# CS104 Project Report : Web Crawler

DHVANIL GHEEWALA
22B0923

14 June 2023

## Introduction

I am Dhvanil gheewala this is my project for CS104 :Web Crawler.
A web crawler, crawler or web spider, is a computer program that's used to search and automatically index website content and other information over the internet. These programs, or bots, are most commonly used to create entries for a search engine index.

## Libraries

The libraries and dependencies used in this project are:

- **BeautifulSoup** : Beautiful Soup is a Python package for parsing HTML and XML documents

- **requests** : Used to get send '$http$' requests from the site

- **argparse** : To receive the arguments passed by user while running the code

- **numpy and matplotlib** : For creativity and data representation[1]

- **warnings**: To remove unnecessary warnings and errors

- **os**: To remove existing specific png

## Usage

To use this program, head to the terminal and type "$python3\ web\_crawler.py$ -u $< site - name >$ -t $< threshold >$ -o $< output - file - name >$"

- **-u**: Enter a valid URL
  for the URL, If not given then it will print an error on the command line.

- **-t**: Enter a positive integer
  for the threshold of recursiveness, must be greater than 0, give an error
  for an invalid threshold

- **-o**: Enter the name of the file where the links will be added (not append).
  It will be created if it doesn't exist

For an output file, If not provided then by default print on the command line.

- **-f**:type *Yes* or *y* (first letter is case insensitive) if you want size of links

# Code Structure

## Function of base code[2]

Using the base code from *https://www.geeksforgeeks.org/*, we can scrape all urls from the given site and it prints them all without any segregation.Such a data is a very back-breaking and tedious job. Figure 1 was taken from [1]

```
http://example.webscraping.com///places/default/user/register?_next=/places/default/index
http://example.webscraping.com///places/default/user/register?_next=/places/default/index/places/default/user/regi
ster
http://example.webscraping.com///places/default/user/register?_next=/places/default/index/places/default/user/regi
ster/places/default/user/register
http://example.webscraping.com///places/default/user/register?_next=/places/default/index/places/default/user/regi
ster/places/default/user/register/places/default/user/register
http://example.webscraping.com///places/default/user/register?_next=/places/default/index/places/default/user/regi
ster/places/default/user/register/places/default/user/register/places/default/user/register
http://example.webscraping.com///places/default/user/register?_next=/places/default/index/places/default/user/regi
ster/places/default/user/register/places/default/user/register/places/default/user/register/places/default/user/re
gister
http://example.webscraping.com///places/default/user/register?_next=/places/default/index/places/default/user/regi
ster/places/default/user/register/places/default/user/register/places/default/user/register/places/default/user/re
gister/places/default/user/register
http://example.webscraping.com///places/default/user/register?_next=/places/default/index/places/default/user/regi
ster/places/default/user/register/places/default/user/register/places/default/user/register/places/default/user/re
gister/places/default/user/register/places/default/user/register
http://example.webscraping.com///places/default/user/register?_next=/places/default/index/places/default/user/regi
ster/places/default/user/register/places/default/user/register/places/default/user/register/places/default/user/re
gister/places/default/user/register/places/default/user/register/places/default/user/register
```

Figure 1: Output of a basic Web crawler

## Functions of modified code and Extra features

Using *web_crawler.py* will provide a detailed segregated output according to whether it is Internal or External link and further divided into file types along with size(if required) of each link

It can save the links crawled into a file of you choice and also saves a detailed plot of :

1. Number of files vs type of file

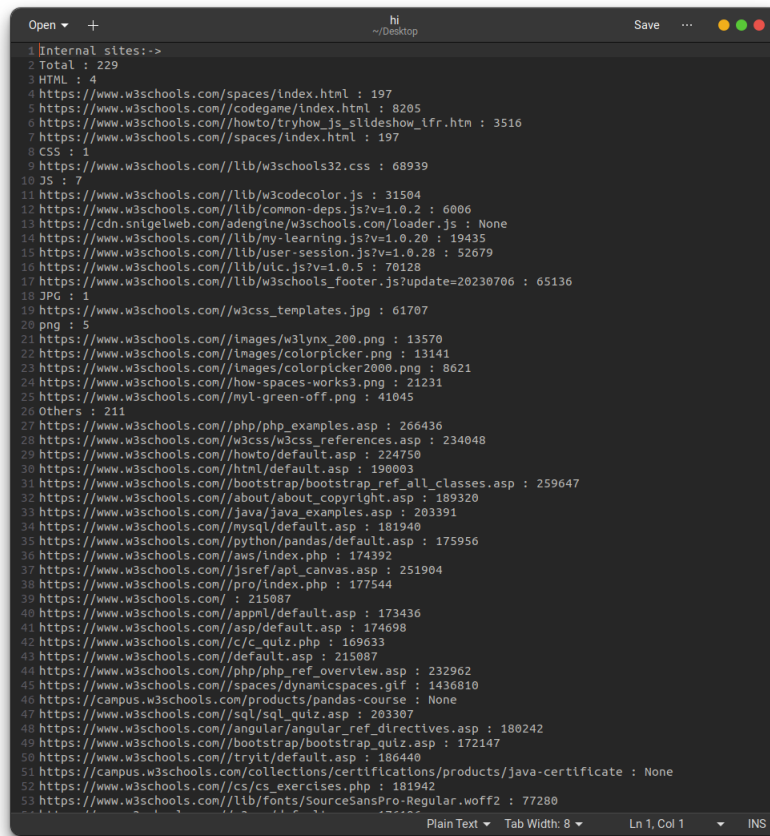2. Total file size vs type of file **(only if *-f* tag is given value *yes*)**

**Note :** It deletes the earlier png files to avoid confusion in some scenarios.

Figure 2: Functioning of web_crawler.py

## Working

- With the help of BeautifulSoup, the code extracts the links in the *href* and *src* tags and then changing them to make valid urls(some are relative ,so the code prefixes the main site etc.).

- The code is based on recursive crawling through the url given by the user.

- The code finally segregates all the files into various categories like *HTML, CSS, JS, JPG, PNG and Others* using a nested dictionary.

Figure 3: Output of web_crawler.py

# References

[1]  URL: https://www.geeksforgeeks.org/bar-plot-in-matplotlib/.

[2]  URL: https://www.geeksforgeeks.org/python-program-to-recursively-scrape-all-the-urls-of-the-website/.