# AI Assisted Literature Review: An Entity Extraction Framework for Scholarly Literature

Presented by:
Moobashara Jawed

# Named Entity Recognition (NER)

## What is Entity Recognition?

Automated identification of key terms (entities) in text into categories like diseases, species or locations.

Recent studies from India `GPE` have shown that mutations in the BRCA1 `GENE` gene are strongly associated with breast cancer `DESEASE` in Homo sapiens `SPECIES` . Researchers at ICMR `ORG` reported this finding in a paper published earlier this year `DATE` , highlighting its potential for improving early diagnosis and treatment strategies.

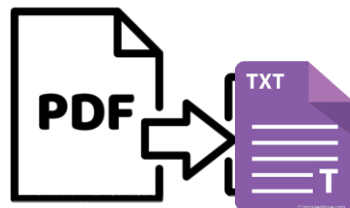| | |
|---|---|
| GPE | GENE |
| DATE | DESEASE |
| ORG | Species |

# NER Workflow

**1** Automated article acquisition via **Pygetpapers**

**2** Document text extraction via PDF parsers

**3** Text tokenization at sentence granularity

"We love NLP!"

Tokenization

"We" "love" "NLP" "!"

**4** 4. Named Entity Recognition (NER)

Recent studies from India GPE have shown that mutations in the BRCA1 GENE gene are strongly associated with breast cancer DESEASE in Homo sapiens SPECIES . Researchers at ICMR ORG reported this finding in a paper published earlier this year DATE , highlighting its potential for improving early diagnosis and treatment strategies.

# Step 1: Automated article acquisition via Pygetpapers`



```
[2]   !pygetpapers --query '"lung cancer"' --xml --pdf --limit 100 --output disease --save_query

      INFO: Total Hits are 581386
      100it [00:00, 191084.46it/s]
      INFO: Saving XML files to /content/disease/*/fulltext.xml
          0% 0/100 [00:00<?, ?it/s]INFO: Wrote the pdf file for PMC12037952
          1% 1/100 [00:03<05:29,  3.33s/it]INFO: Wrote the pdf file for PMC12023508
          2% 2/100 [00:06<04:55,  3.02s/it]INFO: Wrote the pdf file for PMC12022210
          3% 3/100 [00:09<04:56,  3.06s/it]INFO: Wrote the pdf file for PMC12040151
```

# Step 2 & 3: PDF parsing and Text tokenization

2. Parsing is the process of extracting and structuring data from documents, such as PDFs.

3. Tokenization is the process of splitting text into smaller units, such as sentences or words (tokens).



**PDF** → **Text Content**

"We love NLP!"

↓

**Tokenization**

"We"  "love"  "NLP"  "!"

- ➤ PDF Parsing Packages: **fitz (from PyMuPDF)**
- ➤ **Fitz**: Extracts raw texts from pdfs

- ➤ Tokenization: **nltk.tokenize.sent_tokenize**
- ➤ Splits text into sentences

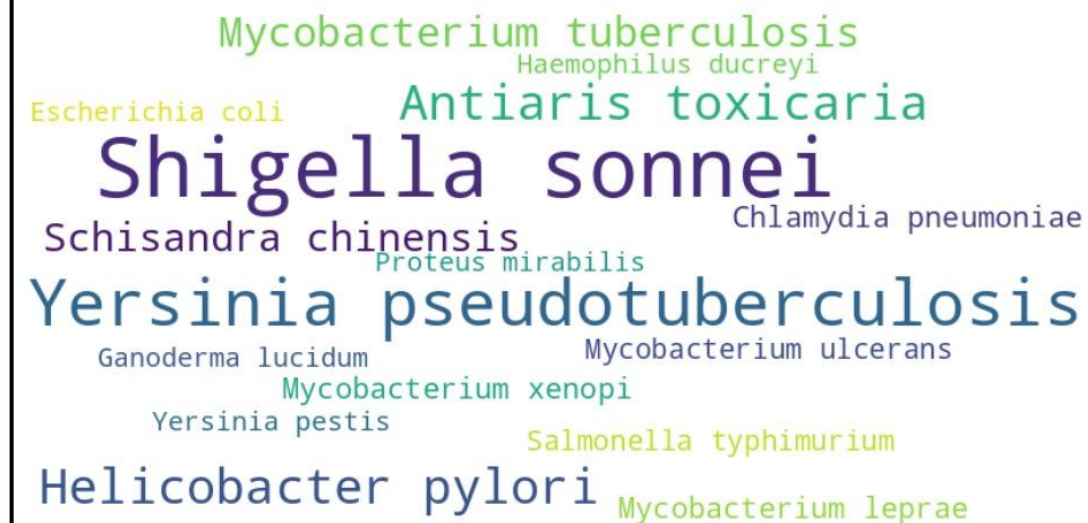Step 4: Named Entity Recognition (NER)

# Species Extraction

| | Species | Frequency | PMC12181914 | PMC12249456 | PMC12261365 |
|---|---|---|---|---|---|
| 1 | Shigella sonnei | 8 | 8 | 0 | 0 |
| 2 | Yersinia pseudotuberculosis | 6 | 6 | 0 | 0 |
| 3 | Helicobacter pylori | 3 | 3 | 0 | 0 |
| 4 | Antiaris toxicaria | 3 | 0 | 0 | 0 |
| 5 | Mycobacterium tuberculosis | 2 | 0 | 2 | 0 |
| 6 | Schisandra chinensis | 2 | 0 | 0 | 2 |
| 7 | Mycobacterium xenopi | 1 | 0 | 1 | 0 |
| 8 | Mycobacterium leprae | 1 | 0 | 1 | 0 |
| 9 | Chlamydia pneumoniae | 1 | 1 | 0 | 0 |
| 10 | Salmonella typhimurium | 1 | 1 | 0 | 0 |

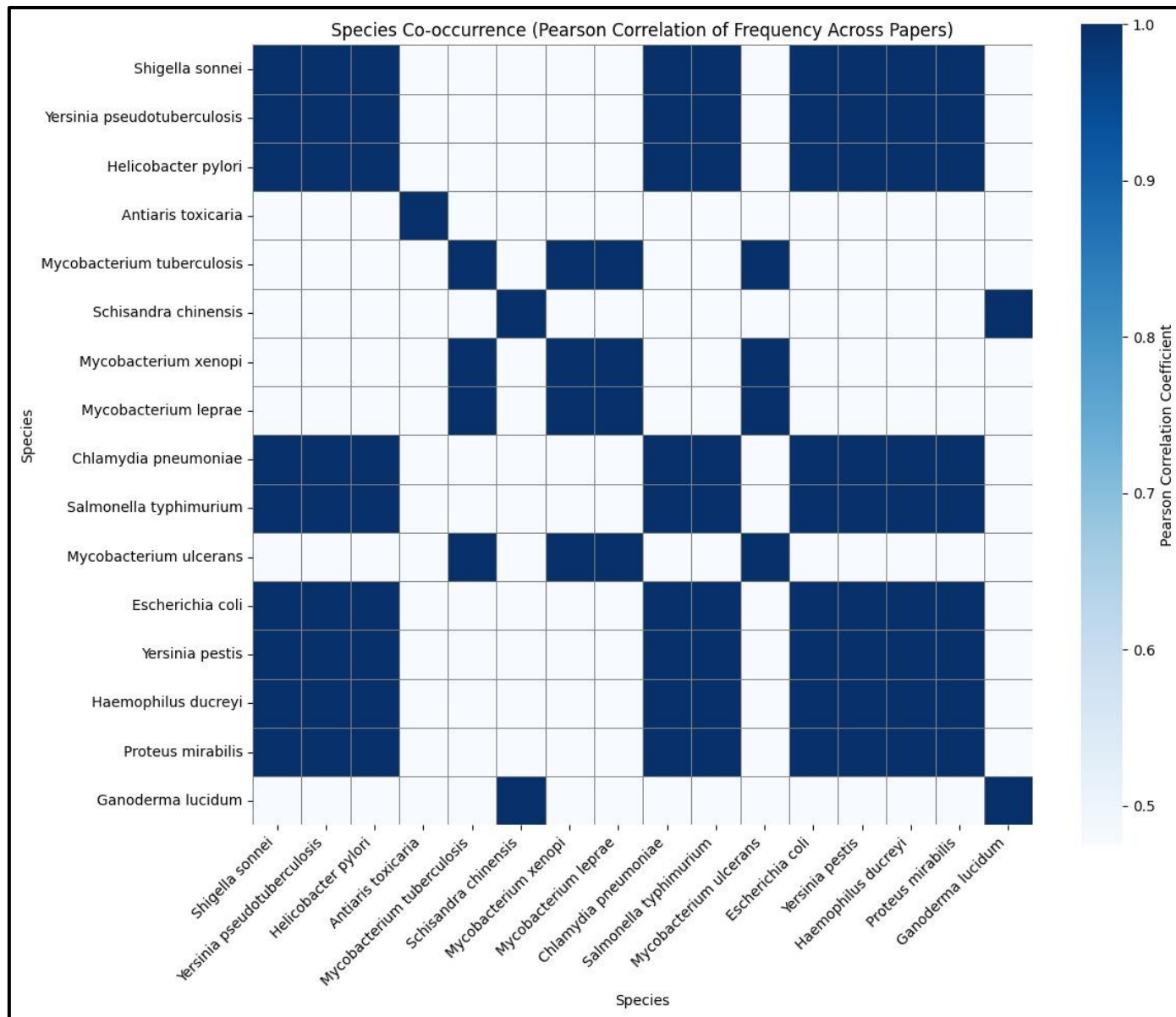**Frequency table showing the top 10 species across 3 PMCIDs out of the 50 papers**

➢ NER Model: **scispaCy en_core_sci_md** (scientific text)
➢ Validation:
   • Regex for binomial nomenclature (Genus species)
   • **GBIF** API verification (EXACT species matches)

# Species Extraction
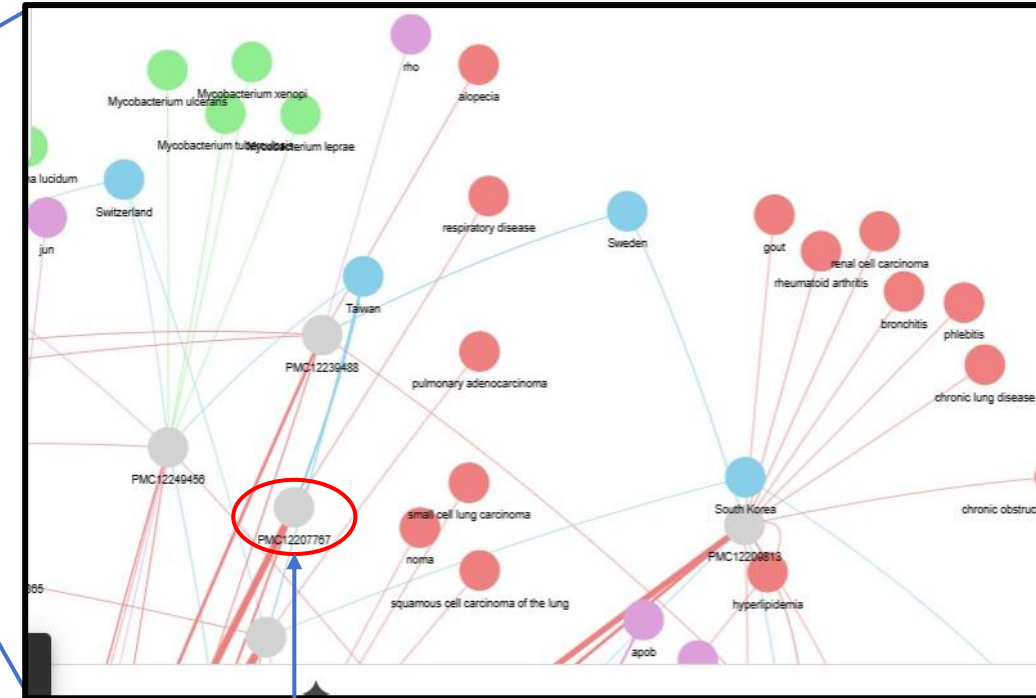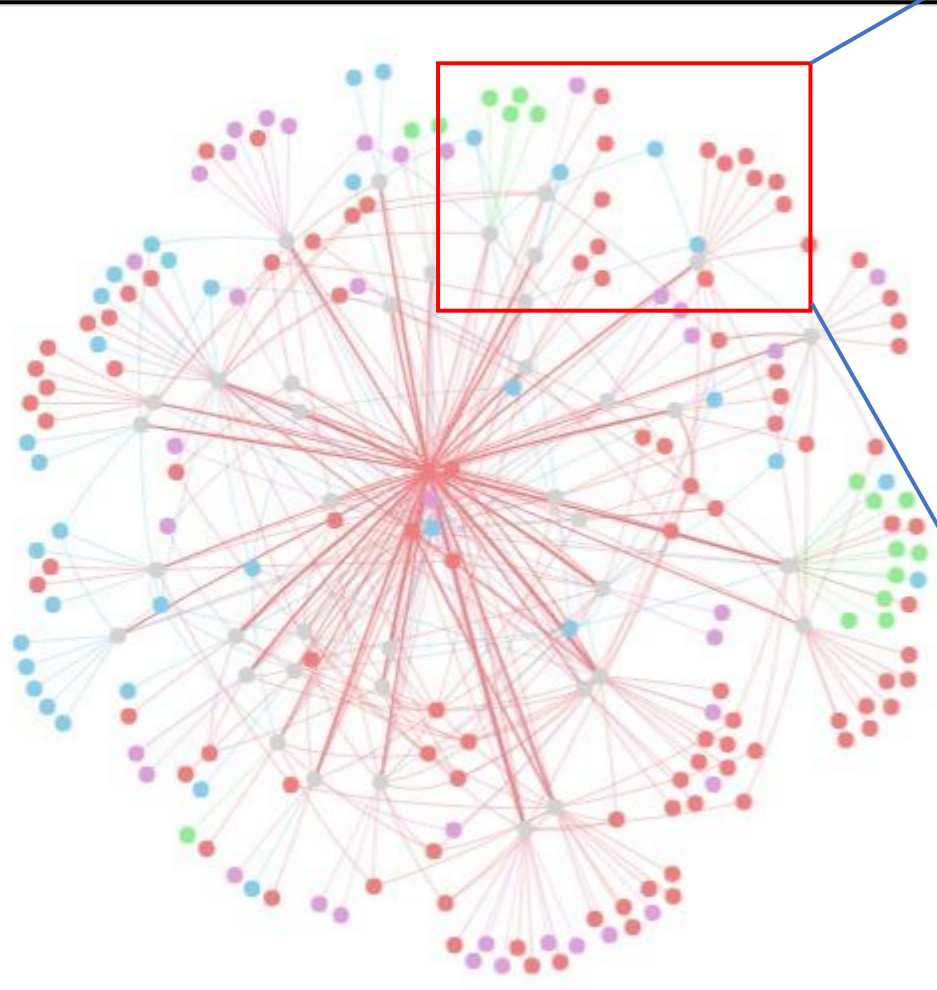


Word cloud of species obtained
from the 50 papers



Co-occurrence plot of species for 50 papers

# Network of the Extracted Entities