# docanalysis demo

Automatically download 1000s of papers, extract scientific data

Ayush Garg,

High School,

Singapore

Shweata N. Hegde,

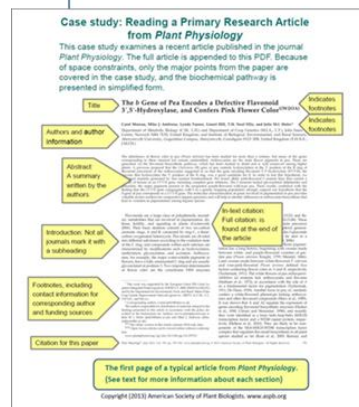3rd year undergrad (plant science & education)
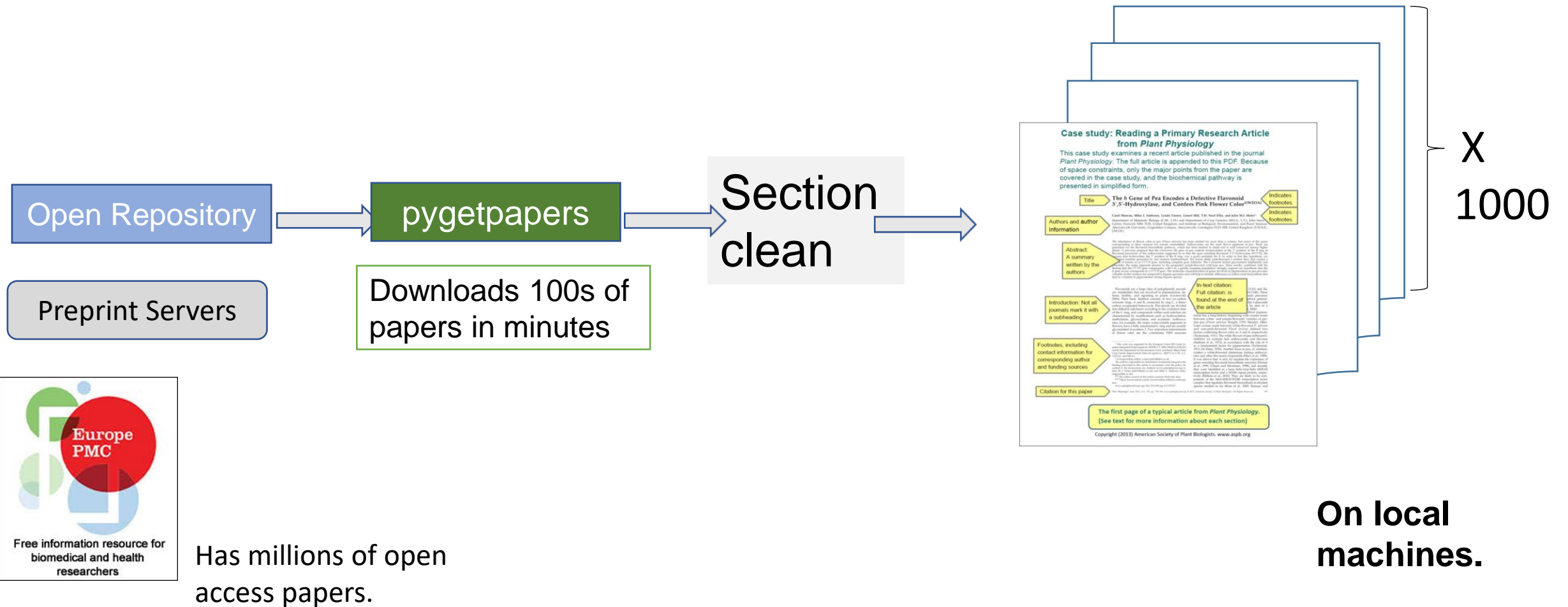
India

Early Career Researchers

MADICES

# docanalysis – a command line tool

Automatically download 1000s of papers, extract scientific data

```
pip install docanalysis
```

# Step 1: Download papers

```
docanalysis --run_pygetpapers -q ""XANES AND EXAFS AND XRD" " -k 10 --project_name xanes_madices
```

Open Repository → pygetpapers → Section clean →

Preprint Servers

Downloads 100s of papers in minutes

X 1000

On local machines.

**Europe PMC**
Free information resource for biomedical and health researchers

Has millions of open access papers.

# Step 2: Section papers
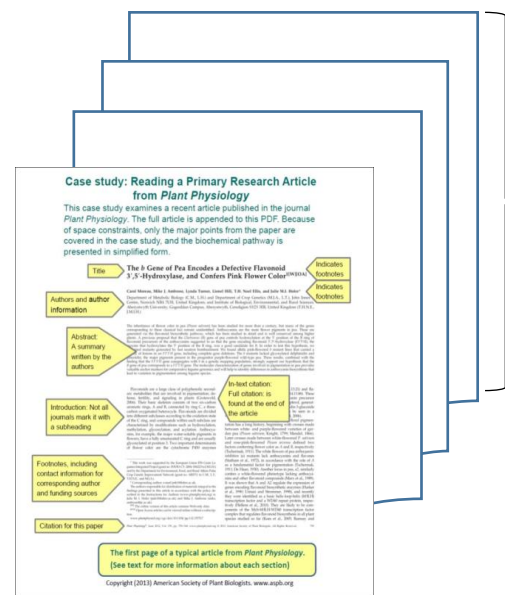
```
docanalysis --project_name xanes_madices --run_sectioning
```

- Article Sections

- <**JATS**> :  Journal Article Tag Suite

# Step 3: Extract Entities

```
docanalysis --project_name xanes_madices --run_sectioning --output entities_202202019.csv
```



**SciSpacy/other NLP tools**

AUTOMATIC
1 paper/s

Background: The purpose of this study is to analyze the surface morphology and elemental composition of zirconia implants before and after photofunctionalization.
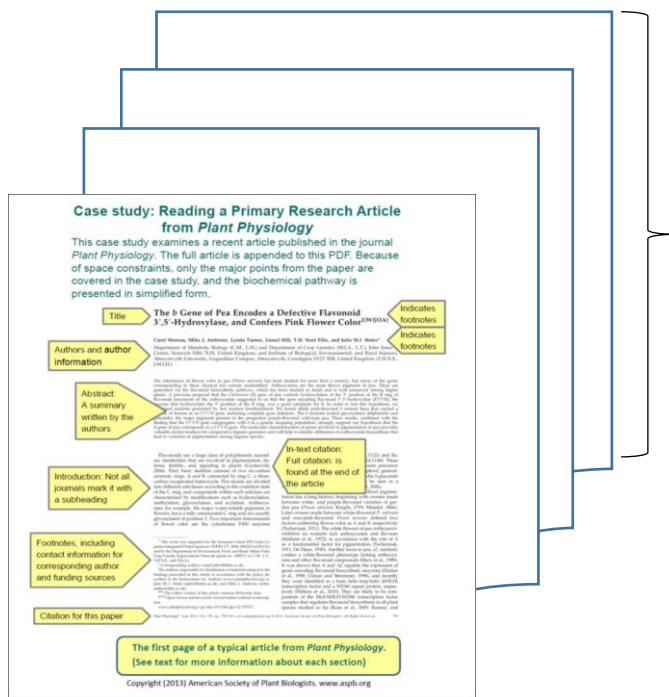
The experimental procedures followed the rules

## 100s of Entities

10 papers on XANES
AND EXAFS AND XRD

Ayush Garg, Shweata Hegde, Daniel Mietchen

https://github.com/petermr/docanalysis

# Step 4: Make dictionary

```
docanalysis--project_name xanes_madices --make_ami_dict entities_20220209
```



10 papers on XANES AND EXAFS AND XRD

```
'surface', 30,
 'redox', 25,
'NIB', 18,
'sXAS', 16,
'mRIXS', 15
'Fe', 13,
'TM', 13,
'DNN', 12,
'MSE', 11,
'capacity', 11,
'surfaces', 10,
'Mn', 9,
'spectra', 7,
'Ti', 7,
'N', 7,
'Al', 7,
'XANES spectra', 6...
```

WIKIDATA

Snippet of the dictionary

Annotate the literature

# docanalysis – Command line Tool

## Automatic!
< 2 secs/paper

| STEP 1 | STEP 2 | STEP 3 | STEP 4 | STEP 5 |
|--------|--------|--------|--------|--------|

| download papers (EPMC) | Section papers | Extract entities | Make dictionary | Search literature |
|---|---|---|---|---|

`--run_pygetpapers`

`--run_sectioning`

`--output`

`--make_ami_dict`



pygetpapers

Ayush Garg

Peter Murray-Rust

https://allenai.git
hub.io/scispacy/

```
pip install docanalysis
```

```
docanalysis --run_pygetpapers -q "XANES AND EXAFS AND
XRD" -k 10 --project_name xanes_exafs_xrd --
run_sectioning --output entities_202202019.csv --
make_ami_dict entities_20220209
```