# Introduction to Web Science

**Assignment 3**

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 Luxembourg

Submission until:  November 16, 2016, 10:00 a.m.
Tutorial on:  November 18, 2016, 12:00 p.m.

The main objective of this assignment is for you understand different concepts that are associated with the "Web". In this assignment we cover two topics: 1) DNS & 2) Internet.

These tasks are not always specific to "Introduction to Web Science". For all the assignment questions that require you to write a 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.
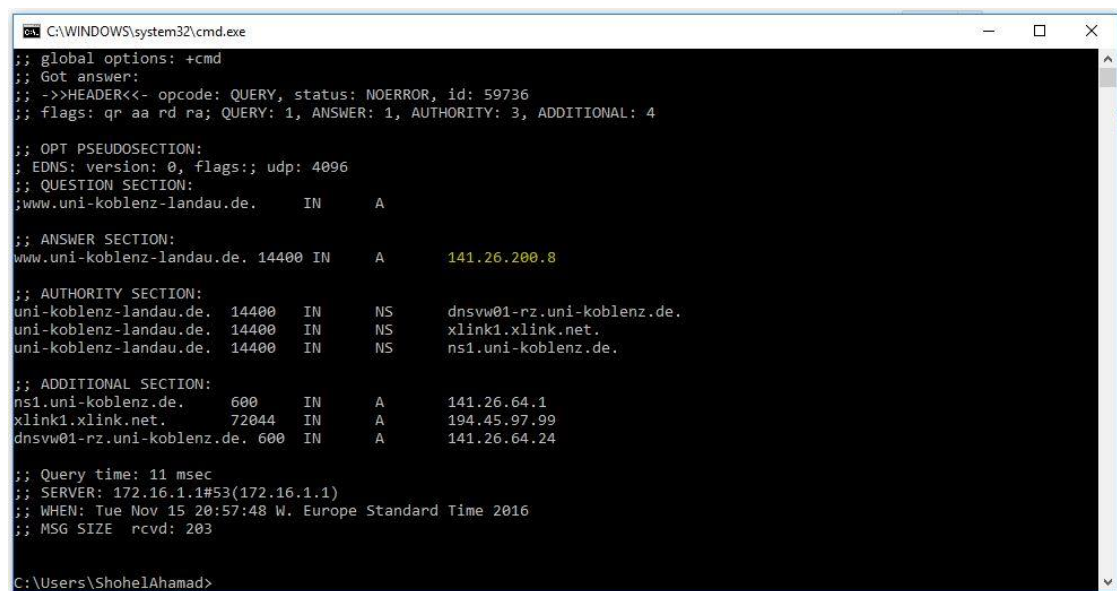
Team Name: mike

# 1 DIG Deeper (5 Points)

Assignment 1 started with you googling certain basic tools and one of them was "*dig*".

1. Now using that dig command, find the IP address of `www.uni-koblenz-landau.de`

2. In the result, you will find "SOA". What is SOA?

3. Copy the SOA record that you find in your answer sheet and explain each of the components of SOA with regards to your find. Merely integrating answers from the internet wont fetch you points.

Try the experiment once from University network and once from Home network and see if you can find any differences and if so, clarify why.

**Answers:**

1. The IP address of `www.uni-koblenz-landau.de` is 141.26.200.8

```
C:\WINDOWS\system32\cmd.exe                                        —    □    ×

;; global options: +cmd
;; Got answer:
;; ->>HEADER<<- opcode: QUERY, status: NOERROR, id: 59736
;; flags: qr aa rd ra; QUERY: 1, ANSWER: 1, AUTHORITY: 3, ADDITIONAL: 4

;; OPT PSEUDOSECTION:
; EDNS: version: 0, flags:; udp: 4096
;; QUESTION SECTION:
;www.uni-koblenz-landau.de.      IN      A

;; ANSWER SECTION:
www.uni-koblenz-landau.de. 14400 IN     A       141.26.200.8

;; AUTHORITY SECTION:
uni-koblenz-landau.de.  14400   IN      NS      dnsvw01-rz.uni-koblenz.de.
uni-koblenz-landau.de.  14400   IN      NS      xlink1.xlink.net.
uni-koblenz-landau.de.  14400   IN      NS      ns1.uni-koblenz.de.

;; ADDITIONAL SECTION:
ns1.uni-koblenz.de.     600     IN      A       141.26.64.1
xlink1.xlink.net.       72044   IN      A       194.45.97.99
dnsvw01-rz.uni-koblenz.de. 600  IN      A       141.26.64.24

;; Query time: 11 msec
;; SERVER: 172.16.1.1#53(172.16.1.1)
;; WHEN: Tue Nov 15 20:57:48 W. Europe Standard Time 2016
;; MSG SIZE  rcvd: 203

C:\Users\ShohelAhamad>
```

2. The SOA record defines the global parameters for the zone (domain) in Domain Name System (DNS) and represents the start of a zone and thus administer zone files . It defines which server is primary nameserver, contact information (E-mail), how secondary nameservers get updated, and the default (minimum) Time-To-Live values.

3. SOA record takens from both, Home network and Uni network show the same results for `www.uni-koblenz-landau.de`

   Uni network

Home network



The SOA record is
uni-koblenz-landau.de. 10371 IN SOA dnsvw01.uni-koblenz-landau.de. root.dnsvw01.uni-koblenz-landau.de. 2016110401 14400 900 604800 14400

a) **uni-koblenz-landau.de.** - is the main name in this zone.

b) **10371** - TTL defines the duration in seconds that the record may be cached by client side programs.In this case this record will be stored almost 3 hours.

c) **IN** - The class IN shows the type of record - IN equates to Internet

d) **SOA** - Shortcuts for Start Of Authority

e) **dnsvw01.uni-koblenz-landau.de.** - Nameserver -Is primary server which holds the zone files. It is an external server in this case, that is why the entire domain name is specified followed by a dot.

f) **root.dnsvw01.uni-koblenz-landau.de.** - This is a DOMAIN NAME that indicates the E-mail address of the person responsible for this zone, and to which email may be sent to report errors or problems. In this case, email is sent to root@dnsvw01.uni-koblenz-landau.de, but written as root.dnsvw01.uni-koblenz-landau.de.

g) **2016110401** - This is a serial number (32-bit integer) that must be incremented on the primary name server whenever a change is made. This number has to increment, whenever any change is made to the Zone file. The standard convention is to use the date of update YYYYMMDDnn, where nn is a revision number in case more than one updates are done in a day. In this case, the last update was made on 04.11.2016. And this was the only update made that day.

h) **14400** - Refresh, This time(in seconds) represents how often a secondary secondary name server gets a copy of the primary zone or sees if the serial number for the zone has increased (so it knows to request a new copy of the data for the zone). In this case, such a request is sent every 4 hours

i) **900** - Retry, 32 bit value in seconds. Defines the time between retries if the secondary name server fails to contact the primary when refresh has expired. In this case the secondary server will wait 15 minutes and then try to re-contact the master server after failing the first time it because it was down. The Retry value (time in seconds) will tell it when to get back.

j) **604800** - Expiry, The number of seconds (32-bit integer) that lets the secondary name server(s) know how long they can hold the information before it is no longer considered valid (authoritative). In this case, the secondary name server will keep the record for 7 days.

k) **14400** - MINIMUM ("Minimum TTL") - The number of seconds that the records in the zone are valid for (time-to-live, or TTL).

# 2 Exploring DNS (10 Points)

In the first part of this assignment you were asked to develop a simple TCP Client Server. Now, using **that** client server setup. This time a url should be send to the server and the server will split the url into the following:

http://www.example.com:80/path/to/myfile.html?key1=value1&key2=value2#InTheDocument

1. Protocol

2. Domain

3. Sub-Domain

4. Port number

5. Path

6. Parameters

7. Fragment

The Protocol for sending the URL will be a string terminated with \r \n.

P.S.: You are **not** allowed to use libraries like `urlparse` for this question. You will also not use "Regular Expressions" for this.

**Answer:**
Client.py

```
 1: import socket
 2:
 3: #creating the socket
 4: PORT = 8080
 5: s = socket.socket()
 6: s.connect(('localhost', PORT))
 7:
 8: #User enters URL until URL ends with  \r\n
 9: url=""
10: while True:
11:     print("Pleas Enter a URL:")
12:     url = input('->')
13:     if (url[len(url) - 5:len(url)]=="\\r \\n"):
14:         break
15:
16: print("The following URL was sent to server: ")
17: print(url)
18: s.send(url.encode('utf-8')) #sending url variable over socket
```

```
Run:   server   client
   ▶  ↑   C:\Users\Slobodan\AppData\Local\Programs\Python\Python36\pythonw.exe C:/Users/Slobodan/client.py
          Pleas Enter a URL:
   ■  ↓   ->http://www.example.com:80/path/to/myfile.html?key1=value1&key2=value2#InTheDocument\r \n
   II  ⇥  The following URL was sent to server:
          http://www.example.com:80/path/to/myfile.html?key1=value1&key2=value2#InTheDocument\r \n

   ⚹  🖶   Process finished with exit code 0
   ✖  🗑
   »
```

Server.py

```
 1: import socket
 2:
 3: #creating socket
 4: PORT = 8080
 5: s = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
 6: s.bind(('localhost', PORT))     #connecting the port with the localhost machine
 7: s.listen(1) # server listens to one client
 8:
 9: conn, addr = s.accept()   # accepting the connection
10: tmp = conn.recv(1024)
11:
12: url = tmp.decode('utf-8')
13:
14: i=0
15: protocol = ""
16: subdomain=""
17: domain=""
18: port=""
19: full_path=""
20: param=""
21: fragment=""
22:
23:
24: if (len(url)>0):
25: # splitting the protocol part of the url
26:     while(url[i:i+3]!= "://"):
27:         protocol = protocol + url[i]
28:         i = i + 1
29:     print("Protocol: \t" + protocol)
30:     i = i + 3
31:
32: #splitting the domain part
33:     array=[]
34:     while (url[i] != "/" and  url[i] != ":"):
35:         domain=domain + url[i]
36:         i=i+1
37:     print("Domain: \t" +domain)
38:
39: #splitting the subdoman part
```

```
40:        sub = domain.split(".")
41:        j=0
42:        while (j<len(sub) -2 ):
43:            subdomain = subdomain + sub[j] + ", "
44:            j=j+1
45:        print("Sub-domain:\t" + subdomain)
46:
47: #splitting the port (if there is any in the url)
48:        if(url[i]==":"):
49:            while(url[i]!= "/"):
50:                port=port+url[i]
51:                i=i+1
52:        else:
53:            port="No port was specified!"
54:            i=i+1
55:        print("Port: \t\t" + port)
56:
57: #splitting the full_path part
58:        while (url[i] != "?" and url[i:i+5] != "\\r \\n"):
59:            full_path = full_path + url[i]
60:            i = i + 1
61:        print("Full path: \t" + full_path)
62:
63: #splitting the parameters
64:        i=i+1
65:        while (url[i] != "#" and url[i-1:i+4] != "\\r \\n"):
66:            if(url[i]=="&"):
67:                param = param + '\n \t\t\t'
68:            else:
69:                param = param + url[i]
70:            i = i + 1
71:        print("Parameters:\t" + param)
72:
73: #splitting the fragments
74:        while (url[i-1:i+4] != "\\r \\n"):
75:            fragment = fragment + url[i]
76:            i = i + 1
77:        print("Fragments: \t" + fragment[1:len(fragment)-1])
78:
79: else:
80:        print("No URL was received!")
```

```
Run:    server    client
  ►  ↑   C:\Users\Slobodan\AppData\Local\Programs\Python\Python36\pythonw.exe C:/Users/Slobodan/server.py
  ■  ↓   Protocol:   http
  ||  ⇄   Domain:     www.example.com
  🖵  ▦   Sub-domain: www,
  ⚲  🖶   Port:       :80
  ✖  🗑   Full path:  /path/to/myfile.html
  ?       Parameters: key1=value1
                      key2=value2
          Fragments:  InTheDocument

          Process finished with exit code 0
```

# 3 DNS Recursive Query Resolving (5 Points)

You have solved the "Routing Table" question in Assignment 2. We updated the routing tables once more. resulting in the following tables creating the following topology

**Table 1:** Routing Table

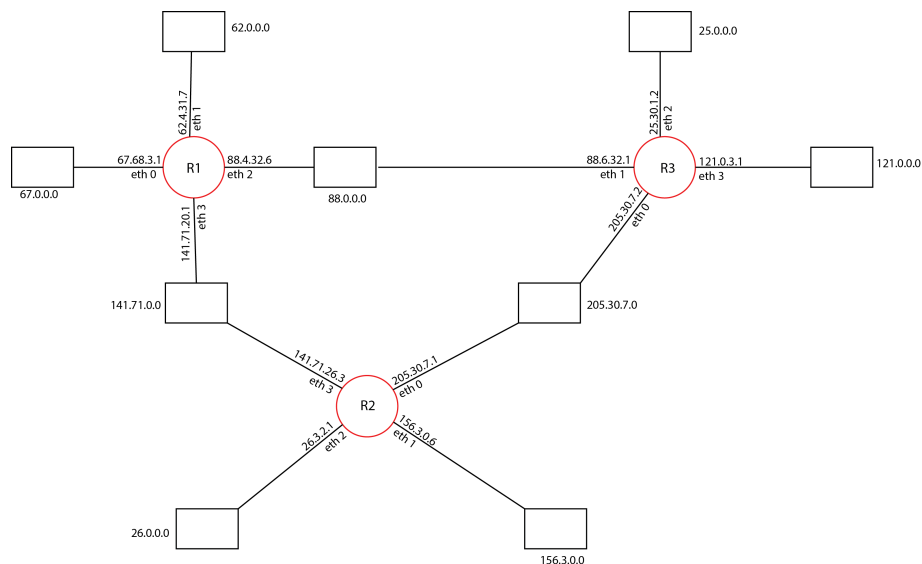| Router1 | | | | Router2 | | | | Router3 | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Destination | Next Hop | Interface | | Destination | Next Hop | Interface | | Destination | Next Hop | Interface |
| 67.0.0.0 | 67.68.3.1 | eth 0 | | 205.30.7.0 | 205.30.7.1 | eth 0 | | 205.30.7.0 | 205.30.7.2 | eth 0 |
| 62.0.0.0 | 62.4.31.7 | eth 1 | | 156.3.0.0 | 156.3.0.6 | eth 1 | | 88.0.0.0 | 88.6.32.1 | eth 1 |
| 88.0.0.0 | 88.4.32.6 | eth 2 | | 26.0.0.0 | 26.3.2.1 | eth 2 | | 25.0.0.0 | 25.03.1.2 | eth 2 |
| 141.71.0.0 | 141.71.20.1 | eth 3 | | 141.71.0.0 | 141.71.26.3 | eth 3 | | 121.0.0.0 | 121.0.3.1 | eth 3 |
| 26.0.0.0 | 141.71.26.3 | eth3 | | 67.0.0.0 | 141.71.20.1 | eth 3 | | 156.3.0.0 | 205.30.7.1 | eth 0 |
| 156.3.0.0 | 88.6.32.1 | eth 2 | | 62.0.0.0 | 141.71.20.1 | eth 3 | | 26.0.0.0 | 205.30.7.1 | eth 0 |
| 205.30.7.0 | 141.71.26.3 | eth 3 | | 88.0.0.0 | 141.71.20.1 | eth 3 | | 141.71.0.0 | 205.30.7.1 | eth 0 |
| 25.0.0.0 | 88.6.32.1 | eth 2 | | 25.0.0.0 | 205.30.7.2 | eth 0 | | 67.0.0.0 | 88.4.32.6 | eth 1 |
| 121.0.0.0 | 88.6.32.1 | eth 2 | | 121.0.0.0 | 205.30.7.2 | eth 0 | | 62.0.0.0 | 88.4.32.6 | eth 1 |



**Figure 1:** DNS Routing Network

Let us asume a client with the following ip address 67.4.5.2 wants to resolve the following domain `subdomain.webscienceexampledomain.com` using the DNS.

You can further assume the root name server has the IP address of 25.8.2.1 and the name-server for `webscienceexampledomain.com` has the IP address 156.3.20.2. Finally the sub-domain is handled by a name server with the IP of 26.155.36.7.

Please explain how the traffic flows through the network in order to resolve the recursive DNS query. You can assume ARP tables are cached so that no ARP-requests have to be made.

**Hint: You can start like this**:

67.4.5.2 creates an IP packet with the source address XXXXXX an destination address YYYYY inside there is the DNS request. This IP packet is send as an ethernet frame to ZZZZZ. ZZZZZ receives the frame and forwards the encapsulated IP packet to ....

Also you can assume the DNS requests and responses will fit inside one IP packet. You also don't have to write down the specific DNS requests and responses in hex.

Answer:

If we assume that Router1 acts as a ISP DNS server, then he is responsible for solving recursive quires. When the client (67.4.5.2) wants to resolve the IP address of subdomain.webscienceexampledomain.com using DNS, then he needs to send the request for finding the IP address of this domain to ISP server (Router1).

Router1 then forwards the request to the top level root server (25.8.2.1) to find the IP address of subdomain.webscienceexampledomain.com. Since this name server (25.8.2.1) is not responsible for a query that is not in his zone, that is why this name server (25.8.2.1) sends a respond, stating the IP address of the DNS server (156.3.20.2) that might be responsible for resolving this dns request.

After Router1 (ISP) receives this respond, he sends the query to (156.3.20.2) webscienceexampledomain.com. asking for the IP address of the domain subdomain.webscienceexampledomain.com. Because this DNS server is also not responsible for resolving the stated dns query, he responds in the same way as the top level root name server. He sends the IP address of a dns server that might be responsible for resolving the domain.

The ISP name server receives the respond from 156.3.20.2 and re- assemble new dns query and send it to 26.155.36.7. Now, this dns server, knows the answer, and sends it back to ISP in which case is then forwarder to the client (67.4.5.2).

Afterwards, ISP server keeps track of this record for the next time when the client is going to issue a resolving the same domain.

# Important Notes

## Submission

- Solutions have to be checked into the github repository. Use the directory name `groupname/assignment3/` 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`.