

Introduction

Having just learned about the Domain Name System (DNS), COMP0023 student Ben Bitdiddle's interest is piqued. In his spare time, Ben runs a small ISP, and of course, he wants to offer his customers a reliable DNS service. Suspicious of recent cache poisoning security exploits in the DNS server implementation he currently runs, he decides to write his own server in Python. With some of the profits from his ISP business, Ben hires you and a fellow student, Alyssa P. Hacker, to help.

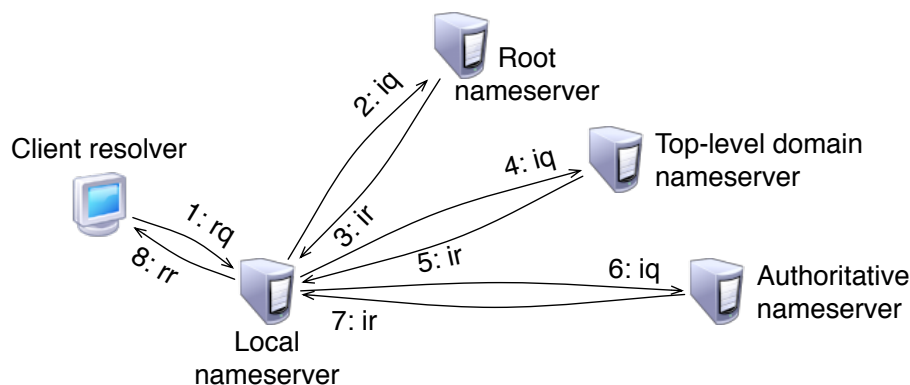


FIGURE 1: A recursive DNS lookup is made up of many iterative lookups, each initiated by the local nameserver.

Recall the overall structure of the DNS from lecture: the local name server is placed close to the pool of clients that it serves, and normally only accepts requests from that pool of clients. When a program running on the client needs to resolve a hostname, a piece of code running on the client called the *resolver* contacts the local nameserver with a *recursive query* for that name (labelled *rq* in Figure 1). This step and the following steps are numbered in consecutive order in the figure. Recall that a recursive query is named as such because it results in the local nameserver making a number of other queries (called *iterative*) on behalf of the recursive query. The local nameserver next makes a series of iterative queries (labelled *iq*) to various other nameservers, which may include a *root* nameserver (authoritative for *.*), a *top-level domain* (TLD) nameserver (authoritative for top-level domains, such as *.edu.*, *.com.*, etc.), and other authoritative nameservers for the domains and subdomains listed in the original recursive query. This sequence of iterative requests and responses, labelled 2–7 in the figure, may be reduced by information present in the local nameserver's cache. Finally, once the local name server has resolved recursive query 1, it replies to the client's resolver with the *recursive response* labelled *rr* in the figure. The resolver then returns the response to the application running on the client.

This is the heart of the DNS nameserver functionality. You can think of the local nameserver as the “workhorse” of the DNS, because it is the one required to implement the recursive lookup algorithm. In fact, most TLD and root nameservers will not honor requests for recursive lookups—they are simply too busy.

This assignment is in two parts. In Part 1, we will understand how recursive queries work in practice by using command line tools and our knowledge about the DNS. In Part 2, we will implement

the recursive query functionality in Ben's DNS server, testing against answers provided by the department's local DNS server. Note that the second part is substantially more challenging than the first—be sure to leave yourself enough time to finish the entire coursework!

Part 1: Manual Recursive Queries

Let's begin by getting a deeper understanding of how DNS queries work in practice. We'll be using `dig`, a command-line UNIX tool that instructs the client's resolver to issue DNS queries to one of the local nameservers listed in `/etc/resolv.conf`, usually `128.16.6.8` (`haig.cs.ucl.ac.uk`) on department systems. `dig` then parses the resulting DNS reply, printing it clearly on the console. For example,

```
$ dig @haig.cs.ucl.ac.uk www.xorp.org
```

Results in output similar to the following:

```
;; Got answer:
;; ->>HEADER<<- opcode: QUERY, status: NOERROR, id: 5220
;; flags: qr rd ra; QUERY: 1, ANSWER: 1, AUTHORITY: 2, ADDITIONAL: 2

;; QUESTION SECTION:
www.xorp.org.                IN      A

;; ANSWER SECTION:
www.xorp.org.    3544      IN      A      208.74.158.171

;; AUTHORITY SECTION:
xorp.org.        3544      IN      NS     ns2.xorp.org.
xorp.org.        3544      IN      NS     ns.xorp.org.

;; ADDITIONAL SECTION:
ns2.xorp.org.    86344     IN      A      193.63.58.145
ns.xorp.org.     86344     IN      A      128.16.70.254

;; Query time: 20 msec
;; SERVER: 128.16.6.8#53(128.16.6.8)
```

From the metadata at the end of the output, we see that `dig` has received an answer from `haig` (`128.16.6.8`). First, the header section tells us that the query was successful (`status: NOERROR`), listing the query's unique identifier, or UID: (`id: 5220`). Second, the header's one-bit flags tell us that this is a response (flag `qr` is set), that `dig` requested a recursive lookup (flag `rd`, for recursion desired, is set), and that `haig` honoured the request for a *recursive* lookup (flag `ra`, for recursion available, is set). Finally, the header tells us to expect one question entry, one answer entry, two authority entries, and two additional entries in the body of the reply. The question entry (`www.xorp.org. IN A`) echoes the question entry that `dig` sent `haig`, querying the Internet (class `IN`) address (`A`) of `www.xorp.org`. The answer entry tells us that `www.xorp.org`'s IP address is `208.74.158.171`. The authority section lists two nameservers authoritative for `xorp.org`, `ns2.xorp.org` and `ns.xorp.org`. The additional section contains two "glue" records specifying the IP addresses of the two authoritative nameservers listed in the authority section.

A key point to note here is that since `dig` requested recursion, `haig` has done the hard part for us. Let's now see how the process works from `haig`'s perspective by making the sequence of iterative queries that `haig` would, using the `+norecurse` flag.

A recursive nameserver needs to be manually configured with the IP addresses of the 13 root nameservers. Some nameservers will tell you the root nameservers they've got configured if you ask nicely.

```
$ dig @haig.cs.ucl.ac.uk . NS
```

haig responds (in part):

```
;; ANSWER SECTION:
```

```
.                459898    IN        NS        l.root-servers.net.
.                459898    IN        NS        g.root-servers.net.
.                459898    IN        NS        c.root-servers.net.
```

```
;; ADDITIONAL SECTION:
```

```
l.root-servers.net.  546298    IN        A        199.7.83.42
g.root-servers.net.  546298    IN        A        192.112.36.4
c.root-servers.net.  546298    IN        A        192.33.4.12
```

So that's the starting point for performing the iterative part of a recursive lookup - we know the IP address of a root nameserver. Lets use `l.root-servers.net` at IP address `199.7.83.42`. Note that `haig` may also tell you an IPv6 address for `l.root-servers.net`, listed in an AAAA record. For simplicity, we'll stick to IPv4 addresses obtained from A records for this coursework.

Supposing that our address caches are empty, let's now use `dig` to make the same series of (non-recursive) requests that `haig` would make given a query for `www.xorp.org` with the `rd` bit set. We start our hypothetical recursive query at `l.root-servers.net` (`199.7.83.42`):

```
$ dig @199.7.83.42 www.xorp.org. +norecurse
```

`l.root-servers.net` responds (in part):

```
;; ->>HEADER<<- opcode: QUERY, status: NOERROR, id: 46135
;; flags: qr; QUERY: 1, ANSWER: 0, AUTHORITY: 6, ADDITIONAL: 12
```

```
;; QUESTION SECTION:
```

```
;www.xorp.org.                IN        A
```

```
;; AUTHORITY SECTION:
```

```
org.                172800    IN        NS        a0.org.afiliias-nst.info.
org.                172800    IN        NS        a2.org.afiliias-nst.info.
org.                172800    IN        NS        b0.org.afiliias-nst.org.
org.                172800    IN        NS        b2.org.afiliias-nst.org.
org.                172800    IN        NS        c0.org.afiliias-nst.info.
org.                172800    IN        NS        d0.org.afiliias-nst.org.
```

```
;; ADDITIONAL SECTION:
```

```
a0.org.afiliias-nst.info. 172800 IN  A    199.19.56.1
a2.org.afiliias-nst.info. 172800 IN  A    199.249.112.1
b0.org.afiliias-nst.org.  172800 IN  A    199.19.54.1
b2.org.afiliias-nst.org.  172800 IN  A    199.249.120.1
c0.org.afiliias-nst.info. 172800 IN  A    199.19.53.1
d0.org.afiliias-nst.org.  172800 IN  A    199.19.57.1
```

The first thing to notice in the above is that we haven't received a direct answer to our query. `l.root-servers.net` has instead observed that our query is for a host in the `org` top-level domain, and referred us to one of several top-level domain nameservers that handle the `org` part of the DNS hierarchy. It's even been kind enough to provide us with the aforementioned glue records that tell us where to continue our query. The next step would be to ask the same question of one of these `org` servers, which should be able to provide information about the `xorp.org` domain. In the following "warmup exercises," you will continue this recursive query in a similar fashion.

Warmup exercises

1. Using `dig` with the `+norecurse` flag always set, targeted at an IP address (e.g. `dig @199.19.56.1` rather than `dig @a0.org.afiliias-nst.info.`), continue the recursive query for `www.xorp.org` on your machine.

Turn in the full list of command lines that you issue as well as the resulting `dig` output. In this and the following exercises, include only the relevant lines in the DIG output in your answer.

[2 marks]

2. Again, using `dig` with the `+norecurse` flag always set, emulate recursive queries in the same manner for the following hostnames. Be sure to start at a root DNS server for each. Note that if you hit a nameserver along the way for which there are no glue records, you may have to go back to the root to lookup the IP address of that nameserver.

`newgate.cs.ucl.ac.uk.`

`www.microsoft.com.`

[2 marks each]

3. What about the name resolution process for `newgate.cs.ucl.ac.uk.` differed from the name resolution process for `www.xorp.org.`? [2 marks]

4. What about the name resolution process for `www.microsoft.com.` differed from the name resolution process for `www.xorp.org.`? [2 marks]

[Total marks for Part 1: 10]

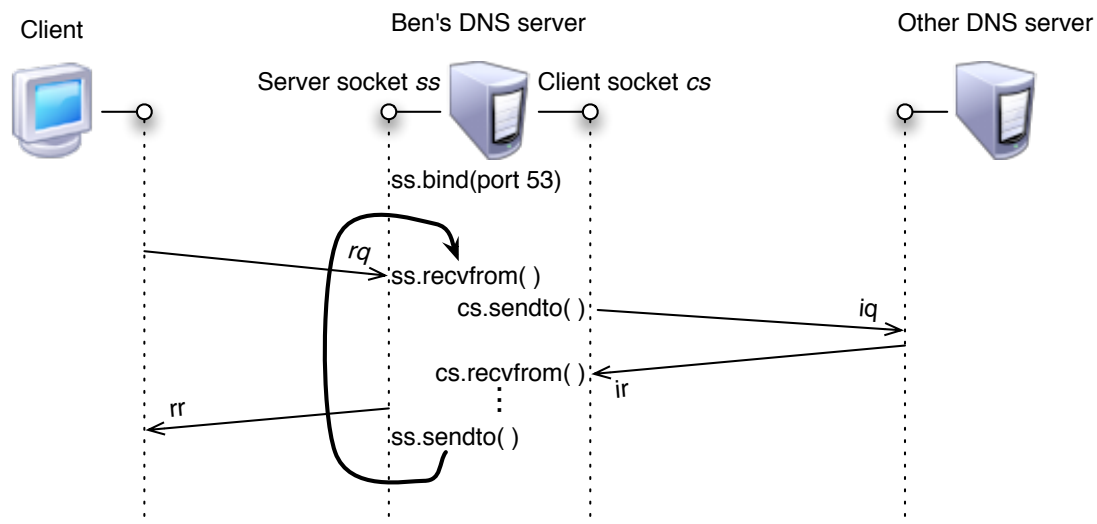


FIGURE 2: Communication flow from client to local nameserver and local nameserver to other DNS servers during one client's request to its local nameserver.

Part 2: Building a Simple Local Nameserver

Now that you are familiar with the process of making recursive queries by hand, you and Alyssa are ready to help Ben build his local nameserver. Ben decides that the nameserver will be a simple “single-threaded” server model that serves one client at a time, completing service of the current client before starting service of the next client. The nameserver we will build will be a “recursive” nameserver, capable of answering the same recursive queries you resolved by hand in Part 1, and will be *caching*, retaining the results of prior queries until they are invalid, using cached data instead of iterative queries whenever possible.

Recall from lecture that normal DNS communication uses UDP, the unreliable datagram delivery service of the Internet, and that applications send and receive UDP datagrams using a software interface called a *socket*. As shown in Figure 2, Ben's server will utilize two sockets: *ss*, a server socket bound to UDP port 53 (the well-known DNS port), and *cs*, a client socket that the server uses to query other servers as a client of those servers. These sockets are shown in Figure 2, along with a time diagram of how communication proceeds in this single-threaded server. First, Ben's local nameserver uses the `bind` call to bind *ss* to UDP port 53. Then, serially, for each incoming client query, the server calls `recvfrom` on the server socket to wait for and receive the recursive query *rq* as shown in the figure. Next, the server makes zero or more iterative queries *iq* using `sendto` and `recvfrom` on the client socket, while processing the results. Finally, the server calls `sendto` on the server socket to send its response to the client.

The preceding steps execute in a continuous loop, as indicated in the figure, with `ss.recvfrom` blocking until the next incoming request.

Getting started

All programming for this coursework should run under Linux on the department's lab machines. We have ensured that the code we give you to build upon works correctly on these lab machines. Note that these machines are accessible over the Internet, so you may work on the coursework either from home (using the CS department's Sun Secure Global Desktop server, accessible on the web at <http://www.cs.ucl.ac.uk/csrw>), or by logging into the machine in person, in lab. You should also be able to ssh into `knuckles.cs.ucl.ac.uk`.

To get started, once you've logged into a lab machine, unpack the distribution tarball

`http://nrg.cs.ucl.ac.uk/ste/comp0023/courseworks/dns-cw.tgz` into your home directory:

```
$ tar xvzf dns-cw.tgz
```

This will create a directory called `dns-cw` which we will refer to as the lab distribution directory.

Source code walkthrough

The main nameserver is located in file `ncsdns.py`, in which Alyssa has coded up the single-threaded structure described above. Your task is to implement the recursive lookup algorithm described above, as well as caching functionality so that your server doesn't make any unnecessary iterative queries.

The first problem you encounter is the need to parse and construct the various DNS packet payload formats. In short order, Alyssa codes up the following functionality to perform much of the necessary work. In package `libs.inetlib`, Alyssa provides you with two helper classes. `DomainName` is a class that represents DNS domain names and contains logic to parse and construct packet payload representations. `InetAddr` provides the same functionality for an IPv4 address.

Parsing and constructing DNS packets

The Python `recvfrom()` library function yields packet payloads as **strings of binary-valued data**, while the python `sendto()` library function expects packet payloads in the same format. We now describe how you may handle this requirement.

Every DNS packet payload is made up of a header, followed by zero or more *question entries*, followed by zero or more *resource records*. For more information about DNS packet formats, see the class lecture notes and the references at the end of this assignment.

Domain names in responses from production DNS servers are often compressed according to RFC1035. Rather than requiring you to implement the string compression and decompression algorithms, package `libs.inetlib` contains a `DomainName` class that implements the decompression algorithm for you. The tricky part, however, is the semantics of how long a domain name is. Suppose your server receives a packet with a domain name inside whose compressed binary representation is 10 bytes long, but whose string length is 13 bytes long. When you construct a `DomainName` object `d` from that compressed binary representation, should `len(d)` evaluate to 10 or 13? The approach we have chosen is to have `len(d)` evaluate to its compressed length, in this case 10, but any copies of `d` that are made by the Python call `copy(d)` get uncompressed and hence have a length of 13. This issue comes up in resource record handling; see below.

The following Python classes in package `libs.dnslib` parse and construct DNS protocol payloads using the classes in `libs.inetlib` discussed above.¹

- `Header` represents a DNS protocol payload header; its constructor creates a new `Header` object from user-specified arguments, while its `fromData` method returns a new `Header` object from user-supplied binary string data.
- `QE` represents the question entry that is usually the first entry after the header. Like `Header`, `QE` has a constructor that creates a `QE` from individual arguments, and a `fromData` method that returns a `QE` object from binary string data.
- `RR` is an abstract class that represents a generic resource record, and contains logic common to the various types of resource records. All resource records contain a domain name, a time-to-live, and a field specifying the length of the data they contain.

¹Note that Python packages reside in subdirectories corresponding to the package name, so the directory path to Alyssa's DNS data format handling code is `libs/dnslib` in the lab distribution directory.

RR's `fromData` method returns a pair (`rr`, `length`) where `rr` is one of the classes below derived from `RR`, which represent resource records of various types, and `length` is the length in bytes of the resource record returned, in the binary string representation supplied to the user's call to `fromData`. Any `DomainName` objects that `RR` or its derivatives contain are copies of the originals, and so have a `length` commensurate with their string representation length. See above for a discussion of string representation `DomainName` length versus binary representation `DomainName` length. Thus, as a user of the `RR` classes, you will need to refer to the second element of the pair that `fromData` returns in order to find out how many bytes a given `RR` takes up in the binary string you give `fromData`.

The classes derived from `RR` are:

- `RR_A` represents an “address” resource record, which contains an IPv4 address.
- `RR_CNAME` represents a canonical name address record, which contains a canonical host-name (a `DomainName`).
- `RR_NS` represents a nameserver resource record, which contains the domain name of another nameserver (a `DomainName`).
- `RR_SOA` represents a start-of-authority resource record, which contains information about a domain.

Debugging and logging

Alyssa provides a `hexdump` function in package `libs.util` to assist you in debugging the contents of packets received or sent from or to other DNS servers. When called with one argument whose value is a packed binary string containing any data, `hexdump` will list the data in human-readable hexadecimal and ASCII form, in two columns, with byte indices in another column to the left. The hexadecimal display is one byte (two hexadecimal characters) at a time. By judicious use of `hexdump` on each packet that you send or receive on the network, combined with cross-referencing the DNS specifications, you can resolve many issues.

A sample invocation of `hexdump` in the context of the provided code is:

```
while 1:
    (data, address,) = ss.recvfrom(512) # DNS limits UDP msgs to 512 bytes
    if not data:
        log.error("client provided no data")
        continue
    else:
        print "Query received from client is:\n", hexdump(data)
        queryheader = Header.fromData(data)
        print "Query header received from client is:\n", hexdump(queryheader.pack())
```

The resulting output visible on the console when the server runs, answering a query for `www.cs.ucl.ac.uk` is:

```
Query received from client is:
0000  38 8E 01 00 00 01 00 00 00 00 00 00 03 77 77 77      8.....www
0010  02 63 73 03 75 63 6C 02 61 63 02 75 6B 00 00 01      .cs.ucl.ac.uk...
0020  00 01                                                ..

Query header received from client is:
0000  38 8E 01 00 00 01 00 00 00 00 00 00 00          8.....
```

Action: To familiarize yourself with DNS payloads (and thus prepare for debugging the results when you begin to construct your own), reconcile the above hexdump with the DNS packet format, understanding where each field's value is in the hexdump.

Alyssa has also provided a logging infrastructure to assist in the debugging of your server. Package `libs.util` sets up logging to the console and file `ncsdns.log`, rotating that log file and its predecessors to files `ncsdns.log.1`, `ncsdns.log.2`, `ncsdns.log.3`, etc., all in the directory from which you run your server.² Then to debug your server, make calls from `ncsdns.py` to the `logger` object with debugging messages you find useful. Each logging call takes two arguments: an importance level, and a string to output to the log. The importance levels and the corresponding call to make in your code are as follows:

Level	Numeric value	How to call
CRITICAL	50	<code>logger.critical(string)</code>
ERROR	40	<code>logger.error(string)</code>
WARNING	30	<code>logger.warning(string)</code>
INFO	20	<code>logger.info(string)</code>
DEBUG	10	<code>logger.debug(string)</code>
DEBUG1	9	<code>logger.log(DEBUG1, string)</code>
DEBUG2	8	<code>logger.log(DEBUG2, string)</code>

Then in file `libs/util.py` you can configure calls to `setLevel` as directed by the comments in order to alter the verbosity level of the debugging output directed to either the console or the log file. When you select a given verbosity level x you will then see only logging calls tagged with a level whose numeric value is x or greater.

Caching

One of the requirements in the next section is that your server must answer queries from an internal *cache* whenever possible, instead of querying another nameserver. Since each resource record that your nameserver receives as a result of making a query on another nameserver contains a *time to live* field (as described in lecture), upon receiving each query, your nameserver can deduce which resource records are still valid, and use the valid records to answer the query.

Requirements and marking scheme

To earn marks on this coursework, your nameserver must generate timely (within the default timeout of the `dig` tool) and valid responses to the queries we issue according to the following requirements, for **every query** that your nameserver responds to.

Following are the requirements for your nameserver. Failure to meet any of these requirements for *any* of the queries we make to your nameserver when we evaluate it will result in a loss of marks.

1. Your nameserver must run with unchanged versions of our library files.
2. Your nameserver must answer recursive queries within 60 seconds without entering an infinite loop or terminating without responding, no matter what information it receives from other DNS servers. Note that this requirement means your nameserver should always reply to the client within this deadline, even if its reply must indicate failure to resolve the query fully.
3. Upon answering a query, your nameserver must wait for a subsequent query, and answer it in turn.
4. Your nameserver must follow all canonical name `CNAME` aliases for a query and return the correct address records for the canonical host name in the query.
5. Your nameserver must return at least one authority record from a subdomain of the canonical host name. That subdomain must be the most highly-qualified subdomain that has at least one authoritative nameserver.

²Note that the rotation happens when `ncsdns` exits, so to read the most recent log, check `ncsdns.log.1`.

For example, if query *a.b.c.d* aliases with CNAME entries to canonical hostname *e.f.g.h*, you must return an authority record of a nameserver for subdomain *f.g.h*, if one exists. Otherwise, you must return an authority record of a nameserver for subdomain *g.h* if one exists, and so on.

6. Your nameserver must return the correct glue records for nameservers it mentions in the authority section of its reply.
7. When possible, your nameserver must use the cache to answer a given query, without generating any network traffic.
8. Your nameserver must cope with not receiving a reply to requests it sends, for any reason (failure of *other* DNS servers, network failures, timeouts due to network congestion, for example). In particular:
 - (a) Your implementation must not crash or exit or hang if it doesn't get a reply to a request it sends.
 - (b) Your implementation must not immediately give up and return an error to the client if it doesn't get a reply to a request it sends.
 - (c) Your implementation must make a robust effort to resolve the client's request successfully even if it receives no reply to one of its own requests in the course of resolving the client's request. Consider that UDP packets aren't delivered reliably on the Internet. So if you don't get a reply, there's a chance the request packet simply wasn't delivered. Even if a particular remote DNS server your implementation tries to contact is down and unreachable entirely, it may be possible in many cases to still resolve the query. (Consider, for example, that there are often multiple NS records for the same domain.)

[40 marks]

Negative requirements. Your nameserver does **not** need to implement the following functions that production DNS servers provide:

1. Your nameserver does not need to respond to queries containing multiple questions (multiple hostnames to resolve, for example).
2. Your nameserver does not need to respond to queries of type other than A (address). In particular, your nameserver does not need to respond to or include IPv6 AAAA records.
3. Your nameserver does not need to respond to iterative queries.
4. You are not required to return time-to-live fields that take into account network propagation delay.

Documentation. No separate design document is required, but you must comment your code thoroughly, to fully explain how it works.

[10 marks]

Testing your nameserver

Python 3.6 is installed on the lab machines. To run and test your server, use the following command:

```
$ python3 ./ncsdns.py
```

The server responds something like:

```
./ncsdns.py: listening on port 40229
```

indicating that it is listening on port 40229. Note that the server listens on a different port every time you run it.

To workaroud the need for “root” access to the machine we are using, our server socket `ss` listens on a high-numbered “ephemeral” port instead of port 53. The server prints the port it is listening on to the standard output when it starts, as above.

For testing purposes, you should send `dig` queries for a variety of DNS hostnames to the port your server is listening on using the `-p [port number]` option. In a different window, use `dig` to query your local nameserver:

```
$ dig @127.0.0.1 -p [port number] [query hostname]
```

To test the server in this example, you would type:

```
$ dig @127.0.0.1 -p 40229 www.xorp.org.
```

Initially, `dig` will time out, not having received a response from your incomplete DNS server.

Once your server responds to queries, use appropriate queries to `dig` to check the responses your server returns against the requirements listed above. Try testing using the following host names, referring to the log as you test for insights as to what the possible bugs in your server are. These are among the host names we will use to mark your server:

```
haig.cs.ucl.ac.uk.  
www.xorp.org.  
lxr.linux.no.  
www.caltech.edu.  
www.microsoft.com.  
www.yahoo.com.
```

Suggestion: For each DNS query, you might want to compare your nameserver’s responses with the ones from a production nameserver (e.g., `haig.cs.ucl.ac.uk`). Our automated tests will indeed check for consistency of your nameserver’s output with the one of a standard `dig` command run on departmental machines.

Note that your results might vary from query to query for domains that use DNS load balancing to return different results in round-robin order or DNS-based content distribution networks that change results based on client identity. We will only mark your server against the requirements listed above, thus taking this into account.

[Total marks for Part 2: 50]

Hand in instructions

Submission of this coursework is electronic using the COMP0023 Moodle web site.

Hand in Part 1 of this coursework by cutting and pasting your work into a plain ASCII text file named `dns-cw-part1.txt` and submitting only this text file using Moodle.

For Part 2, hand in the file `ncsdns.py` completed with your implementation of the nameserver.

The teaching staff is not responsible for coursework submissions not made in accordance with the above instructions—in particular, we do not accept coursework submissions via email.

[Total marks for this coursework: 60]

References

The following references on the Domain Name System may be useful in completing this assignment:

1. RFC1034: Domain names—Concepts and facilities (web)
2. RFC1035: Domain names—Implementation and specification (web)
3. RFC1536: Common DNS implementation errors and suggested fixes (web)
4. Python documentation
<http://docs.python.org>
<http://docs.python.org/tutorial/index.html>
5. The Python library reference is available at:
<http://docs.python.org/library/index.html>

Announcements

Please monitor the class Piazza discussion forum during the period between now and the due date for the coursework. Any announcements (*e.g.*, helpful tips on how to work around unexpected problems encountered by others) will be sent there.

Academic honesty

This coursework is an individual coursework. You have to complete it alone; you may not consult with other students in the course (or with other people who are not taking the course) about what algorithm to use or how to write your code. You may not look at other students' code (from this year or past years), nor show them yours. You are free to read reference materials found on the Internet (and any other reference materials). You may of course use the code we have given you. All other code you submit must be written by you alone. Copying of code from student to student is a serious infraction; it will result in automatic awarding of zero marks to all students involved, and is viewed by the UCL administration as cheating under the regulations concerning Examination Irregularities (normally resulting in exclusion from all further examinations at UCL). The course staff use extremely accurate plagiarism detection software to compare code submitted by all students and identify instances of copying of code; this software sees through attempted obfuscations such as renaming of variables and reformatting, and compares the actual parse trees of the code. Rest assured that it is far more work to modify someone else's code to evade the plagiarism detector than to write code for the assignment yourself!