



TECHNICAL UNIVERSITY OF MOLDOVA  
FACULTY OF COMPUTERS, INFORMATICS AND MICROELECTRONICS  
DEPARTMENT OF SOFTWARE ENGINEERING AND AUTOMATION

WEB PROGRAMMING  
LABORATORY WORK #2

---

## HTTP, caching and content negotiation

---

*Author:*  
Popa EUGENIU  
std. gr. FAF-202

*Supervisor:*  
Alexei ȘERȘUN

Chișinău 2023

# 1 Task

- Write a command line program, using go2web executable as a starting point.
- The program should implement at least the following CLI:

```
1 go2web -u <URL> # make an HTTP request to the specified URL and print the
  response
2 go2web -s <search-term> # make an HTTP request to search the term using your
  favorite search engine and print top 10 results
3 go2web -h # show this help
```

- The responses from request should be human-readable (e.g. no HTML tags in the output).

# 2 Results

- I have written the command line program in Python. [1].
- The program implements the following CLI:

```
1 go2web -u <URL>
2 go2web -s <search-term>
3 go2web -h
```

These commands can be used in this way: `python go2web.py [option] [argument]`

- The response from the request is completely human readable.

Function to access the given URL:

```
1 import re
2 import socket
3 import sys
4 from urllib.parse import parse_qs, unquote, urlparse
```

Import the required modules, socket for creating a TCP/IP socket and connecting to the web server, and BeautifulSoup for parsing the HTML response.

```
1 s = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
2 s.settimeout(5)
```

This creates a new TCP/IP socket object using the socket module, with the AF\_INET address family (IPv4) and SOCK\_STREAM socket type. The settimeout method sets a timeout period of 5 seconds on the socket, meaning that if a connection cannot be established within that time, a socket.timeout exception will be raised.

```
1 start_pos = url.find("//") + 2
2 end_pos = url.find("/", start_pos)
3 if end_pos == -1:
4     host = url[start_pos:]
5     path = ""
6 else:
7     host = url[start_pos:end_pos]
8     end_pos_query = url.find("?")
9     if end_pos_query == -1:
10         end_pos_frag = url.find("#")
```

```

11         if end_pos_frag == -1:
12             end_pos_frag = len(url)
13             path = url[end_pos:end_pos_frag]
14         else:
15             path = url[end_pos:end_pos_query]

```

These lines parse the given URL to extract the host and path components of the HTTP GET request. It finds the position of the double slashes // in the URL and adds 2 to get the start position of the host. It also finds the first forward slash / after the double slashes to determine the end of the host component and the beginning of the path component. If the forward slash is not found, it means that the URL does not contain a path component, so the path variable is set to an empty string. Otherwise, it extracts the path component from the URL, which may contain query parameters and/or a fragment identifier.

```

1     port = 80
2     try:
3         s.connect((host, port))
4     except socket.timeout:
5         print("Connection timed out.")
6         s.close()
7         exit()

```

This sets the default port number for HTTP requests to 80 and establishes a connection to the web server using the connect() method of the socket object. It tries to connect to the host and port specified, and if the connection times out (due to the settimeout value), it will print an error message, close the socket, and exit the program.

```

1     print("Sending request to", host, "at path", path, "\n")
2
3     request = "GET " + path + " HTTP/1.1\r\nHost: " + host + "\r\n\r\n"
4     s.sendall(request.encode())

```

This prints a message indicating the host and path of the HTTP GET request and then creates a GET request header with HTTP/1.1 protocol version, host name, and a blank line to signal the end of the header. It sends the request using the sendall() method of the socket object, which sends the entire message in one go.

```

1     response = b""
2     try:
3         while True:
4             buffer_size = 3072
5             data = s.recv(buffer_size)
6             if not data:
7                 break
8             response += data
9     except socket.timeout:
10         pass

```

This receives the response from the server and stores it in a response variable. It uses a loop to receive data in chunks of 3072 bytes until there is no more data to receive. The recv() method of the socket object is used to receive data. If the method returns an empty byte string, it means that there is no more data to receive, and the loop breaks. If the method times out (due to the settimeout value), the loop will also break.

```

1     content = response.decode()
2     body_start = content.find('\r\n\r\n') + 4
3     body = content[body_start:]
4

```

```

5 soup = BeautifulSoup(body, 'html.parser')
6 body = soup.get_text()
7 body = ' '.join(body.split())
8
9 print('\nContent:')
10 print(body)

```

This decodes the response from bytes to a string and extracts the body of the response by finding the start of the body after the header. It then uses BeautifulSoup to parse the HTML response and extract the text content of the HTML. Finally, it cleans up the text by removing any extra whitespace and prints the resulting text to the console.

Overall, this function takes a URL as input, connects to the web server using a TCP/IP socket, sends an HTTP GET request, receives the response, extracts the body content of the response, and prints it to the console.

Output example:

```

1 python go2web.py -u https://quotes.toscrape.com/
2 Sending request to quotes.toscrape.com at path /
3
4 Content:
5 Quotes to Scrape Quotes to Scrape Login "The world as we have created it is a process
  of our thinking. It cannot be changed without changing our thinking". by Albert
  Einstein (about) Tags: change deep-thoughts thinking world "It is our choices,
  Harry, that show what we truly are, far more than our abilities". by J.K. Rowling
  (about) Tags: abilities choices "There are only two ways to live your life. One is
  as though nothing is a miracle. The other is as though everything is a miracle".
  by Albert Einstein (about) Tags: inspirational life live miracle miracles "The
  person, be it gentleman or lady, who has not pleasure in a good novel, must be
  intolerably stupid". by Jane Austen (about) Tags: aliteracy books classic humor
  "Imperfection is beauty, madness is genius and it's better to be absolutely
  ridiculous than absolutely boring". by Marilyn Monroe (about) Tags: be-yourself
  inspirational "Try not to become a man of success. Rather become a man of value".
  by Albert Einstein (about) Tags: adulthood success value "It is better to be hated
  for what you are than to be loved for what you are not". by André Gide (about)
  Tags: life love "I have not failed. I've just found 10,000 ways that won't work".
  by Thomas A. Edison (about) Tags: edison failure inspirational paraphrased "A
  woman is like a tea bag; you never know how strong it is until it's in hot water".
  by Eleanor Roosevelt (about) Tags: misattributed-eleanor-roosevelt "A day without
  sunshine is like, you know, night". by Steve Martin (about) Tags: humor obvious
  simile Next → Top Ten tags love inspirational life humor books reading friendship
  friends truth simile Quotes by: GoodReads.com Made with by Scrapinghub

```

Function to print the top 10 results of the search term:

```

1 import re
2 import socket
3 import sys
4 from urllib.parse import parse_qs, unquote, urlparse

```

Import the required modules.

```

1 s = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
2 s.settimeout(5)
3
4 host = "www.google.com"
5 port = 80
6
7 try:
8     s.connect((host, port))

```

```

9     except socket.timeout:
10         print("Connection timed out.")
11         s.close()
12         exit()

```

This creates a TCP/IP socket and connects to the Google search page. If a connection cannot be established within 5 seconds, the function prints an error message and exits.

```

1     encoded_search_term = search_term.replace(" ", "+")
2     request = f'GET /search?q={encoded_search_term} HTTP/1.1\r\nHost: www.google.com\r\n\r\n'
3     s.sendall(request.encode())

```

This encodes the search term by replacing any spaces with "+" characters and sends an HTTP GET request to the Google search page. The request includes the encoded search term as a query parameter.

```

1     response = b""
2     try:
3         while True:
4             buffer_size = 3072
5             data = s.recv(buffer_size)
6             if not data:
7                 break
8             response += data
9     except socket.timeout:
10         pass

```

This receives the response from the server in chunks of 3072 bytes until there is no more data to receive. If the method times out (due to the settimeout value), the loop will break.

```

1     print("Top 10 results from Google:")
2     print()
3
4     soup = BeautifulSoup(response, 'html.parser')
5     links = soup.find_all(class_="egMi0 kCrYT")
6
7     for i, elem in enumerate(links[:10]):
8         title = elem.find('h3').get_text()
9         url_elem = elem.find('a')
10        url = unquote(url_elem['href'])
11        url_match = re.search('(?P<url>https?:/[^\&\s]+)', url)
12        url = url_match.group("url")
13        url = unquote(url)
14        url = re.sub('&sa=.*', '', url)
15        url_parts = urlparse(url)
16        url_params = parse_qs(url_parts.query)
17
18        if 'url' in url_params:
19            url = url_params['url'][0]
20
21        print(f'{{int(i+1)}}. {{title}}')
22        print(f'Access the link: {url}')
23        print()

```

This parses the HTML response using BeautifulSoup and extracts the top 10 search results. For each result, it extracts the title and URL from the HTML, decodes the URL using unquote(), removes any unnecessary query parameters, and prints the result to the console.

Output example:

```

1 python go2web.py -s "hello kitty"
2 Top 10 results from Google:
3
4 1. The Official Home of Hello Kitty & Friends
5 Access the link: https://www.sanrio.com/
6
7 2. HELLO KITTY — Sanrio
8 Access the link: https://www.sanrio.com/collections/hello-kitty
9
10 3. Hello Kitty — Wikipedia
11 Access the link: https://en.wikipedia.org/wiki/Hello_Kitty
12
13 4. Produse Hello Kitty — eMAG.ro
14 Access the link: https://www.emag.ro/brands/brand/hello-kitty
15
16 5. Hello Kitty — Facebook
17 Access the link: https://www.facebook.com/hellokitty/
18
19 6. Avril Lavigne — Hello Kitty (Official Video) — YouTube
20 Access the link: https://www.youtube.com/watch?v=LiaYDPRedWQ
21
22 7. Hello Kitty in romana — Seria 1 — YouTube
23 Access the link: https://www.youtube.com/watch?v=lyNWq5oXFZY
24
25 8. Hello Kitty (@hellokitty) / Twitter
26 Access the link: https://twitter.com/hellokitty
27
28 9. Hello Kitty (@hellokitty) • Instagram photos and videos
29 Access the link: https://www.instagram.com/hellokitty/
30
31 10. Cumpara autocolante auto Hello Kitty AK47 in Chisinau — Autoshina
32 Access the link: https://autoshina.md/ro/autocolante-auto-hello-kitty-ak47/

```

Function to print the help menu:

```

1 def help():
2     help_message = """Available CLI:
3     -u <URL>           # make an HTTP request to the specified URL and print the
4     response
5     -s <search-term> # make an HTTP request to search the term using Google and
6     print top 10 results
7     -h               # show this help
8 Usage: python go2web.py [option] [argument]"""
9
10    print(help_message)

```

### 3 Conclusion

This laboratory work was quite challenging but also rewarding. It required a good understanding of HTTP requests and responses, as well as the ability to work with sockets and parse HTML content. One of the most challenging aspects was making HTTP requests without using any third-party libraries. This required a lot of low-level programming using the socket module in Python. However, it also helped to gain a better understanding of how HTTP works under the hood and how web browsers interact with web servers. Overall, this laboratory work provided a great opportunity to learn more about web development and HTTP, and to develop some practical skills in making HTTP requests and parsing HTML content. It also helped to understand the importance of caching and content negotiation in web development.

## References

[1] Python

## 4 GitHub

<https://github.com/eugencic/utm-pw/tree/main/lab2>