

Udacity CS101: Building a Search Engine

Unit 4: Building an Index

[Introduction \(Introduction\)](#)

[Q-1: Quiz \(Data Structures\)](#)

[Add to Index \(Add to Index\)](#)

[Q-2: Quiz \(Add to Index\)](#)

[Lookup Procedure \(Lookup\)](#)

[Q-3: Quiz \(Lookup\)](#)

[Building the Web Index \(Building the Web Index\)](#)

[Q-4: Quiz \(Add Page to Index\)](#)

[Finishing The Web Crawler \(Finishing the Web Crawler\)](#)

[Q-5: Quiz \(Finishing the Web Crawler\)](#)

[Startup \(Startup\)](#)

[The Internet \(The Internet\)](#)

[Networks \(Networks\)](#)

[Q-6 Quiz \(Networks\)](#)

[Smoke Signals \(Smoke Signals\)](#)

[Latency \(Latency\)](#)

[Q-7 Quiz \(Latency\)](#)

[Bandwidth \(Bandwidth\)](#)

[Bits \(Bits\)](#)

[Q-8 Quiz \(Bits\)](#)

[Measuring Bandwidth \(What Is Your Bandwidth\)](#)

[Q-9: Quiz \(What Is Your Bandwidth\)](#)

[Measuring Latency \(Traceroute\)](#)

[Q-10: Quiz \(Traveling Data\)](#)

[Making a Network \(Making a Network\)](#)

[Protocols \(Protocols\)](#)

[Conclusion](#)

[Answer Key](#)

Introduction (Introduction)

In unit 4 you are going to learn how to finish the code for your search engine and how to respond to a query when someone wants the given web pages that correspond to a given keyword. You will also learn about how networks and the world wide web work to understand more about how you can build up your search index.

The main new computer science idea you will learn is how to build complex data structures. You will learn how to design a structure that you can use so that you can respond to queries without needing to rescan all the web pages every time you want to respond to a query. The structure you will build for this is called an index. The goal of an **index** is to map a keyword and where that keyword is found. For example, in the index of a book you can see a page number which serves as a map to where a term or concept can be found. The key ideas in index will allow us to find references to what we want. With a search engine the index gives you a way for a keyword to map to a list of web pages, which are the urls where those particular web pages appear. Once you have done the work of building an index, then the look-ups are really fast.

Q-1: Quiz (Data Structures)

Deciding on data structures is one of the most important parts of building software. As long as you pick the right data structure, the rest of the code will be a lot easier to write.

Which of these data structures would be a good way to represent the index for your search engine?

- a. [*<keyword₁>*, *<url_{1, 1}>*, *<url_{1, 2}>*, *<keyword₂>*, ...]
- b. [[*<keyword₁>*, *<url_{1, 1}>*, *<url_{1, 2}>*], [*<keyword₂>*, *<url_{2, 1}>*, ...]]
- c. [[*<url_{1, 1}>*, [*<keyword₁>*, *<keyword₂₃>*, ...]], [*<url_{2, 1}>*, [*<keyword₂>*, ...]]]
- d. [[*<keyword₁>*, *<url_{1, 1}>*, *<url_{1, 2}>*]], [*<keyword₂>*, [*<url_{2, 1}>*]], ...]

[Answer to Q-1](#)

Add to Index (Add to Index)

Q-2: Quiz (Add to Index)

Define a procedure, `add_to_index`, that takes three inputs:

- an index `[[<keyword>, [<url>, ...]], ...]`
- a keyword string
- a url string

If the keyword is already in the index, add the url to the list of urls associated with that keyword.

If the keyword is not in the index, add an entry to to the index: `[keyword, [url]]`

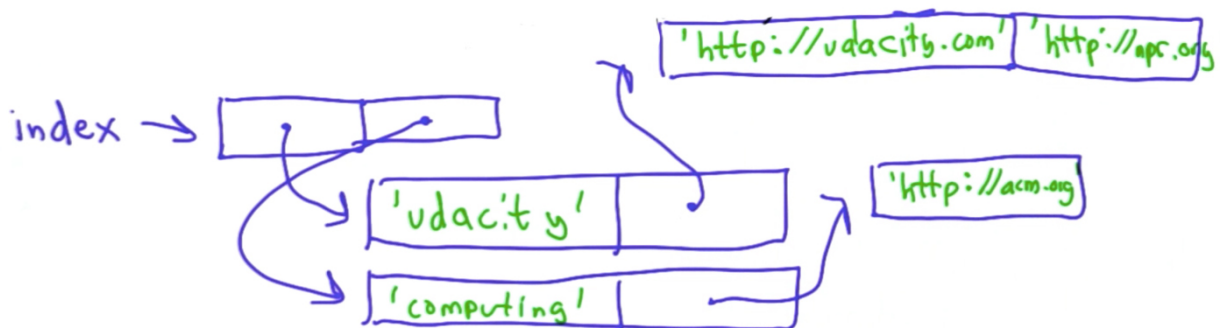
For example:

```
index = []
add_to_index(index, 'udacity', 'http://udacity.com')
add_to_index(index, 'computing', 'http://acm.org')
```

This code starts with the empty list `index`. After the two lines of code the empty list will contain two lists beginning with the keywords, `udacity` and `computing`.

```
index = []
add_to_index(index, 'udacity', 'http://udacity.com')
add_to_index(index, 'computing', 'http://acm.org')
add_to_index(index, 'udacity', 'http://npr.org')
```

In this code, `udacity` is already in the index and you don't want to add a new entry to the index itself. Since `udacity` is already in the index, what you want to do is add the new url to the list already associated with that keyword.



[Answer to Q-2](#)

Lookup Procedure (Lookup)

Q-3: Quiz (Lookup)

Define a procedure, lookup, that takes two inputs:

- An index: A list where each element of the list is a list containing a keyword and a list as its second element. The second list element is a list of urls where that keyword appears.

-The keyword to lookup

The output should be a list of the urls associated with the keyword. If the keyword is not in the index, the output should be an empty list.

For example:

```
lookup('udacity') → ['http://udacity.com', 'http://npr.org']
```

[Answer to Q-3](#)

Building the Web Index (Building the Web Index)

To build your web index you want to find a way to separate all the words on a web page. It is possible to use the concepts you've already seen to build a procedure to do this, however, Python has a built-in operation that will make this much simpler.

Split. When you invoke the split operation on a string the output is a list of the words in the string.

```
<string>.split()
[<word>, <word>, ... ]
```

For example,

```
quote = "In Washington, it's dog eat dog. In academia, it's exactly the
opposite. --- Robert Reich"
print quote.split()
['In', 'Washington,', "'it's", 'dog', 'eat', 'dog.', 'In', 'academia,',
 "'it's", 'exactly', 'the', 'opposite.', '---', 'Robert', 'Reich']
```

This operation does a pretty good job of separating out the words in the list so that they will be useful. However, in the case of **'Washington,'** which was followed by a comma in the quote, for the keyword you would not want to include the comma. While this isn't perfect, it is going to good enough for now.

Here is another example of how **split** works, using triple quotes ("""). Using the triple quotes you can define one string over several lines:

```
quote = """
    There's no buisness like show business,
    but there are several businesses like accounting.
    (David Leterman)
    """
print quote.split()
["There's", 'no', 'buisness', 'like', 'show', 'business,', 'but',
 'there', 'are', 'several', 'businesses', 'like', 'accounting.'
 '(David', 'Leterman)']
```

This still has similar problems to the first example, where the parentheses are included in the word **'(David'**.

Q-4: Quiz (Add Page to Index)

Define a procedure, `add_page_to_index`, that takes three inputs:

- index
- url (string)
- content (string)

It should update the index to include all of the word occurrences found in the page content by adding the url to the word's associated url list.

For example:

```
index = []

add_page_to_index(index, 'fake.test', "This is a test")
print index
[['This', ['fake.test']], ['is', ['fake.test']], ['a', ['fake.test']],
 ['test', ['fake.test']]]
```

Printing at `index[0]`:

```
index = []

add_page_to_index(index, 'fake.test', "This is a test")
print index[0]
['This', ['fake.test']]
```

Printing at `index[1]`:

```
index = []

add_page_to_index(index, 'fake.test', "This is a test")
print index[1]
['is', ['fake.test']]
```


Now, add a page called **real.test**:

```
index = []

add_page_to_index(index, 'fake.test', "This is a test")
add_page_to_index(index, 'real.test', "This is not a test")

print index
[['This', ['fake.test', 'real.test']], ['is', ['fake.test',
'real.test']], ['a', ['fake.test']], ['test', ['fake.test',
'real.test']], ['not', ['real.test']]]
```

Have a look at the entries when you **index[1]**:

```
index = []

add_page_to_index(index, 'fake.test', "This is a test")
add_page_to_index(index, 'real.test', "This is not a test")

print index
print index[1]
[['This', ['fake.test', 'real.test']], ['is', ['fake.test',
'real.test']], ['a', ['fake.test']], ['test', ['fake.test',
'real.test']], ['not', ['real.test']]]
['is', ['fake.test', 'real.test']]
```

Have a look at the entries when you `index[4]`:

```
index = []

add_page_to_index(index, 'fake.test', "This is a test")
add_page_to_index(index, 'real.test', "This is not a test")

print index

print index[1]

print index[4]
[['This', ['fake.test', 'real.test']], ['is', ['fake.test',
'real.test']], ['a', ['fake.test']], ['test', ['fake.test',
'real.test']], ['not', ['real.test']]]

['is', ['fake.test', 'real.test']]

['not', ['real.test']]
```

You have already defined a procedure for responding to a query, so check out if it works on this index, searching for the keyword 'is':

```
index = []

add_page_to_index(index, 'fake.test', "This is a test")
add_page_to_index(index, 'real.test', "This is not a test")

print lookup(index, 'is') # you should expect to see that this keyword
appears on both pages
['fake.test', 'real.test']
```

Now try searching for a keyword, **'udacity'**, that is not in either of the urls:

```
index = []

add_page_to_index(index, 'fake.test', "This is a test")
add_page_to_index(index, 'real.test', "This is not a test")

print lookup(index, 'udacity') # you should expect to see an empty list
[]
```

Well, that was successful! Try to define this procedure.

[Answer to Q-4](#)

Finishing The Web Crawler (Finishing the Web Crawler)

Returning to the code you wrote before for crawling the web, make some modifications to include the code you've just written.

First, a quick recap on how the `crawl_web` code below works before incorporating the indexing:

```
def crawl_web(seed):
    tocrawl = [seed]
    crawled = []
    while tocrawl:
        page = tocrawl.pop()
        if page not in crawled:
            union(tocrawl, get_all_links(get_page(page)))
            crawled.append(page)
    return crawled
```

First, you defined two variables `tocrawl` and `crawled`. Starting with the `seed` page, `tocrawl` keeps track of the pages left to crawl, whereas `crawled` keeps track of the pages that have already been crawled. When there are still pages left to crawl, remove the last page from `tocrawl` using the `pop` method. If that page has not been crawled yet, get all the links from the page and add the new ones to `tocrawl`. Then, add the page that was just crawled to the list of `crawled` links. When there are no more pages to crawl, return the list of `crawled` pages.

Adapt the code so that you can use the information found on the pages crawled. The changed code is below. First, add the variable `index`, to keep track of the content on the pages along with their associated urls. Since you are really interested in the `index`, this is what we will return.

It is possible to return both `crawled` and `index`, but to keep it simple just `return index`. Next, add a variable, `content` to replace `get_page(page)`. This variable will be used twice, once in the code already there and once in the code to be filled in for the quiz. The procedure `get_page(page)` is expensive as it requires a web call, so we don't want to call it more often than is necessary. Using the variable `content` means that the call to the procedure `get_page(page)` only needs to be performed once and then the result is stored and can be used over and over without having to go through the expensive call again.

Q-5: Quiz (Finishing the Web Crawler)

Fill in the missing line using the variable **content**.

```
def crawl_web(seed):
    tocrawl = [seed]
    crawled = []
    index = []
    while tocrawl:
        page = tocrawl.pop()
        if page not in crawled:
            content = get_page(page)
            # FILL IN THE MISSING LINE BELOW HERE

            union(tocrawl, get_all_links(content))
            crawled.append(page)
    return index
```

[Answer to Q-5](#)

Startup (Startup)

You now have a functioning web crawler! From a seed you can find a set of pages; for each of these pages you can add the content from that page to an index, and return that index. Additionally, since you have already written the code you can do the look-up that will return the pages for that keyword.

But you're not completely done yet. In unit 5, you will see how to make a search engine faster and in unit 6 you will learn how to find the best page for a given query rather than returning all the pages.

Before then, you need to understand more about how the Internet works and what happens when you request a page on the world wide web.

The Internet (The Internet)

So far, you've been using the `get_page` function where you pass in a `url` and it passes out the content of the page. This is the python code that does that:

```
def get_page(url):
    try:
        import urllib
        return urllib.urlopen(url).read()
    except:
        return ""
```

The table below shows what each line of the code does.

Code	Explanation
<code>urllib.urlopen(url)</code>	This opens the web page at url
<code>urllib.urlopen(url).read()</code>	and reads the page requested which is a string.
<code>return urllib.urlopen(url).read()</code>	That string is then returned.
<code>import urllib</code>	The code which contains <code>urlopen</code> is in the library <code>urllib</code> , so we have to import that.
<code>try: except: return ""</code>	This is an exception handler. It's called a try block. We try these things but they might not always work. There might be an error. The page might not be returned, or the url might be bad or the page times out. If we request a url which we can't be loaded, it jumps to the <code>except:</code> block and returns an empty string instead of producing an error.

Looking at the code for `get_page` doesn't help you understand what is going on when you request a page on the web because it is all hidden in the Python library, `urllib`. To understand more, read about how the Internet works.

Networks (Networks)

The **Internet** is a particular type of network. So what is a network?

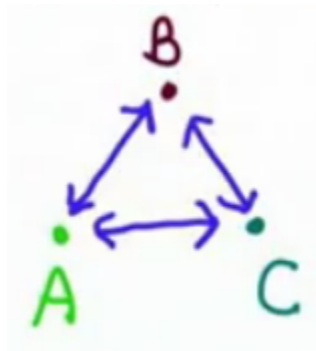
A **network** is a group of entities that can communicate, even though they are not directly connected.

Meet Alice and Bob. They can communicate directly, so this is not a network by our definition:



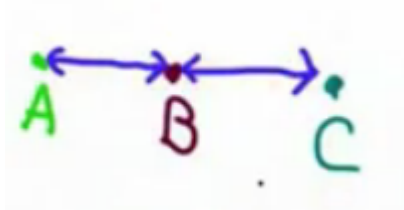
Not a Network

Next, Carol comes along. Alice, Bob and Carol can all communicate directly so this is still not a network.



Not a Network

In this example, Alice and Carol can not communicate directly but they can still communicate via Bob, so this is a network.



Network

Many people would consider all three examples of a network, but you want to think about a more precise definition. A message that has to go through two hoops, so that people who are not directly connected can still communicate, is a network. It must have at least three people.

Q-6 Quiz (Networks)

This quiz tests your understanding of what a network is, and your knowledge of world history!

How old is the idea of a network?

- a. 10 years old
- b. 30 years old
- c. 100 years old
- d. 1000 years old
- e. over 3000 years old

[Answer to Q-6](#)

Smoke Signals (Smoke Signals)

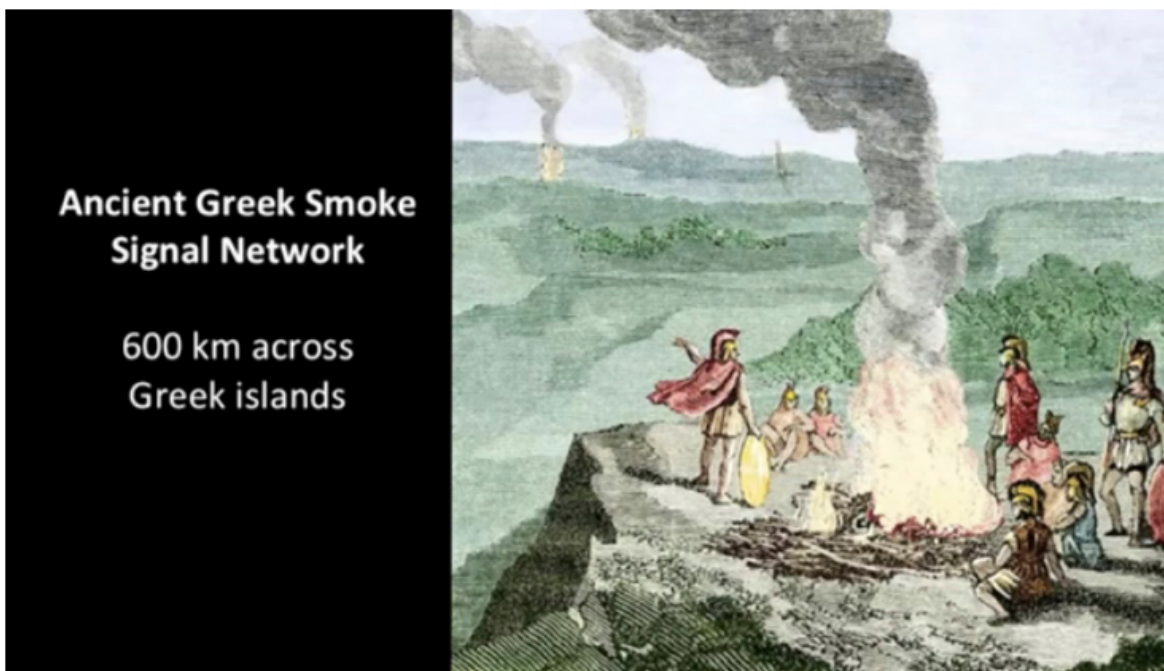
Humans have been using networks since at least as long as we can record. They might have just been networks where people just talked to each other and passed on message but even 3000 years ago, there were sophisticated networks using technology. In the Illiad, written about 800 BC, Homer describes a network from approximately 1200 B.C.

Thus, from some far-away beleaguered island, where all day long the men have fought a desperate battle from their city walls, the smoke goes up to heaven; but no sooner has the sun gone down than the light from the line of beacons blazes up and shoots into the sky to warn the neighbouring islanders and bring them to the rescue in their ships.

Iliad, Homer

written ~800 BC describing events ~1200BC

The Ancient Greeks had a line of beacons spread over the Greek Islands, which could be lit to announce an oncoming attack. When the soldiers at one of those beacons saw another beacon lit they would light their own beacon, spreading the message across the islands.



What does it take to make a network like this work? Suppose a message is to be sent from Rhodes to Sparta. There are many points along the way where smoke signals can be sent. First, the beacon is lit at Rhodes to be seen by the places near Rhodes. To reach Sparta, Naxos needs to resend the message by lighting its own beacon, which can be seen at Melos, Cydonia and Athens. When Melos lights its beacon, Sparta receives the message.



Several different things are needed in order to make a network like this work:

- a way to encode and interpret messages - Greeks: "Agamennon is arriving" specific smoke signal
- a way to route messages, that is, a way to indicate which nodes should forward the message. Most likely the Ancient Greeks did this by lighting all the beacons, which is very wasteful. For the Greeks, a message such as 'the enemies are arriving' would be fine to spread throughout. This is called "flooding the network." To send a message between two points is a harder problem.
- rules for deciding who gets to use resources -- who gets to send a message. Greeks: Generals have priority.

All these things are needed on the Internet too.

Latency (Latency)

Latency is the time it takes a message to get from the source to the destination.

Latency is measured in seconds, or for a fast network today, in milliseconds. There are 1000 milliseconds in 1 second. (If you care about your ping in online games, it's latency that matters.)

Q-7 Quiz (Latency)

Zeus, the all powerful Greek god, wants to send a message from Rhodes to Sparta. He thinks that the latency of the smoke signalling system is too high, which means it takes too long for a message to be sent.

How could Zeus (remember, he's all powerful) reduce the latency between Rhodes and Sparta? (Choose one or more answers.)

- a. Make the signalling nodes further apart, so it takes fewer hops.
- b. Threaten the soldiers to make the signalling fires quicker.
- c. Add colours to the smoke so there are more messages per signal.
- d. Increase the speed of light.

[Answer to Q-7](#)

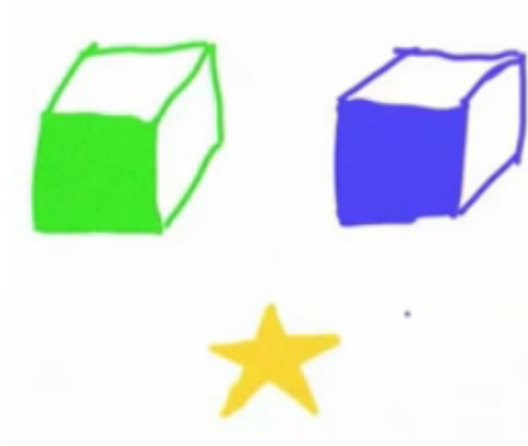
Bandwidth (Bandwidth)

Bandwidth is the amount of information that can be transmitted per unit of time. Bandwidth does not consider the start up time (that's the latency) to get the first bit across, but the amount of information that can be got across once the first bit has arrived. In other words, it's the rate at which information is transferred.

Bandwidth is measured in units of [information/time], such as bits per second. On the Internet it is often measured in Mbps which is megabits per second.

Bits (Bits)

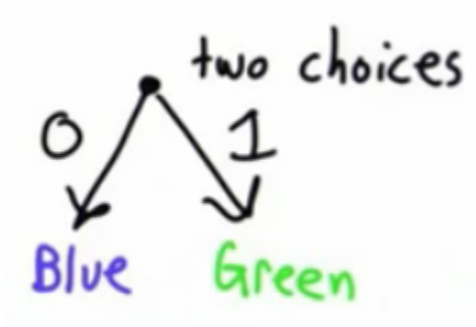
A **bit** is the smallest unit of information.



For example, if you ask, "Is the gold star in the green box?" you get one bit of information back. This bit of information takes you from having two choices to just one. If the answer is yes, you know the green box is the box you should open. If the answer is no, you know it's the blue box.

In computing, rather than talking about bits in terms of yes or no, you talk about them as 0 and 1, where 0 more readily maps to no and 1 to yes.

The diagram below shows how the question, “Is the gold star in the green box?” reduces the two choices to one with just one bit of information. If the answer is yes, that is 1, then the choice labelled with that arrow leads to the green box. If the answer is no, which is 0, the arrow leads to the blue box.

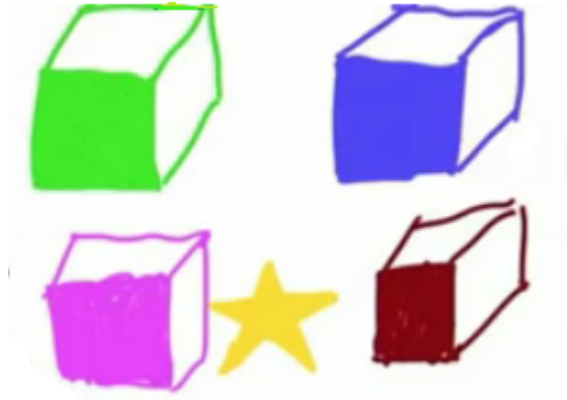


Therefore, if there are two choices then knowing which to choose is one bit of information. For the Ancient Greeks' network, they could only transmit one bit of information by lighting the fire. If the fire was lit, then yes (1) or not lit (0). If that were the case, they would only have been able to send one message - the enemy is coming.

If all that can be sent is 0s and 1s, can anything more interesting be sent? Can just one thing be sent?

Yes! Anything you want can be sent!

Suppose that instead of two boxes, there are four boxes: green, blue, purple and red. There is still only one gold star.



You could ask:

- Is the star in the green box?
- Is the star in the purple box?
- Is the star in the blue box?

As long as you know the gold star is in one of the boxes, you don't need to ask "Is the star in the red box?" because a negative answer to the other three questions already tells you it is. In order to be sure about where the gold star is, you need to ask at most, three questions, which takes in three bits. Is it possible to do better?

Q-8 Quiz (Bits)

How many bits are needed to find the gold star?

- 1
- 2
- 3

[Answer to Q-8](#)

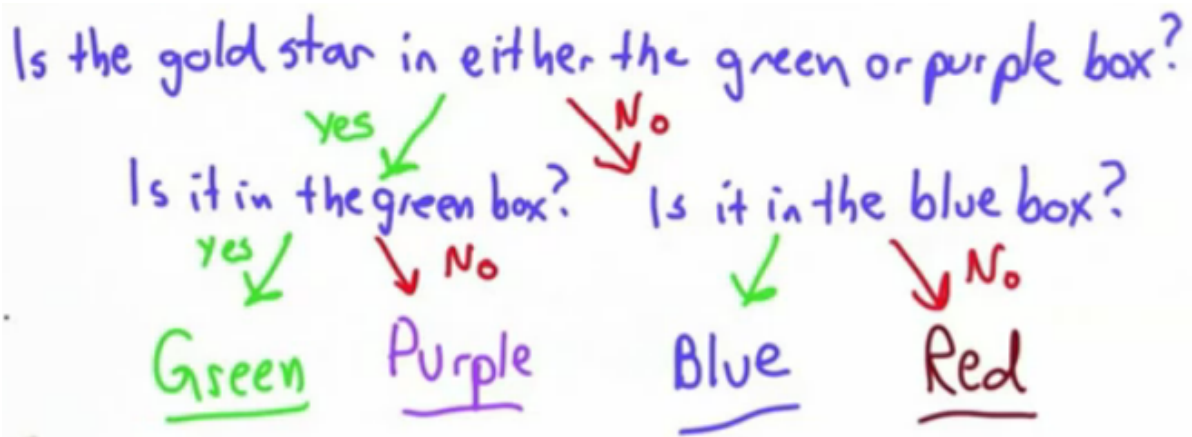
Before, these were the questions:

- Is the star in the green box?
- Is the star in the purple box?
- Is the star in the blue box?

Three questions might be needed, but two questions are enough if the answers to the questions give more information. The problem with asking “Is the star in the green box?” is that the answer is “No” three quarters of the time and “Yes” only one quarter. A guess of “No” would be right more than half the time. This means that the question does not provide a full bit’s worth of information.

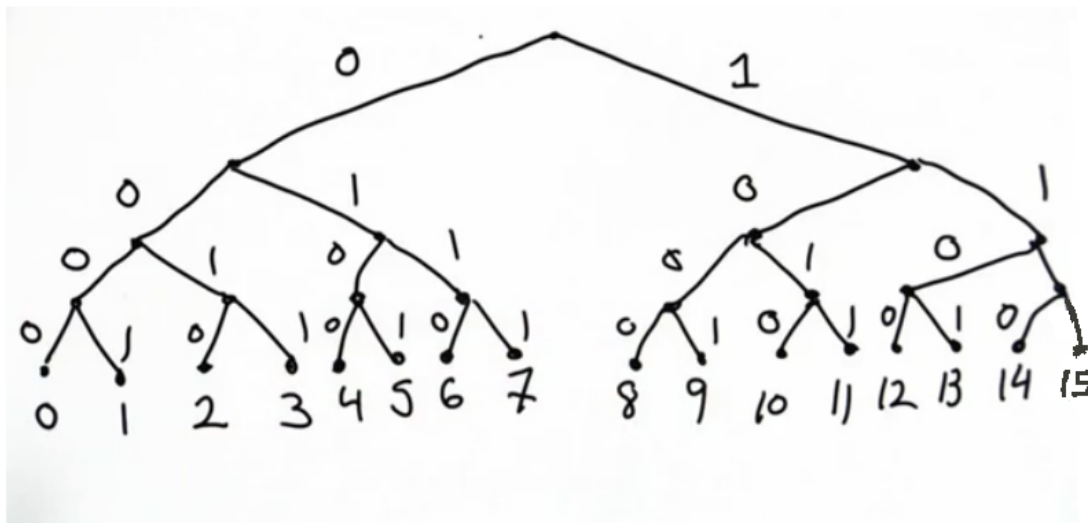
Instead, you would be better off asking a question that has an answer equally likely to be “Yes” as it is to be “No”. So, an initial question could be “Is the gold star in either the green or purple box?” If the answer is yes, then the number of choices is cut in half and the gold star is in the green or purple box. If the answer is no, again, the number of choices has halved and the gold star is in the blue or the purple box. After that question, only one more question is needed to determine exactly which box the star is in.

The diagram below illustrates that only a total of two questions is needed:

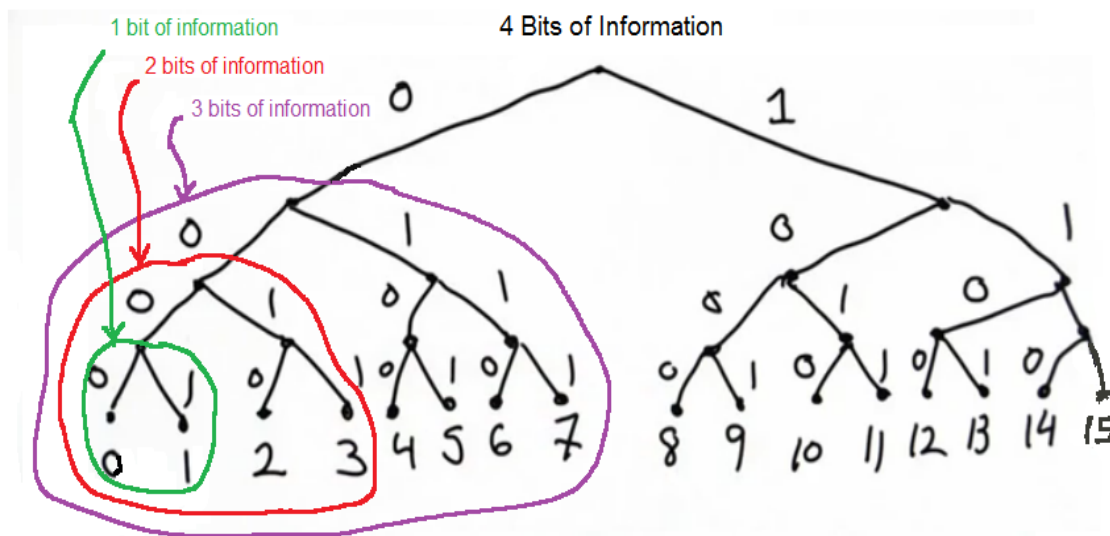


This example shows that you can encode four things in bits. There are four colours, and two bits are needed to find an answer.

Any number of things can be encoded with more bits. In the diagram below, sixteen numbers are encoded using four bits.



For every extra bit added, double the number of things can be encoded as there is one more yes/no question.



This diagram tells you:

For one bit, 2 numbers can be encoded.

For two bits, it's double that number so, $2 * 2 = 2^2 = 4$ numbers.

For three bits, it's double again, which means a total of $2 * 2 * 2 = 2^3 = 8$ numbers encoded.

For four bits, as in the example above, $2 * 2 * 2 * 2 = 2^4 = 16$ numbers are encoded.

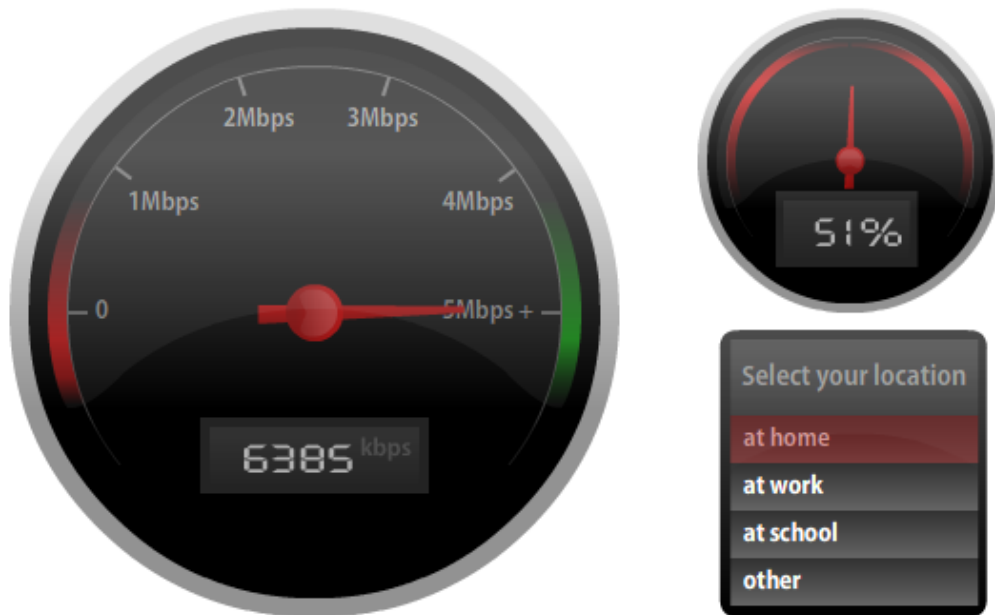
In general, the number of things which can be encoded in n bits is 2^n .

Once a number of things can be distinguished, information can be sent. Anything that is discrete, such as strings or lists, can be converted into a number and that number can then be transmitted. Once bits can be sent, anything can be sent. The number of bits needed is a measure of the amount of information that can be sent.

Bits are directly related to bandwidth, which is a measure of the amount of information that can be sent. Bandwidth is measured in bits per second, each bit being a yes/no question -- 0s and 1s. Those 0s and 1s could be encoding a string of text, like in a web page, or an image.

Measuring Bandwidth (What Is Your Bandwidth)

There are different ways to measure the bandwidth over your Internet connection. One way is to use <http://www.cnet.co.uk/broadband-speedtest/>. This site is not necessarily that accurate. It will try sending messages to figure out your bandwidth and depends on your location. In the image below, the bandwidth shown is 6385 kbps (kilobits per second) which is 6.385 Mbps (megabits per second). If you do the test several times, you'll find the values vary.



Bandwidth does vary since you are sharing the net with other people who may be doing different things. There are many reasons why bandwidth varies.

Q-9: Quiz (What Is Your Bandwidth)

This is more of a survey than a quiz.

What is your bandwidth?

- less than 1 kbps
- 1-10 kbps
- 10-100 kbps
- 100-1000 kbps 1000 kbps = 1Mbps
- 1-10 Mbps
- 10-20 Mbps
- 20-40 Mbps
- 40-60 Mbps
- 60-100 Mbps
- 100-200 Mbps

[Answer to Q-9](#)

Measuring Latency (Traceroute)

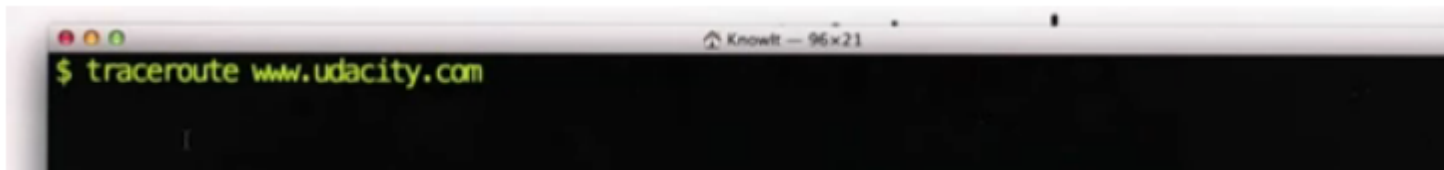
You can learn a lot about the Internet by measuring latency to different destinations. For the map of Greece, there were hops on the way from Rhodes to Sparta.

Rhodes ->Naox -> Melos ->Sparta

Using an application called **traceroute**, the hops on the Internet can be seen between the machine it is run on and the destination.

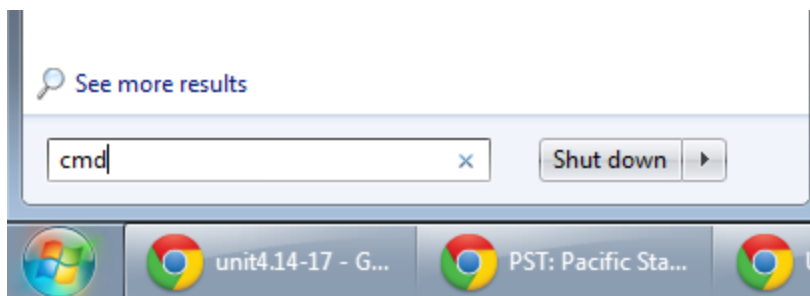
Mac:

On a mac, you can create a shell and run **traceroute** directly.



Windows:

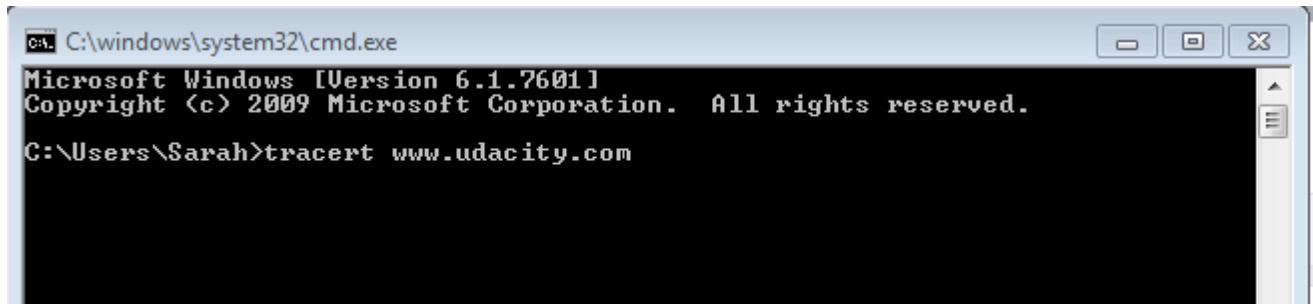
7 and Vista: Click on the start menu. Type cmd in the search bar and hit enter.



Earlier versions of Windows:

Click on the start menu and select run. In the open: box, type cmd and hit enter.

You'll get up a command prompt window.

A screenshot of a Windows command prompt window. The title bar shows the path C:\windows\system32\cmd.exe. The window content displays the following text: Microsoft Windows [Version 6.1.7601] Copyright (c) 2009 Microsoft Corporation. All rights reserved. C:\Users\Sarah>tracert www.udacity.com

Note that the command in Windows is **tracert** and not **traceroute**.

Unix:

Run **traceroute** from a terminal window.

Once you've typed **traceroute** (or **tracert** for Windows), pick where you want the destination to be. For the first example, **traceroute www.udacity.com**, you can see all the hops along the way. **Traceroute** is sending packets across the network looking at all the intermediate hops to figure out the route it takes to get a packet from the current location to www.udacity.com. A **packet** is a piece of information sent over the Internet. It contains a header and a message. The header contains information about, among other things, the packet source, destination and its "hop limit". Hop limit is the amount of hops the packet is allowed to go through. This is so that a packet can not get stuck in an infinite loop, jumping back and forth between the same places. The hop limit is used in **traceroute** to send packets out different distances on the route to the destination. The hop limit is only 1 byte long, which is 8 bits. From the previous discussion on bits, you saw that the amount of information which can be encoded by 8 bits is $2^8 = 256$ so the maximum number of hops is 255 (since it goes from 0 to 255 inclusive.) This puts a maximum distance a packet can travel on the Internet at 255 hops. Fortunately, everywhere on the Internet is connected by far fewer hops than that!

```
tracert to ghs.l.google.com (74.125.127.121), 64 hops max, 52 byte packets
 1  192.168.1.1 (192.168.1.1)  1.570 ms  1.284 ms  1.152 ms
 2  10.1.10.1 (10.1.10.1)  3.306 ms  2.843 ms  1.601 ms
 3  96.157.108.1 (96.157.108.1)  40.692 ms  79.370 ms  34.974 ms
 4  te-4-2-ur03.santaclara.ca.sfba.comcast.net (68.85.191.9)  37.454 ms  13.988 ms  11.822 ms
 5  te-0-4-0-3-ar01.sfsutro.ca.sfba.comcast.net (68.87.226.122)  13.332 ms  17.473 ms  16.512 ms
 6  pos-3-0-0-0-cr01.sanjose.ca.ibone.comcast.net (68.86.90.93)  18.350 ms  19.463 ms  22.402 ms
 7  pos-0-1-0-0-pe01.529bryant.ca.ibone.comcast.net (68.86.87.2)  17.168 ms  15.704 ms  16.234 m
s
 8  66.208.228.226 (66.208.228.226)  16.080 ms  13.738 ms  15.677 ms
 9  72.14.232.136 (72.14.232.136)  16.542 ms  16.395 ms
    72.14.232.138 (72.14.232.138)  18.251 ms
10  209.85.250.66 (209.85.250.66)  19.343 ms
    209.85.250.60 (209.85.250.60)  31.734 ms
    209.85.250.63 (209.85.250.63)  17.711 ms
11  216.239.47.186 (216.239.47.186)  44.509 ms  39.825 ms  31.657 ms
12  72.14.233.138 (72.14.233.138)  35.075 ms  38.356 ms
    72.14.233.140 (72.14.233.140)  34.252 ms
13  64.233.174.127 (64.233.174.127)  34.072 ms
    64.233.174.97 (64.233.174.97)  33.404 ms  39.034 ms
14  * 216.239.46.6 (216.239.46.6)  57.757 ms
    216.239.46.22 (216.239.46.22)  36.094 ms
15  pz-in-f121.1e100.net (74.125.127.121)  37.314 ms  39.523 ms  40.487 ms
$
```

Here, you can see each hop. It took fifteen steps to get to www.udacity.com and the total time was about 39 milliseconds. You can see there are several different times (see circled figures) because the application is doing multiple tests and the time might vary from test to test.

Although the **tracert** is to www.udacity.com, from the first line of the **tracert** you can see that the site is being served by a server at google.com.

```
tracert to ghs.l.google.com (74.125.127.121), 64 hops max, 52 byte packets
```

The first hop shows the ip 192.168.3.1. This is an Internet address that refers to your router. You will always start there.

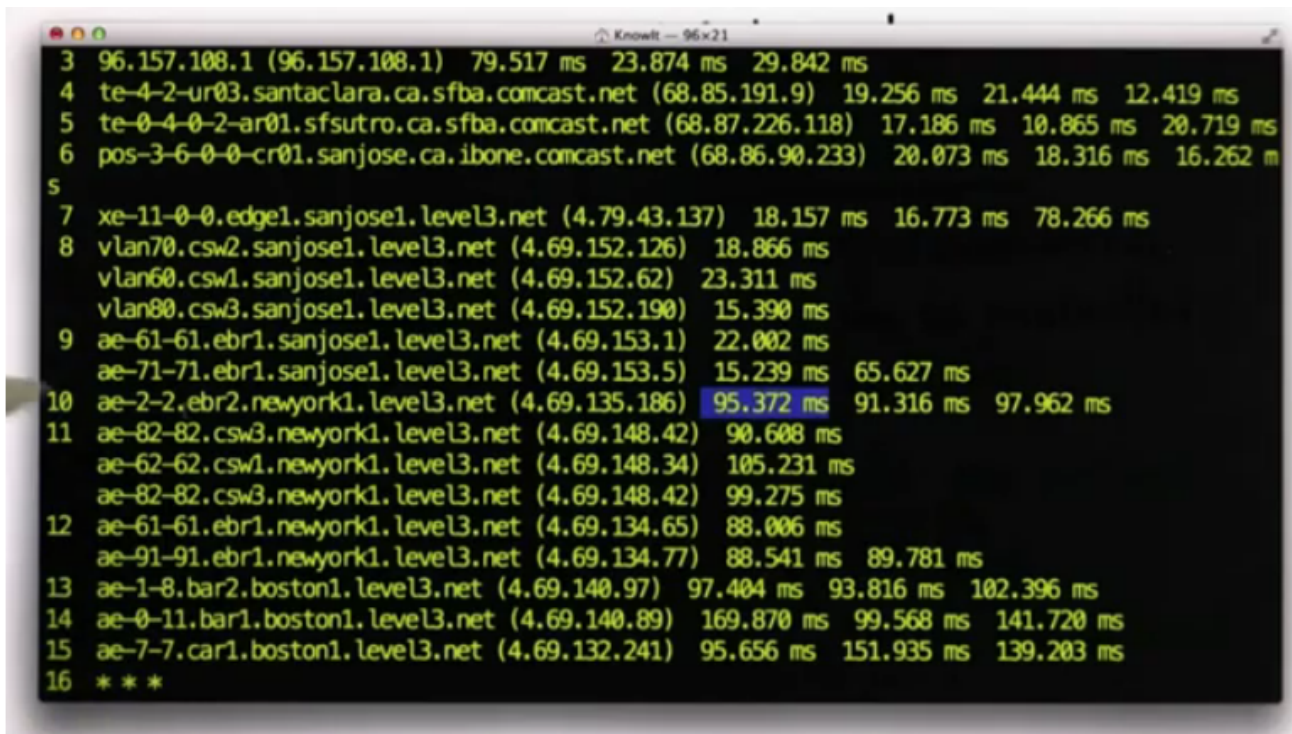
```
1  192.168.1.1 (192.168.1.1)  1.570 ms  1.284 ms  1.152 ms
```


After that you can see all the time amounts it takes to get to all the different sites. For example, to get to the Comcast site in Santa Clara took between 11.8 and 37.5 milliseconds.

```
4 te-4-2-ur03.santaclara.ca.sfba.comcast.net (68.85.191.9) 37.454 ms 13.988 ms 11.822 ms
```

Each hop on the way took a few milliseconds, taking a about 40 milliseconds for 15 hops.

The next hop is somewhere further away at MIT, which is in Boston, Massachusettes, the other side of the country from Palo Alto, California. To do this, the command is `tracert www.mit.edu`.



```
3 96.157.108.1 (96.157.108.1) 79.517 ms 23.874 ms 29.842 ms
4 te-4-2-ur03.santaclara.ca.sfba.comcast.net (68.85.191.9) 19.256 ms 21.444 ms 12.419 ms
5 te-0-4-0-2-ar01.sfsutro.ca.sfba.comcast.net (68.87.226.118) 17.186 ms 10.865 ms 20.719 ms
6 pos-3-6-0-0-cr01.sanjose.ca.ibone.comcast.net (68.86.90.233) 20.073 ms 18.316 ms 16.262 ms
7 xe-11-0-0.edge1.sanjose1.level3.net (4.79.43.137) 18.157 ms 16.773 ms 78.266 ms
8 vlan70.csw2.sanjose1.level3.net (4.69.152.126) 18.866 ms
  vlan60.csw1.sanjose1.level3.net (4.69.152.62) 23.311 ms
  vlan80.csw3.sanjose1.level3.net (4.69.152.190) 15.390 ms
9 ae-61-61.ebr1.sanjose1.level3.net (4.69.153.1) 22.002 ms
  ae-71-71.ebr1.sanjose1.level3.net (4.69.153.5) 15.239 ms 65.627 ms
10 ae-2-2.ebr2.newyork1.level3.net (4.69.135.186) 95.372 ms 91.316 ms 97.962 ms
11 ae-82-82.csw3.newyork1.level3.net (4.69.148.42) 90.608 ms
  ae-62-62.csw1.newyork1.level3.net (4.69.148.34) 105.231 ms
  ae-82-82.csw3.newyork1.level3.net (4.69.148.42) 99.275 ms
12 ae-61-61.ebr1.newyork1.level3.net (4.69.134.65) 88.006 ms
  ae-91-91.ebr1.newyork1.level3.net (4.69.134.77) 88.541 ms 89.781 ms
13 ae-1-8.bar2.boston1.level3.net (4.69.140.97) 97.404 ms 93.816 ms 102.396 ms
14 ae-0-11.bar1.boston1.level3.net (4.69.140.89) 169.870 ms 99.568 ms 141.720 ms
15 ae-7-7.car1.boston1.level3.net (4.69.132.241) 95.656 ms 151.935 ms 139.203 ms
16 * * *
```

Getting to udacity.com was not very far, geographically. To get to MIT, data travels through Santa Clara, to San Jose and then across the country to New York and finally to Boston. The locations are marked on this Google map: <http://g.co/maps/u9ssk>.

You can see that Palo Alto, Santa Clara and San Jose are all on the West Coast of the USA and then there is a huge jump across the country to New York and Boston on the East Coast.

Sometimes from the host names you can guess the locations and sometimes you cannot. Unlike with the Ancient Greek fires, where the distance between hops was limited, the distance between hops on the Internet can be thousands of miles. There is probably a fiber optic cable between San Jose and New York and therefore no need for any routing directions between them -- the data just goes straight through. The big time distance is from San Jose to New York, as highlighted on the screenshot above.

In total it took the packets about 100 ms to get across the country.

The three asterisks (*) at the bottom of the **traceroute** application's output indicate the packet is not actually getting to the final destination because there is no response from the web server at MIT. When this happens, you can try to repeat the **traceroute** with different time out options. If you want to terminate the **traceroute** before it's finished, you can use **ctrl-c**.

It takes longer to run **traceroute** than it does to send a message since the application is sending many messages and trying to find out all the points on the way.

About the furthest place from Palo Alto is Madagascar. Below is the output from **traceroute** www.mairie-antananarivo.mg

The packet starts in Palo Alto, then goes from Santa Clara to Dallas, and then through many servers from the same company.

```
Knowit - 96x21
1 192.168.1.1 (192.168.1.1) 1.490 ms 1.304 ms 1.144 ms
2 10.1.10.1 (10.1.10.1) 1.751 ms 1.774 ms 1.606 ms
3 96.157.108.1 (96.157.108.1) 37.883 ms 32.633 ms 30.144 ms
4 te-4-2-ur03.santaclara.ca.sfba.comcast.net (68.85.191.9) 17.342 ms 15.213 ms 13.077 ms
5 te-1-10-0-5-ar01.sfsutro.ca.sfba.comcast.net (68.85.155.58) 18.533 ms 17.581 ms 19.159 ms
6 te-4-0-0-4-cr01.dallas.tx.ibone.comcast.net (68.86.90.149) 19.326 ms 14.483 ms 16.329 ms
7 pe-0-3-0-0-pe01.529bryant.ca.ibone.comcast.net (68.86.87.142) 18.952 ms 16.333 ms 16.132 ms
...
8 te4-7.mpd01.sjc04.atlas.cogentco.com (154.54.11.185) 216.503 ms 203.360 ms 211.722 ms
9 te0-0-0-6.mpd22.sfo01.atlas.cogentco.com (154.54.28.81) 16.552 ms
  te0-1-0-7.mpd21.sfo01.atlas.cogentco.com (154.54.5.37) 15.274 ms
  te0-0-0-6.mpd21.sfo01.atlas.cogentco.com (154.54.2.165) 17.892 ms
10 te0-4-0-3.mpd22.mci01.atlas.cogentco.com (154.54.45.85) 63.339 ms
  te0-3-0-3.mpd22.mci01.atlas.cogentco.com (154.54.7.225) 59.774 ms
  te0-1-0-3.mpd22.mci01.atlas.cogentco.com (154.54.24.106) 60.102 ms
11 te0-4-0-3.mpd22.ord01.atlas.cogentco.com (154.54.30.177) 77.186 ms
  te0-3-0-2.mpd21.ord01.atlas.cogentco.com (154.54.7.137) 78.872 ms
  te0-2-0-3.mpd21.ord01.atlas.cogentco.com (154.54.30.106) 78.425 ms
12 te0-5-0-6.ccr21.bos01.atlas.cogentco.com (154.54.45.250) 168.038 ms
  te0-2-0-5.ccr21.bos01.atlas.cogentco.com (154.54.42.250) 186.106 ms
  te0-4-0-9.ccr21.bos01.atlas.cogentco.com (154.54.44.250) 186.106 ms
```

```
KnowIt - 96x21
12 te0-5-0-6.ccr21.bos01.atlas.cogentco.com (154.54.45.250) 168.038 ms
   te0-2-0-5.ccr21.bos01.atlas.cogentco.com (154.54.42.250) 186.106 ms
   te0-0-0-5.ccr21.bos01.atlas.cogentco.com (154.54.28.114) 182.315 ms
13 te0-1-0-2.ccr21.lon13.atlas.cogentco.com (154.54.1.94) 168.321 ms
   te0-0-0-2.ccr21.lon13.atlas.cogentco.com (154.54.5.162) 176.444 ms
   te0-4-0-3.ccr21.lpl01.atlas.cogentco.com (154.54.44.194) 169.850 ms
14 154.54.60.101 (154.54.60.101) 170.542 ms
   154.54.60.45 (154.54.60.45) 198.435 ms
   te0-4-0-4.mpd21.lon13.atlas.cogentco.com (154.54.37.174) 170.146 ms
15 te0-4-0-0.ccr21.mrs01.atlas.cogentco.com (130.117.0.153) 184.718 ms
   te0-3-0-5.ccr21.mrs01.atlas.cogentco.com (154.54.59.238) 186.853 ms
   te0-2-0-5.ccr21.mrs01.atlas.cogentco.com (154.54.59.234) 185.577 ms
16 te1-1.ccr01.nce02.atlas.cogentco.com (130.117.3.158) 180.841 ms
   te0-3-0-0.ccr21.mrs01.atlas.cogentco.com (154.54.59.226) 183.394 ms
   te0-2-0-0.ccr21.mrs01.atlas.cogentco.com (154.54.36.185) 179.543 ms
17 te1-1.ccr01.nce02.atlas.cogentco.com (130.117.3.158) 180.494 ms 184.437 ms
   te2-2.ccr01.nce03.atlas.cogentco.com (154.54.59.46) 191.302 ms
18 149.6.56.30 (149.6.56.30) 195.841 ms 193.635 ms 197.468 ms
19 149.6.56.30 (149.6.56.30) 195.605 ms 196.506 ms 195.620 ms
20 int-bd2.abconnect.net (217.16.0.2) 195.313 ms * *
```

Finally, it escapes from Texas at hop 18 and arrives in Washington DC (149.6.56.30), where it leaves the USA to head for France (217.16.0.2) at hop 19. It already takes nearly twice the time it took to get to MIT without having reached Madagascar yet. (Note: the information about the locations come from <http://whatismyipaddress.com> and may not be accurate.)

Q-10: Quiz (Traveling Data)

We observed that the latency between Palo Alto and Cambridge, MA is 100 milliseconds. The distance between Palo Alto, CA and Cambridge, MA is 4300 kilometers. At what fraction of the speed of light ($\sim 300\,000$ km/s) did my data travel?

a. $\frac{1}{7000}$

b. $\frac{1}{70}$

c. $\frac{1}{3}$

d. $\frac{1}{700}$

e. $\frac{1}{7}$

f. 0.8

[Answer to Q-10](#)

Making a Network (Making a Network)

Several different things are needed in order to make a network like this work.

1. A Way to encode and interpret messages

Greeks: “Agamennon is arriving” specific smoke signal

Internet: message bits electrons/photons

Any message can be encoded in bits and then the bits can be encoded on the wire. How that encoding actually works is pretty complicated and it’s not something that will be covered in this class. There are lots of different ways to do it.

2. A way to route messages

Greeks: directing smoke signals

Internet: routers figure out next hops

For all the routers along the path, a message comes in and the router has to decide where to send it on. Maybe the router has a table saying where to send it next. Nirvada is on the way to Boston from California, so it could send the message on to Nirvada from California, but that wasn’t what happened in the traceroute. The message was sent to San Jose where there is a big strong pipe straight across the country to New York and from there the message was sent on to Boston.

Routing is a challenging problem which won’t be covered in more detail in this class.

3. Rules for deciding who gets to use resources

Greeks: generals have priority

Internet: best effort service

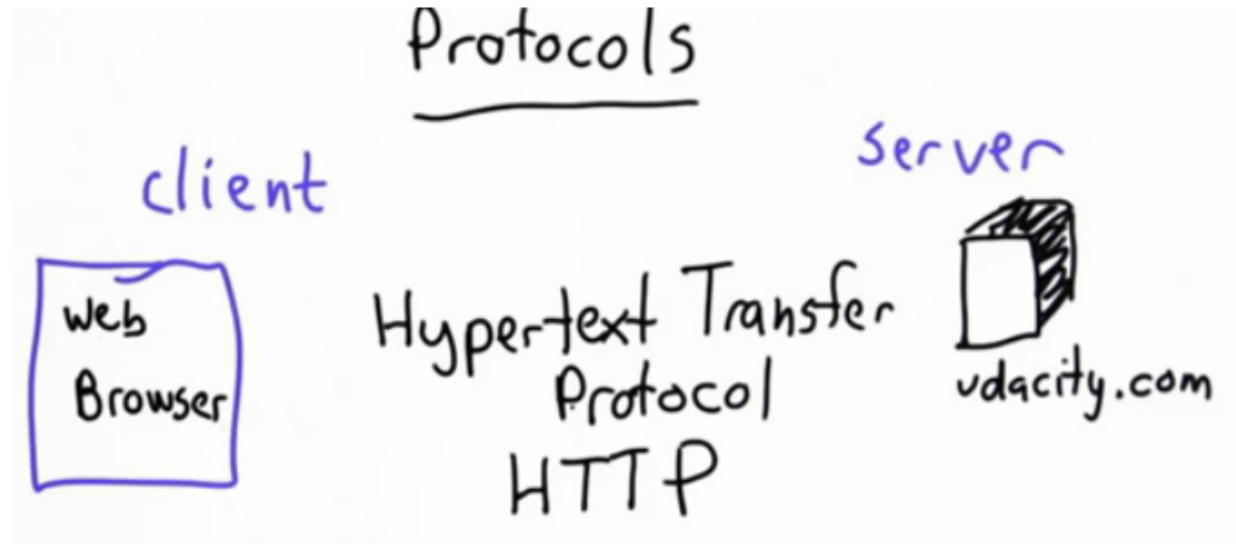
Unlike with the Greeks, there are no real rules on the Internet for who gets resources. It’s much more of a wild west where everywhere along the network gets to decide on its own what rules to apply. It’s called a best effort service. If two messages need to be sent by a router at the same time, the router decides which one to send on. This means that your package might get dropped. There is no guarantee that a package will reach its destination on the Internet.

This class won’t cover anything more on routing and rules, so you are encouraged to take a future networking class to learn more about these. However, we will look a little bit more about how messages work for the web.

Protocols (Protocols)

A **protocol** is a set of rules, that people agree to, which determines how two entities can talk to each other.

For the web, the protocol gives rules about how a client and a server talk to each other. The client is the web browser and the server is the web server.

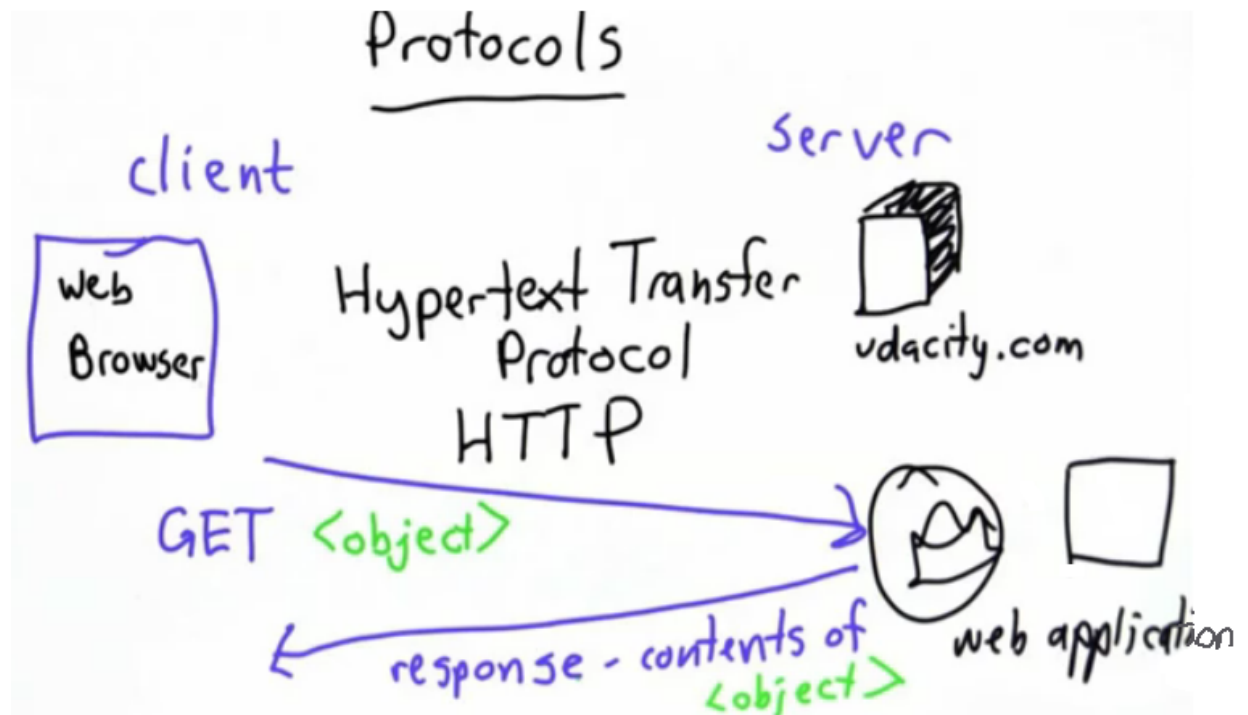


The protocol used on the web is called Hypertext Transfer Protocol which is abbreviated as HTTP. When you look in your browser, almost all the urls that you use start with http. That indicates that when the page is requested, the protocol to be used to talk to the server is this Hypertext Transfer Protocol. It's a very simple protocol, and only has two main messages. One of those messages is **GET**. The client can send a message to the server which says **GET** followed by the name of the object you want to get, **GET <object>**. That's all the client does. The python code for `get_page`,

```
def get_page(url):
    try:
        import urllib
        return urllib.urlopen(url).read()
    except:
        return ""
```

calls a library function `urllib.urlopen(url)` which actually does this.

After the client sends the message, the server will receive it. The server runs some code on it, finds the file that was requested, perhaps runs some more code, and then sends back a response with the contents of the requested **<object>**. That's the whole protocol.



If you click on a link in your browser, your web browser works out which url you are requesting and sends a **GET** message to the correct web server specified by that url. When it gets a response, it processes and then renders it.

If you'd like to learn more about what the server does, take the Web App course, [CS253 Web Application Engineering](#) and if you're interesting in how a web browser works, take the course [CS262: Programming Languages](#).

Conclusion

Hopefully you understand at a high level what a web browser does when requesting data over the Internet. There is nothing magic about it. The process is all about sending messages across the Internet and receiving responses which are text. That text is processed by a browser, or even by the web crawler you've programmed.

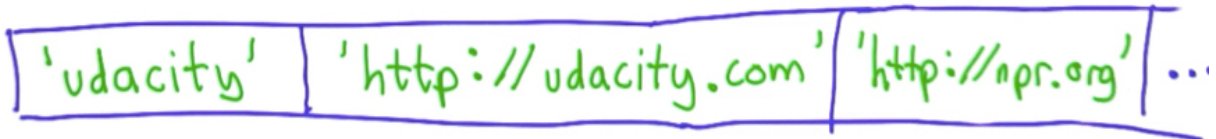
The search engine so far works but it isn't fast or smart. In unit 5, you'll look at how to make the search engine scale and respond to queries faster. In unit 6, you'll see learn how to find the best response to a query, that is, to respond with the best page rather than all the pages.

Answer Key

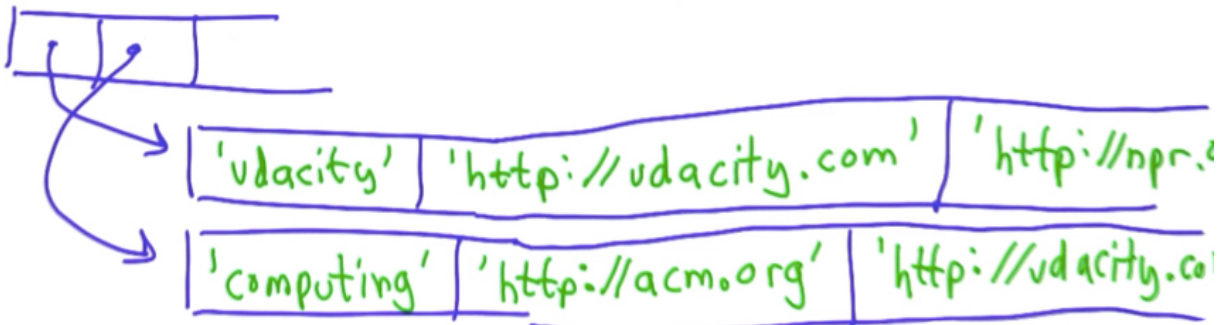
A-1: Answer (Data Structures)

The best choice is d, and b is a close runner up. Choices a and c would be really difficult. Here's why:

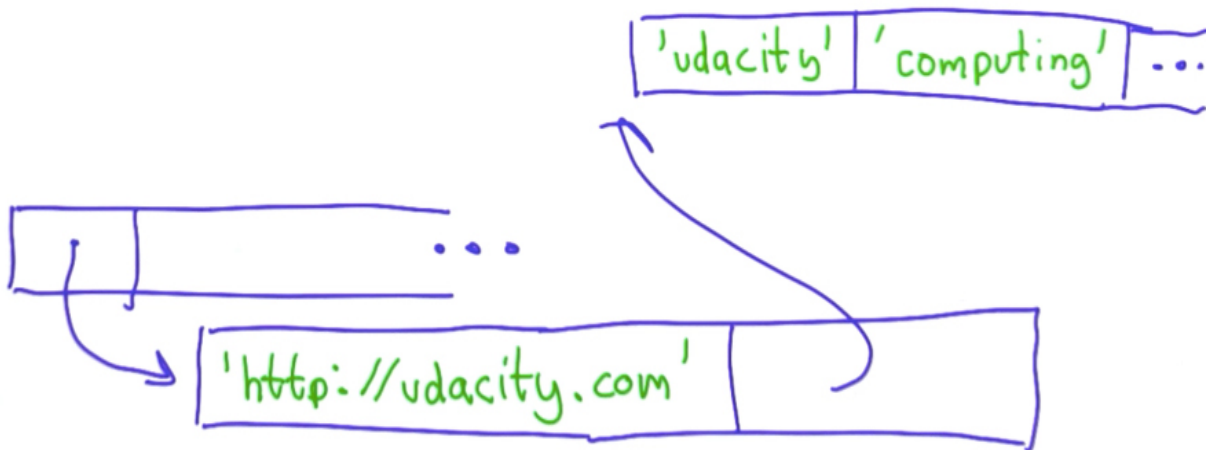
a - Hard to read because everything is all in one list. This structure will also make it hard to loop over the keywords.



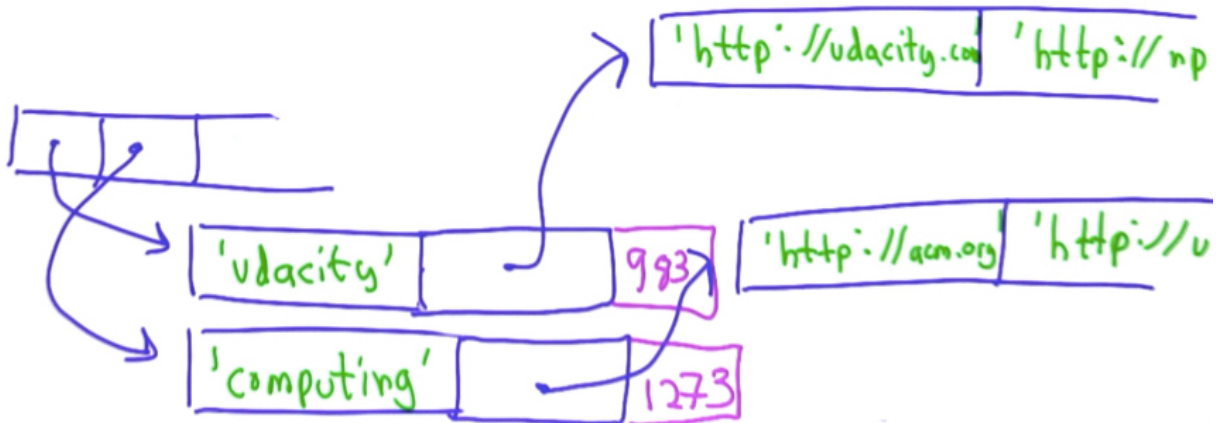
b - This option is okay, but not as good as d. The structure of b is such that there is a list, where each element of the list is a list, and the element lists are themselves lists. The big advantage to this structure is that it is easy to tell the keywords from the urls, the keyword is always the first element of the list. It is also easier to go through the keywords because for each list you just have to look to the first element, check to see if it is the keyword you are looking for, and if not, move on to the next one. The only downfall to this structure is that the distinction between the keywords and the urls could be even more clear.



c - While this has more structure, but it does not make it easy to look up the pages where a keyword appears. to look for a particular keyword you would have to look in each entry, look in the second part of the entry, scan it to see if the keyword appears, if it does then you want the url in your result, and if it doesn't then that url is not in the result. This option is not the best because you would have to scan through all the pages, which would take too much time.



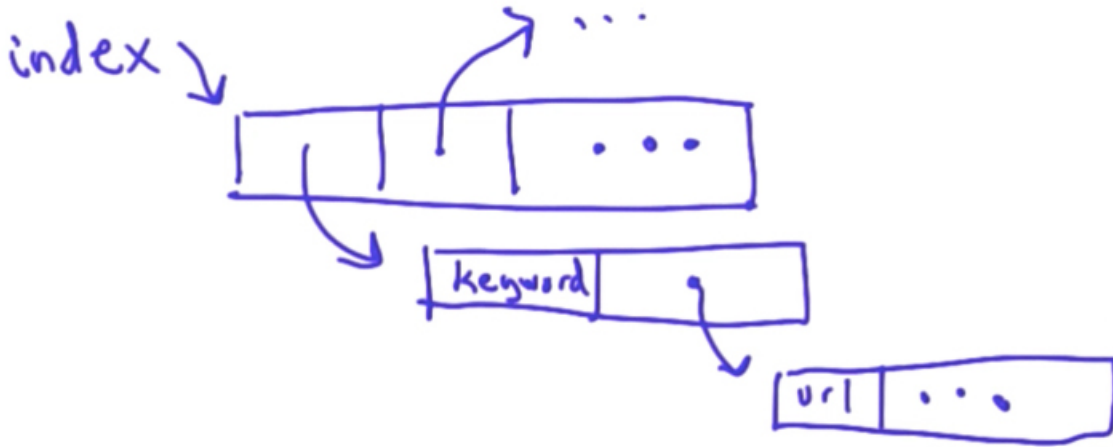
d- In this structure there are just two elements in the inner list; the keyword followed by a list of urls. This is the best option because it makes a very clear separation between the keyword and the list of urls. It also means that if you decide you want to keep track of something else, like the number of times someone searches for each keyword, you can easily do that by adding an extra element. This is something that would be more difficult to do with the structure of option



b.

A-2: Answer (Add to Index)

Here is an image of the data structure that you are looking to write a procedure for:



```
def add_to_index(index, keyword, url):
    for entry in index: # loop through the elements of index, giving
        each one the name entry
        if entry [0] == keyword: # test to see if the value at position
            0 of entry identical to the keyword that's passed in
            entry[1].append(url) # if it is equal you want to append the
            url to the list of urls associated with that entry
            return # stop the loop
    #not found, add new entry
```

If the keywords are not found, you want to add a new entry. The new entry is going to have as its value a list containing two elements, the keyword and the second element will be a list containing the urls that were found that have that keyword. So far, there is only one url that was passed in to **index**. To do this, add to your code:

```
def add_to_index(index, keyword, url):
    for entry in index:
        if entry [0] == keyword:
            entry[1].append(url)
            return # stop the loop
    index.append([keyword, [url]])
```

Let's test this procedure, starting with an empty list for the **index**.

```
index = []
add_to_index('udacity', 'http://udacity.com')

print index
Traceback (most recent call last):
  File "input/test.py", line 11, in <module>
    add_to_index('udacity', 'http://udacity.com')
TypeError: add_to_index() takes exactly 3 arguments (2 given)
```

According to the error message, you have to pass in three parameters, index, keyword and url. So, go ahead and add index in there.

```
index = []
add_to_index(index, 'udacity', 'http://udacity.com')

print index
[['udacity', ['http://udacity.com']]]
```

Try the next one:

```
index = []
add_to_index(index, 'udacity', 'http://udacity.com')
add_to_index(index, 'computing', 'http://acm.org')

print index
[['udacity', ['http://udacity.com']], ['computing', [http://acm.org]]]
```

Now, try adding a new url to an already existing keyword:

```
index = []
add_to_index(index, 'udacity', 'http://udacity.com')
add_to_index(index, 'computing', 'http://acm.org')
add_to_index(index, 'udacity', 'http://npr.org')

print index
[['udacity', ['http://udacity.com', 'http://npr.org']], ['computing',
http://acm.org]]
```

A-3: Answer (Lookup)

```
def lookup(index, word): # two parameters
    for entry in index: # loop through the entries in the index
        if entry [0] == keyword:
            return entry [1] # if found then return the urls associated
    with that entry
    return [] # if no keyword is found, then return empty list
```

Try this in the interpreter:

```
def lookup(index, word):
    for entry in index:
        if entry [0] == keyword:
            return entry [1]
    return []

index = []
add_to_index(index, 'udacity', 'http://udacity.com')
add_to_index(index, 'computing', 'http://acm.org')
add_to_index(index, 'udacity', 'http://npr.org')

print lookup(index, 'udacity')
['http://udacity.com', 'http://npr.org']
```

Try looking up something that is not in the list:

```
print lookup(index, 'audacity')
[]
```

A-4: Answer (Add Page to Index)

The goal is to define a procedure, `add_page_to_index`, which takes in three inputs, the `index`, the `url` and the content at that url. This takes two steps:

1. split the page into its component words,
2. then add each word along with the url to the index.

The `split` method divides the content into its component words, while the procedure, `add_to_index` adds a word and a url to the index. However, this is not sufficient to satisfy the second step. You still need to do this for each of the words found during the first step. A `for` loop will take care of this. Remember that the point of this procedure is to modify the index, so there will be no return results.

This is how you can put the code together using this structure. Use the `split` method to divide the content into its component words and store them in the `words` variable.

```
def add_page_to_index(index,url,content):  
    words=content.split()
```

Go through each of the words in `words`, which we can do using a `for` loop, and naming the variable `word`:

```
def add_page_to_index(index,url,content):  
    words=content.split()  
    for word in words:
```

adding each of the words to the index by calling `add_to_index` and passing in the `index`, the `word` and the `url`.

```
def add_page_to_index(index,url,content):  
    words=content.split()  
    for word in words:  
        add_to_index(index, word, url)
```

Note that if a word occurs more than once on the same page, we're going to keep adding it to the index each time. This means there might be more than one occurrence of the same url in the list of urls associated with a keyword. Depending on what we want our search engine to do and how we want our it to respond to queries, this may or may not be a good thing. It will be discussed further in a later class.

To try this code in the python interpreter, you'll need the three procedures from earlier. These are printed below, followed by a reminder of what each procedure does and then some examples you can try out for yourself.

```
def add_to_index(index,keyword,url):
    for entry in index:
        if entry[0] == keyword:
            entry[1].append(url)
            return
    #not found, add a new entry
    index.append([keyword,[url]])

def lookup(index,keyword):
    for entry in index:
        if entry[0]==keyword:
            return entry[1]
    return []

def add_page_to_index(index,url,content):
    words=content.split()
    for word in words:
        add_to_index(index, word,url)
```

First, the **add_to_index** procedure takes in the **index**, a **keyword** and a **url**. It goes through the entries in the **index**, checking to see if it already contains the **keyword**. If it does, it appends the **url** to the entry that matches the **keyword**. If the **keyword** is not found, the procedure adds a new entry to the index that is a list of the **keyword** and the single **url**: **[keyword, [url]]**.

Next, the **lookup** procedure takes an **index** and a **keyword**. It looks at each entry **[<keyword>, <list of urls>]** in the index to see if the **keyword** you're looking for is at the first position of the entry. If it is, it returns the list of urls which corresponds to that **keyword**.

Finally, add on the **add_page_to_index** procedure you just defined.

Examples:

1. Here is the same example from before:

```
index=[]
add_page_to_index(index,'fake.test',"This is a test")
print index

index=[]
add_page_to_index(index,'not.test',"This is not a test")
print index
[['This', ['fake.test']], ['is', ['fake.test']], ['a', ['fake.test']],
 ['test', ['fake.test']]]
[['This', ['not.test']], ['is', ['not.test']], ['not', ['not.test']],
 ['a', ['not.test']], ['test', ['not.test']]]
```

2. To convince yourself that the code is working, here's a more complex example using a quotation from Scott Adams on dilbert.com, followed by one from [Randy Pausch](http://randy.pausch.com).

```
index=[]
add_page_to_index(index, 'http://dilbert.com',
    """
    Another strategy is to ignore the fact that you
    are slowly killing yourself by not sleeping and
    exercising enough. That frees up several hours a
    day. The only downside is that you get fat and die.
    --- Scott Adams on Time Management
    """)
add_page_to_index(index, 'http://randy.pausch',
    """
    Good judgement comes from experience, experience
    comes from bad judgement. If things aren't going
    well it probably means you are learning a lot and
    things will go better later.
    --- Randy Pausch
    """)
print index
[['Another', ['http://dilbert.com']], ['strategy', ['http://dilbert.com']],
 ['is', ['http://dilbert.com', 'http://dilbert.com']], ['to', ['http://dilbert.com']],
 ['ignore', ['http://dilbert.com']], ['the', ['http://dilbert.com']],
 ['fact', ['http://dilbert.com']], ['that', ['http://dilbert.com', 'http://dilbert.com']],
 ['you', ['http://dilbert.com', 'http://dilbert.com', 'http://randy.pausch']],
 ['are', ['http://dilbert.com', 'http://randy.pausch']], ['slowly', ['http://dilbert.com']],
 ['killing', ['http://dilbert.com']], ['yourself', ['http://dilbert.com']],
 ['by', ['http://dilbert.com']], ['not', ['http://dilbert.com']],
 ['sleeping', ['http://dilbert.com']], ['and', ['http://dilbert.com', 'http://randy.pausch']],
 ['exercising', ['http://dilbert.com']], ['enough.', ['http://dilbert.com']],
 ['That', ['http://dilbert.com']], ['frees', ['http://dilbert.com']],
 ['up', ['http://dilbert.com']], ['several', ['http://dilbert.com']],
 ['hours', ['http://dilbert.com']], ['a', ['http://dilbert.com', 'http://randy.pausch']],
 ['day.', ['http://dilbert.com']], ['The', ['http://dilbert.com']], ['only', ['http://dilbert.com']],
 ['downside', ['http://dilbert.com']], ['get', ['http://dilbert.com']],
 ['fat', ['http://dilbert.com']], ['die.', ['http://dilbert.com']], ['---',
 ['http://dilbert.com', 'http://randy.pausch']], ['Scott', ['http://dilbert.com']],
 ['Adams', ['http://dilbert.com']], ['on', ['http://dilbert.com']],
 ['Time', ['http://dilbert.com']], ['Management', ['http://dilbert.com']],
 ['Good', ['http://randy.pausch']], ['judgement', ['http://randy.pausch']],
 ['comes', ['http://randy.pausch', 'http://randy.pausch']],
 ['from', ['http://randy.pausch', 'http://randy.pausch']],
 ['experience,', ['http://randy.pausch']], ['experience', ['http://randy.pausch']],
 ['bad', ['http://randy.pausch']], ['judgement.', ['http://randy.pausch']],
 ['If', ['http://randy.pausch']], ['things',
```

```
['http://randy.pausch', 'http://randy.pausch']], ["aren't", ['http://randy.pausch']], ['going', ['http://randy.pausch']], ['well', ['http://randy.pausch']], ['it', ['http://randy.pausch']], ['probably', ['http://randy.pausch']], ['means', ['http://randy.pausch']], ['learning', ['http://randy.pausch']], ['lot', ['http://randy.pausch']], ['will', ['http://randy.pausch']], ['go', ['http://randy.pausch']], ['better', ['http://randy.pausch']], ['later.', ['http://randy.pausch']], ['Randy', ['http://randy.pausch']], ['Pausch', ['http://randy.pausch']]]
```

It's quite big as there were a lot of words in the two quotes. Note that some words appear in both lists and some in just one.

To take a closer look at what is going on, you can try the following queries (It's probably best to remove the **print index** first so you can see the results more clearly):

```
print lookup(index, 'you')
['http://dilbert.com', 'http://dilbert.com', 'http://randy.pausch']
```

In this result you see that **'you'** occurred twice in the dilbert.com quote, but just once in the Randy Pausch quote.

```
print lookup(index, 'good')
[]
```

The word **'good'** does not appear in either quote.

```
print lookup(index, 'bad')
['http://randy.pausch']
```

Bad appears in just the Randy Pausch quote.

Using this code you can look up words in your index and get the urls where they are found. You can add pages to your index and record all the words in that page in the location they occur. The one thing left to do is to connect the code here with the code for crawling the web.

A-5 Answer (Finishing the Web Crawler)

```
add_page_to_index(index, page, content)
```

A-6: Answer (Networks)

The answer is d, over 3000 years old

A-7: Answer (Latency)

Three of the four answers will reduce latency.

a. CORRECT: Make the signalling nodes further apart, so it takes fewer hops. This would reduce latency as at each hop, there is a time delay as the troops have to see the fire and make their own fire. Fewer hops would mean the message could get across faster.

b. CORRECT: Threaten the soldiers to make the signalling fires quicker. Making the troops work faster would reduce the delay between the receiving and sending on of the messages, which would in turn increase the speed across the network, and hence reduce the latency.

c. INCORRECT: Add colours to the smoke so there are more messages per signal. This will increase the amount of information that can be sent across a network, but it will not increase the speed at which that information is sent. The amount of information which can be sent is called the bandwidth.

d. CORRECT: Increase the speed of light.

This would increase the speed that information travels between hops. It wouldn't make much difference overall as the time spent travelling between hops is very small compared with the delays at hops.

A-8: Answer (Bits)

2

A-9: Answer (What Is Your Bandwidth)

There is no correct (or wrong) answer to this quiz.

A-10: Answer (Traveling Data)

The answer is e, 17, which is pretty good. There aren't many things you can think of in terms of a fraction of the speed of light of over 10%.

The data travelled across the country at about 43,000 km/s. Most of that time taken is not the time on the wire. The speed through optical fiber is about 50% slower than the speed light travels in a vacuum. What takes up most of the time is all the routers the data has to go through. The traceroute showed that it went through about 20 routers. Each router had to take a packet in, figure out where to send it, and then send it on. All this happened in the 100 ms it took for the data to travel across the country.

The python code to calculate this is

```
distance = 4300. #km between Palo Alto and Cambridge
c=300000        #km/s (approximate speed of light)
time = 0.1      # 100 ms = 0.1

light_time = distance / c #the time in seconds it would take
                    #light to travel across the country

print light_time
print time/light_time
0.014333333333333
6.97674418605
```

The output 6.97674418605 printed above is the fraction (time for the data to travel) / (travel time at the speed of light). The question asked for (travel time at the speed of light) / (time it took for the data to travel). This is simply one over the value calculated, giving the answer 17.

With a direct vacuum tunnel from Palo Alto to Cambridge, the data would have travelled at about seven times the speed of the current Internet, taking only around 0.014 seconds, which is $0.014 \times 1000 = 14$ ms. Unfortunately, not many people can afford a direct vacuum tunnel between the two places they want to send data!

With optic fiber and no routers along the route, the journey would have taken about 28 ms (since it takes twice as long as light in a vacuum). This means that the time lost at routers is about 72 ms.



