# Introduction to Web Science

**Assignment 4**

Prof. Dr. Steffen Staab                René Pickhardt

staab@uni-koblenz.de                rpickhardt@uni-koblenz.de

Korok Sengupta

koroksengupta@uni-koblenz.de

Institute of Web Science and Technologies
Department of Computer Science
University of Koblenz-Landau

Submission until:   November 23, 2016, 10:00 a.m.
Tutorial on:   November 25, 2016, 12:00 p.m.

In this assignment we cover two topics: 1) **HTTP** & 2) **Web Content**

For all the assignment questions that require you to write code, make sure to include the code in the answer sheet, along with a separate python file. Where screen shots are required, please add them in the answers directly and not as separate files.

Delta Group Members

Oana Dumitrasc
odumitrasc@uni-koblenz.de

Alisa Becker
alisabecker@uni-koblenz.de

Omar Aly
oaly@uni-koblenz.de

# 1 Implementing a simplified HTTP GET Request (15 Points)

The goal of this exercise is to review the hyptertext transfer protocol and gain a better understanding of how it works.

Your task is to use the python programming language to create an HTTP client (http-client.py) that takes a URL as a command line argument and is able to download an arbitrary file from the World Wide Web and store it on your hard drive (in the same directory as your python code is running). The program should also print out the complete HTTP header of the response and store the header in a seperated file.

Your programm should only use the socket library so that you can open a TCP socket and and sys library to do command line parsing. You can either use urlparse lib or your code from assignment 3 in order to process the url which should be retrieved.

Your programm should be able to sucessfully download at least the following files:

1. http://west.uni-koblenz.de/en/studying/courses/ws1617/introduction-to-web-science

2. http://west.uni-koblenz.de/sites/default/files/styles/personen_bild/public/_IMG0076-Bearbeitet_03.jpg

**Use of libraries like `httplib`, `urllib`, etc are not allowed in this task.**

## 1.1 Hints:

There will be quite some challenges in order to finnish the task

- Your program only has to be able to process HTTP-responses with status 200 OK.

- Make sure you receive the full response from your TCP socket. (create a function handling this task)

- Sperated the HTTP header from the body (again create a function to do this)

- If a binary file is requested make sure it is not stored in a corrupted way

## 1.2 Example

```
1: python httpclient.py http://west.uni-koblenz.de/index.php
2:
3: HTTP/1.1 200 OK
4: Date: Wed, 16 Nov 2016 13:19:19 GMT
5: Server: Apache/2.4.7 (Ubuntu)
6: X-Powered-By: PHP/5.5.9-1ubuntu4.20
7: X-Drupal-Cache: HIT
8: Etag: "1479302344-0"
9: Content-Language: de
```

```
10: X-Frame-Options: SAMEORIGIN
11: X-UA-Compatible: IE=edge,chrome=1
12: X-Generator: Drupal 7 (http://drupal.org)
13: Link: <http://west.uni-koblenz.de/de>; rel="canonical",<http://west.uni-koblenz.de
14: Cache-Control: public, max-age=0
15: Last-Modified: Wed, 16 Nov 2016 13:19:04 GMT
16: Expires: Sun, 19 Nov 1978 05:00:00 GMT
17: Vary: Cookie,Accept-Encoding
18: Connection: close
19: Content-Type: text/html; charset=utf-8
```

The header will be printed and stored in index.php.header. The retrieved html document will be stored in index.php

**HTTP Client - GET Request**

```
 1: import socket
 2: from urllib.parse import urlparse
 3: import time
 4:
 5:
 6: # function to separate header from body & print header to console
 7: # also save body in the appropriate file format according to the
 8: # content type value of the header
 9: def getHeaderAndBody(receivedData, fPath):
10:     # find the starting part of body
11:     startOfBody = receivedData.find(b"\r\n\r\n")
12:     # creating the header part
13:     header = receivedData[:startOfBody]
14:     # creating the body part
15:     # the constant 4 is just to escape the \r\n\r\n characters
16:     body = receivedData[(startOfBody + 4):]
17:
18:     # if the HTTP status is 200 OK we process header & body
19:     # else we just print out the header
20:     if (str(header)).find("200 OK") >= 0:
21:         # creating the header dictionary
22:         # to know the content type
23:         headerDictionary = createHeaderDict(header.decode("utf-8"))
24:         # extracting the content type of the body
25:         contentType = headerDictionary["Content-Type"]
26:         # getting the image name by splitting on the url path
27:         splittedPath = fPath.split("/")
28:         fileName = splittedPath[-1]
29:         # according to the content type
30:         # save the body in appropriate format
31:         if contentType.find("image") >= 0:
32:             # writing to the body file
33:             # which will be created if it doesn't exist
```

```
34:             # using 'wb' to write bytes to the file and not strings
35:             bodyFile = open(fileName, "wb")
36:             bodyFile.write(body)
37:             # writing to the header file
38:             # which will be created if it doesn't exist
39:             headerFile = open("image_header.txt", 'w')
40:             headerFile.write(header.decode("utf-8"))
41:         elif contentType.find("text/html") >= 0:
42:             # if the file name doesn't have extension we add it
43:             # using the content type value
44:             content = (((contentType.split(";"))[0]).split("/"))[1]
45:             fileExtension = "" if fileName.find(".") >= 0 else content
46:             # writing to the body file
47:             # which will be created if it doesn't exist
48:             # using 'w' to write string to the file
49:             bodyFile = open(fileName + "." + fileExtension, 'w')
50:             bodyFile.write(body.decode("utf-8"))
51:             # writing to the header file
52:             # which will be created if it doesn't exist
53:             headerFile = open("html_header.txt", 'w')
54:             headerFile.write(header.decode("utf-8"))
55:     # showing the header file content on the console
56:     print(header.decode("utf-8"))
57:
58:
59: # method to create a dictionary for header content from header string
60: def createHeaderDict(header):
61:     headerDict = dict()
62:     headerLines = header.split("\r\n")
63:     for i in range(1, len(headerLines)):
64:         keyValue = headerLines[i].split(": ")
65:         headerDict[keyValue[0]] = keyValue[1]
66:
67:     return headerDict
68:
69:
70: # creating function to keep socket open till we receive all data
71: # or connection timeout(default 3 seconds)
72: def recvFullResponse(socket, timeout=3.0):
73:     # using non blocking sockets
74:     socket.setblocking(0)
75:
76:     # variables to hold our partial and final completed data
77:     # received from the server
78:     allData = list()
79:     partialData = ''
80:
81:     # begin counting time for timeout
82:     startTime = time.time()
```

```
83:
84:     # variable to determine how much data we want to receive
85:     bytesSize = 1024
86:
87:     while True:
88:         # we don't need to wait if we didn't receive any further data
89:         # than what we already received
90:         if allData and time.time() - startTime > timeout:
91:             break
92:         # wait double timeout if we didn't get any data at all
93:         # then break with a readable message to the user
94:         elif time.time() - startTime > timeout * 2:
95:             return "Time out and No Data Recieved"
96:         else:
97:             # try catch block
98:             # trying to receive data if nothing received we just do
99:             # another loop till the server responded with some data
100:            try:
101:                # receiving data
102:                partialData = socket.recv(bytesSize)
103:                # check if the get request result status is 200 OK
104:                # the bytesSize variable is used to not recheck
105:                # for each and every next part of the data
106:                if bytesSize == 1024 and (str(partialData)).find("200 OK") < 0:
107:                    return ''.join(partialData)
108:                else:
109:                    bytesSize = 4096
110:                    if partialData:
111:                        # append new data to already received data
112:                        allData.append(partialData)
113:                        # for sure reset start time if we get data
114:                        # to be sure we didn't get timeout before
115:                        # data receiving completion
116:                        startTime = time.time()
117:                    # wait for a while before doing another loop
118:                    else:
119:                        time.sleep(0.1)
120:
121:            except:
122:                pass
123:
124:     # use join with empty byte to return data
125:     # in a one byte sequence, for later extracting header purposes
126:     return b''.join(allData)
127:
128:
129: try:
130:     # create TCP socket
131:     sckt = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
```

```
132: except:
133:     print("Socket Creation Failed")
134:
135: # accepting user inputs from command line
136: userInput = input("URL: ")
137: # parsing url using urlparse library
138: parsedURL = urlparse(userInput)
139: filePath = parsedURL.path
140: # creating server address using the server name from
141: # parseURL.netloc parameter and setting the port to
142: serverAdd = (parsedURL.netloc, 80)
143: # creating get request using the path extracted from the given URL
144: getRequest = "GET " + filePath + " HTTP/1.0\r\nHost: " + parsedURL.netloc + "\r\n
145:
146: try:
147:     # connecting to the server
148:     sckt.connect(serverAdd)
149: except socket.error:
150:     print("Connection to the server failed")
151:
152: try:
153:     # sending request to the server
154:     sckt.sendall(bytes(getRequest, "UTF-8"))
155:     # storing received data in a variable
156:     receivedData = recvFullResponse(sckt)
157:     # doing the final work
158:     getHeaderAndBody(receivedData, filePath)
159: except:
160:     print("Sending Message Failed")
161:
162: # closing socket
163: sckt.close()
```

## 2 Download Everything (15 Points)

If you have successfully managed to solve the previous exercise you are able to download a web page from any url. Unfortionately in order to successfully render that very webpage the browser might need to download all the included images

In this exercise you should create a python file (downloadEverything.py) which takes two arguments. The first argument should be a name of a locally stored html file. The second argument is the url from which this file was downloaded.

Your program should

1. be able to find a list of urls the images that need to be downloaded for successful rendering the html file.

2. print the list of URLs to the console.

3. call the program from task 1 (or if you couldn't complete task 1 you can call wget or use any python lib to fulfill the http request) to download all the necessary images and store them on your hard drive.

**To finnish the task you are allowed to use the 're' library for regualar expressions and everything that you have been allowed to use in task 1.**

### 2.1 Hints

1. If you couldn't finnish the last task you can simulate the relevant behavior by using the program wget which is available in almost any UNIX shell.

2. Some files mentioned in the html file might use relative or absolut paths and not fully qualified urls. Those should be fixed to the correct full urls.

3. In case you run problems with constructing urls from relative or absult file paths you can always check with your web browser how the url is dereferenced.

**Download Everything**

```
 1: import socket
 2: from urllib.parse import urlparse
 3: import re
 4:
 5: # importing another python file
 6: # this syntax it is used in order to avoid
 7: # the problems caused by the hyphen in the name
 8: taskOne = __import__("http-client")
 9:
10: def getUrls(fileName):
```

```
11:      # opening the file from which we want to get
12:      # the urls for the images
13:      file = open(fileName, 'r')
14:      fileContent = file.read()
15:      # find all the image tags in order to get the correct urls
16:      urls = re.findall('<img [^>]*src=\"([^\"]+)', fileContent)
17:      prsdURL = urlparse(inputUrl)
18:      # looping over the urls
19:      # if the first character is /
20:      # then this is a relative path
21:      # and we have to make it absolute
22:      for url in urls:
23:          if url[0] == '/':
24:              urls[urls.index(url)] = prsdURL.scheme + "://" + prsdURL.netloc + url
25:      return urls
26:
27:
28: # accepting user inputs from command line
29: inputUrl = input("URL: ")
30: localFileName = input("File Name: ")
31: imagesUrls = getUrls(localFileName)
32: # printing out the urls
33: print(imagesUrls)
34: for imgUrl in imagesUrls:
35:
36:      try:
37:          # create TCP socket
38:          sckt = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
39:      except:
40:          print("Socket Creation Failed")
41:
42:      # parsing url using urlparse library
43:      parseURL = urlparse(imgUrl)
44:      flPath = parseURL.path
45:      # creating server address using the server name from
46:      # parseURL.netloc parameter and setting the port to
47:      serverAdd = (parseURL.netloc, 80)
48:      # creating get request using the path extracted from the given URL
49:      getRequest = "GET " + flPath + " HTTP/1.0\r\nHost: " + parseURL.netloc + "\r\n
50:
51:      try:
52:          # connecting to the server
53:          sckt.connect(serverAdd)
54:      except socket.error:
55:          print("Connection to the server failed")
56:
57:      try:
58:          # sending request to the server
59:          sckt.sendall(bytes(getRequest, "UTF-8"))
```

```
60:            # using the methods from the first task in order to download the images
61:            taskOne.getHeaderAndBody(taskOne.recvFullResponse(sckt), flPath)
62:        except:
63:            print("Sending Message Failed")
64:
65:        sckt.close()
```

## Important Notes

### Submission

- Solutions have to be checked into the github repository. Use the directory name `groupname/assignment4/` in your group's repository.

- The name of the group and the names of all participating students must be listed on each submission.

- Solution format: all solutions as *one* PDF document. Programming code has to be submitted as Python code to the github repository. Upload *all* `.py` files of your program! Use `UTF-8` as the file encoding. *Other encodings will not be taken into account!*

- Check that your code compiles without errors.

- Make sure your code is formatted to be easy to read.
    - Make sure you code has consistent indentation.
    - Make sure you comment and document your code adequately in English.
    - Choose consistent and intuitive names for your identifiers.

- Do *not* use any accents, spaces or special characters in your filenames.

### Acknowledgment

This latex template was created by Lukas Schmelzeisen for the tutorials of "Web Information Retrieval".

### LaTeX

Currently the code can only be build using LuaLaTeX, so make sure you have that installed. If on Overleaf, there's an error, go to settings and change the LaTeXengine to `LuaLaTeX`.