# Python Web Crawler

Ankeet Saha

June 10, 2023

## 0.1 Introduction

This report describes a Python web crawler developed as part of the DIC project. The web crawler takes a website URL as input and returns the links found on the website. It provides various command line arguments for customization.

## 0.2 Modules Used

- argparse

- requests

- urllib

- bs4

- lxml

- colorama

- urllib3

  To ensure that the web crawler functions correctly in your system you can use the command `pip3 install <names of the Modules mentioned above>`

## 0.3 Usage

To run the web crawler, use the following command line arguments:
`finalProject.py -u <url> -t <threshold> -o <output file> -H`
The `-t` and `-o` arguments are optional, while `-H` is an optional flag.
-u is a compulsory argument

### 0.3.1 Command Line Arguments/Flags

- `-t`: Specifies the recursive depth up to which the website will be crawled. If not provided, the crawler will crawl until there are no more links.

- `-o`: Specifies the output file where the results will be saved. If not provided, the contents are printed to the terminal.

- `-H`: A flag that enables error handling from the server side.If not provided, all links found are printed, regardless of encountering errors or not. It handles the request exception errors.

- `-u`: This is a compulsory argument which is used to provide a url to crawl.

## 0.4 Crawler Characteristics

The Python web crawler developed for this project has the following characteristics:

1. Takes a website URL as input and returns the links found on the website.

2. Allows specifying the recursive depth for crawling using the `-t` command line argument.

3. Supports saving the results to an output file using the `-o` command line argument.

4. Provides error handling from the client side using the `-H` flag.

## 0.5 Output Format

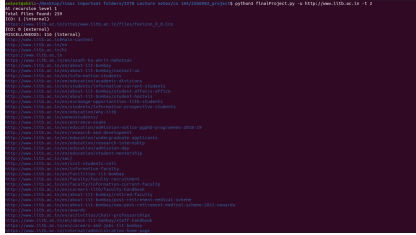The web crawler generates the following output:

- File names in each recursive depth are printed.

- The number of files printed in each recursive depth is displayed.

- The files are segregated based on file types.

- The number of files for each file type in each recursive depth is shown.

- Each file type is further categorized as internal or external.
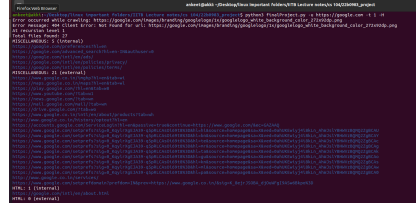
## 0.6 Customizations

The web crawler includes the following customizations:

1. Valid links are colorized in blue to differentiate them from other information.

2. Files are categorized as internal or external.

3. Error handling is implemented for the following cases:

   (a) Timeout error
   (b) SSL certificate verification error
   (c) Request exception error from the client side

## 0.7  Screenshots



(a) Example 1



(b) Example 2

# Bibliography

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

[2] URL: https://stackoverflow.com/questions/50270232/
scrape-all-of-sublinks-of-a-website-recursively-in-python-using-beautiful-soup.