# Exuberanter: A Versatile Python Tool for Literature Data Extraction

Email: wenne@chalmers.se    Teo Lovmar[1], Johan Bengtsson-Palme[1, 2, 3] and Marcus Wenne[1, 2]
Phone: +46 725 039 395

**CHALMERS**

## Motivation

- **Efficiency Need:** Literature reviews are time-consuming and demand substantial manual effort.

- **Reproducibility Concern:** Manual literature review outcomes can vary and may lack reproducibility.

- **Knowledge Expansion:** The rapid growth of scientific literature amplifies these challenges.

- **Our Solution:** We provide a tool for efficient, reproducible literature reviews, with semi-automatic data extraction capabilities and an intuitive user interface.

## Features

- **Efficient Literature Review:** Optimize your literature surveying with our semi-automated approach.

- **Downloading articles:** Simplify the task of obtaining, and assessing research articles from PubMed.

- **User customization:** Specify what type of information you would like to extract from the different sections of the articles.

- **Data Extraction:** Extract data using regex and plot data extraction tools in a semi-automatic fashion.

- **User Interface:** Navigate and interact with the processed data in a straightforward, user-friendly manner.

## Workflow

### Pre defined search query

First the user specifies one or multiple query terms which will be used to search PubMed.

eg. '"antibiotic resistance genes" AND q-PCR'

### Download articles and their metadata

Using the NCBI API Efetch, Exuberanter downloads all available open access articles, and their accompanying metadata.

Articles not available through the API can be manually downloaded and accessed by Exuberanter.
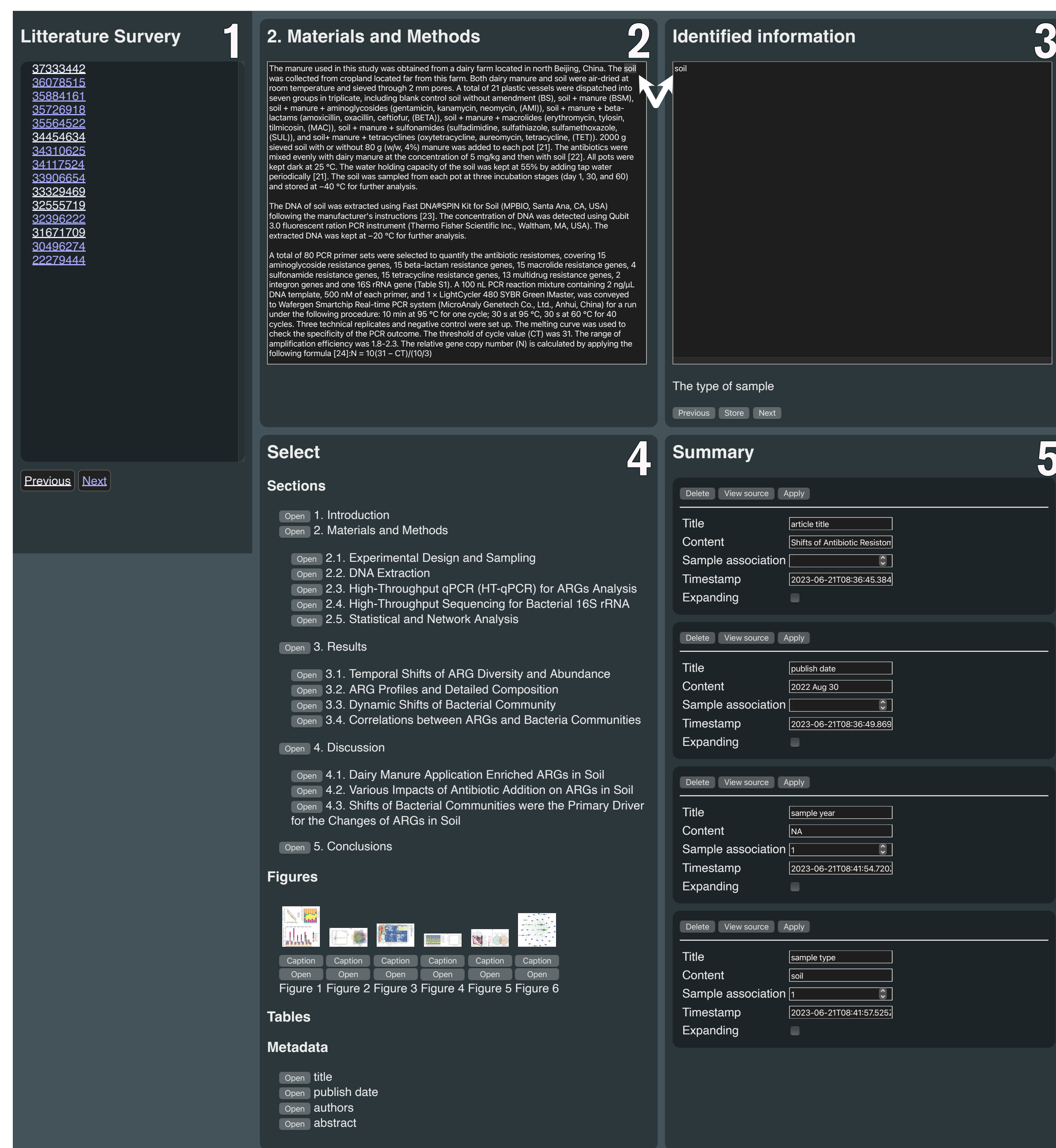
### Extract pre specified data from articles

The user specifies what type of information Exuberanter should extract using regex or an NLM API (to come). All data is stored in a JSON file.

### Verify or modify the extracted data via a UI

Via the UI the user can manually verify the pre-extracted data, or manually change it. The tool has WebPlotDigitizer (*https://automeris.io/WebPlotDigitizer*) integrated which facilitates plot data extraction.

## Interface



**1.** Navigation between individual articles (PID).

**2.** Article text where data points identified by the tool are highlighted (left arrow).

**3.** Either automatically detected, or manually specify values for each datapoint (right arrow).

**4.** The various article sections, figures or tables can be selected here.

**5.** Individual datapoints collected during the literature survey.

## Code examples

```
1   # The publishing date of the article. Extracted from the article metadata
2   'publish date': {
3       'targets': { 'metadata': ['publish date'] },
4       'regex': re.compile(r'.+', re.DOTALL),
5       'description': {
6           'info': 'The publish date of the article',
7           'data': 'any date format',
8       }
9   },
10
11  # The type of sample collected. Can be one of multiple categories
12  'sample type': {
13      'sample associated': True,
14      'targets': { 'sections': ['method'] },
15      'categories': {
16          'soil': re.compile('|'.join(['soil', 'soils', 'Clay']), re.IGNORECASE),
17          'manure': re.compile('|'.join(['manure', 'guano', 'feces']), re.IGNORECASE),
18          'sewage': re.compile('|'.join(['sewage', 'waste water']), re.IGNORECASE),
19      },
20      'description': {
21          'info': 'The type of sample',
22          'data': '"soil", "manure or "sewage"',
23      }
24  },
```

## GitHub

https://github.com/mnyt-aqw/Exuberanter/

## Help!

We are currently having problems **extracting data from tables** due to the large variability in how journals format them.

Tips on how to solve this are greatly appreciated!

1. **Division of Systems and Synthetic Biology, Department of Life Sciences, SciLifeLab, Chalmers University of Technology, Gothenburg, Swede**
2. **Centre for Antibiotic Resistance Research in Gothenburg (CARe), Gothenburg, Sweden**
3. **Department of Infectious Diseases, Institute of Biomedicine, The Sahlgrenska Academy, University of Gothenburg, Gothenburg, Sweden**