# DA5402 - Assignment 1

Anirudh Rao be21b004

## Module 1

To scrape Google News, the `selenium` library in Python was used. This was chosen as it has the ability to scrape lazy loading websites. It uses a web driver for a web browser to interact with websites and scrape information from them. For this project, the web driver for Google Chrome was used.

The Chrome web driver, powered by `selenium`, opens the homepage of Google News. The URL for this is not hard-coded and is instead configurable via the attached `config.json` file.

To account for lazy loading, a sleep time is incorporated to allow the website to load sufficient content for scraping. This sleep time is also configurable using `config.json`.
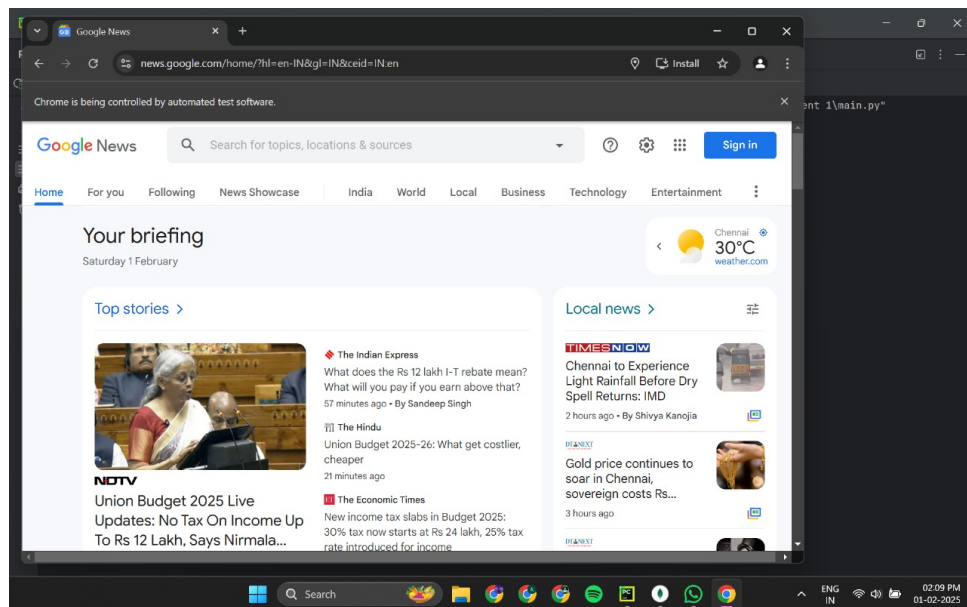


Figure 1: `selenium` opening the Google News homepage

## Module 2

Next, `selenium` clicks on the 'Top stories' link to navigate to that page. This string is not hard-coded and is instead configurable via `config.json`.
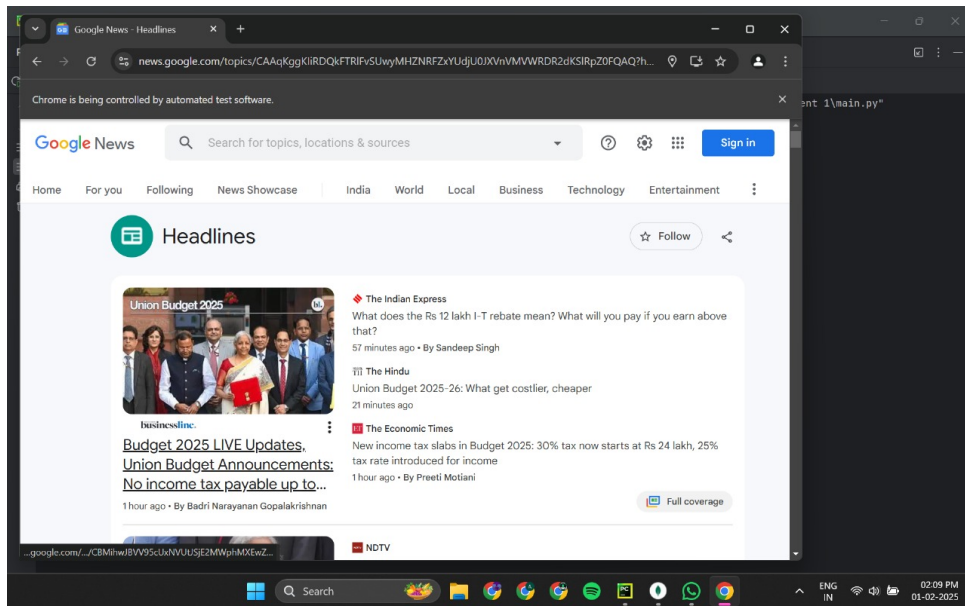
Figure 2: `selenium` navigating to the 'Top stories' page

## Module 3

`selenium` scrolls through the 'Top stories' page until the height of the webpage stays constant. This scrolling is performed using the configured sleep time, keeping in mind the lazy loading of the webpage.
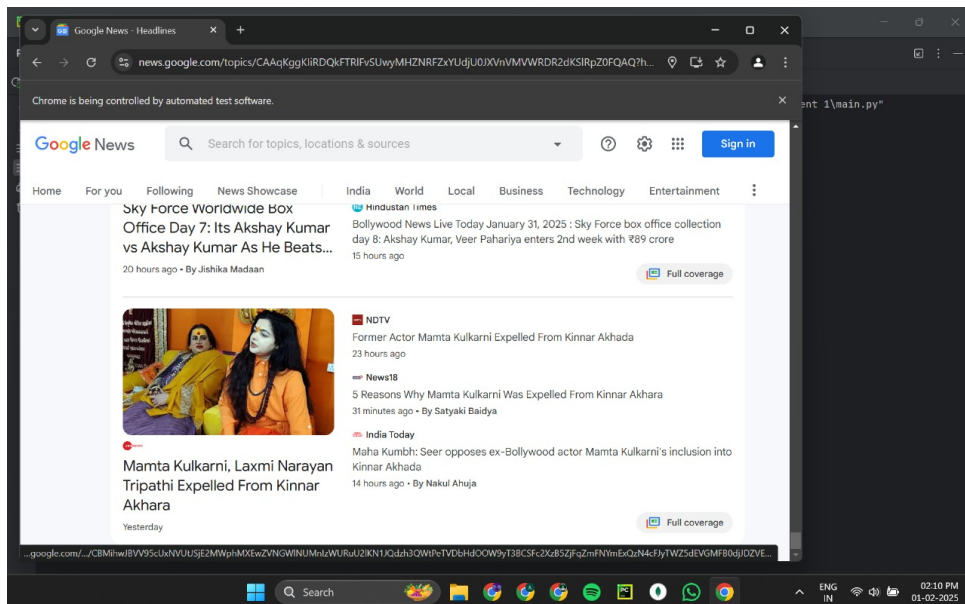


Figure 3: `selenium` navigating to the bottom of the 'Top stories' page

Finally, `selenium` uses CSS selectors to scrape the headline, thumbnail, newspaper, article URL, and article date. To account for Google News changing its layout, these CSS selectors are not hard-coded and are instead configurable via `config.json`.
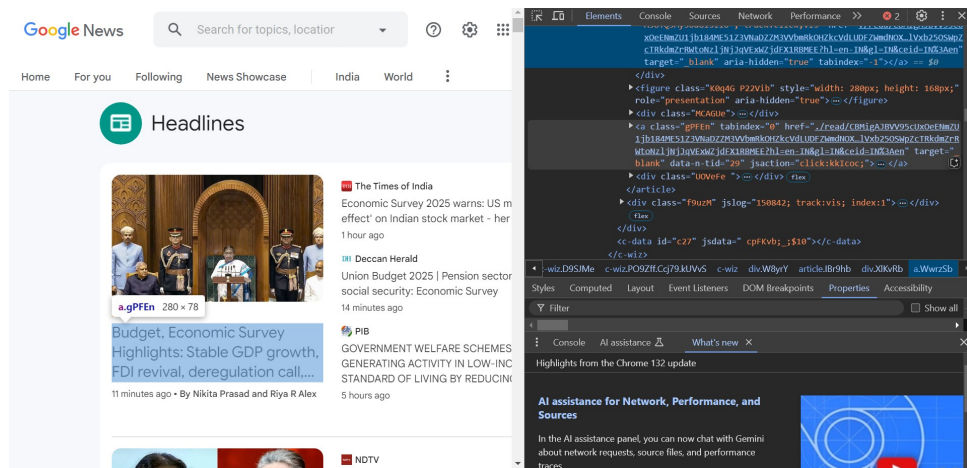
Figure 4: Example of a CSS selector using the class (`gPFEn`) of the headline

## Module 4

MongoDB was chosen as the database to store the scraped data. The library `pymongo` is used to interact with MongoDB programmatically. The database is stored locally and not on the cloud.

Connecting to the MongoDB database requires some configurations, which are not hard-coded and are specified in the `config.json` file. In addition to this, MongoDB Compass must be installed locally on the system being used.

The scraped data and metadata, along with the scrape time, are added by the script to the database as a 'document', which is equivalent to a Python dictionary. The thumbnail is stored in a binary format that is friendly to MongoDB. The original image can be obtained from the binary data using the `PIL` library.
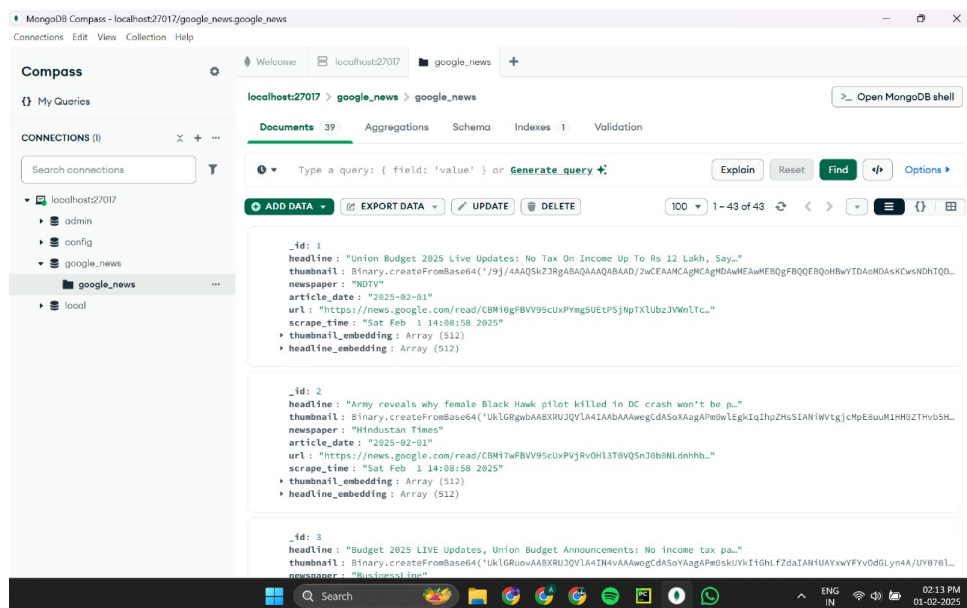


Figure 5: Scraped data in the MongoDB database

# Module 5

Before adding data to the database, it is important to check whether the same, or similar, articles are already present in the database. Two articles may have semantically similar headlines or visually similar thumbnails. Having such articles in our image-caption dataset is redundant and not useful.

To compare the similarities of the headlines and thumbnails of articles, we use embeddings from OpenAI's CLIP model. Both the headline embedding and the thumbnail embedding (vectors of size 512) of each article are stored in the MongoDB database. If the headline embeddings or thumbnail embeddings of two articles have a cosine similarity above a threshold that is configured in `config.json`, they are deemed duplicates. If one of these articles is already present in the database, the second article is not added to the database.
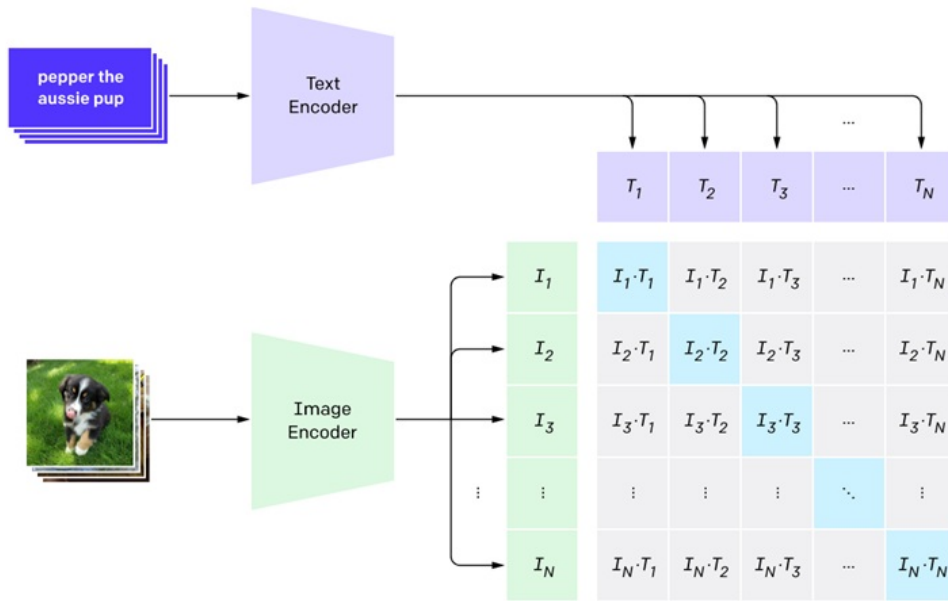


Figure 6: CLIP embeddings



Figure 7: Comparing image similarity using CLIP

# Module 6

All modules are orchestrated in a cascaded manner in the `main.py` script.
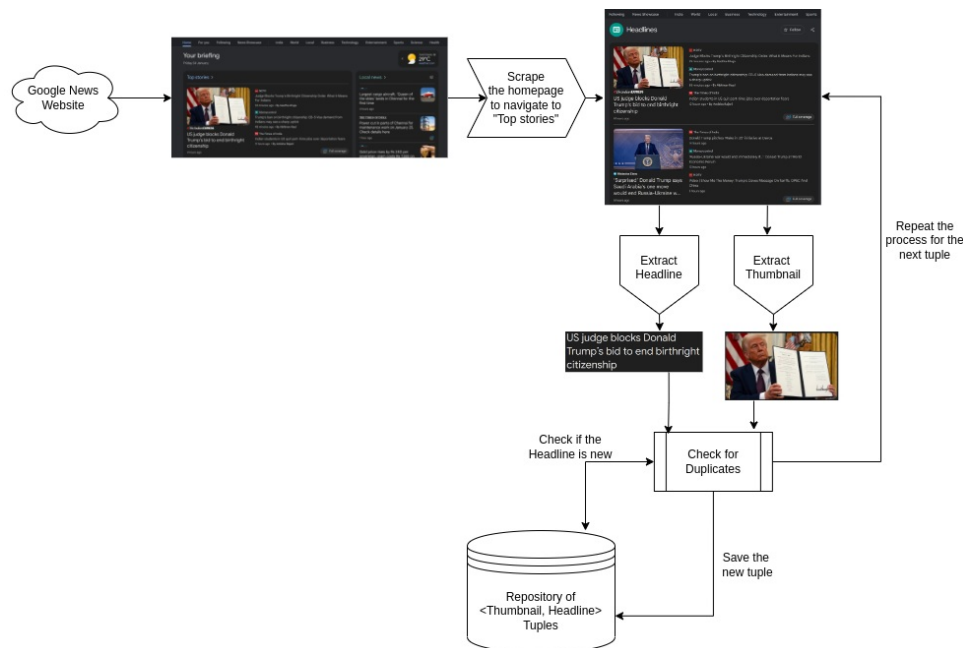


Figure 8: The flow of the script

The `logging` library is used to log time stamps of invocation, error statuses, and other relevant details in the `script.log` file. The script takes about 10 minutes to run completely.
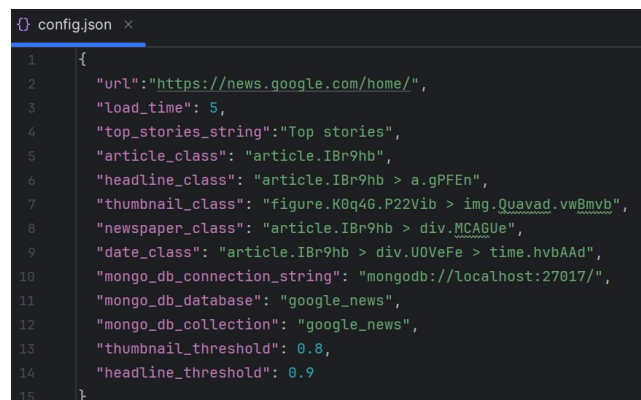


Figure 9: The `script.log` file after running the script twice

A `requirements.txt` file is also provided to install the libraries necessary to run this script.

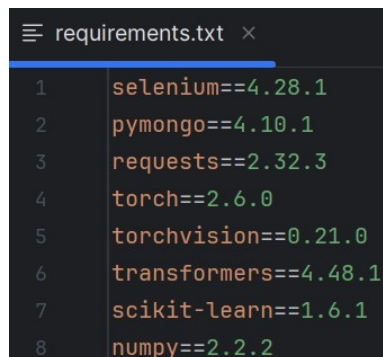The script is made friendly to be executed as a CronJob.

- There are no `print` statements in the script.

- Logging and error handling are performed diligently.

- The parameters are configurable through `config.json`.

- The shebang line is provided at the top of `main.py` for the Python interpreter.

The script can be scheduled as a CronJob using the `crontab` after granting the necessary permissions.

```json
{} config.json  ×
1    {
2        "url":"https://news.google.com/home/",
3        "load_time": 5,
4        "top_stories_string":"Top stories",
5        "article_class": "article.IBr9hb",
6        "headline_class": "article.IBr9hb > a.gPFEn",
7        "thumbnail_class": "figure.K0q4G.P22Vib > img.Quavad.vwBmvb",
8        "newspaper_class": "article.IBr9hb > div.MCAGUe",
9        "date_class": "article.IBr9hb > div.UOVeFe > time.hvbAAd",
10       "mongo_db_connection_string": "mongodb://localhost:27017/",
11       "mongo_db_database": "google_news",
12       "mongo_db_collection": "google_news",
13       "thumbnail_threshold": 0.8,
14       "headline_threshold": 0.9
15   }
```

Figure 10: The `config.json` file

```
≡ requirements.txt  ×
1    selenium==4.28.1
2    pymongo==4.10.1
3    requests==2.32.3
4    torch==2.6.0
5    torchvision==0.21.0
6    transformers==4.48.1
7    scikit-learn==1.6.1
8    numpy==2.2.2
```

Figure 11: The `requirements.txt` file