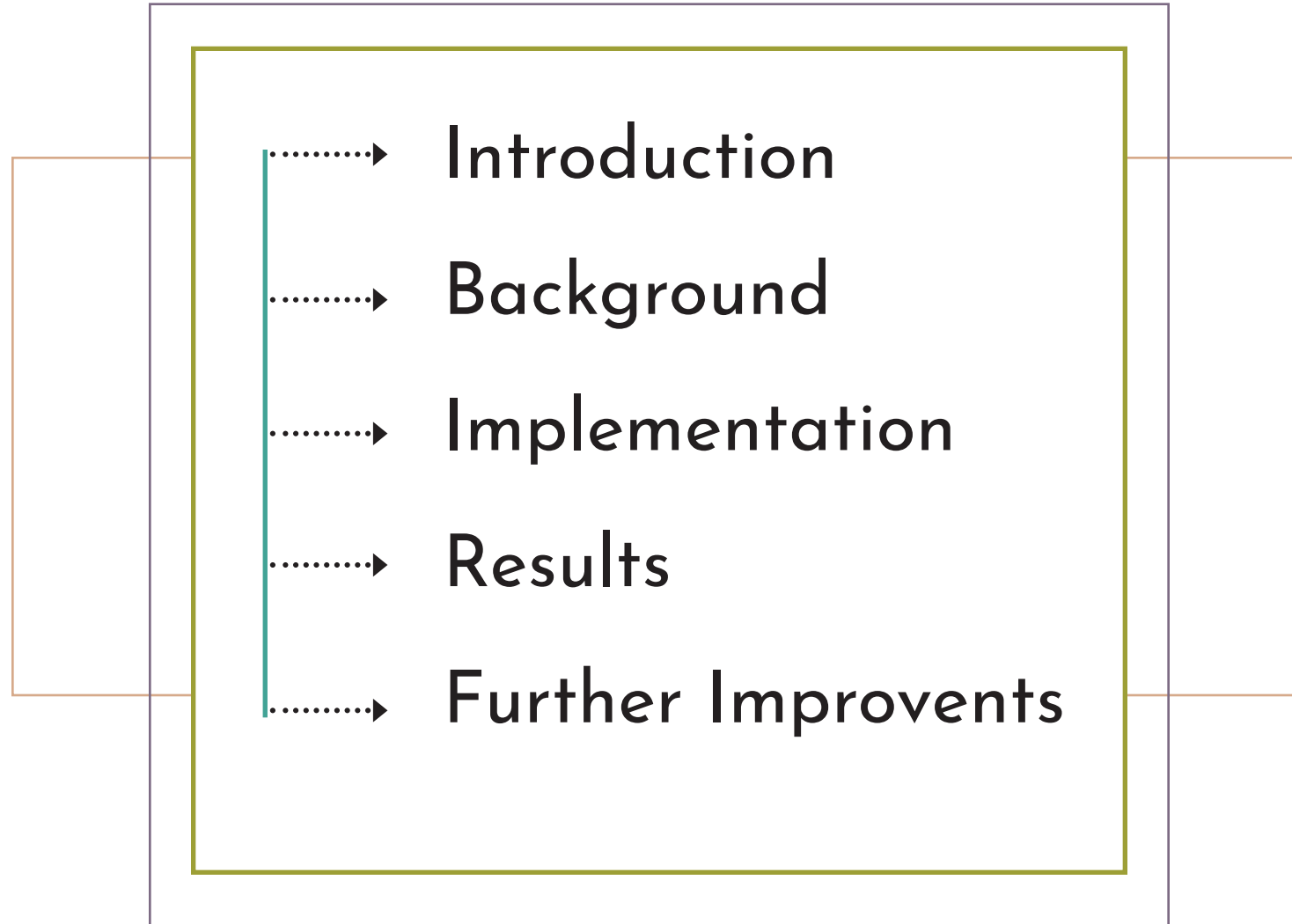


# Abbreviation detection for biomedical articles

by Sonja Kenari

# Agenda

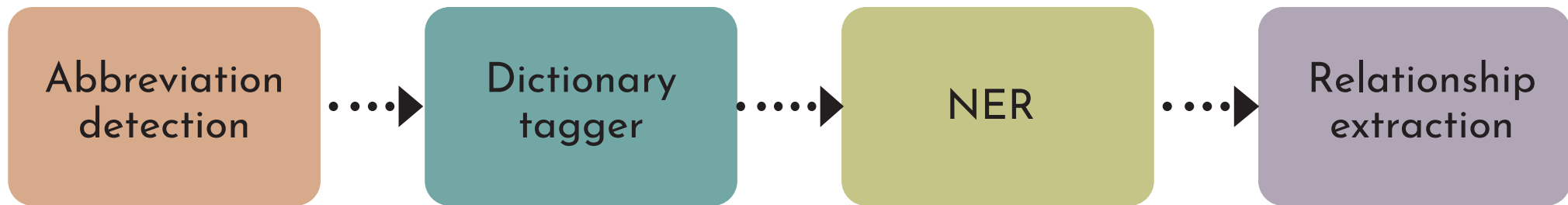


# Introduction

Full project description

## COVID-19 Open Research Dataset Challenge (CORD-19):

.....→ *What do we know about vaccines and therapeutics?*



# Introduction

Abbreviation Detection

spaCy ..... Python library for NLP

## Abbreviation detection

**Makes it easier to:**



Find articles of interest faster



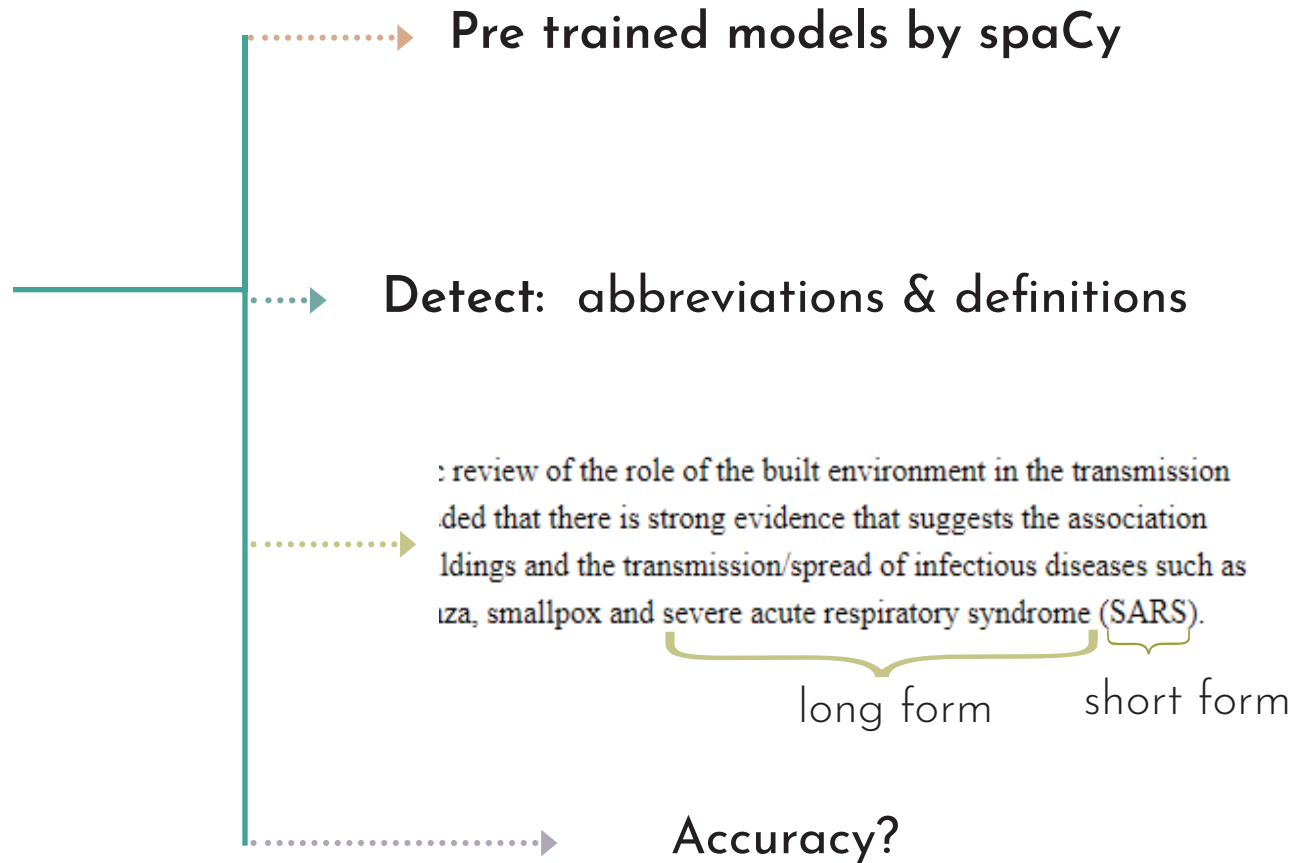
Keep up with the amount of new abbreviations

# Background

Abbreviation Detection

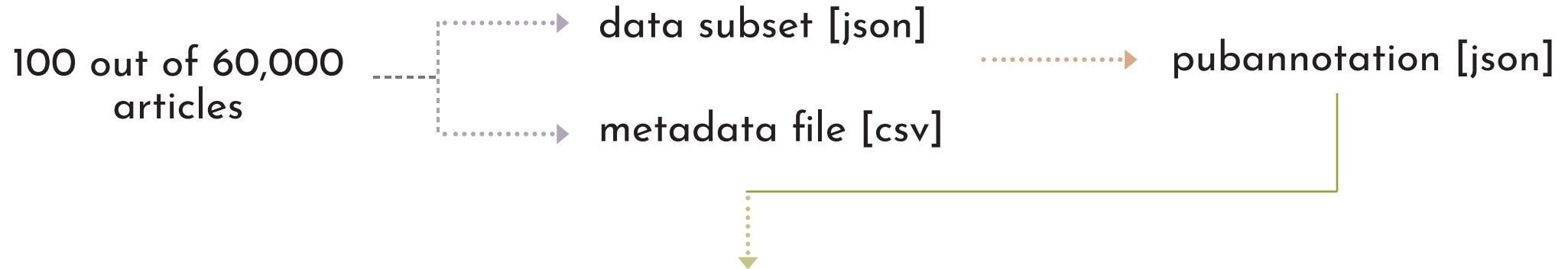
scispaCy:

## AbbreviationDetector

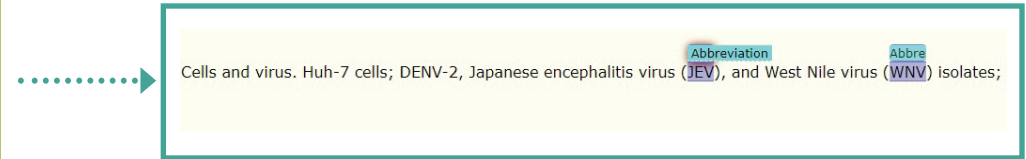


# Implementation

Generate Pubannotations



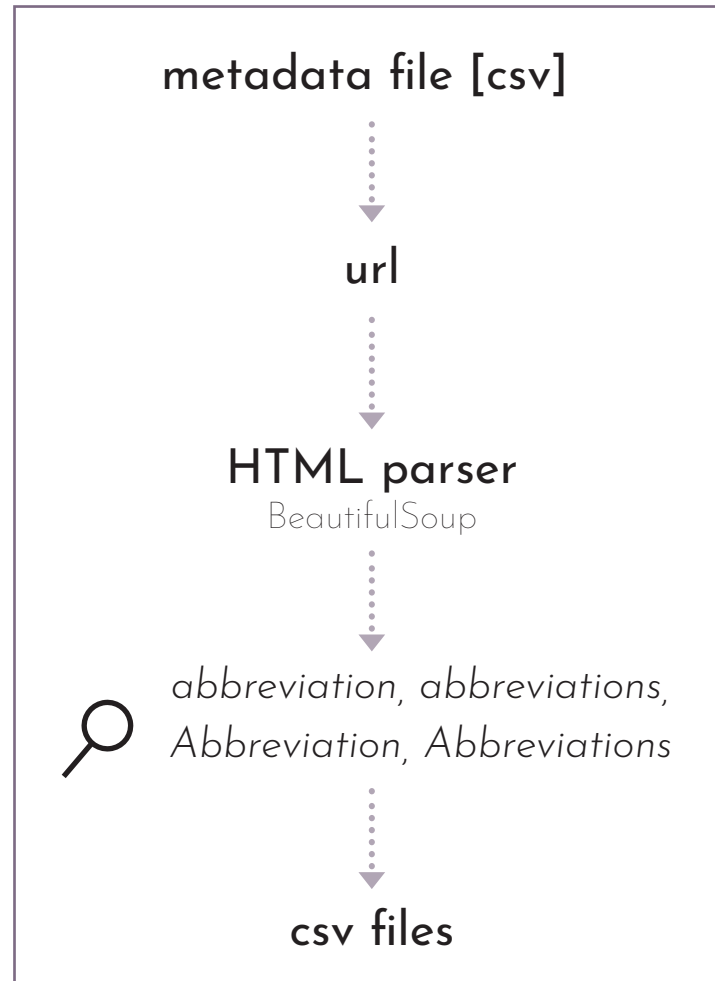
```
1 {
2   "cord_uid": "0nm7wgf5",
3   "source_x": "PMC",
4   "pmcid": "PMC5075077",
5   "divid": "3",
6   "text": "Cells and virus. Huh-7 cells; DENV-2, Japanese encephalitis virus (JEV), and West Nile virus (WNV) isolates;
7   "project": "cdlai_CORD-19",
8   "denotations": [
9     {
10      "id": "S-scispacy-abbr_T1",
11      "span": {
12        "begin": 67,
13        "end": 70
14      },
15      "obj": "Abbreviation"
16    },
17    {
18      "id": "S-scispacy-abbr_T2",
19      "span": {
20        "begin": 94,
21        "end": 97
22      },
23      "obj": "Abbreviation"
24    }
25  ]
26 }
```



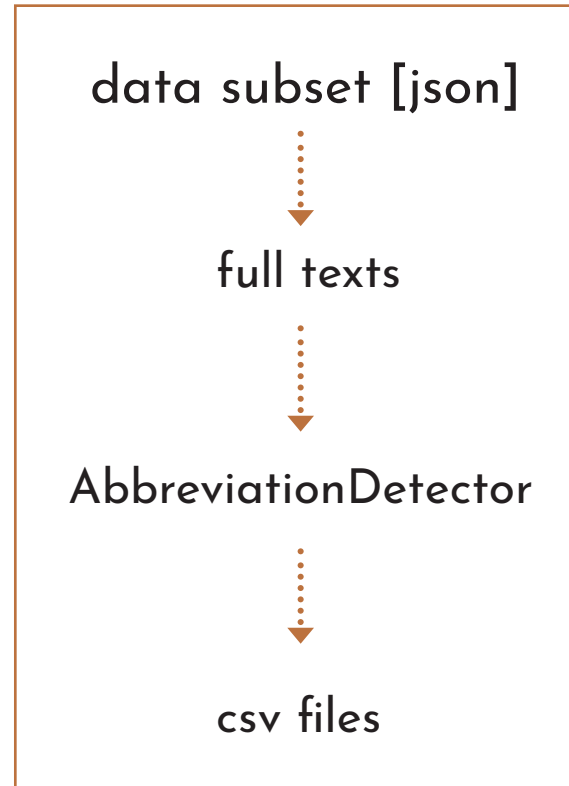
# Implementation

Generating files of abbreviations

web scraping



scispaCy



output file format

| Short form | Long form                                     |
|------------|---|
| CRP        | C-reactive protein                            |
| FTD        | Fast-track Diagnostics                        |
| IQR        | interquartile range                           |
| OR         | odds ratio                                    |
| PCV        | pneumococcal conjugate vaccination            |
| PCT        | procalcitonin                                 |
| RSV        | respiratory syncytial virus                   |
| RT-PCR     | real-time multiplex polymerase chain reaction |
| WHO        | World Health Organization                     |

# Implementation

Evaluation

Compare the **2** {  
detected abbreviations with spaCy [csv]  
detected abbreviations with web scraping [csv]

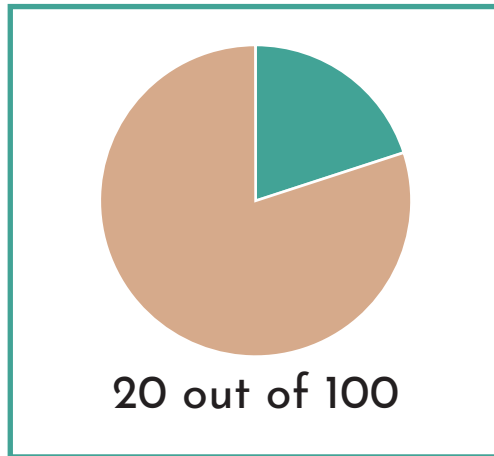
$$\frac{\text{Number unique short forms detected by spaCy}}{\text{Number short forms detected by web scraping}} = (\%)$$

$$\frac{\text{Number unique long forms detected by spaCy}}{\text{Number long forms detected by web scraping}} = (\%)$$



# Result

Abbreviation lists in



short forms hit rate

Highest: 87.5%  
Lowest: 25%

long forms hit rate

Highest: 52.6%  
Lowest: 0%

notable faults

- spaCy weak on long form
- text from json files not updated after url articles
- faults in denotation extraction

# Further Improvements

## **spaCy**

Improve the results

## **Optimize programs**

Make more time efficient

## **web scraper**

Update data

## **Extract from Pubannotations**

Instead of full text extraction

**Thank you  
for listening!**

---

Questions...?

Sonja Kenari

nat14sta@student.lu.se

2020-05-29