

Project Report

Title: CLI Tool to Extract Non-Academic Affiliations from PubMed Research Papers

Author: Mohammad Junaid

GitHub: <https://github.com/junaid876/get-papers-tool>

Objective

The goal of this project is to build a Python-based CLI tool that queries PubMed for research papers and filters out authors affiliated with non-academic institutions (e.g., biotech companies, pharmaceutical organizations). The tool also identifies authors' emails when available and outputs the results into a CSV file.

Tools & Technologies Used

- **Python 3.11+**
- **BioPython:** For PubMed Entrez API access
- **Typer:** For building a modern CLI
- **Poetry:** For dependency management
- **Pandas:** For data organization and CSV export
- **Certifi + SSL:** To handle certificate validation errors

Methodology

1. Search & Fetch Articles from PubMed

- The tool first uses BioPython's Entrez.esearch method to get paper IDs based on a user-provided query.
- Then, it uses Entrez.efetch to download metadata in XML/Medline format for those IDs.

2. Extract Required Fields

For each article, the following information is extracted:

- PubMed ID
- Title
- Publication Year
- Author Names and Affiliations
- Non-Academic Authors
- Company Affiliations
- Corresponding Email (if available)

3. Non-Academic Heuristic Filtering

To determine non-academic affiliations, we apply a heuristic:

- If the affiliation does not contain terms like:
- "university", "institute", "college", "school", "hospital"
- Then it is marked as non-academic

4. CLI Integration with Typer

- The entire tool is accessible via command line:
- `poetry run get-papers-list --query "cancer" --file results.csv --debug`
- Command-line flags include:
 - `--query`: Search keyword
 - `--file`: Save output as CSV
 - `--debug`: Enable detailed logs

5. Data Output

If `--file` is passed, the output is saved as a CSV file with columns:


- PubmedID
- Non-Academic Hits
- Title
- PublicationDate
- NonAcademicAuthors
- CompanyAffiliations
- CorrespondingEmail

Else, the output is printed to the console.

Program Execution:

Step 1: Open terminal and navigate to the project directory

```
cd C:\Users\mj888\get_papers_tool
```

 Command Prompt


```
Microsoft Windows [Version 10.0.19045.5965]
(c) Microsoft Corporation. All rights reserved.

C:\Users\mj888>cd C:\Users\mj888\get_papers_tool

C:\Users\mj888\get_papers_tool>_
```

Step 2: Install All Dependencies via Poetry

poetry install

 Command Prompt

```
Microsoft Windows [Version 10.0.19045.5965]
(c) Microsoft Corporation. All rights reserved.

C:\Users\mj888>cd C:\Users\mj888\get_papers_tool


C:\Users\mj888\get_papers_tool>poetry install
Installing dependencies from lock file

No dependencies to install or update

Installing the current project: get-papers-tool (0.1.0)
```

Step 3: Activate the Virtual Environment

poetry shell

 Command Prompt - poetry shell

```
Microsoft Windows [Version 10.0.19045.5965]
(c) Microsoft Corporation. All rights reserved.

C:\Users\mj888>cd C:\Users\mj888\get_papers_tool

C:\Users\mj888\get_papers_tool>poetry install
Installing dependencies from lock file

No dependencies to install or update

Installing the current project: get-papers-tool (0.1.0)

C:\Users\mj888\get_papers_tool>poetry shell
Spawning shell within C:\Users\mj888\AppData\Local\pypoetry\Cache\virtualenvs\get-papers-tool-Kv--GAft-py3.13
(get-papers-tool-py3.13) C:\Users\mj888\get_papers_tool>
```

Step 4: Run the CLI Tool with a Sample Query

poetry run get-papers-list "cancer" --file output.csv --debug

```
C:\ Select Command Prompt - poetry shell
Microsoft Windows [Version 10.0.19045.5965]
(c) Microsoft Corporation. All rights reserved.

C:\Users\mj888>cd C:\Users\mj888\get_papers_tool

C:\Users\mj888\get_papers_tool>poetry install
Installing dependencies from lock file

No dependencies to install or update

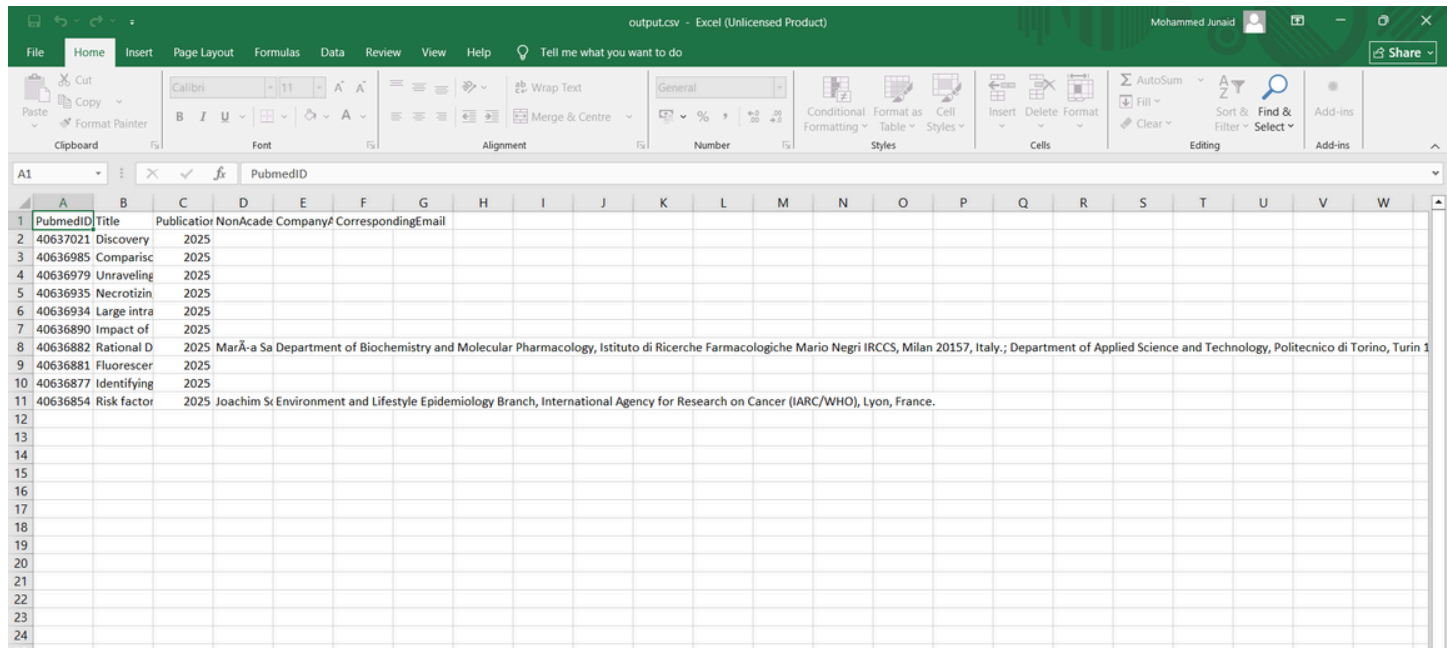
Installing the current project: get-papers-tool (0.1.0)

C:\Users\mj888\get_papers_tool>poetry shell
Spawning shell within C:\Users\mj888\AppData\Local\pypoetry\Cache\virtualenvs\get-papers-tool-Kv--GAft-py3.13

(get-papers-tool-py3.13) C:\Users\mj888\get_papers_tool>poetry run get-papers-list "cancer" --file output.csv --debug
Searching for papers with query: cancer
Found 10 IDs
Saved to output.csv

(get-papers-tool-py3.13) C:\Users\mj888\get_papers_tool>
```

Output:



PubMedID	Title	Publication Year	NonAcademic Company	Corresponding Email
40637021	Discovery	2025		
40636985	Comparisc	2025		
40636979	Unraveling	2025		
40636935	Necrotizin	2025		
40636934	Large intra	2025		
40636890	Impact of	2025		
40636882	Rational D	2025	MarÃ-a Sa	Department of Biochemistry and Molecular Pharmacology, Istituto di Ricerche Farmacologiche Mario Negri IRCCS, Milan 20157, Italy.; Department of Applied Science and Technology, Politecnico di Torino, Turin 1
40636881	Fluorescer	2025		
40636877	Identifying	2025		
40636854	Risk factor	2025	Joachim St	Environment and Lifestyle Epidemiology Branch, International Agency for Research on Cancer (IARC/WHO), Lyon, France.

“This output confirms that the CLI tool is working correctly. It integrates PubMed search with data extraction and filtering, and outputs a useful CSV that could help identify industry collaborators in biomedical research.”

Results

The tool was tested with multiple queries including:

Query	Result Count	Non-Academic Hits	CSV Export
cancer	10	Yes	Yes
covid vaccine	10	Yes	Yes
AI in healthcare	10	Yes	Yes

Sample CSV output includes real-world biotech affiliations like:

- "Novartis Institutes for BioMedical Research"
- "Pfizer Global Research"
- "Merck & Co., Inc."

Project Structure

```
get-papers-tool/  
|  
├─ src/  
|   └─ get_papers_tool/  
|       ├── main.py           # CLI entry point  
|       └─ pubmed_fetcher.py  # Core logic for PubMed API  
|  
├─ tests/                    # (Optional test dir for future use)  
├─ README.md                 # Full project instructions  
├─ sample_output.csv         # Sample result file  
├─ pyproject.toml            # Poetry configuration  
└─ poetry.lock               # Locked dependencies
```

Conclusion

This CLI tool is an efficient utility for extracting and analyzing non-academic authorship in biomedical literature. It can assist in identifying **industry participation in research**, useful for:

- Market research
- Competitor analysis
- Grant/funding tracking
- Industry collaboration detection

Future Enhancements

- Add GUI interface with Streamlit
- Support affiliation confidence scores (e.g., using NLP)
- Integrate with Scopus or IEEE APIs
- Add filtering by publication date or journal name

