# Data Preparation Process
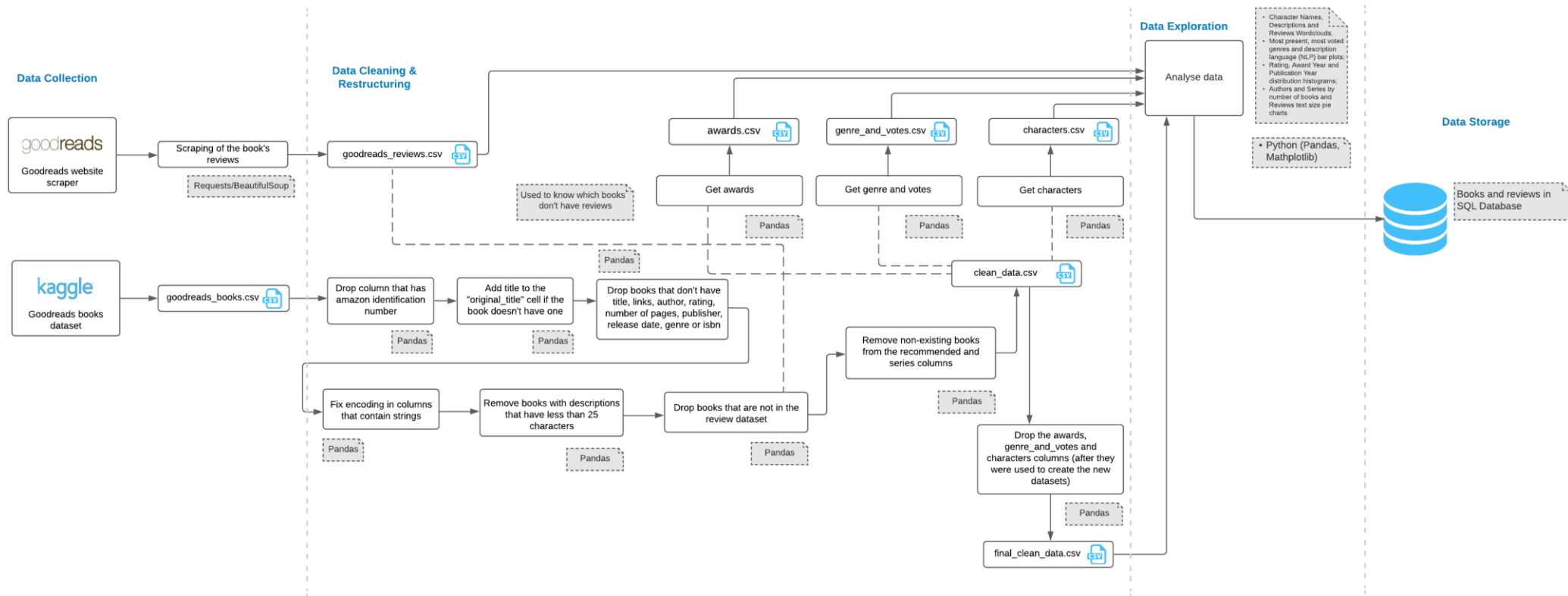
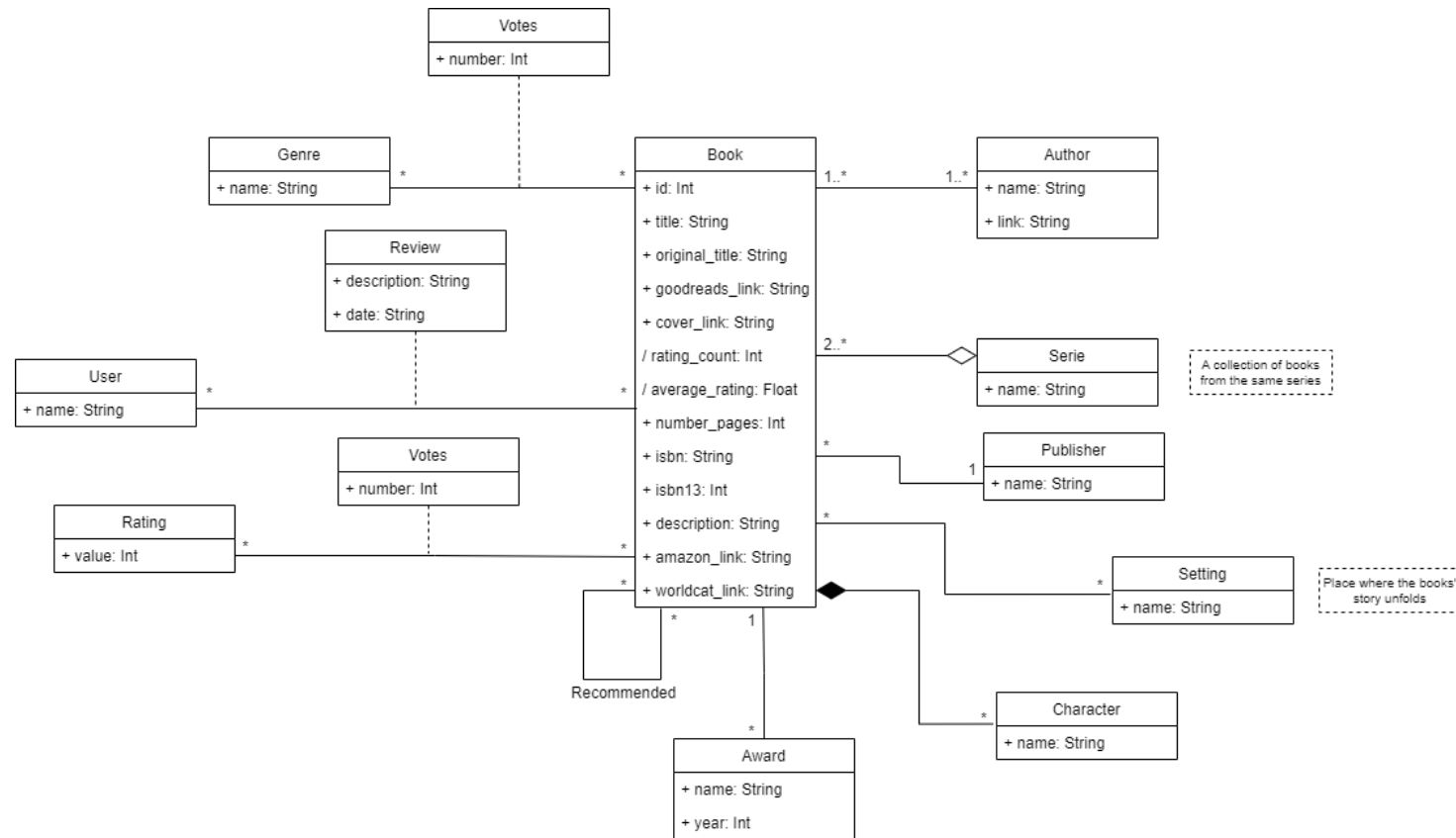## GOODREADS' BOOKS AND REVIEWS

Diogo Almeida (up201806630)
Pedro Queirós (up201806329)

# Data Extraction and Preparation
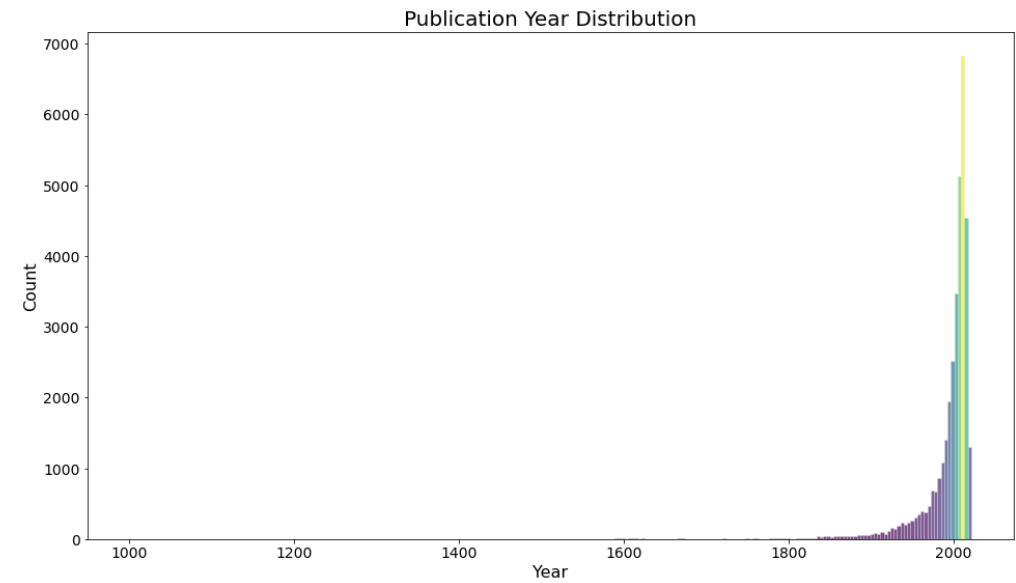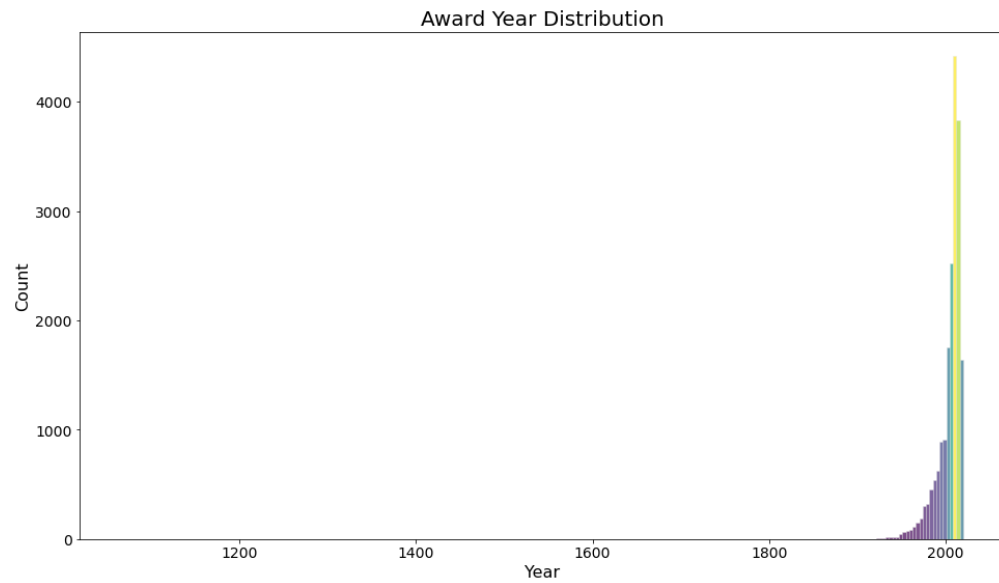
# Domain Conceptual Model

# Data Characterization

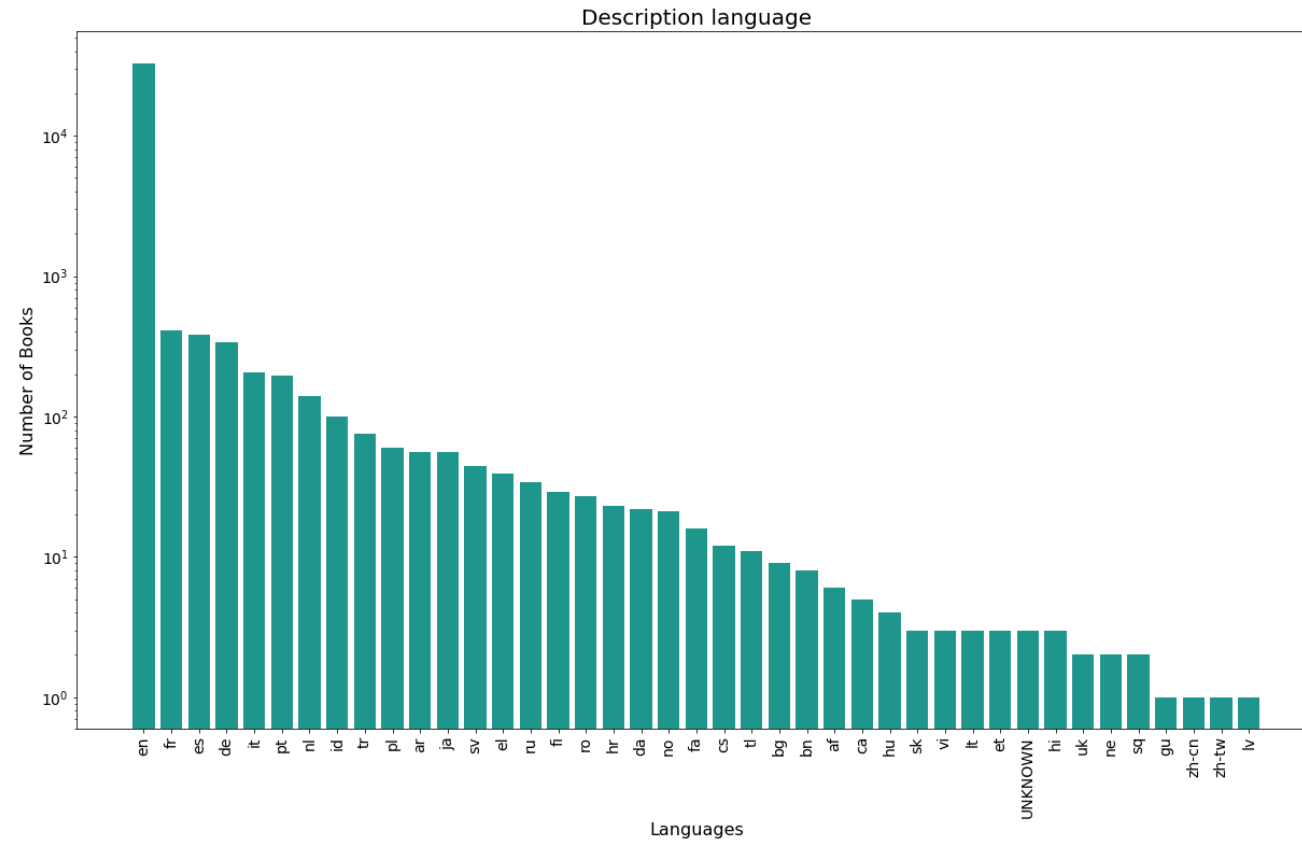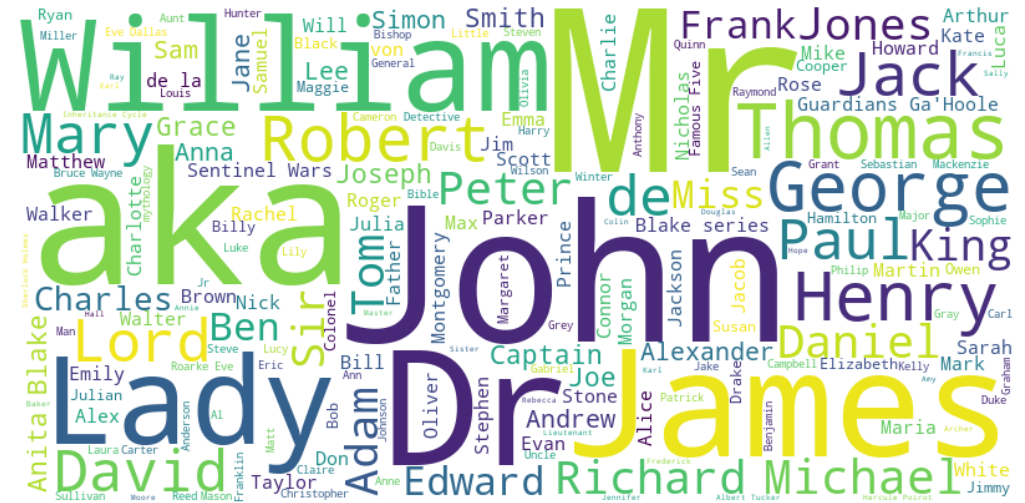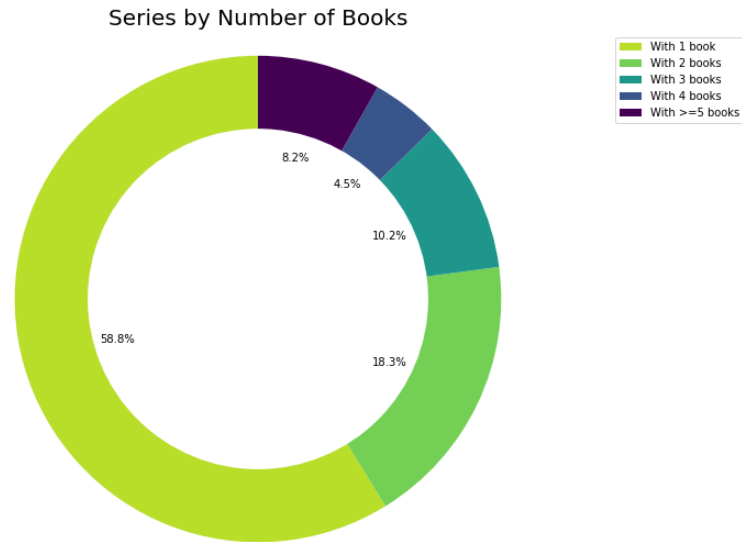# Data Characterization

# Data Characterization



Series by Number of Books

- With 1 book
- With 2 books
- With 3 books
- With 4 books
- With >=5 books

58.8%
18.3%
10.2%
4.5%
8.2%

# Conclusions and Future Work

For the extraction and preparation of the information the following steps were made:

❖Gathering data and enrich it with data from diferentes sources

❖Data cleaning

❖Analyse the final data

Regarding future work, some of the implemented search tasks will be:

❖Search books by author

❖Search books based on their language

❖Search books by genre, awards or even by their reviews or rating

# References

❖ "Goodreads Books – 31 features" dataset:
- https://www.kaggle.com/austinreese/goodreads-books

❖ Goodreads website:
- https://www.goodreads.com

❖ Python libraries: pandas, requests, beautifulsoup and spacy:
- https://pandas.pydata.org/
- https://docs.python-requests.org/en/latest
- https://beautiful-soup-4.readthedocs.io/en/latest
- https://spacy.io/universe/project/spacy-langdetect