Analysis report for
**Web and Social Computing Assignment 2**

By Sampat Kr Ghosh
Roll No. 192IT020

**Introduction:**
This reports explains the implementation of web crawling using icrawler package. Performance analysis is doen on number of threads used in multi-threaded crawling.

**Performance comparison of sequential and multithreaded crawler:**



Above figure shows sequential crawling using icrawler



Above figure shows multi-threaded crawling.

From the time taken, we can say that as the number if threads increases, the time decreases for crawling.

Below are the screenshot where
Input:  Starting URL, list of words, max number of pages to be download.

Try each one of the below variants, while designing crawling strategies.
a. Crawl only from pages whose body text includes one of the words in the list.
b. Crawl only from pages whose title includes one of the words in the list.
c. Give priority to page with most words from list (either most different words or most occurrences of words)
Output: Search results page

```
Console ☒        Task List    Outline    Problems
<terminated> MultipleCrawlerController [Java Application] C:\Pro
Enter url:
http://infotech.nitk.ac.in/
Enter words space separated:

david reshma
Enter max number of pages to fetch:
1000
[david, reshma]
http://infotech.nitk.ac.in/
```

```
=======Time taken: 151 seocnds
====================Search Result========================
People | Department of Information Technology
url: http://infotech.nitk.ac.in/node/118
weight: 3
-----------------------------------------------------
Sidney David Rosario | Department of Information Technology
url: http://infotech.nitk.ac.in/faculty/sidney-david-rosario
weight: 3
-----------------------------------------------------
People | Department of Information Technology
url: http://infotech.nitk.ac.in/people?order=field_phdcategory&sort=asc
weight: 3
-----------------------------------------------------|
Reshma U | Department of Information Technology
url: http://infotech.nitk.ac.in/research-scholars/reshma-u
weight: 4
-----------------------------------------------------
Reshma U | Department of Information Technology
url: http://infotech.nitk.ac.in/research-scholars/reshma-u
weight: 4
-----------------------------------------------------
People | Department of Information Technology
url: http://infotech.nitk.ac.in/people?qt-professorquicktab=3
weight: 3
-----------------------------------------------------
People | Department of Information Technology
url: http://infotech.nitk.ac.in/people?qt-professorquicktab=2
weight: 3
-----------------------------------------------------
People | Department of Information Technology
url: http://infotech.nitk.ac.in/people?qt-professorquicktab=5
weight: 3
-----------------------------------------------------
People | Department of Information Technology
url: http://infotech.nitk.ac.in/people?qt-professorquicktab=4
weight: 3
-----------------------------------------------------
Reshma | Department of Information Technology
url: http://infotech.nitk.ac.in/content/reshma
```