networks (Homework)

In this task we consider so-called *Long Fat Networks*. This refers to connections that have a high transmission rate but, in particular, also a high delay. Examples include satellite links which as a result of the long distance have high propagation delays. In particular, we want to investigate the impact on TCP congestion control.

a)* For TCP, the send window is selected depending on the receive window as well as the congestion control window. What is the exact relationship between the windows?

$$w_s = \min(w_r, w_c)$$

Let two users be connected to the Internet at high transmission rates via a geostationary satellite. The RTT between the two users is 800 ms, the transmission rate be $r = 24$ Mbit/s.

b)* How large should the transmission window (measured in bytes) be selected so that continuous transmission is possible?

The first ACK can arrive after one RTT at the earliest, provided that serialization times are neglected. This results in the following for the transmission window

$$w_s \geq \text{RTT} \cdot r = 800 \cdot 10^{-3} \text{ s} \cdot 24 \cdot 10^6 \frac{\text{bit}}{\text{s}} = 100 \cdot 24 \cdot 10^3 \text{ B/s} = 2.4 \text{ MB.}$$

c)* Why is the situation in subproblem b) a problem for TCP flow control? Hint: Have a look at the TCP header.

Since the transmit window is chosen as a minimum of receive and congestion windows, and the receiver informs the transmitter of its receive window via the receive window field, which is limited to 16 bit, the transmit window is also limited to a maximum value of $(2^{16} - 1) \text{ B} = 65535 \text{ B}$. However, according to subproblem b), we would need a transmit window of size $2.4 \cdot 10^6 \text{ B}$.

d)* Read Section 2 of RFC 1323 (<http://www.ietf.org/rfc/rfc1323.txt>, see Appendix). Describe the solution to the problem from subproblem c).

We need the *TCP-Window-Scaling* option, which ensures that the receive window is scaled by 2^x . The „shift.cnt“ field of the TCP window scaling option specifies the x exponent.

e) Determine the minimum value for the shift.cnt field of the TCP window scaling option.

$$\begin{aligned}
 (2^{16} \cdot 2^x) - 1 &\geq 2,4 \cdot 10^6 \\
 x &\geq \text{ld} \left(\frac{2,4 \cdot 10^6 + 1}{2^{16}} \right) \\
 &= \text{ld}(36,62) \approx 5,19 \Rightarrow x = 6
 \end{aligned}$$

Explanation: We look for the smallest exponent x such that the maximum receive-window is greater than the value of $2,4 \cdot 10^6$ B calculated in subproblem b). The receive window is 16 bit wide, so it can take the value $0xffff = 2^{16} - 1$ at most. So this value has to be scaled by 2^x .

A quick look at the size of TCP's sequence number space ($|S| = 2^{32}$, since SEQ and ACK numbers are 32 bit long fields) shows that we don't get a problem like in the "sliding window protocols" task.

f) Specify the header of the first TCP SYN packet that establishes the connection. To do this, use the concrete numerical values from the specification. A TCP header is shown again in Figure 3.1 as a reminder. There you will also find two forms for the solution.

Note: It is not necessary to fill the header in binary. However, please make it clear whether the numbers are represented in hexadecimal, decimal, or binary.

Assume that the size of the traffic control window is currently half of the value calculated in subproblem b). The MSS is 1200 B and the TCP connection is currently in the congestion-avoidance phase.

g) How long does it take for the window to fully utilize the line?

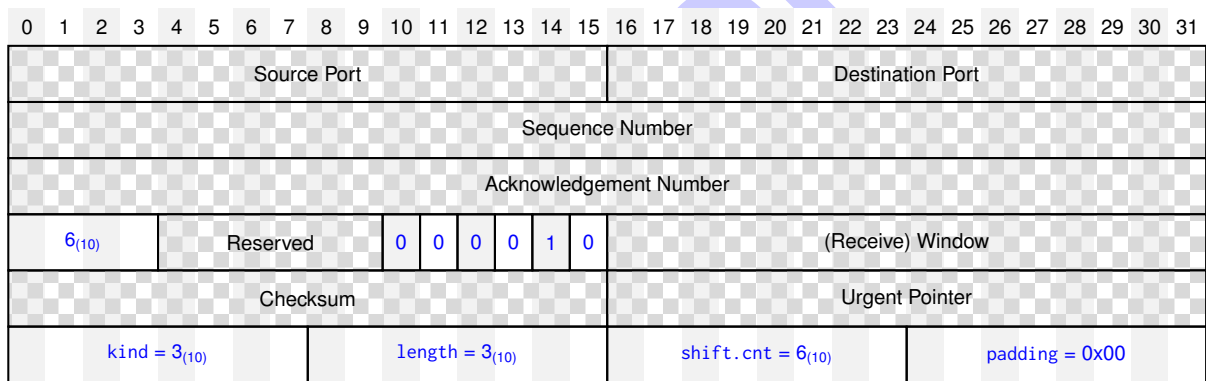
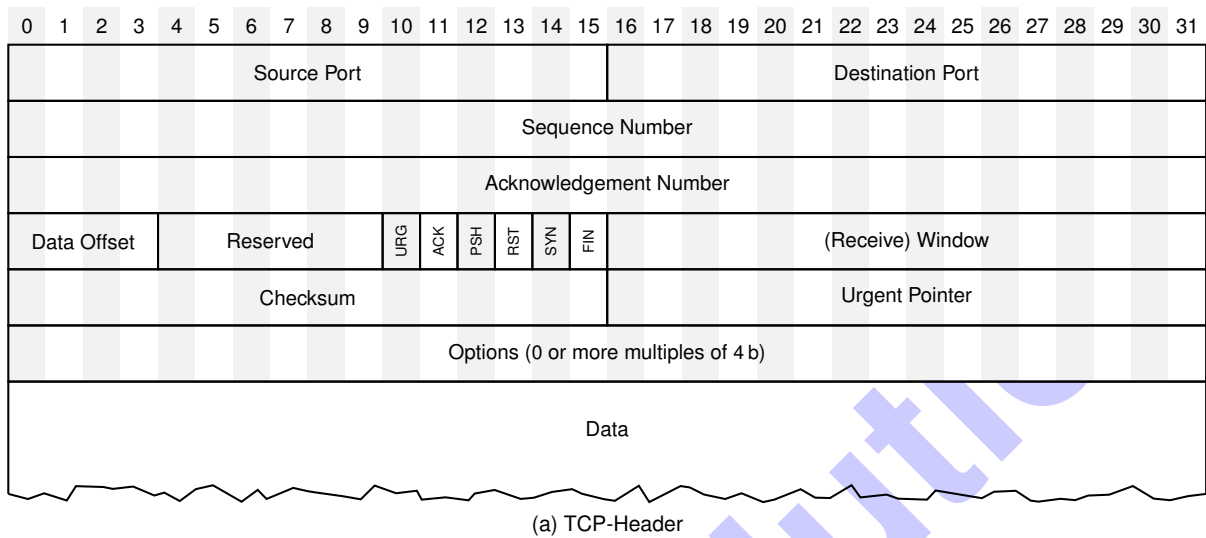
Note: The congestion control window is not directly affected by TCP window scaling.

The window is enlarged by 1 MSS per RTT. Consequently it takes

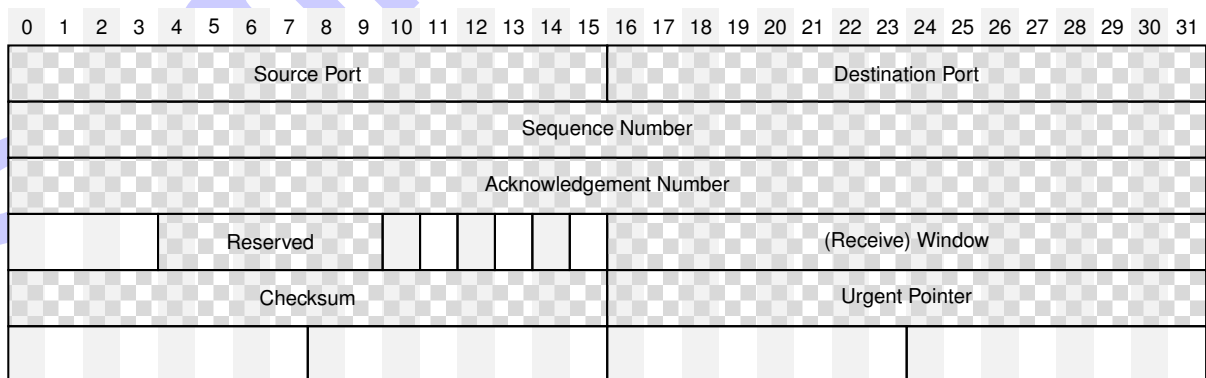
$$\frac{1,2 \cdot 10^6 \text{ B}}{1200 \text{ B}} \cdot 0,8 \text{ s} = \frac{96 \cdot 10^2}{12} \text{ s} = 800 \text{ s}$$

h) Does the result of subproblem g) result in a problem?

Yes. It takes more than 10 min for TCP to fully utilize the receive window again — way too long.



(b) Preprint



(c) Another preprint

Figure 3.1: TCP headers for solving task 3

2. TCP WINDOW SCALE OPTION

2.1 Introduction

The window scale extension expands the definition of the TCP window to 32 bits and then uses a scale factor to carry this 32-bit value in the 16-bit Window field of the TCP header (SEG.WND in RFC-793). The scale factor is carried in a new TCP option, Window Scale. This option is sent only in a SYN segment (a segment with the SYN bit on), hence the window scale is fixed in each direction when a connection is opened. (Another design choice would be to specify the window scale in every TCP segment. It would be incorrect to send a window scale option only when the scale factor changed, since a TCP option in an acknowledgement segment will not be delivered reliably (unless the ACK happens to be piggy-backed on data in the other direction). Fixing the scale when the connection is opened has the advantage of lower overhead but the disadvantage that the scale factor cannot be changed during the connection.)

The maximum receive window, and therefore the scale factor, is determined by the maximum receive buffer space. In a typical modern implementation, this maximum buffer space is set by default but can be overridden by a user program before a TCP connection is opened. This determines the scale factor, and therefore no new user interface is needed for window scaling.

2.2 Window Scale Option

The three-byte Window Scale option may be sent in a SYN segment by a TCP. It has two purposes: (1) indicate that the TCP is prepared to do both send and receive window scaling, and (2) communicate a scale factor to be applied to its receive window. Thus, a TCP that is prepared to scale windows should send the option, even if its own scale factor is 1. The scale factor is limited to a power of two and encoded logarithmically, so it may be implemented by binary shift operations.

TCP Window Scale Option (WSopt):

```
Kind: 3 Length: 3 bytes
+-----+-----+-----+
| Kind=3 |Length=3 |shift.cnt|
+-----+-----+-----+
```

This option is an offer, not a promise; both sides must send Window Scale options in their SYN segments to enable window scaling in either direction. If window scaling is enabled, then the TCP that sent this option will right-shift its true receive-window values by 'shift.cnt' bits for transmission in SEG.WND. The value 'shift.cnt' may be zero (offering to scale, while applying a scale factor of 1 to the receive window).

This option may be sent in an initial <SYN> segment (i.e., a segment with the SYN bit on and the ACK bit off). It may also be sent in a <SYN,ACK> segment, but only if a Window Scale option was received in the initial <SYN> segment. A Window Scale option in a segment without a SYN bit should be ignored.

The Window field in a SYN (i.e., a <SYN> or <SYN,ACK>) segment itself is never scaled.