



1/6/2025

PubMed Paper Fetcher

Aganitha - Python & DevOps



Rakesh D

B. TECH CSE

PES UNIVERSITY

PubMed Paper Fetcher

PubMed Paper Fetcher: Approach, Methodology, and Results

1. Introduction

This report describes the approach, methodology, and results of a Python-based program that fetches research papers from the PubMed API based on a user-defined query. The program identifies papers where at least one author is affiliated with a pharmaceutical or biotech company and saves the results in a CSV file.

2. Approach

The project follows a systematic approach:

1. **Understanding the Requirements:** The task involves querying PubMed, filtering relevant papers, extracting required details, and saving the data.
2. **Using the PubMed API:** The **Entrez Programming Utilities (E-utilities)** provided by NCBI are used to fetch paper details.
3. **Processing Data:** Extract relevant metadata such as PubMed ID, title, publication date, author affiliations, and corresponding author emails.
4. **Saving Data:** The final extracted data is formatted and stored in a CSV file for easy analysis.
5. **Command-line Interface:** The script is built as a command-line tool supporting query input, debugging, and file output options.

3. Methodology

3.1 Tools and Technologies Used

- **Python:** Main programming language.
- **Requests Library:** Used for making HTTP requests to the PubMed API.
- **Pandas:** Used for data processing and exporting results to CSV.
- **Poetry:** Used for dependency management and packaging.
- **Git & GitHub:** Used for version control and project management.

3.2 Implementation Steps

Step 1: Setting Up the Project Environment

- Installed necessary dependencies using Poetry.
- Created a Python script (fetch_papers.py) to fetch and process data.

Step 2: Fetching Papers from PubMed

- Constructed an API request using requests.
- Used the `esearch.fcgi` endpoint to search for papers based on user queries.

- Retrieved a list of relevant PubMed IDs.

Step 3: Fetching Paper Details

- Used the esummary.fcgi API to retrieve metadata for the obtained PubMed IDs.
- Extracted the title, publication date, and other relevant details.

Step 4: Filtering Author Affiliations

- Attempted to extract non-academic authors and company affiliations.
- Due to API limitations, advanced NLP techniques can be used in future iterations.

Step 5: Saving Results

- Used Pandas to store extracted data in a structured CSV format.
- Implemented a --file option to allow users to specify output filenames.

Step 6: Command-line Functionality

- Implemented argparse to accept user queries.
- Added debugging mode (--debug) to print additional execution details.

Step 7: Packaging and Version Control

- Configured pyproject.toml for Poetry packaging.
- Added GitHub integration for code sharing and collaboration.

4. Results

4.1 Execution Example

To fetch papers related to "cancer research" and save results in papers.csv:

```
poetry run get-papers-list "cancer research" -f papers.csv
```

4.2 Sample Output (CSV Format)

PubmedID	Title	Publication Date	Non-academic Author(s)	Company Affiliation(s)	Corresponding Author Email
12345678	Cancer Study 1	2023-05-10	Dr. John Doe	ABC Biotech Inc.	john.doe@example.com
87654321	Immunotherapy Advances	2022-11-22	Dr. Jane Smith	XYZ Pharma Ltd.	jane.smith@example.com

5. Challenges and Future Improvements

5.1 Challenges

- Limited author affiliation data:** The PubMed API does not directly provide detailed author affiliations, requiring additional processing.

- **Rate Limits:** The API enforces request limits, requiring optimization in querying.

5.2 Future Improvements

- **Use NLP for Affiliation Extraction:** Implement Named Entity Recognition (NER) to better identify pharmaceutical/biotech affiliations.
- **Database Storage:** Store results in a relational database for further analysis.
- **Web-based Interface:** Develop a user-friendly web interface to run queries and view results interactively.

6. Conclusion

The PubMed Paper Fetcher successfully retrieves research papers based on user queries, processes relevant details, and saves them in CSV format. While certain API limitations exist, the script provides a solid foundation for automating research paper retrieval and analysis. Future enhancements can improve accuracy and usability.