## CS432 Assignment3

1. While researching the how to handle large files of text I came across <a href="http://stackoverflow.com/questions/1131220/get-md5-hash-of-big-files-in-python">http://stackoverflow.com/questions/1131220/get-md5-hash-of-big-files-in-python</a>. Initially I didn't know of the hashlib or md5, but upon researching those I came across <a href="https://docs.python.org/2/library/hashlib.html">https://docs.python.org/2/library/hashlib.html</a> which further informed me of how to work with the library and its resources. I constructed a python scrip by merging my original draft to process the raw data(shown below) along with a block of code I found from the second link(also shown below).

## Original code

```
import sys
import fileinput
import commands

import os

rinf(sys.version)

RAW DATA
input = open('test.txt', 'r')

for line in input:

url = line

print(line)

counter_1 = 0

for line in input:

filename = str(counter_1)
print str(line)

curling = 'curl ' + '"' + str(line).rstrip('\n') + '"' + '> ./rawData/' + filename

print curling
output_1 = commands.getoutput(curling)

input.close()
```

## Code from link

```
>>> import hashlib
>>> m = hashlib.md5()
>>> m.update("Nobody inspects")
>>> m.update(" the spammish repetition")
>>> m.digest()
'\xbbd\x9c\x83\xdd\x1e\xa5\xc9\xd9\xde\xc9\xa1\x8d\xf0\xff\xe9'
>>> m.digest_size
16
>>> m.block_size
64
```

## CS432 Assignment3

The end result is-

```
imput = open("URIfinal.txt", 'r')

for line in input:
    url = line
    url = url.replace('\n', '')

# referenced from http://stackoverflow.com/questions/1131220/get-md5-hash-of-big-files-in-python
def convert_to_md5(content):
    m = hashlib.md5()
    m.update(content)
    return m.hexdigest()

output = convert_to_md5(url)
print output

# % wget - 0 www.cnn.com http://www.cnn.com/provided on assignment prompt
os.system(" wget -0 /bcarey/HTML" + output + ".txt" + url)
```

This document simply downloads the raw content of each URL.

I simply copied this file into a separate .py file, and altered the final line with the lynx command provided on the prompt as well to strip away everything except for the text file

```
% lynx -dump -force_html www.cnn.com > www.cnn.com.processed
```

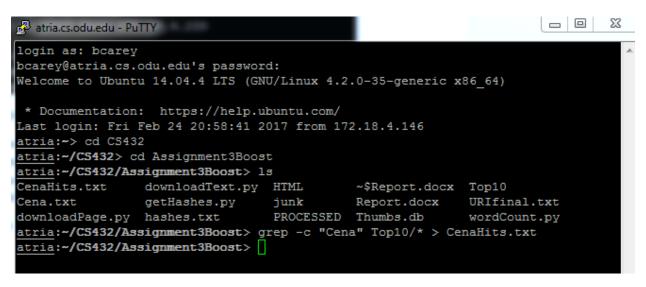
Both files(page\_content.py & text\_only.py) were ran through Putty. It took roughly 30 minutes for both processes to be completed.

2. For question two I had to review grep, as I had not used it since CS252. I found tutorials on grep at <a href="http://tldp.org/LDP/Bash-Beginners-Guide/html/sect\_04\_02.html">http://tldp.org/LDP/Bash-Beginners-Guide/html/sect\_04\_02.html</a> & <a href="http://www.computerhope.com/unix/ugrep.htm">http://www.computerhope.com/unix/ugrep.htm</a> . Once I did that it became fairly easy.

```
-I,
--files-with-matches

Instead of the normal output, print the name of each input file from which output would normally have been printed. The scanning will stop on the first match.
```

I used the previous command from commputerhope.com and dumped all matches into a text file. My query term was "Cena" as in John Cena. I then used grep again to yield the individual word counts for Cena. I placed all files in my Top 10 folder, and grep –c for the ones with the most hits. I removed all other files from my Top10 folder, all that remains are the 10 pages with the highest occurrence of Cena and the Top10.txt with their names. This was done by-



All generated .txt files can be found inside of my Q2 folder.

3.	For	part three I chose to use http://www.prchecker.info/check page rank.php					
		Web Page URL:	http://www.muscleandfitness.com/athletes-celebrities/news/john-cena- no-reason-bring-back-attitude-era				
		The Page Rank:	0/10				
		(the page rank value is <b>0</b> from 10 possible points)					
		Web Page UR	L: http://wwe.com				
		The Page Ran	6/10				
		(the page rank value is 6 from 10 possible points)					
		Web Page UR	L: http://bleacherreport.com/wwe				
		The Page Ran	k: 4/10				
		(the page ran	k value is <b>4</b> from 10 possible points)				

Web Page URL:	http://www.skysports.com/wwe/news/14203/10760037/ww smackdown-john-cena-beats-randy-orton-thanks-to-luke-har
The Page Rank:	0/10
(the page rank va	alue is <b>0</b> from 10 possible points)
Web Page URL:	http://www.cagesideseats.com/
The Page Rank:	6/10
(the page rank v	value is <b>6</b> from 10 possible points)
Web Page URL:	http://www.wrestlinginc.com/wi/news/2015/0505/594002/tripl comments-on-sami-zayn-vs-john-cena/
The Page Rank:	0/10
(the page rank val	ue is <b>0</b> from 10 possible points)
Web Page URL:	http://www.foxsports.com/wwe/story/john-cena-promises-fanditching-wwe-hollywood-122716
The Page Rank:	0/10
(the page rank val	lue is <b>0</b> from 10 possible points)
Web Page URL:	http://hollywoodlife.com/2016/08/21/aj-styles-beats-john-cesummerslam-2016-wwe-results/#!
_	

Web Page URL:	http://shop.wwe.com/						
The Page Rank:		5/10					
(the page rank va	alue is <b>5</b> from 10 possible points)						
Web Page URL:	https://www.washingtonpost.com/ne lead/wp/2016/04/04/shaquille-oneal appearances-at-wrestlemania-32/?ut	-john-cena-make-surprise-ring-					
The Page Rank:		0/10					
(the page rank value is <b>0</b> from 10 possible points)							
Web Page URL: http://www.wrestlinginc.com/wi/news/wwe-news/							
The Page Rank:		0/10					
(the page rank value is <b>0</b> from 10 possible points)							