

Delete scripts/.gitkeep

Update README.md

Add files via upload

Add files via upload

Update LICENSE

Scholarly_Paper_Crawler

scripts

LICENSE

inputs.txt

README.md

webcrawler_paper_search.bat

7 months ago

7 months ago

4 months ago

7 months ago

7 months ago

Scholarly_Paper_Crawler

Author: Suk Jin Mun

Version: 1.0.0 Year: 2024 License: MIT

Email: rOysJmUN[at]gMail.com (all lower case letters)

This is an automated tool designed to streamline the process of searching, retrieving, and analyzing academic papers from Google Scholar. It automates the collection of research data, including paper details, PDF downloads, and compound identification from abstracts.

For more information or to report issues, please contact the author.

Table of contents

- 1.0 About
- 2.0 Important Notice
- 3.0 Directory Structure
 - 3.1 html_parsing
 - 3.2 arXiv_xml
 - 3.3 pdf_first_100_sentences
- 4.0 Pipeline Manual
 - 4.1 Setup
 - 4.2 Running the Pipeline
 - 4.3 Checking Output
- 5.0 Features
 - 5.1 Spin Detection
 - 5.2 Keyword Permutation
- 6.0 Files Explained
 - 6.1 inputs.txt
 - o 6.2 webcrawler_paper_search.bat
 - 6.3 pdf_files
 - 6.4 csv_files
 - 6.5 scripts
- 7.0 API Management
 - 7.1 Free Trial Accounts
 - 7.2 Paid Accounts
 - 7.3 Usage Monitoring
 - 7.4 Best Practices
- <u>8.0 Notes</u>
- 9.0 Known Issues and Future Updates
- 10.0 License
- 11.0 Credits

1.0 About

This pipeline automates the search and retrieval of research paper metadata based on specified keywords. It uses Google Scholar to gather information on research papers, including links and author details, and attempts to download PDFs from Sci-Hub. Designed primarily for studying and researching materials, this pipeline is especially useful for material science purposes. However, it is versatile enough to support other research areas as well, as it retrieves PDFs based on the keywords provided. At the end of the process, the pipeline organizes and stores the information in a separate XLSX file for easy reference.

2.0 Important Notice

This pipeline performs automated web scraping that may be flagged as bot-like behavior by academic servers. Running it on institutional networks (universities or research institutes) risks IP bans that could affect all users on that network and disrupt legitimate research activities. Instead, use a private environment (home network, personal VPN) to protect both your institution and ensure successful data collection. This approach prevents potential institutional policy violations, avoids administrative issues with IT departments, and maintains uninterrupted access to academic resources for all users. For optimal performance and security, it is strongly recommended to use a reliable VPN service with IP rotation capabilities - this provides an additional layer of protection for both your personal IP and API usage. This helps preventing potential blocks from Google Scholar and other academic servers while maintaining consistent access to the ScraperAPI service. Users must thoroughly review this documentation before using the pipeline, as a complete understanding is essential for successful operation of this pipeline.

3.0 Directory Structure

Directory Structure

```
root directory/
├─ inputs.txt
                                     # Configuration file for input parameters
— webcrawler_paper_search.bat
                                     # Batch file to execute the Python script
  pdf_files/
                                     # Directory to store downloaded PDF files
                                    # Subfolder for each search session, named based on input
   └─ [search_session_name]/
   pdf_first_100_sentences/
       LastName[Year]_first100.txt # Text files named after their source PDFs
 — csv_files/
                                    # Directory to store generated result files
  [search_session_name].xlsx
                                   # Excel file for each search session, named based on input
├─ html_parsing/
                                    # Directory to store raw HTML from Google Scholar searches
  # Directory to store raw HIML From Google Scholar Searches

[search_session_name]/ # Subfolder containing raw HTML files for each search result
                                    # Directory to store XML responses from arXiv API
├─ arXiv_xml/
   | Subfolder containing XML files from arXiv queries | # Subfolder containing XML files from arXiv queries
   scripts/
                                     # Directory for Python scripts
   webcrawler_paper_search.py # Main Python script to run the pipeline
```

3.1 html parsing

- Stores raw HTML content from Google Scholar search results
- Each search session has its own subfolder
- Helps with debugging and analysis of search results
- Files are named as paper_N_raw.html where N is the result number

3.2 arXiv xml

- Stores XML responses from arXiv API queries
- Each search session has its own subfolder
- Contains detailed metadata about papers found on arXiv
- Useful for tracking and debugging arXiv search results

3.3 pdf_first_100_sentences

- A directory that stores text files containing the first 100 sentences from each PDF
- Each search session has its own subfolder named after the session

- Text files are named to match their source PDFs (e.g., Duan2021_first100.txt)
- Used for guick content preview and text analysis without opening PDFs
- Files contain extracted and cleaned text from the beginning of each paper

4.0 Pipeline Manual

4.1 Setup

First, go to https://www.scraperapi.com/signup and create a free account. They offer 5,000 free API requests per month. After signing up, you will be directed to https://dashboard.scraperapi.com/. Retrieve your API Key in the API Key section. Then, Enter search parameters in the inputs.txt file. Format example:

```
O> Input API Key
'API KEY YOU OBTAINED FROM SCRAPERAPI.COM'

1> Specify the name for this search session (e.g., "1D_AFM_material_search"):
1D_AFM_material_search

2> Specify the maximum number of searches to perform for each keyword combination (e.g., "10"):
3

3> Provide keywords for the search (e.g., "1D antiferromagnetic chain system"):
one dimensional antiferromagnetic chain system, 1D antiferromagnetic chain system, quasi-1d AFM

4> Search for compound name information? (Y/N) (e.g., "Y"):
Y
```

Note: Processing 45 papers (15 keyword combinations with 3 searches each) typically takes approximately 5 hours due to mandatory waiting periods and rate limiting

4.2 Running the Pipeline

- Execute webcrawler_paper_search.bat to install required dependencies and launch the script automatically.
- The script reads from inputs.txt, performs the search on Google Scholar, attempts to download PDFs from Sci-Hub, and saves results into an Excel file within csv_files/.

4.3 Checking Output

- The Excel file will contain titles, authors, publication year, keywords, Google Scholar links, and PDF filenames.
- You can explore and analyze the results from the csv_files/ directory.
- Downloaded PDFs can be found in the pdf_files/[search_session_name]/ directory.

5.0 Features

5.1 Spin Detection

The pipeline includes automated detection of quantum mechanical spin values from papers. This feature:

- Extracts physically valid spin values following quantum mechanical principles
- Supports both integer spins (S=0, S=1, S=2) and half-integer spins (S=1/2, S=3/2, S=5/2)
- Validates spin values against quantum mechanical constraints:
 - Only accepts standard quantum spin values (0, 1/2, 1, 3/2, 2, 5/2)
 - Enforces proper denominator (only /2 for half-integer spins)
 - o Rejects non-physical spin values
 - Limits maximum spin value to 2 (as per Standard Model constraints)
- Searches for spin values in both paper titles and PDF content
- Reports "No valid quantum spin value found" when no valid spin is detected

The spin values are extracted in formats:

- S = X/Y (with space)
- S=X/Y (without space)
- S=X (integer values) Where X and Y are integer numbers following quantum mechanical constraints.

5.2 Keyword Permutation

The pipeline implements comprehensive keyword combination searching that examines all possible permutations of the input keywords. This feature:

- · Generates all possible combinations of input keywords, including individual terms and their permutations
- · Searches each combination separately in Google Scholar
- Performs the specified number of searches (max_results) for EACH combination
- Maintains a wait period between combinations to comply with rate limits
- · Automatically skips duplicate papers found across different combinations

For example, if keywords "A, B, C" are provided (line 3> of inputs.txt), the pipeline will search for:

- 1. Individual keywords: "A", "B", "C"
- 2. Pairs with all permutations: "A B", "A C", "B A", "B C", "C A", "C B"
- 3. Triples with all permutations: "A B C", "A C B", "B A C", "B C A", "C A B", "C B A"

The total number of combinations grows factorially with the number of keywords:

```
Number of combinations for different keyword counts:

2 keywords: 3 combinations

3 keywords: 15 combinations

4 keywords: 64 combinations

5 keywords: 325 combinations
```

Example of search behavior:

Input: 2> Specify the maximum number of searches to perform: 10 3> Provide keywords: A, B, C

Result:

• Will perform 10 searches for "A"

6 keywords: 1956 combinations

- Will perform 10 searches for "B"
- Will perform 10 searches for "C"
- Will perform 10 searches for "A B"
- Will perform 10 searches for "A C"
- Will perform 10 searches for "B A" ...and so on for all 15 combinations. In such case, a total of 15x10 papers should be processed.

Each combination is treated as a separate search query, with the specified maximum number of results collected for each. This ensures thorough coverage of the literature while avoiding duplicate entries. The script automatically handles:

- Tracking results per combination
- · Removing duplicate papers across combinations
- · Maintaining proper delays between searches
- · Recording which search terms found each paper

6.0 Files Explained

6.1 inputs.txt

- This file contains the parameters for the search pipeline.
- Parameters include the search session name, the maximum number of search results to retrieve, and keywords.

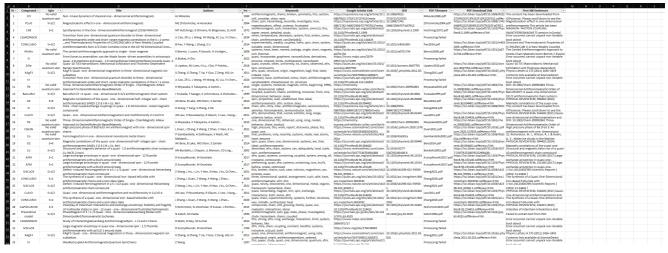
6.2 webcrawler_paper_search.bat

- A batch file that automates the execution of the webcrawler_paper_search.py script.
- It installs required Python packages and runs the main script.

6.3 pdf_files\

- · A directory that stores PDFs associated with the search results.
- Each search session creates a subfolder based on the session name.
- The script attempts to download PDFs from Sci-Hub.
- PDFs are named using the format: LastName[Year].pdf (e.g., Duan2021.pdf)

6.4 csv_files\



- This directory contains Excel files generated from each search session.
- Each session creates a file named after the session with detailed metadata including:
 - o Paper titles, authors, and year
 - Extracted compound information
 - Quantum spin values (following physical constraints)
 - Keywords and links
 - First 100 sentences from PDFs

6.5 scripts\

· The main Python script that handles reading inputs, querying Google Scholar, downloading PDFs from Sci-Hub, and







7.0 API Management

7.1 Free Trial Accounts

- Initial free credits: 5,000 requests
- Trial duration: 7 days from account creation
- · Credits are one-time only and do not renew
- After 7 days or using all credits (whichever comes first):
 - o Remaining credits expire
 - API key stops working
 - Must upgrade to paid plan to continue

7.2 Paid Accounts

- API keys remain active as long as account is in good standing
- · Credits reset monthly based on subscription plan

• The key itself doesn't expire unless manually revoked

7.3 Usage Monitoring

- 1. Through ScraperAPI Dashboard (https://www.scraperapi.com/dashboard):
 - View current credit balance
 - Monitor daily/monthly usage statistics
 - o Check plan status and renewal date
 - Track remaining trial days (for free accounts)
- 2. Through Console Messages: Common API Status Messages:

```
"Error in API request: 401 Client Error: Unauthorized"

→ API key is invalid or credits depleted
"Error in API request: 429 Too Many Requests"

→ Temporary rate limit, script will automatically retry
"Success"

→ API key is valid and has available credits
```

9

7.4 Best Practices

- Each Google Scholar search page uses 1 credit
- Typically 10 results per page
- Example: 100 papers might use 10-20 credits

For Free Trial:

- Plan searches carefully within 7-day window
- · Start with smaller searches to test the system
- Save important searches for when familiar with tool

For All Users:

- · Check credit balance before starting large search sessions
- Monitor console output for API-related error messages
- · Keep track of credit usage for planning
- · Consider upgrading if you need regular access
- Save API usage statistics for future planning

8.0 Notes

- The script saves raw HTML from Google Scholar searches for analysis
- ArXiv is used as an alternative source when Sci-Hub download fails
- The script implements exponential backoff and retry mechanisms for robust web scraping
- Added better error handling and logging for debugging purposes
- The script uses cloudscraper to bypass DDoS protection
- · Downloads have mandatory delays (60 seconds minimum) between requests to comply with rate limits
- Spin values are validated against quantum mechanical principles before being included in the output
- Due to the factorial growth of combinations, users should be mindful when using more than 4 keywords, as this can lead to very long execution times. Each combination requires its own set of searches with appropriate waiting periods to avoid overloading the search servers.
- This project is licensed under the MIT License see the LICENSE file for details.

9.0 Known Issues and Future Updates

Ver 1.0.0

5/28/25, 1:34 AM SukjinMun/Scholarly_Paper_Crawler: An automated scholarly literature pipeline that systematically searches, ...

Issues

- Current compound name and spin value extraction algorithms exhibit reduced accuracy. This indicates the pipeline requires further refinement of the extraction logic.
- The pipeline currently lacks the capability to extract and classify material properties such as crystal structure (single crystal/powder form), synthesis methods, and dimensional characteristics.
- PDF retrieval success rate is inconsistent due to intermittent connectivity issues with Sci-Hub and arXiv APIs, which requires more robust fallback mechanisms.

Future Updates

Releases

No releases published

Packages

No packages published

Languages

Python 98.8%Batchfile 1.2%