# Introduction to Web Science

**Assignment 5**

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 30, 2016, 10:00 a.m.
Tutorial on:   December 2, 2016, 12:00 p.m.

Please look at the lessons 1) **Dynamic Web Content** & 2) **How big is the Web?**

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.

Team Name: Tango

1. Mariya Chkalova
   mchkalova@uni-koblenz.de

2. Arsenii Smyrnov
   smyrnov@uni-koblenz.de

3. Simon Schauß
   sschauss@uni-koblenz.de

4. Lukas Härtel
   lukashaertel@uni-koblenz.de

# 1 Creative use of the Hyptertext Transfer Protocol (10 Points)

HTTP is a request response protocol. In that spirit a client opens a TCP socket to the server, makes a request, and the server replies with a response. The server will just listen on its open socket but cannot initiate a conversation with the client on its own.

However you might have seen some interactive websites which notify you as soon as something happens on the server. An example would be Twitter. Without the need for you to refresh the page (and thus triggering a new HTTP request) they let you know that there are new tweets available for you. In this exercise we want you to make sense of that behaviour and try to reproduce it by creative use of the HTTP protocol.

Have a look at `server.py`[1] and `webclient.html` (which we provide). Extend both files in a way that after `webclient.html` is servered to the user the person controlling the server has the chance to make some input at its commandline. This input should then be send to the client and displayed automatically in the browser without requiring a reload. For that the user should not have to interact with the webpage any further.

## 1.1 webclient.html

```
 1: <html>
 2: <head>
 3:         <title>Abusing the HTTP protocol - Example</title>
 4: </head>
 5: <body>
 6:         <h1>Display data from the Server</h1>
 7:         The following line changes on the servers command line
 8:         input: <br>
 9:         <span id="response" style="color:red">
10:                 This will be replaced by messages from the server
11:         </span>
12: </body>
13: </html>
```

## 1.2 Hints:

- This exercise is more like a riddle. Try to focuse on how TCP sockets and HTTP work and how you could make use of that to achieve the expected behaviour. Once you have an idea the programming should be straight-forward.

- The Javascript code that you need for this exercise was almost completely shown in one of the videos and is available on Wikiversity.

---

[1]you could store the code from http://blog.wachowicz.eu/?p=256 in a file called `server.py`

- In that sense we only ask for a "proof of concept" nothing that would be stable out in the wilde.

  - In particular, don't worry about making the server uses multithreading. It is ok to be blocking for the sake of this exercise.

- Without use of any additional libraries or AJAX framework we have been able to solve this with 19 lines of Javascript and 11 lines of Python code (we provide this information just as a way for you to estimate the complexity of the problem, don't worry about how many lines your solution uses).

**Listing 1** Task 1

```python
1: #!/usr/bin/python
2:
3: import socket  # Networking support
4: import signal  # Signal support (server shutdown on signal receive
       ↪ )
5: import time  # Current time
6:
7:
8: class Server:
9:     """ Class describing a simple HTTP server objects."""
10:
11:     def __init__(self, port=80):
12:         """ Constructor """
13:         self.host = ''  # <-- works on all avaivable network
              ↪ interfaces
14:         self.port = port
15:         self.www_dir = 'www'  # Directory where webpage files are
              ↪ stored
16:
17:     def activate_server(self):
18:         """ Attempts to aquire the socket and launch the server
              ↪ """
19:         self.socket = socket.socket(socket.AF_INET, socket.
              ↪ SOCK_STREAM)
20:         try:  # user provided in the __init__() port may be
              ↪ unavaivable
21:             print("Launching HTTP server on ", self.host, ":",
                  ↪ self.port)
22:             self.socket.bind((self.host, self.port))
23:
24:         except Exception as e:
25:             print("Warning: Could not aquite port:", self.port, "\
                  ↪ n")
26:             print("I will try a higher port")
27:             # store to user provideed port locally for later (in
                  ↪ case 8080 fails)
28:             user_port = self.port
```

```
29:                    self.port = 8080
30:
31:                try:
32:                    print("Launching HTTP server on ", self.host, ":",
                          ↪  self.port)
33:                    self.socket.bind((self.host, self.port))
34:
35:                except Exception as e:
36:                    print("ERROR: Failed to acquire sockets for ports
                          ↪  ", user_port, " and 8080. ")
37:                    print("Try running the Server in a privileged user
                          ↪   mode.")
38:                    self.shutdown()
39:                    import sys
40:                    sys.exit(1)
41:
42:            print("Server successfully acquired the socket with port:"
                  ↪  , self.port)
43:            print("Press Ctrl+C to shut down the server and exit.")
44:            self._wait_for_connections()
45:
46:        def shutdown(self):
47:            """ Shut down the server """
48:            try:
49:                print("Shutting down the server")
50:                s.socket.shutdown(socket.SHUT_RDWR)
51:
52:            except Exception as e:
53:                print("Warning: could not shut down the socket. Maybe
                      ↪  it was already closed?", e)
54:
55:        def _gen_headers(self, code, poll=False):
56:            """ Generates HTTP response Headers. Ommits the first line
                  ↪  ! """
57:
58:            # determine response code
59:            h = ''
60:            if (code == 200):
61:                h = 'HTTP/1.1 200 OK\n'
62:            elif (code == 404):
63:                h = 'HTTP/1.1 404 Not Found\n'
64:
65:            # write further headers
66:            current_date = time.strftime("%a, %d %b %Y %H:%M:%S", time
                  ↪  .localtime())
67:            h += 'Date: ' + current_date + '\n'
68:            h += 'Server: Simple-Python-HTTP-Server\n'
69:            if poll:
70:                h += 'Content-Type: application/octet-stream\n'
```

```
71:                    h += 'Cache-Control: no-cache\n'
72:               h += 'Connection: close\n\n'  # signal that the conection
                      ↪ wil be closed after complting the request
73:
74:               return h
75:
76:        def _wait_for_connections(self):
77:               """ Main loop awaiting connections """
78:               while True:
79:                      print("Awaiting New connection")
80:                      self.socket.listen(3)  # maximum number of queued
                           ↪ connections
81:
82:                      conn, addr = self.socket.accept()
83:                      # conn - socket to client
84:                      # addr - clients address
85:
86:                      print("Got connection from:", addr)
87:
88:                      data = conn.recv(1024)  # receive data from client
89:                      string = bytes.decode(data)  # decode it to string
90:
91:                      # determine request method  (HEAD and GET are
                           ↪ supported)
92:                      request_method = string.split(' ')[0]
93:                      print("Method: ", request_method)
94:                      print("Request body: ", string)
95:
96:                      # if string[0:3] == 'GET':
97:                      if (request_method == 'GET') | (request_method == '
                           ↪ HEAD'):
98:                       # file_requested = string[4:]
99:
100:                      # split on space "GET /file.html" -into-> ('GET','
                            ↪ file.html',...)
101:                      file_requested = string.split(' ')
102:                      file_requested = file_requested[1]  # get 2nd
                            ↪ element
103:
104:                      # Check for URL arguments. Disregard them
105:                      file_requested = file_requested.split('?')[0]  #
                            ↪ disregard anything after '?'
106:
107:                      if (file_requested == '/'):  # in case no file is
                            ↪ specified by the browser
108:                          file_requested = '/webclient.html'  # load
                                ↪ index.html by default
109:                      if (file_requested == '/poll'):
110:                          header = self._gen_headers(200, poll=True).
```

5

```
                                    ↪ encode()
111:                    conn.sendall(header)
112:                    while True:
113:                        msg = input('msg: ').encode()
114:                        conn.sendall(msg)
115:
116:                file_requested = self.www_dir + file_requested
117:                print("Serving web page [", file_requested, "]")
118:
119:                ## Load file content
120:                try:
121:                    file_handler = open(file_requested, 'rb')
122:                    if (request_method == 'GET'):  # only read the
                            ↪  file when GET
123:                        response_content = file_handler.read()  #
                                ↪ read file content
124:                    file_handler.close()
125:
126:                    response_headers = self._gen_headers(200)
127:
128:                except Exception as e:  # in case file was not
                        ↪ found, generate 404 page
129:                    print("Warning, file not found. Serving
                            ↪ response code 404\n", e)
130:                    response_headers = self._gen_headers(404)
131:
132:                    if (request_method == 'GET'):
133:                        response_content = b"<html><body><p>Error
                                ↪ 404: File not found</p><p>Python
                                ↪ HTTP server</p></body></html>"
134:
135:                server_response = response_headers.encode()  #
                        ↪ return headers for GET and HEAD
136:                if (request_method == 'GET'):
137:                    server_response += response_content  # return
                            ↪ additional conten for GET only
138:
139:                conn.send(server_response)
140:                print("Closing connection with client")
141:                conn.close()
142:
143:            else:
144:                print("Unknown HTTP request method:",
                        ↪ request_method)
145:
146:
147: def graceful_shutdown(sig, dummy):
148:     """ This function shuts down the server. It's triggered
149:     by SIGINT signal """
```

```
150:        s.shutdown()   # shut down the server
151:        import sys
152:        sys.exit(1)
153:
154:
155:    ##############################################################
156:    # shut down on ctrl+c
157:    signal.signal(signal.SIGINT, graceful_shutdown)
158:
159:    print("Starting web server")
160:    s = Server(80)   # construct server object
161:    s.activate_server()   # aquire the socket
```

## 2 Web Crawler (10 Points)

Your task in this exercise is to "crawl" the `Simple English Wikipedia`. In order to execute this task, we provide you with a mirror of the Simple English Wikipedia at `141.26.208.82`.

You can start crawling from `http://141.26.208.82/articles/g/e/r/Germany.html` and you can use the `urllib` or `doGetRequest` function from the last week's assignment.

Given below is the strategy that you might adopt to complete this assignment:

1. Download `http://141.26.208.82/articles/g/e/r/Germany.html` and store the page on your file system.

2. Open the file in python and extract the local links. (Links within the same domain.)

3. Store the file to your file system.

4. Follow all the links and repeat steps 1 to 3.

5. Repeat step 4 until you have downloaded and saved all pages.

### 2.1 Hints:

- Before you start this exercise, please have a look at Exercise 3.

- Make really sure your crawler doesn't follow external urls to domains other than `http://141.26.208.82`. In that case you would start crawling the entire web

- Expect the crawler to run about 60 Minutes if you start it from the university network. From home your runtime will most certainly be even longer.

- It might be useful for you to have some output on the crawlers commandline depicting which URL is currently being fetched and how many URLs have been fetched so far and how many are currently on the queue.

- You can (but don't have to) make use of breadth-first search.

- It probably makes sense to take over the full paths from the pages of the Simple English Wikipedia and use the same folder structure when you save the html documents.

- You can (but you don't have to) speed up the crawler significantly if you use multithreading. However you should not use more than 10 threads in order for our mirror of Simple English Wikipedia to stay alive.

**Listing 2** Task 1

```
1: import curses
2: import json
3: import re
```

```
 4: import os
 5: import time
 6: from threading import Thread
 7: from urllib.parse import urlparse, urljoin, urldefrag
 8: from urllib.request import urlopen
 9:
10:
11: class Crawler:
12:     POOL_SIZE = 10
13:
14:     def __init__(self, start_url):
15:         self.scr = curses.initscr()
16:         curses.cbreak()
17:         curses.noecho()
18:         curses.curs_set(0)
19:         self.scr.keypad(1)
20:         self.start_url = urlparse(start_url)
21:         self.stack = set()
22:         self.url_to_urls = dict()
23:         self.error_urls = set()
24:         self.worker_id_to_working = dict()
25:         self.compiled_regex = re.compile('<a[^>]+href="([^">]+)"')
26:
27:     @staticmethod
28:     def get(url):
29:         with urlopen(url) as connection:
30:             return connection.read()
31:
32:     def extract_urls(self, payload):
33:         return self.compiled_regex.findall(payload)
34:
35:     def crawl(self, worker_id):
36:         while True:
37:             stack_length = len(self.stack)
38:             if len(self.stack) == 0 and True not in self.
                  ↪ worker_id_to_working.values():
39:                 break
40:             elif stack_length > 0:
41:                 self.worker_id_to_working[worker_id] = True
42:                 url = self.stack.pop()
43:                 try:
44:                     document = self.get(url)
45:                     relative_urls = self.extract_urls(document.
                        ↪ decode())
46:                     absolute_urls = list(map(lambda href: urljoin(
                        ↪ url, href), relative_urls))
47:                     for u in absolute_urls:
48:                         defraged_u = urldefrag(u).url
49:                         if defraged_u not in self.url_to_urls and
```

```
                                          ↪ urlparse(u).netloc == self.start_url
                                          ↪ .netloc:
50:                                 self.stack.add(defraged_u)
51:                     self.url_to_urls.update({url: absolute_urls})
52:                     self.write_document(url, document)
53:                 except Exception as e:
54:                     self.url_to_urls.update({url: str(e)})
55:                 finally:
56:                     self.worker_id_to_working[worker_id] = False
57:
58:     @staticmethod
59:     def write_document(url, document):
60:         path = './output%s' % urlparse(url).path
61:         os.makedirs(os.path.dirname(path), exist_ok=True)
62:         with open(path, 'wb') as handle:
63:             handle.write(document)
64:
65:     def write_statistics(self):
66:         with open('./statistics.json', 'w') as handle:
67:             json.dump(self.url_to_urls, handle)
68:
69:     def log(self):
70:         start = time.time()
71:         while True:
72:             if len(self.stack) == 0 and True not in self.
                    ↪ worker_id_to_working.values():
73:                 break
74:             c_urls_len = len(self.url_to_urls)
75:             c_stack_len = len(self.stack)
76:             end = time.time()
77:             start_end_diff = end - start
78:             stack_delta_per_second = c_stack_len / start_end_diff
79:             urls_per_second = c_urls_len / start_end_diff
80:             self.scr.addstr(0, 0, 'stack size: %i' % c_stack_len)
81:             self.scr.addstr(1, 0, 'mean stack delta second: %i' %
                    ↪ int(stack_delta_per_second))
82:             self.scr.addstr(2, 0, 'unique urls: %i' % c_urls_len)
83:             self.scr.addstr(3, 0, 'mean urls per second: %i' % int
                    ↪ (urls_per_second))
84:             self.scr.refresh()
85:             time.sleep(1 / 10)
86:
87:     def start(self):
88:         try:
89:             self.stack.add(self.start_url.geturl())
90:             workers = []
91:             for worker_id in range(0, self.POOL_SIZE):
92:                 thread = Thread(target=self.crawl, args=[worker_id
                        ↪ ])
```

```
 93:                    self.worker_id_to_working.update({worker_id: False
                    ↪ })
 94:                    workers.append(thread)
 95:                    thread.start()
 96:               logger = Thread(target=self.log)
 97:               logger.start()
 98:               for thread in workers:
 99:                    thread.join()
100:               logger.join()
101:               self.write_statistics()
102:          finally:
103:               curses.nocbreak()
104:               self.scr.keypad(0)
105:               curses.echo()
106:               curses.endwin()
107:
108:
109: if __name__ == '__main__':
110:     crawler = Crawler('http://141.26.208.82/articles/g/e/r/Germany
             ↪ .html')
111:     # crawler = Crawler(b'http://localhost:8000/index.html')
112:     crawler.start()
```

# 3 Web Crawl Statistics (10 Points)

If you have successfully completed the first exercise of this assignment, then please provide the following details. You may have to tweak your code in the above exercise for some of the results.

## 3.1 Phase I

1. Total Number of *webpages* you found.

2. Total number of links that you encountered in the complete process of crawling.

3. Average and median number of links per web page.

4. Create a *histogram* showing the distribution of links on the crawled web pages. You can use a bin size of 5 and scale the axis from 0-150.
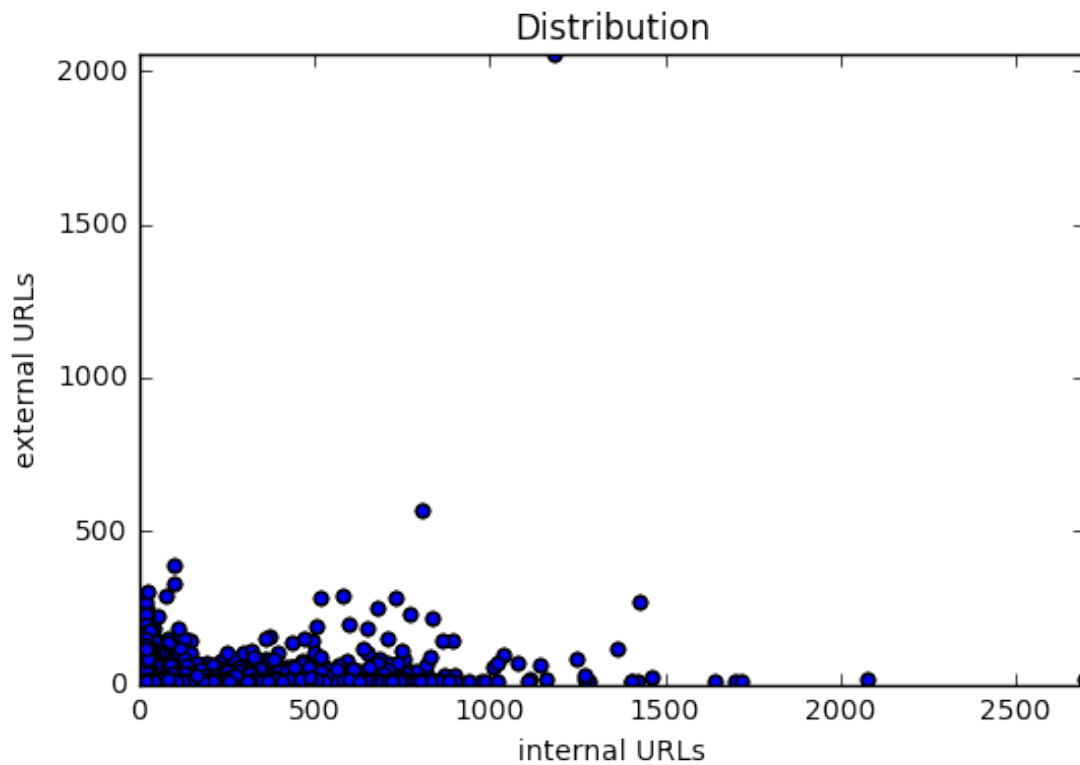
## 3.2 Phase II

1. For every page that you have downloaded, count the number of internal links and external links.

2. Provide a *scatter plot* with number of internal links on the X axis and number of external links on the Y axis.

total number of webpages: 83895

total number of links 6042511

average number of links 72.0246856189

mean number of links 57.0

**Listing 3** Task 1

```
 1: import json
 2: import pandas
 3: import numpy as np
 4: from matplotlib import pyplot as plot
 5: from urllib.parse import urlparse
 6:
 7: statistics_path = './statistics.json'
 8: url = urlparse('http://141.26.208.82/articles/g/e/r/Germany.html')
 9:
10: def load_statistics():
11:     with open(statistics_path, 'r') as f:
12:         return json.load(f)
13:
14: def filter_accessible_links(stats):
15:     return {k: [urlparse(i) for i in v] for k, v in stats.items()
          ↪ if isinstance(v, list)}
16:
17: def filter_dead_links(stats):
18:     return {k: v for k, v in stats.items() if isinstance(v, str)}
19:
20: statistics = load_statistics()
21:
```

```
22: accessible_links = filter_accessible_links(statistics)
23: dead_links = filter_dead_links(statistics)
24:
25: print('total number of webpages:', len(accessible_links))
26:
27: links_per_link = [(len(v)) for k, v in accessible_links.items()]
28:
29: total_links = sum(links_per_link)
30: print('total number of links', total_links)
31:
32: average_per_webpage = np.mean(links_per_link)
33: print('average number of links', average_per_webpage)
34:
35: median_per_webpage = np.median(links_per_link)
36: print('average number of links', median_per_webpage)
37:
38: internal_to_external = [(len([i for i in v if i.netloc == url.
        ↪ netloc]),
39:                         len([i for i in v if i.netloc != url.
                                ↪ netloc])) for k, v in
                                ↪ accessible_links.items()]
40: internal_urls = [i for i, e in internal_to_external]
41: external_urls = [e for i, e in internal_to_external]
42:
43: print(sum(internal_urls))
44: print(sum(external_urls))
45:
46: plot.scatter(internal_urls, external_urls)
47: plot.xlabel("internal URLs")
48: plot.ylabel("external URLs")
49: plot.title("Distribution")
50: plot.xlim((0, max(internal_urls)))
51: plot.ylim((0, max(external_urls)))
52: plot.show()
```

## Important Notes

### Submission

- Solutions have to be checked into the github repository. Use the directory name `groupname/assignment5/` 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`.