Matthew Gomez CS 432 – Web Science Assignment 1 1/28/16

Program Narrative

For the first assignment of this course, we were tasked to write a program written in Python, which should be able to read any URI from and able to extra the links present on the page. Redirections are also required to be shown via command line if such a task is present. The program is also required to pull all PDF files laying around in the links alongside their content length; bytes size of the files.

The program consists of two functions, first function titled "LinkExtractor", extracts the links directly from the page the specific URI directs to. The second function titled "LinkHandler" serves the purpose of getting a request and response while also handling the pdf extraction. As well as reporting their individual byte sizes.

PART 1

1. <u>Curl command</u>: <u>curl</u> --data "q=vt hokie" http://search.vt.edu/search/pages.html

```
🚱 linux.cs.odu.edu - PuTTY
                                   "q=vt hokie" http://search.vt.edu/search/pages.html
<html lang="en">
 chead>
cmeta http-equiv="Content-Type" content="text/html; charset=utf-8" />
cmeta http-equiv="X-UA-Compatible" content="IE=edge" />
cmeta name="viewport" content="width=device-width, initial-scale=1.0" />
ctitle>Virginia Tech | Search - Web</title>
clink rel="shortcut icon" type="image/x-icon" href="//www.assets.cms.vt.edu/image/s/favicon.ico" />
clink rel="stylesheet" href="/search/assets/css/base.css" type="text/css" media=
"screen" />
  screen" />
link rel="stylesheet" helf-"/search/assets/css/baset.css" type-"text/css" media-
screen" />
link rel="stylesheet" href="/search/assets/css/enhanced.css" type="text/css" me
ia-"screen" />
script type="text/javascript"
src="//www.assets.cms.vt.edu/jquery/archives/jquery-1.10.latest.min.js">

 script type="text/javascript" src="/search/assets/js/search_utils.js"></script>
<a href="#vt-skipto-search">Skip to Search</a><a href="#vt-skipto-results">Skip to Results</a>
            <div id="container">
                         <div id="header_container">
    <!-- BEGIN HEADER -->
<a href="http://www.givingto.vt.edu">Giving</a><a href="http://www.lib.vt.edu">Libraries</a><a href="http://maps.vt.edu">Maps & amp;
Locations</a><a href="http://www.vt.edu/az_index.html">A to</a>
    Index</a>
                        /

//ul>

/ul id="vt_student_tools">

/li>/a
                                                 href="https://banweb.banner.vt.edu/ssb/prod/twbk
                                     Spa</a><a href="https://scholar.vt.edu/portal">Scholar</a><
                        >
```

<!DOCTYPE html>

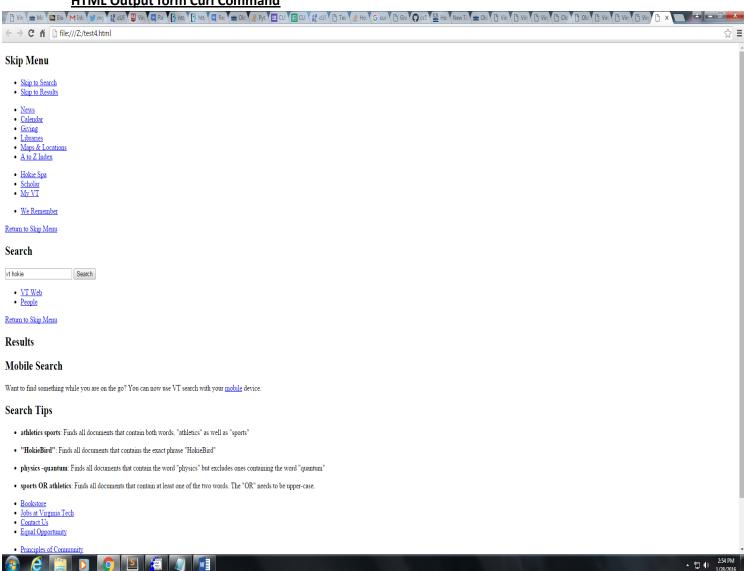
```
<script type="text/javascript" src="/search/assets/js/search utils.js"></script>
<script type="text/javascript" src="/search/assets/js/search_pages.js"></script>
</head>
<body>
<div class="vt_skip">
   <h2><a id="vt-skipto-menu">Skip Menu</a></h2>
   ul>
       <a href="#vt-skipto-search">Skip to Search</a>
       <a href="#vt-skipto-results">Skip to Results</a>
   </div>
    <div id="container">
       <div id="header container">
           <!-- BEGIN HEADER -->
<div id="header">
    <div id="vt_logo">
       <a id="vt_home_btn" title="Virginia Tech" href="http://www.vt.edu"></a>
   </div>
   <div id="vt utilities">
       ul id="vt toplinks">
           <a href="http://www.vtnews.vt.edu">News</a>
           <a href="http://www.calendar.vt.edu">Calendar</a>
           <a href="http://www.givingto.vt.edu">Giving</a>
           <a href="http://www.lib.vt.edu">Libraries</a>
           <a href="http://maps.vt.edu">Maps & amp;</a>
                   Locations</a>
           <a href="http://www.vt.edu/az_index/index.html">A to Z Index</a>
       ul id="vt student tools">
           <a
href="https://banweb.banner.vt.edu/ssb/prod/twbkwbis.P_WWWLogin">Hokie
                   Spa</a>
           <a href="https://scholar.vt.edu/portal">Scholar</a>
           <a href="https://my.vt.edu/">My VT</a>
       ul id="vt_we_remember">
           <a href="http://www.weremember.vt.edu">We Remember</a>
       </div>
</div>
<!-- END HEADER -->
       </div>
       <div id="content_container">
```

```
<div id="content">
<div class="vt_skip">
    <a href="#vt-skipto-menu">Return to Skip Menu</a>
    <h2><a id="vt-skipto-search">Search</a></h2>
</div>
                <div id="vt_search_block">
                    <form action="#" onSubmit="return executeQuery()" method="get"
name="vt_search_form" id="vt_header_search_form">
                        <input type="text" maxlength="50" placeholder="Search pages and
people" name="q"
                             value="vt hokie" id="vt_search_box" autocomplete="off"/>
                        <button id="vt_go_button">
                             <span class="vt skip">Search</span>
                        </button>
                    </form>
                </div>
                <div id="navigation">
                    ul>
                        class="current"><a href="#">VT Web</a>
                    <
                         <a
href="people.html;isessionid=0A02594B4BAC966B5EA2A1520BB386C6.mt-prod-4?q=vt+hokie"
id="vt-people-nav">People</a>
                    </div>
<div class="vt_skip">
    <a href="#vt-skipto-menu">Return to Skip Menu</a>
    <h2><a id="vt-skipto-results">Results</a></h2>
</div>
                <div id="results">
                    <div id="vt_gcse_script">
            <noscript>
            <div class="noscript">
                    It looks like you have JavaScript turned off. See search results
href="http://www.google.com?cx=012042020361247179657:wmrvw9b99ug&cof=FORID:11&ie
=UTF-8&q=vt hokie">here</a>.
                </div>
            </noscript>
            <div id="vt_gcse_results" class="gcse-searchresults-only" data-resultsetsize="7"</pre>
data-gname="vt_gcse_results">
            </div>
```

```
</div>
               </div>
               <div id="rb_content">
                   <h2>Mobile Search</h2>
                   >
                       Want to find something while you are on the go? You can now use VT
                       search with your <a href="/search/m">mobile</a> device.
                   <h2>Search Tips</h2>
                   ul>
                       >
                           >
                               <strong>athletics sports</strong>: Finds all documents that
                               contain both words, "athletics" as well as "sports"
                           >
                               <strong>"HokieBird"</strong>: Finds all documents that
contains
                               the exact phrase "HokieBird"
                           >
                               <strong>physics -quantum</strong>: Finds all documents
that
                               contain the word "physics" but excludes ones containing the
word
                               "quantum"
                           >
                           >
                               <strong>sports OR athletics</strong>: Finds all documents
that
                               contain at least one of the two words. The "OR" needs to be
                               upper-case.
                           </div>
```

```
</div>
        </div>
        <div id="footer container">
<!-- BEGIN FOOTER -->
<div id="footer">
    <l
        <a href="http://www.bookstore.vt.edu">Bookstore</a>
        <a href="http://www.jobs.vt.edu/">Jobs at Virginia Tech</a>
        <a href="http://www.vt.edu/contacts/">Contact Us</a>
        <a href="http://www.vt.edu/about/equal-opportunity.html">Equal</a>
                Opportunity</a>
   <a
            href="http://www.vt.edu/diversity/principles-of-community.html">Principles
                of Community</a>
        <a href="http://www.vt.edu/about/privacy.html">Privacy</a>
                Statement</a>
        <a href="http://www.vt.edu/about/acceptable-use.html">Acceptable</a>
                Use Policy</a>
        <a href="http://www.vt.edu/about/accessibility.html">Accessibility</a>
        © 2016 Virginia Polytechnic Institute and State University. <span id="version-
number">VT Search: 2.4.3</span>
</div>
<!-- END FOOTER -->
        </div>
   </div>
<script type="text/javascript">
var gaJsHost = (("https:" == document.location.protocol) ? "https://ssl." : "http://www.");
document.write(unescape("%3Cscript src="" + gaJsHost + "google-analytics.com/ga.js"
type='text/javascript'%3E%3C/script%3E"));
</script>
<script type="text/javascript">
var pageTracker = _gat._getTracker("UA-5217491-2");
pageTracker._trackPageview();
</script>
</body>
```

HTML Output form Curl Command



PART 2- URIS

http://www.cs.odu.edu/~mln/

```
linux.cs.odu.edu - PuTTY
  BeautifulSoup([your markup], "html5lib")
markup_type=markup_type))
Traceback (most recent call last):
   File "432Assignment1_Gomez.py", line 78, in <module>
      res=linkHandler(links.get('href'))
   File "432Assignment1_Gomez.py", line 44, in linkHandler response = urllib2.urlopen(request)
File "/usr/lib/python2.7/urllib2.py", line 127, in urlopen return _opener.open(url, data, timeout)
File "/usr/lib/python2.7/urllib2.py", line 396, in open
protocol = req.get_type()

File "/usr/lib/python2.7/urllib2.py", line 258, in get_type
raise ValueError, "unknown url type: %s" % self.__original
ValueError: unknown url type:
ValueError: unknown uri type:
sirius:~> python 432Assignment1_Gomez.py http://www.cs.odu.edu/~mln/
['432Assignment1_Gomez.py', 'http://www.cs.odu.edu/~mln/']
/usr/local/lib/python2.7/dist-packages/bs4/__init__.py:166: UserWarning: No parser was explicitly specifi
on another system, or in a different virtual environment, it may use a different parser and behave diffe
To get rid of this warning, change this:
  BeautifulSoup([your markup])
to this:
  BeautifulSoup([your markup], "html5lib")
markup_type=markup_type))
http://www.cs.odu.edu/~mln/ None
http://www.odu.edu None
http://www.cs.odu.edu/~mln/ None
http://www.cs.odu.edu/~mln/pubs/ None
http://www.cs.odu.edu/~mln/teaching/ None
http://www.cs.odu.edu/~mln/service/
http://www.larc.nasa.gov/
            http://www.nasa.gov/centers/langley/home/index.html
http://sils.unc.edu/ None
http://www.openarchives.org/pmh/ 6028
http://www.openarchives.org/ore/ 7265
http://www.mementoweb.org/guide/rfc/ID/ 207674
http://www.openarchives.org/rs/toc 4590
http://www.nsf.gov/awardsearch/showAward.do?AwardNumber=0643784
            http://www.nsf.gov/awardsearch/showAward?AWD ID=0643784
http://www.cs.odu.edu/~mln/nsf-cv-2014.pdf 88700
http://www.cs.odu.edu/~mln/lineage.html None
http://www.cs.odu.edu/~mln/travel.html None
http://www.cs.odu.edu/~mln/mln-ad.pdf 92868
https://newsle.com/MichaelLNelson 85647
https://storify.com/michaelnelson/coverage-of-ws-dl-members-and-research None
http://ws-dl.blogspot.com/ None
http://twitter.com/phonedude_mln
            https://twitter.com/phonedude_mln
             277803
https://twitter.com/phonedude mln 277803
sirius:~>
```

```
http://www.nsf.gov/awardsearch/showAward.do?AwardNumber=0643784
        http://www.nsf.gov/awardsearch/showAward?AWD_ID=0643784
http://www.cs.odu.edu/~mln/cv.pdf 291564
http://www.cs.odu.edu/~mln/nsf-cv-2014.pdf 88700
http://www.cs.odu.edu/~mln/lineage.html None
http://www.cs.odu.edu/~mln/travel.html None
http://www.cs.odu.edu/~mln/mln-ad.pdf 92868
https://newsle.com/MichaelLNelson 85647
https://storify.com/michaelnelson/coverage-of-ws-dl-members-and-research None
http://ws-dl.blogspot.com/ None
http://twitter.com/phonedude mln
        https://twitter.com/phonedude mln
https://twitter.com/phonedude mln 277803
sirius:~> python 432Assignment1_Gomez.py http://www.cs.odu.edu/~mln/teaching/cs532-s16/test/pdfs.html
['432Assignment1_Gomez.py', 'http://www.cs.odu.edu/~mln/teaching/cs532-s16/test/pdfs.html']
/usr/local/lib/python2.7/dist-packages/bs4/_init__spy:166: UserWarning: No parser was explicitly specified, so I'm us on another system, or in a different virtual environment, it may use a different parser and behave differently.
To get rid of this warning, change this:
 BeautifulSoup([your markup])
to this:
 BeautifulSoup([your markup], "html5lib")
  markup_type=markup_type))
http://twitter.com/webscidl
        https://twitter.com/webscidl
http://www.dlib.org/dlib/november15/vandesompel/11vandesompel.html 62959
http://arxiv.org/abs/1508.02315 None
http://arxiv.org/abs/1508.02315 None
http://www.cs.odu.edu/~mln/pubs/ht-2015/hypertext-2015-temporal-violations.pdf 2184076
http://www.cs.odu.edu/~mln/pubs/tpdl-2015/tpdl-2015-annotations.pdf 622981
http://arxiv.org/pdf/1512.06195
        http://arxiv.org/pdf/1512.06195.pdf
        1748961
http://www.cs.odu.edu/~mln/pubs/tpdl-2015/tpdl-2015-off-topic.pdf 4308768
http://www.cs.odu.edu/~mln/pubs/tpdl-2015/tpdl-2015-stories.pdf 1274604
http://www.cs.odu.edu/~mln/pubs/tpdl-2015/tpdl-2015-profiling.pdf 639001
http://dx.doi.org/10.1007/s00799-015-0150-6
        http://link.springer.com/article/10.1007%2Fs00799-015-0150-6
        None
http://arxiv.org/abs/1506.06279 None
        http://link.springer.com/article/10.1007 \$ 2Fs 00799-015-0155-1
        None
http://bit.ly/1ZDatNK
        http://www.cs.odu.edu/~mln/pubs/jcdl-2015/jcdl-2015-temporal-intention.pdf
http://www.cs.odu.edu/~mln/pubs/jcdl-2015/jcdl-2015-mink.pdf 1254605
http://www.cs.odu.edu/~mln/pubs/jcdl-2015/jcdl-2015-arabic-sites.pdf 709420
http://www.cs.odu.edu/~mln/pubs/jcdl-2015/jcdl-2015-dictionary.pdf 2350603
http://bit.ly/jcdl-pdf
        http://www.cs.odu.edu/~mln/teaching/cs532-s16/test/pdfs.html
http://dx.doi.org/10.1007/s00799-015-0140-8
        http://link.springer.com/article/10.1007%2Fs00799-015-0140-8
        None
sirius:~>
```

PART 3

Bow-Tie Graph Values:

IN: O, P, M

SCC: A, B, C, G

OUT: D, H

Tendrils: I, J, L

Tubes: N

Disconnected: E, F, K

