# 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.

Group Name : Yankee
Candidates :

1. Sabin Bhattarai - 216203590
   sbhattarai@uni-koblenz.de

2. Biplov K.C. - 216203865
   biplov@uni-koblenz.de

3. Syed Salman Ali. - 216203923
   salmanali@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.d
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

**Answer 1**



**Figure 1:** Obtained output showing the printed header file

```
 1: import socket
 2: import sys
 3: from urllib.parse import urlparse
 4:
 5: ENCODING = 'utf-8'
 6:
 7: #takes command line arguments and performs downloading files
 8: def mainFunction(commandArguments , totalArg):
 9:     if (totalArg > 1):
10:         downloadArbitaryFileWEB(commandArguments[1])
11:     else:
12:         print ("\n Invaild URL -- Default URL shown below will be used : " +
13:         "\n http://west.uni-koblenz.de/en/studying/courses/ws1617/" +
14:         "introduction-to-web-science \n\n")
15:         downloadArbitaryFileWEB("http://west.uni-koblenz.de/en/studying/" +
16:         "courses/ws1617/introduction-to-web-science")
17:
18:
19: #Connects to the server and receives response. writes to a file
20: def downloadArbitaryFileWEB(url , content_typeLength=14):
21:     #Connecting to Server and obtaining responses in Byte format
22:     responseInBytes = _clientConnectAndResponse(url)
23:
24:     headerBytes = responseInBytes[:responseInBytes.find(b'\r\n\r\n')]
25:     remainingBytes = responseInBytes[responseInBytes.find(b'\r\n\r\n'):]\
26:                                 .strip(b'\r\n\r\n')
27:
28:     #Here we check if HTTP responses with 200 OK
29:     if (headerBytes.find(b'\r\n')) > 0 :
30:         if headerBytes[:headerBytes.find(b'\r\n')].decode(ENCODING).lower()\
31:                     != "http/1.1 200 ok":
32:             print ("Cannot determine if HTTP 200 OK. Further Processing halts")
33:             return
34:
35:     #Display Header
36:     print('\n')
37:     print (headerBytes.decode(ENCODING))
38:
39:     #checking if the URL is for an image
40:     isImage = _checkIfImage(headerBytes , content_typeLength)
41:     headerFile = open('index.php.header' , 'w', newline='\n')
42:     headerFile.write(headerBytes.decode(ENCODING))
43:
44:     if isImage:
45:         imageToFile = open(urlparse(url).path.split('/')[-1], 'wb')
46:         imageToFile.write(remainingBytes)
47:         imageToFile.close()
48:     else:
49:         htmlToFile = open('index.php', 'wb')
```

```
50:            htmlToFile.write(remainingBytes)
51:            htmlToFile.close()
52:
53:
54: #Helper function that handles socket connection and also gives response in bytes
55: def _clientConnectAndResponse(url, timeout=10, receive_buffer=8096):
56:     parsed = urlparse(url)
57:     try:
58:         host, port = parsed.netloc.split(':')
59:     except ValueError:
60:         host, port = parsed.netloc, 80
61:
62:     clientSocket = socket.create_connection((host, port), timeout)
63:
64:     method  = 'GET %s HTTP/1.0\r\n\r\n' % parsed.path
65:     clientSocket.sendall(bytes(method, ENCODING))
66:
67:     response = [clientSocket.recv(receive_buffer)]
68:     while response[-1]:
69:         response.append(clientSocket.recv(receive_buffer))
70:
71:     responseInBytes = b''.join(r for r in response)
72:     return responseInBytes
73:
74: #Takes headerBytes and checks if response from server was an image file or not
75: def _checkIfImage (headerBytes, content_typeLength):
76:     if (headerBytes.find(b'Content-Type:') >= 0):
77:         headerBytesAfterContentType = headerBytes[headerBytes.find\
78:                                             (b'Content-Type:'):]
79:     else:
80:         print ("DOESNOT HAVE A CONTENT TYPE")
81:         raise ValueError('No Content Type in header File')
82:     if (headerBytesAfterContentType.find(b';') >= 0):
83:         fileFormat = headerBytesAfterContentType[:headerBytesAfterContentType\
84:                                             .find(b';')]
85:     else:
86:         fileFormat = headerBytesAfterContentType
87:     #print (fileFormat.decode(ENCODING))
88:     fileFormatStrip = fileFormat.decode(ENCODING)[content_typeLength:]
89:
90:     if fileFormatStrip.lower().find('image/') >= 0:
91:         return True
92:     return False
93:
94: if __name__ == "__main__":
95:      #Main Function call alongside command line Argument Handling
96:      totalArg = len(sys.argv)
97:      commandArguments = sys.argv # Get the arguments list
98:      mainFunction(commandArguments , totalArg)
```

## 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.

**Answer 2**

```
 1: import re
 2: import sys
 3: from urllib.parse import urlparse
 4: from urllib.parse import urljoin
 5: from yankee_http_client import downloadArbitaryFileWEB
 6:
 7: ENCODING = 'utf-8'
 8:
 9: def mainDownloadEverythingFunction(commandArguments , totalArg):
10:     if (totalArg > 1):
11:         parseHTML(commandArguments[1] , commandArguments[2])
12:     else:
13:         print ("Please give valid arguments")
14:         return
15:
16: def parseHTML(filename , url, content_typeLength=14):
17:     with open(filename , 'r') as html:
18:         content = html.read()
19:         pat = re.compile (r'<img [^>]*src="([^"]+)')
20:         # Here we also check\ from <link href and images inside them
21:         pat2 = re.compile (r'<link [^>]*href="([^"]+)')
22:
23:         moreimages = pat2.findall(content)
24:         moreimagesFiltered = list(filter(lambda x:\
25:                                     x.find('.ico') >= 0  or\
26:                                     x.find('.png') >= 0  or\
27:                                     x.find('.jpeg') >= 0  , moreimages))
28:
29:         img = pat.findall(content)
30:
31:     allImagesURL = moreimagesFiltered + img
32:     toAppendRelativeImageURL = _getPartToAppendToRelativeImageURL(url)
33:
34:     for imageurl in allImagesURL:
35:         if(imageurl.find('://') > 0):
36:             downloadArbitaryFileWEB(imageurl)
37:             #perform download for that url
38:         else:
39:             link = urljoin(toAppendRelativeImageURL , imageurl)
40:             downloadArbitaryFileWEB(link) #perform download for that url
41:
42:
43: # Helper function that returns the front part of the url that is to be appened
44: #to relative urls for valid image construction
45: def _getPartToAppendToRelativeImageURL(url):
46:     parsed = urlparse(url)
47:     try:
```

```
48:          host, port = parsed.netloc.split(':')
49:      except ValueError:
50:          host = parsed.netloc
51:      return parsed.scheme +'://'+ host
52:
53:
54: #Main Function call alongside command line Argument Handling
55: totalArg = len(sys.argv)
56: commandArguments = sys.argv # Get the arguments list
57: mainDownloadEverythingFunction(commandArguments , totalArg)
```
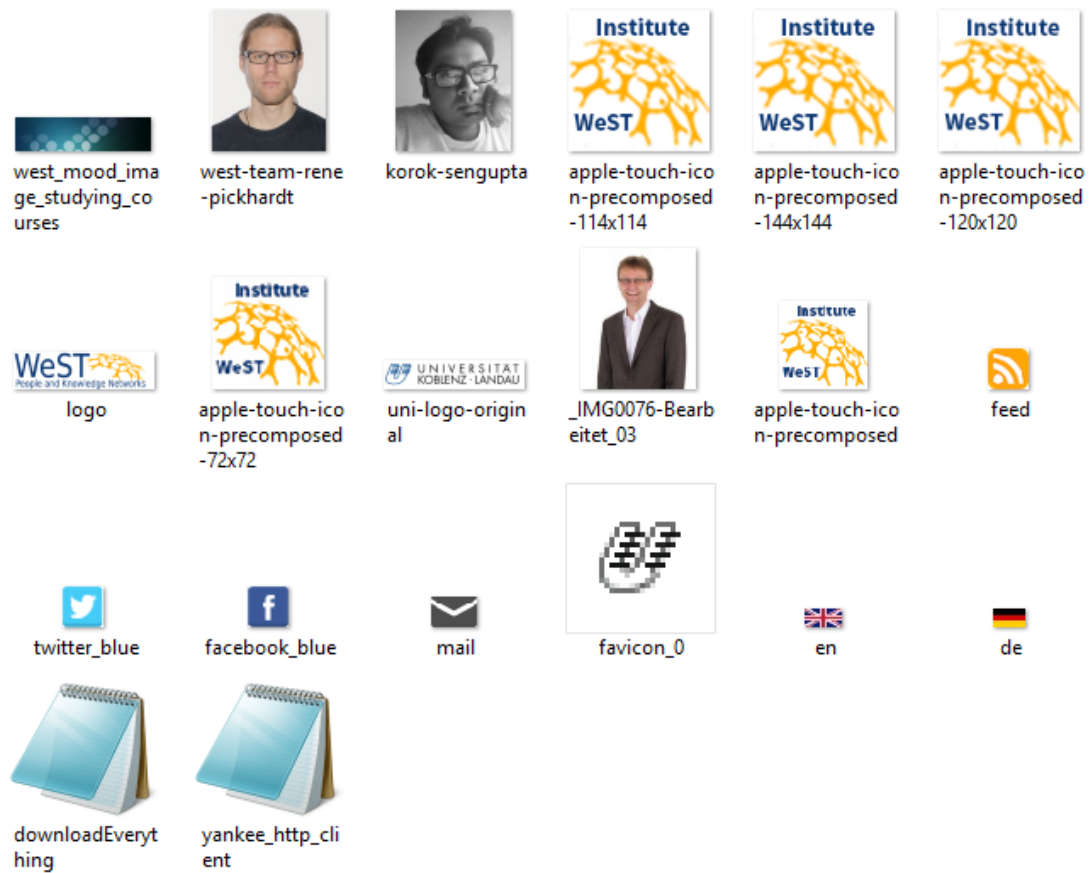


**Figure 2:** Screen-shot of folder containing all the downloaded images

## 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`.