**Web and Social Computing (IT752)**
**Lab Assignment 3**

**Submitted to : Dr. Sowmya Kamath**
**Submitted by : Shaikh Sahil Ahmed**
**Roll number: 192IT022**

## 1. Crawler

A Web crawler, sometimes called a spider or spiderbot and often shortened to crawler, is an Internet bot that systematically browses the World Wide Web, typically for the purpose of Web indexing.
Library used as spider for web crawling. The basic web crawler is uploaded under crawler directory. The main.py file contains the crawler program where it is crawling the website URL: 'http://infotech.nitk.ac.in/'.

## 2. Single Threaded and Multi Threaded Crawler.

The package used is icrawler to implement crawler and imported googleImageCrawler. For single threaded, I have chosen the thread value as 1 and searched for the keyword "Donald Trump", so related to the keyword the pages are crawled within the range of date 1 Jan 2016 to 31 Jan 2016.



So, execution time taken by the single threaded crawler is 10.02 seconds.
Now I go for the Multi Threaded Crawler. Same as above I have chosen the thread value as 5 and searched for the keyword "Donald Trump", so related to the keyword the pages are crawled within the range of date 1 Jan 2016 to 31 Jan 2016.



So, execution time taken by the single threaded crawler is 10.02 seconds. Multi thread is more efficient.

**2. Data structure used for indexing (optimal for searching, insertion... )**

In the case of a crawler though, it would not be needed because starting from the root, you can maintain a "list" of visited URLs and every time you are about to follow a link, check if it has been encountered before. If it has not been encountered then add it to the list and follow it.

It doesn't have to literally be a list (i.e. array) though, it can be a dictionary or other data structure that would help in speeding up the search.
It could also be an SQL database or something like a key-value storage like redis. If you use something like this then all the indexing and querying will be handed for you by the database system with which you could communicate through a standard method (SQL, special API, other).

**3. BFS and DFS crawler**

BFS and DFS crawler is implemented and uploaded under the "bfs dfs" folder.
Giving an input, crawler crawls on the basis of bfs and dfs.