# TRUMP DAILY

Trump daily is a crawler that crawls the top 25 news contents that contains the top-featured man on the planet, US President Donald Trump. Currently, the crawler crawls from

- CNN
- Twitter

## Installation

The crawler is built using python 3.6. So you would need to set up Python 3.6. The crawlers can be run as scheduled task (cron job) or can be run manually.

- Create a mysql/mariadb database named 'trumpdaily' in your localhost/hosted server
- Upload the SQL file 'trumpdaily.sql' to create the tables
- Install the required python3.6 libraries
- Run mainCrawler.py to start the crawler
- Connect your website or web application to query the feed from the database and display on the site

### Required libraries

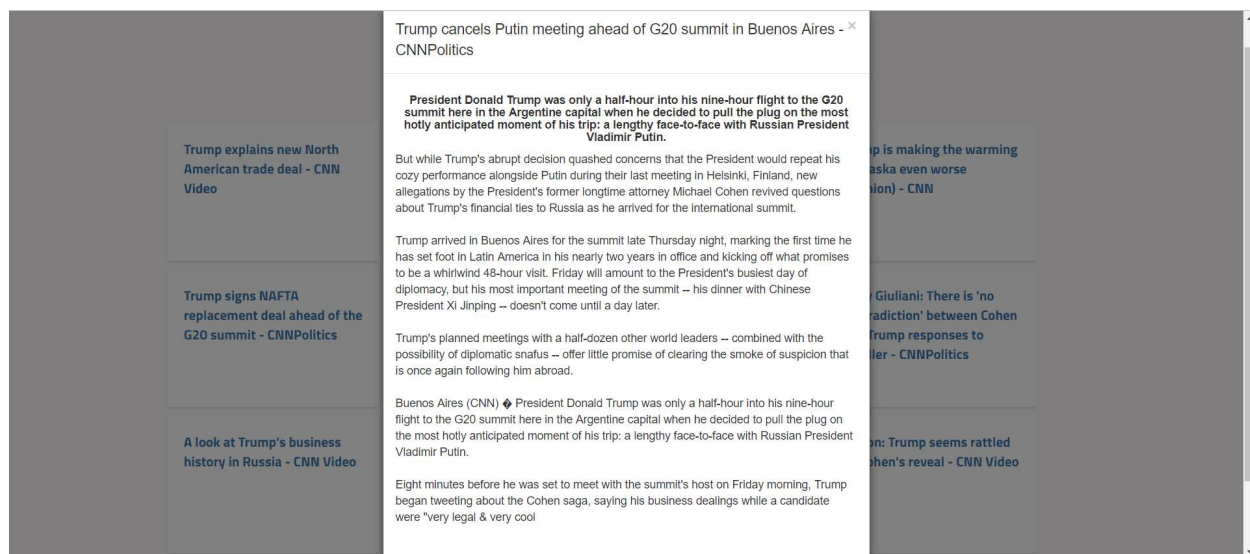The python libraries and the additional tools that are required to set this project are:

| Libraries/Tools | Purpose | Reference |
|---|---|---|
| MySQL | To save the collected feeds | https://www.mysql.com/ |
| Github | Git repository | plugins/github/README.md |
| Selenium 3.141.0 | To render dynamic web | https://pypi.org/project/selenium/ |

| | pages of CNN | |
|---|---|---|
| Beautifulsoup4 4.6.3 | To build effective crawler to parse through contents | https://pypi.org/project/beautifulsoup4/ |
| Configparser 3.5.0 | To parse configuration files | https://pypi.org/project/configparser/ |
| time 1.0.0 | Default python library to add delays into twitter scrapping | https://pypi.org/project/time/ |

# THE WEBSITE

The url to the website is trumpdaily.deciphertechglobal.com

## How to use the website:

- The website is divided into two sections, CNN and Twitter
- If you click CNN logo, it will take you to the CNN section
- If you click the Twitter logo, it will take you to the Twitter section
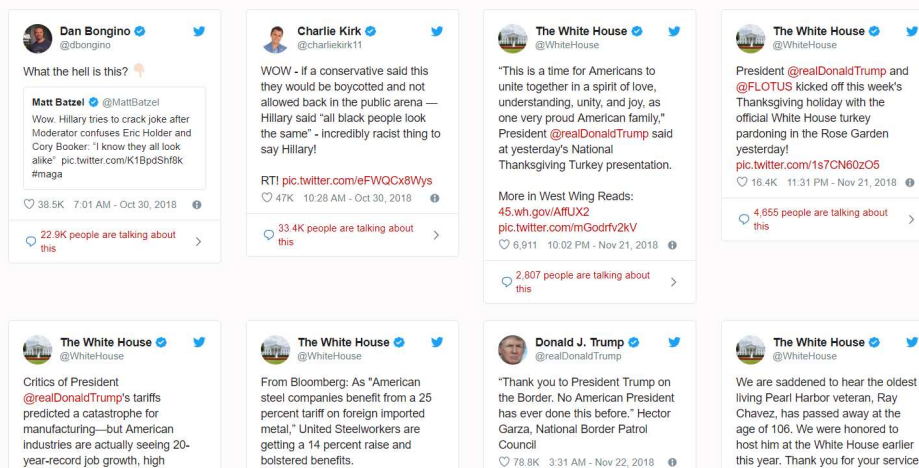
## CNN Section:

- Designed in a modal form
- If you click on any news title it will pop up the news and its details like this



## Twitter section:

- Feed collected from the Twitter tweet IDs
- The IDs are then used to generate Twitter widget using https://platform.twitter.com/widgets.js

DONALD TRUMP TWITTER UPDATES

# CNN CRAWLING

1. To crawl CNN and find news about Trump, the starting point was chosen as the home page (www.cnn.com)
2. Using search engine analytics similarweb (www.similarweb.com), two primary category sections of CNN was chosen which has the highest probability to contain top news about Trump, they are "politics" and "us".
3. Different crawlers were made for each of the pages to crawl the urls of the news finding and scraping from the element named "cd__headline"
4. After getting the URLs, the URLs were passed to a class named "singleCNNNews", which crawls each URL and extract title (element parsed= title), short summary (element parsed=meta description) and description (multiple methods used for different templates used by CNN) (if available)
5. Accumulated data is then returned in the form of JSON and then stored in the database. Error handling has been added to handle any error in between this process

# TWITTER CRAWLING

The ideology behind twitter crawling to get Trump tweets is to get the Tweet IDs from the two official accounts of Donald Trump: @ POTUS and @ realDonaldTrump.

1. At first the page is rendered and scrolled down automatedly 5 times to make larger collection of tweets in the page source
2. Then the data-stream-id, which contains the Tweet ID, is collected from all the 'js-stream-tweet' elements of the page
3. The IDs are then sent to the database