

SELENIUM AND GOOGLE_IMAGES_DOWNLOAD

Web Scraping Images from Google



Anand Suresh

Mar 19 · 4 min read

. . .

Hello Readers! I am Anand, a very enthusiastic developer and an undergrad student.

Recently, i've been trying to figure out different ways to web scrape images from google.com based on search query using python and I've stumbled upon more than one method to do so. Not too sure to say which method is actually better, but to my usage I've incurred that It depends on the query or the kind of images we want to download.

So for this Blog I've decided to split the content into different snippets for ease of use later and add the final codes in the end for those who want to skip the process.

. . .

Method 1: Using "google_images_download" library

Installation:

```
Using PyPi - $ pip3 install google_images_download
```

```
-----  
Using CLI - (cd into working directory)
```

```
$ git clone https://github.com/hardikvasa/google-images-download.git
```

```
$ cd google-images-download
```

```
$ python setup.py install
```

Usage:

Given below is a code snippet as a function call whose argument will be the query as a string. So all we've got to do when embedding this function into a program would be to import the library and call the function with a set of queries.

```
# importing google_images_download module
from google_images_download import google_images_download
```

```
1  def downloadimages(query):
2      response = google_images_download.googleimagesdownload()
3      # keywords is the search query
4      # format is the image file format
5      # limit is the number of images to be downloaded
6      # print urs is to print the image file url
7      # size is the image size which can
8      # be specified manually ("large, medium, icon")
9      # aspect ratio denotes the height width ratio
10     # of images to download. ("tall, square, wide, panoramic")
11     arguments = {"keywords": query,
12                 #"format": "jpg",
13                 "limit":100,
14                 "print_urls":True}
15                 #"size": "large"}
16                 #"aspect_ratio": "panoramic" }
17     '''
18     I've played around with different parameters
19     some yeild resuts and others don't. It can also depend on the searched query
20     '''
21     try:
22         response.download(arguments)
23
24     # Handling File NotFound Error
25     except FileNotFoundError:
26         arguments = {"keywords": query,
27                     "format": "jpg",
28                     "limit":4,
29                     "print_urls":True,
30                     "size": "medium"}
31
32         # Providing arguments for the searched query
33         try:
34             # Downloading the photos based
35             # on the given arguments
36             response.download(arguments)
```

The above code snippet automatically saves the downloaded images within it's respective folder.

. . .

Method 2: Using Selenium

Installation:

To install Selenium using PyPi

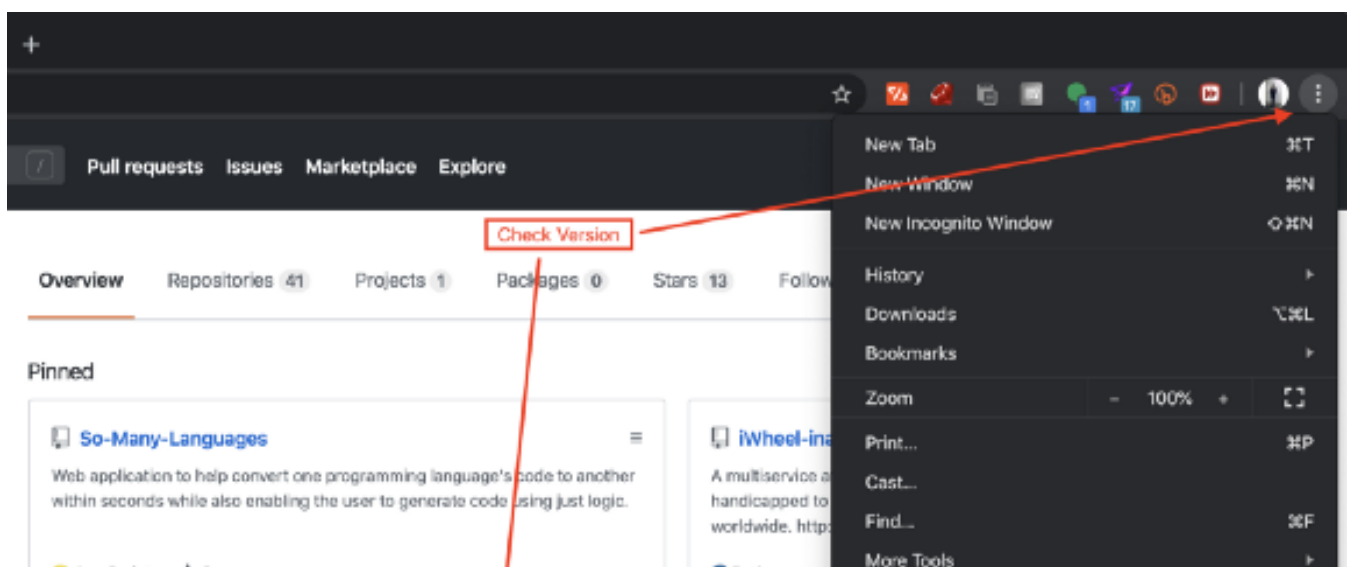
```
$ pip3 install selenium
```

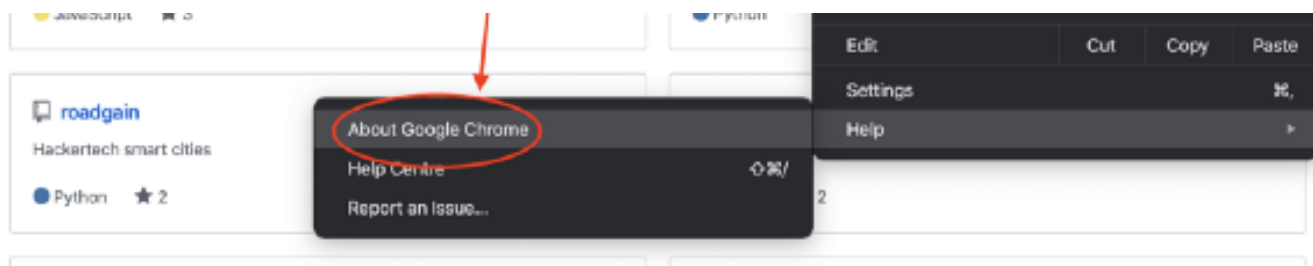
Usage:

Consider Selenium to be this package that lets your program do what a human user would normally do on the browser, like searching for something, going through videos on youtube, downloading images from xyz website and such can all be possible by using selenium to simulate a user's actions.

For this tutorial we are going to pair up Selenium with Google Chrome and automate it to fetch the images we require. To do that we are going to need the Chrome driver. There is a specific chrome-driver to every version of chrome, My chrome version is V80, hence i use the ChromeDriver version 80.

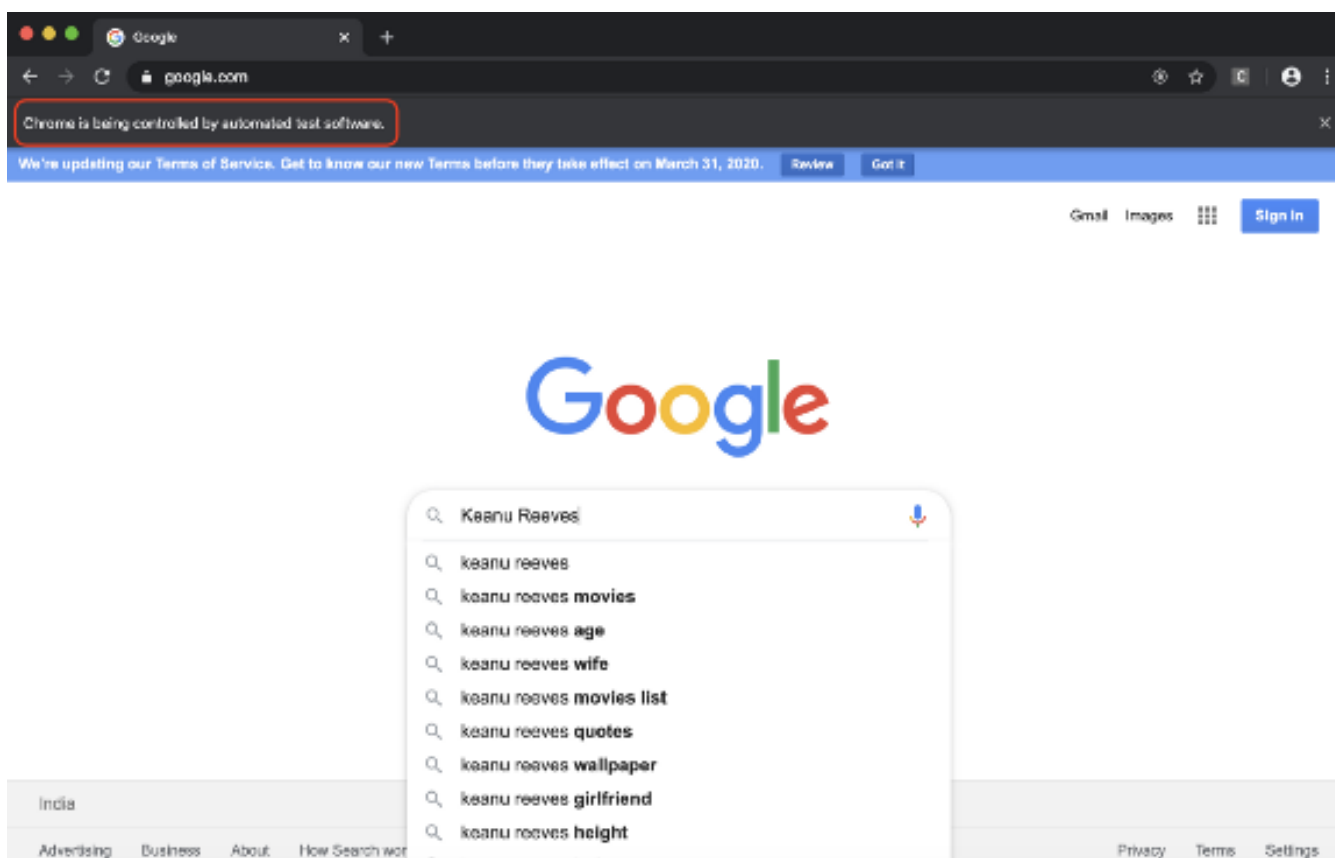
You can find the version of your chrome in the "About Google Chrome" section and download the respective version of chrome-driver from [HERE](#) .





Check Chrome Version

First, we will start by running a simple program which will open up the chrome window and automate a search of the query of our choice from the code.



Automated Searching

```

1  import selenium
2  from selenium import webdriver
3  # This is the path I use
4  # DRIVER_PATH = '/Users/anand/Desktop/chromedriver'
5  # Put the path for your ChromeDriver here
6  DRIVER_PATH = 'Enter/your/path/here'
7  wd = webdriver.Chrome(executable_path=DRIVER_PATH)
8  wd.get('https://google.com')
9  search_box = wd.find_element_by_css_selector('input.gLFyf') #input box selector
10 search_box.send_keys('Any query')
11 #wd.quit()

```

creates a window

Second phase of this program would be to search for the query, move into the images section and get the image links. To achieve the following just like the previous code, we'll be using the chrome-driver to navigate using selectors and pages.

```
fetch_image_urls(query:str, max_links_to_fetch:int, wd:webdriver,
sleep_between_interactions:int=1)
```

The function `fetch_image_urls` takes 4 input arguments but the 4th argument has a default value of 1 (sleep interval) so, using only 3 arguments will suffice. The function returns the image-urls.

1. 'query' is the string input in the textbox from the previous snippet.
2. 'max_links_to_fetch' is an integer input that defines the number of images required after scraping.
3. 'wd' is the path to the chromedriver installed on your PC or mac.

```
1  def fetch_image_urls(query:str, max_links_to_fetch:int, wd:webdriver, sleep_between_interactions:int=1):
2      def scroll_to_end(wd):
3          wd.execute_script("window.scrollTo(0, document.body.scrollHeight);")
4          time.sleep(sleep_between_interactions)
5
6      # build the google query
7      search_url = "https://www.google.com/search?safe=off&site=&tbm=isch&source=hp&q={q}&oq={q}&oeq={q}"
8
9      # load the page
10     wd.get(search_url.format(q=query))
11
12     image_urls = set()
13     image_count = 0
14     results_start = 0
15     while image_count < max_links_to_fetch:
16         scroll_to_end(wd)
17
18         # get all image thumbnail results
19         thumbnail_results = wd.find_elements_by_css_selector("img.Q4LuWd")
20         number_results = len(thumbnail_results)
21
22         print(f"Found: {number_results} search results. Extracting links from {results_start}:{number_results}")
23
24         for img in thumbnail_results[results_start:number_results]:
```

```

25     # try to click every thumbnail such that we can get the real image behind it
26     try:
27         img.click()
28         time.sleep(sleep_between_interactions)
29     except Exception:
30         continue
31
32     # extract image urls
33     actual_images = wd.find_elements_by_css_selector('img.n3VNCb')
34     for actual_image in actual_images:
35         if actual_image.get_attribute('src') and 'http' in actual_image.get_attribute('src'):
36             image_urls.add(actual_image.get_attribute('src'))
37
38     image_count = len(image_urls)
39
40     if len(image_urls) >= max_links_to_fetch:
41         print(f"Found: {len(image_urls)} image links, done!")
42         break
43     else:
44         print("Found:", len(image_urls), "image links, looking for more ...")
45         time.sleep(30)
46         return
47     load_more_button = wd.find_element_by_css_selector(".mye4qd")
48     if load_more_button:
49         wd.execute_script("document.querySelector('.mye4qd').click();")
50
51     # move the result startpoint further down

```

system. Hence, for that we are going to download images using PILLOW module.

```
$ pip3 install pillow
```

```
-----
Program call:
```

```
import PIL
```

once we download the images we need to save the downloaded images in an organised folder of our choice and for that we will be using simple OS commands from the python file.

```
1 def persist_image(folder_path:str,file_name:str,url:str):
```

```

2     try:
3         image_content = requests.get(url).content
4
5     except Exception as e:
6         print(f"ERROR - Could not download {url} - {e}")
7
8     try:
9         image_file = io.BytesIO(image_content)
10        image = Image.open(image_file).convert('RGB')
11        folder_path = os.path.join(folder_path, file_name)
12        if os.path.exists(folder_path):
13            file_path = os.path.join(folder_path, hashlib.sha1(image_content).hexdigest()[:10])
14        else:
15            os.mkdir(folder_path)
16            file_path = os.path.join(folder_path, hashlib.sha1(image_content).hexdigest()[:10])
17        with open(file_path, 'wb') as f:
18            image.save(f, "JPEG", quality=85)
19        print(f"SUCCESS - saved {url} - as {file_path}")
20    except Exception as e:
21        print(f"ERROR - Could not save {url} - {e}")

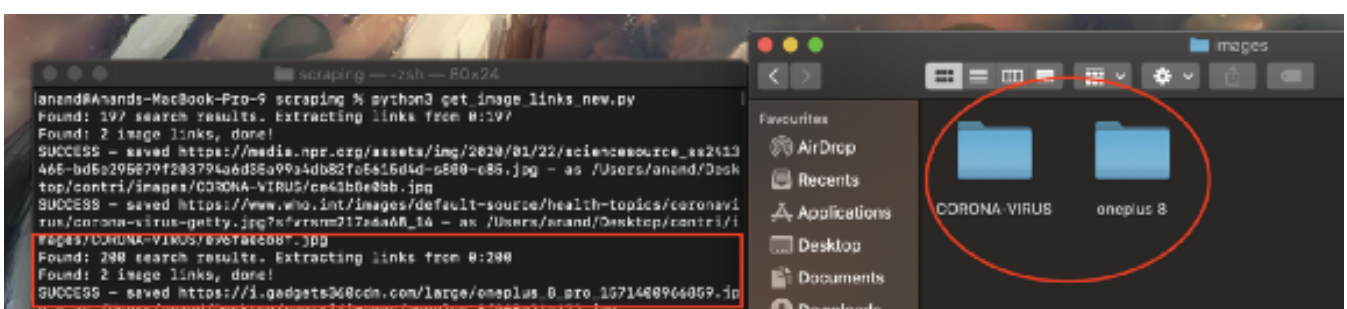
```

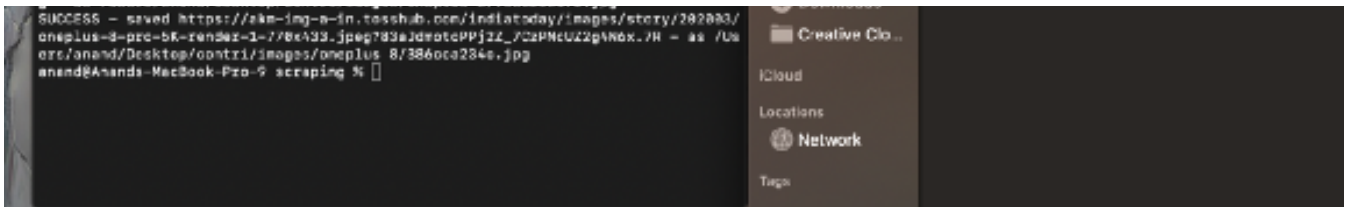
```
persist_image(folder_path:str, file_name:str, url:str)
```

persist_image function has 3 arguments which are all mandatory for the function call.

1. 'folder_path' will be the path to a common folder to where you want to save all the images, mine is "/Users/anand/Desktop/contri/images".
2. 'file_name' this is a str value which is of your choice, but personally i prefer passing file_name = query, because it'll create separate folders with the query as title.
3. 'url' is a string input which we get from the returned values of 'fetch_image_urls' function.

I ran the whole program with all snippets put together with two queries, "coronavirus" and "OnePlus 8". The result is shown below





Created Folders with Respective query images

. . .

Thank You for reading!

Just run the whole code and watch the magic happen :)

libraries to import:

```
import selenium from selenium
import webdriver
import time
import requests
import os
from PIL import Image
import io
import hashlib
```

execution:

```
$ python3 complete_program.py
```

```
1  import selenium
2  from selenium import webdriver
3  import time
4  import requests
5  import os
6  from PIL import Image
7  import io
8  import hashlib
9  # This is the path I use
10 #DRIVER_PATH = '/Users/anand/Desktop/chromedriver'
11 # Put the path for your ChromeDriver here
12 DRIVER_PATH = <enter_your_path>
13
14
15 def fetch_image_urls(query:str, max_links_to_fetch:int, wd:webdriver, sleep_between_interactions:int):
16     def scroll_to_end(wd):
17         wd.execute_script("window.scrollTo(0, document.body.scrollHeight);")
18         time.sleep(sleep_between_interactions)
19
```



```
17
20 # build the google query
21 search_url = "https://www.google.com/search?safe=off&site=&tbm=isch&source=hp&q={q}&oq={q}"
22
23 # load the page
24 wd.get(search_url.format(q=query))
25
26 image_urls = set()
27 image_count = 0
28 results_start = 0
29 while image_count < max_links_to_fetch:
30     scroll_to_end(wd)
31
32     # get all image thumbnail results
33     thumbnail_results = wd.find_elements_by_css_selector("img.Q4LuWd")
34     number_results = len(thumbnail_results)
35
36     print(f"Found: {number_results} search results. Extracting links from {results_start}:")
37
38     for img in thumbnail_results[results_start:number_results]:
39         # try to click every thumbnail such that we can get the real image behind it
40         try:
41             img.click()
42             time.sleep(sleep_between_interactions)
43         except Exception:
44             continue
45
46         # extract image urls
47         actual_images = wd.find_elements_by_css_selector('img.n3VNCb')
48         for actual_image in actual_images:
49             if actual_image.get_attribute('src') and 'http' in actual_image.get_attribute('src'):
50                 image_urls.add(actual_image.get_attribute('src'))
51
52         image_count = len(image_urls)
53
54         if len(image_urls) >= max_links_to_fetch:
55             print(f"Found: {len(image_urls)} image links, done!")
56             break
57     else:
58         print("Found:", len(image_urls), "image links, looking for more ...")
59         time.sleep(30)
60         return
61         load_more_button = wd.find_element_by_css_selector(".mye4qd")
62         if load_more_button:
63             wd.execute_script("document.querySelector('.mye4qd').click();")
64
65     # move the result startpoint further down
66     results_start = len(thumbnail_results)
```

```
67
68     return image_urls
69
70 def persist_image(folder_path:str,file_name:str,url:str):
71     try:
72         image_content = requests.get(url).content
73
74     except Exception as e:
75         print(f"ERROR - Could not download {url} - {e}")
76
77     try:
78         image_file = io.BytesIO(image_content)
79         image = Image.open(image_file).convert('RGB')
80         folder_path = os.path.join(folder_path,file_name)
81         if os.path.exists(folder_path):
82             file_path = os.path.join(folder_path,hashlib.sha1(image_content).hexdigest()[:10])
83         else:
84             os.mkdir(folder_path)
85             file_path = os.path.join(folder_path,hashlib.sha1(image_content).hexdigest()[:10])
86         with open(file_path, 'wb') as f:
87             image.save(f, "JPEG", quality=85)
88         print(f"SUCCESS - saved {url} - as {file_path}")
89     except Exception as e:
90         print(f"ERROR - Could not save {url} - {e}")
91
92 if __name__ == '__main__':
93     wd = webdriver.Chrome(executable_path=DRIVER_PATH)
94     queries = ["CORONA-VIRUS","oneplus 8"] #change your set of queries here
95     for query in queries:
96         wd.get('https://google.com')
```

Get the Medium app

