# Web Scraping to create Image Captioning Dataset

R.Rohith
EP21B030

## 1 CAUTION

Ensure that you have Google Chrome and MongoDB installed. Additionally, verify that all packages listed in the `requirements.txt` file are installed. Simply run `main.py` once you are set.

## 2 Config

The configuration JSON file includes the following three parameters:

- *home_page_url*: The URL of the Google News home page.

- *keywords*: A list of potential keywords to search for the top stories page, containing strings that could represent **Top stories**.

- *image_dir*: The directory, relative to the project directory, where images are saved locally for future reference.

## 3 Module 1

Module 1's goal is to scrape the HTML content from the home page of Google News. The `fetch_page.py` Python file accomplishes this task. Given the *home_page_url*, it returns the HTML content of the home page using the `requests` library.

## 4 Module 2

Module 2 aims to retrieve the URL of the top stories page from the HTML content of the home page. This is handled by the `fetch_topstories_url.py` file, which, when given the HTML content, searches through the text of all $< a >$ tags. It returns the `href` attribute (the URL of the top stories page, relative to the home page) when the text matches a keyword in the *keywords* list.

## 5 Module 3

Module 3 focuses on scraping images and captions from the top stories page. It uses the `selenium` package to handle lazy loading by scraping the HTML content of the top stories page through the Chrome browser. The script then iterates over all `article` tags in the page that contain figures to gather the images, captions, and publication date-time.

## 6 Module 4 & Module 5

Module 4 receives data from Module 3 (including the caption, image URL, date-time, and image data) and passes it through Module 5. Module 5 checks whether an entry with the same caption and date-time already exists in the database. The data is then stored in the **image_captioning database**. Images are saved in the **images collection**, and the captions with metadata are stored in the **captions collection**. The images are associated with the unique ID of their metadata in the captions collection, ensuring a one-to-one mapping between the entries in both collections. Both are implemented in the `save_db.py` file.

## 7   Module 6

Module 6 orchestrates all the previous modules and handles logging of actions along with their timestamps. The logs are saved in the directory named `logs`.



Figure 1: A snapshot of the log file



Figure 2: A sample of scraped image along with the local image directory



Figure 3: Few entries from the **captions collection** in MongoDB Database