

# **Configuration Manual**

MSc Research Project Artificial Intelligence

Praneeth Bandi Student ID: 22183922

School of Computing National College of Ireland

Supervisor:

Rejwanul Haque

#### National College of Ireland





#### **School of Computing**

Student Name:	Praneeth Bandi				
Student ID:	22183922				
Programme:	MSc in Artificial Intelligence	Year:	2023		
Module:	MSc Research Method				
Lecturer:	Rejwanul Haque				
Date:	05/01/2023				
Project Title:	Advanced Google Scholar Scraper: A Content-Based Filtering Approach for Literature Recommendation Using BERT Embeddings				

#### Word Count: 1438 Page Count: 10

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

<u>ALL</u> internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

Signature: Praneeth Bandi

**Date:** 05/01/2024

#### PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST

Attach a completed copy of this sheet to each project (including multiple copies)	
Attach a Moodle submission receipt of the online project submission, to each project (including multiple copies).	
You must ensure that you retain a HARD COPY of the project, both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer.	

Assignments that are submitted to the Programme Coordinator Office must be placed into the assignment box located outside the office.

Office Use Only	
Signature:	
Date:	
Penalty Applied (if applicable):	

## Configuration Manual: Advanced Google Scholar Scraper: A Content-Based Filtering Approach for Literature Recommendation Using BERT Embeddings

Praneeth Bandi Student ID: 22183922

### **1** Introduction

Welcome to the user guide for Advanced Google Scholar Scraper. This guide is intended to help users install and get comfortable with our advanced tool for effective literature search. A complete literature recommendation system The Advanced Google Scholar Scraper combines web scraping, natural language processing (NLP), and an easy-to-use graphical interface into a single streamlined package. However, with the growing amount of scholarly content scattered about the Internet, our scraper equips researchers and students everywhere to quickly cut through Google Scholar's gargantuan heap.

This manual will guide you through installation, system requirements and configuration settings. Setting the program up couldn't be easier. How to wield the power of BERT embeddings and cosine-based similarity metrics, so you can receive highly accurate and relevant literature recommendations. Tkinter-based GUI: explore its user friendly intuitive interface.

You can use this application to study Natural Language Processing, Machine Learning, BERT Models and more. Its cross-discipline robustness means it is equally at home in all such fields. Jump into the guide to make full use of Advanced Google Scholar Scraper, and keep an eye open for further developments as we continue to add functionality.

### 2 System Requirements

Before setting up the Advanced Google Scholar Scraper, ensure that your system meets the following requirements for optimal performance:

Operating System	Windows 10 or Later
Programming Language	Python
Integrated Development Environment (IDE)	PyCharm
Memory (RAM)	4GB+
Graphical Processing Unit (GPU)	If Available
Web Application	Chrome
Processor	Minimum I3

Operating System:

• Windows 10 or later.

Programming Language:

• Python: The scraper is developed in Python, and the system should have Python installed. Download the latest version from python.org.

Integrated Development Environment (IDE):

• PyCharm: We recommend using PyCharm as the preferred IDE for running and managing the scraper. Download the latest version from JetBrains.

Memory (RAM):

• 4GB or higher: Ensure your system has a minimum of 4GB RAM for smooth execution. Higher RAM capacity will enhance performance, especially when dealing with large datasets.

Graphics Processing Unit (GPU):

• If available: While not mandatory, having a GPU can potentially accelerate certain processes, providing a more efficient experience.

Web Application:

• Google: The scraper interacts with Google Scholar, so a stable internet connection and access to Google services are essential.

Processor:

• Intel Core i3 or equivalent: The processor should be at least an Intel Core i3 or equivalent for effective execution of the scraping tasks.

Ensure that your system components align with these specifications to fully leverage the capabilities of the Advanced Google Scholar Scraper.

### **3** Installation Instructions

Follow these steps to set up the Advanced Google Scholar Scraper on your system:

1. Install Python:

Ensure Python is installed on your system. Download the latest version from python.org.

- 2. Install PyCharm: Download and install PyCharm, the recommended Integrated Development Environment (IDE), from JetBrains.
- 3. Install Required Libraries: Open a terminal or command prompt and run the following command to install the necessary Python libraries:

#### pip install pandas requests torch beautifulsoup4 selenium scikit-learn transformers

- 4. Download ChromeDriver: Download the ChromeDriver executable suitable for your Chrome browser version from <u>ChromeDriver Downloads</u>.
- 5. Extract and Add ChromeDriver to PATH: Extract the downloaded ChromeDriver executable and add its location to your system's PATH environment variable.

- 6. Open PyCharm: Open PyCharm and create a new Python project.
- Copy and Run the Code: Copy the provided Python code into a Python file within your PyCharm project. Run the script to ensure it executes without errors.
- 8. Install Missing Libraries: If any library is missing, use the following command in PyCharm's terminal to install it:

#### pip install library\_name

- 9. Configure WebDriver: Ensure the webdriver.Chrome() call is compatible with your Chrome browser version.
- 10. Run the Scraper:

Run the script, and the Tkinter GUI for the Advanced Google Scholar Scraper will appear.

- 11. Enter Query and Settings: Input the desired article title or keyword, specify the number of pages to scrape, and choose an output folder. Optionally, browse for an output folder.
- 12. Click "Extract Data": Click the "Extract Data" button to initiate the scraping process. Monitor the progress through the displayed progress bar.
- 13. Completion Message: Once the extraction is complete, a status message will indicate success. The extracted data will be saved in a CSV file.
- 14. Close the Application: Close the application when finished.

Now you have successfully set up and used the Advanced Google Scholar Scraper on your system. Adjust the code as needed for your specific requirements.

Note: Ensure you have a stable internet connection and access to Google services for the scraper to function properly.

### 4 Installation Demo

#### **1. PyCharm Installation**



Figure 1



Figure 3

🖺 PyCharm Edu Setup		-	
	Choose Install Location Choose the folder in which to insta	ill PyCharm Ed	u.
Setup will install PyCharm E Browse and select another	du in the following folder. To install folder. Click Next to continue.	in a different f	folder, click
Destination Folder			
C:\Program Files\JetBra	ains\PyCharm Edu 2022.2.2	Bro	wse
Space required: 1.4 GB Space available: 21.9 GB			
	< Back	Next >	Cancel
	Figure 4		
🚯 PyCharm Edu Setup		-	
<b>(</b>	Installing Please wait while PyCharm Edu is	being installe	ed.
Extract: app.jar 69%			
Show details			

Figure 5

Click on the .exe file and complete the Installation as in Figure 3, 4 and 5.After the Setup is completed, open the folder of the Code and import the Scholar Scraping Code.

Elle Edit View Navigate Code VCS Help		
	t — 🚯 tinat.py =	
Praneeth E\Current Projects\Praneeth		
Y in foutput		
Bert Model.csv		
Machine Learning.csv		
Natural Language Processing.csv		
R Denal Present & doors		
Im External Libraries		
Scratches and Consoles		
	10 from selenium import webdriver	
	11 ferm selenium.webdriver.common.by import By	
	<pre>16 tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')</pre>	
	17 model = BertModel.from.pretrained('bert-base-uncased')	
	<pre>19 Ddef get_bert_embedding(text):</pre>	
	input ids = tokenizer.encode(text. return.tensori="ot", truncation=True)	
	21 with torch no grad():	
	22 autouts = model(input ids)	
	21 return outputs last hidden state.mean(dim=1).soupeze().tolist()	
	def cosine stallarity score/embedding1 embedding2).	
	nting option similarity (Lashaddiani) Lashaddiani)10161	
	A same to exclantioner stationifier. filoneral-	
	and because of the second seco	
	atth open(itemane, w., Meetine , encounty of a covite.	
	30 Tietunanes = [ ifte, Autors, Can, Content, Salturity ] 31 moltan a sex Distanticonfilm [ independent filment]	
	anter - csv.otcunter(csvire, ris(chines=ris(chines))	
🕑 Version Control 🛛 🛛 Problems 🖾 Terminal 🍩 P	ython Packages 🔹 👁 Python Console 💽 Services	
Download pre-built shared indexes: Reduce the indexing	time and CPU load with pre-built Python packages shared indexes // Always dow (moments a Discovering binary modules	12:29 LF UTF-8 4 spaces Python 3.10 💁 🚭

Figure 6

After Importing the Window will look like as in Figure 6.



Figure 7

Now open the Terminal, Either by Alt+F12 or Clicking on the terminal in bottom left as in Figure 7.



Figure 8

After opening the Terminal, Copy and Paste the pip installation command as in 3.3 Install Required Libraries or Just by simply copy, pasting below command as in Figure 8 **pip install pandas requests torch beautifulsoup4 selenium scikit-learn transformers** And hit enter.

Now All the Installation work is completed.

### 5 How to Run

The Advanced Google Scholar Scraper gives users the option of fine-tuning many different configuration items to suit their own taste. Follow these guidelines to customize the configuration settings:

To Run:

Open Terminal

#### Figure 9

• Type python filename.py such as python final.py and hit enter as in Figure 9.

#### Query Configuration:

멍	📻 Project 🔻								: .
Proj	<ul> <li>Praneeth E\Current Projects\/</li> <li>foutput</li> <li>Bert Model.csv</li> <li>Machine Learning.csv</li> </ul>		4 from tkin 5 6 import pa 7 import re	t <mark>er import</mark> ndas as pd <b>quests</b>	filedialog, messagebox, ttk				Notifications
>	ig Natural Language Proc i‰ final.py i‰ final Report 1.docx i Illi External Libraries i Scratches and Consoles		8 import to 9 from bs4 10 from sele 11 from sele 12 from skle 13 Ø Advanced G	rch import Beau nium import nium.webdr; arn.metrics pogle Scholar Sc	utifulSoup t webdriver iver.common.by import By s.pairwise import cosine_simil caper —	arity			
			14 15 Article Title or 16 Number of 17 18 Output Folder	Keyword: Pages: optional):		Browse	)		
			19 20 21 22		Extract Data		truncation=True)		
т	erminal: Local × 🕂 🗸								
s R f R e	(from sympy->torch) (1. equirement already satis rom cffi>=1.14->trio~=0. equirement already satis s (from wsproto>=0.14->t	3.0) fied: pycparser i 17->selenium) (2. fied: h11<1,>=0.9 rio-websocket~=0.	n c:\us 21) .0 in c:\users\shou 9->selenium) (0.14)	t\appdata∖ 0)	local\packages\pythonsoftware	Foundation.	in.3.10_qbz5n2kfra8p0\localcache\local-packages\p .python.3.10_qbz5n2kfra8p0\localcache\local-packag		
[ [ P	notice] A new release of notice] To update, run: S E:\Current Projects\Pr								
_		🛛 Terminal 📚 Python P	Packages 🛛 🏶 Python Consol	Services					10 2 🔿

#### Figure 10

- Locate the entry box labeled Article Title or Keyword in the GUI as in Figure 10.
- Type in your preferred article title or key word to formulate the search term.
- Modify the query to suit your own needs.

Number of Pages to Scrape:

- In the GUI, find the "Number of Pages" entry box.
- Select the number of pages to scrape from Google Scholar.
- Tune the number to the extent of literature search you require.

**Output Folder Specification:** 

- In the GUI, find the entry box for "Output Folder (optional)."
- You can optionally specify the folder in which the extracted data will be stored.
- Click the "Browse" button to browse and choose an output folder.
- If there is no output folder specified, then the script will save the data in current working directory.

Through these simple steps, users can successfully customize the Advanced Google Scholar Scraper to meet their research needs. An easy-to-use GUI and clearly written code that is well documented make customization and adjustment simple, even for technically illiterate users.

### References

Python Software Foundation. (n.d.). Python.org. Retrieved from https://www.python.org/

JetBrains.(n.d.).PyCharm:Download.Retrievedfromhttps://www.jetbrains.com/pycharm/download/?section=windows

Chromium Projects. (n.d.). ChromeDriver - WebDriver for Chrome. Retrieved from <u>https://chromedriver.chromium.org/downloads</u>