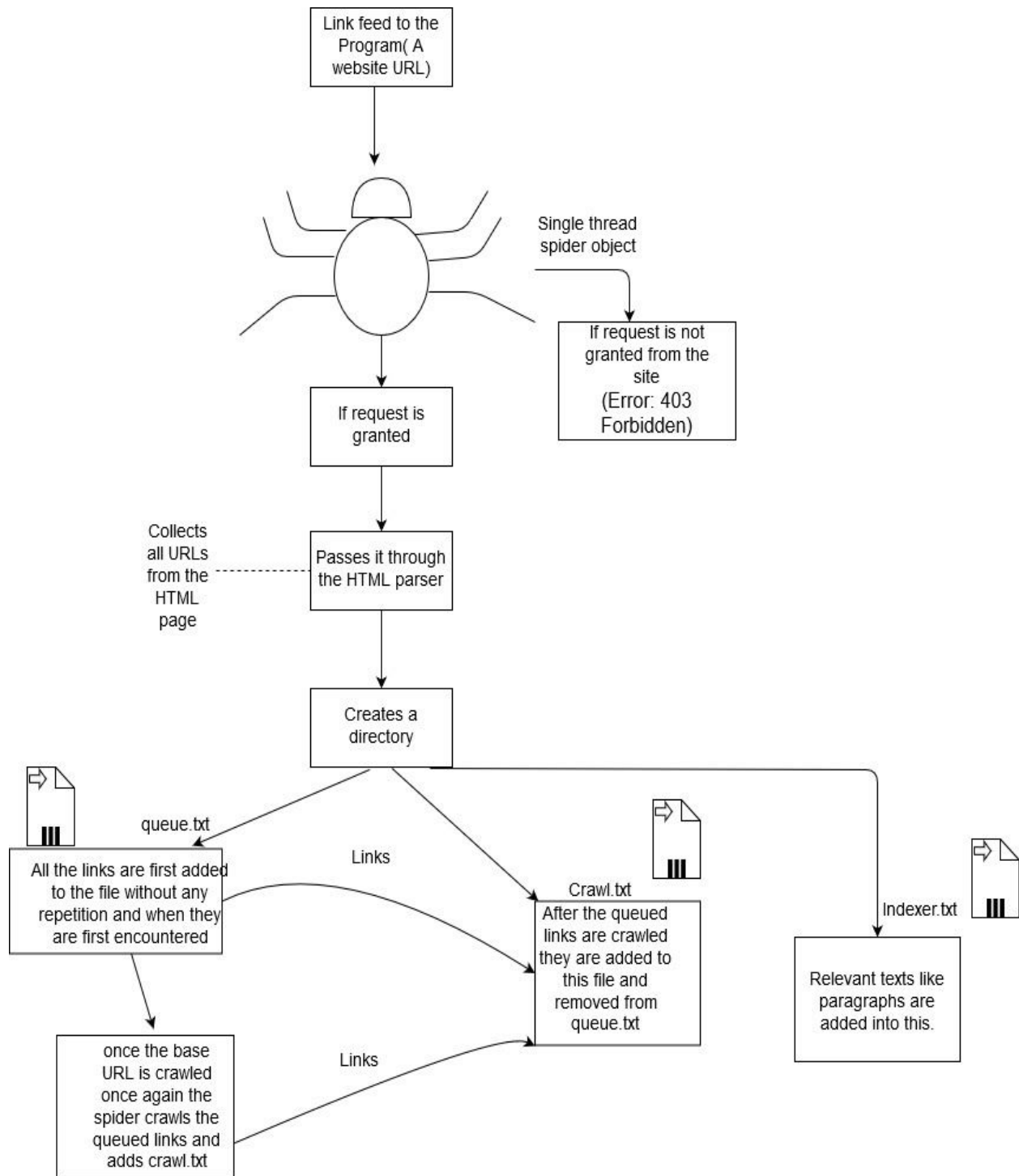# WEB CRAWLER

## Introduction:

- **<u>Definition:</u>** A Web crawler, sometimes called a spider or spider bot and often shortened to crawler, is an Internet bot that systematically browses the World Wide Web, typically for the purpose of Web indexing.

- **<u>Uses:</u>** Web crawlers are mainly used to create a copy of all the visited pages for later processing by a search engine, that will index the downloaded pages to provide fast searches. Crawlers can also be used for automating maintenance tasks on a Web site, such as checking links or validating HTML code.

- **<u>Features:</u>** The key features of Web Crawler implemented in the project are:

  1. The contents of the web page are the corpus from which we retrieve the data later.

  2. The user requests the web page to provide relevant links and text paragraphs and collects them as response.

  3. In this crawler we are crawling through its HTML page and collecting all the relevant links and paragraphs and storing them in different files.

Link feed to the Program( A website URL)

Single thread spider object

If request is not granted from the site
(Error: 403 Forbidden)

If request is granted

Collects all URLs from the HTML page

Passes it through the HTML parser

Creates a directory

queue.txt

Links

Crawl.txt

Indexer.txt

All the links are first added to the file without any repetition and when they are first encountered

After the queued links are crawled they are added to this file and removed from queue.txt

Relevant texts like paragraphs are added into this.

once the base URL is crawled once again the spider crawls the queued links and adds crawl.txt

Links

# Approach and Workflow-

Crawlers can be of many types,we have implemented the one which would a website URL as the user input and give us a file collecting all the links in it and relevant text inside each page of that particular webpage.
We have implemented a two level crawler.In the first level it gets the url and crawls it and after getting all the links in a file in the second level we are again crawling every individual link.
We have used the python programming language to implement the crawler.
Files in our Code
        1.main.py
        2.functions.py
        3.htmlpars.py
        4.crawler.py
        5.textinfo.py

**main.py** is for the user reference where one will give the URL of a particular website as the input to the program.

**crawler.py** has the spider class that is responsible for collecting the user fed link, parse it using htmlparser.py and use the functions.py to get all the possible URLs.

**functions.py** contains some user defined functions to create a directory,create queue and crawl text files,appending and deleting from files,files to set and set file conversions, getting project names and domain and subdomain names.

**htmlpars.py** will create an object of find_link( ) class which is a child class of HTMLParser,By getting a link it will parse its HTML page and get all the **<a>** tags with attribute type **"href"=** and return it in the form of set( ).These are the hidden links in a website in the from of buttons, texts, images which are navigable redirects.

**textinfo.py** is used to extract all the **<p>** tag values so as to create a corpus of text data that would be later indexed.

Finally we will get a Project directory consisting of three files queue.txt,crawl.txt and indexer.txt.

## User Input:



**1. queue.txt** contains all the URLs that we get after crawling the base URL which can be further crawled.

**2. crawl.text** is the file having all the URLs which are crawled.



**3. indexer.txt** contains the all relevant text inside a particular link

## Shortcoming of our web crawler:

The web crawler is not able to crawl certain websites due to **403 Forbidden error,** we have implemented an exception handling for that.
The program being single threaded takes some time to crawl the website.
It will not work for files such as pdfs, text files,images to extract links in that.

By:
**Rohit Kumar Behera**
**Saanidhya Vats**
**Saidatta Sahu**
**Samprita Sahoo**