

Lextract: An Python pipeline for the automated extraction of European Commission market definitions

Shriyan S. Yamali ¹

¹ Newark Charter High School, United States

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

Software

- [Review](#) 
- [Repository](#) 
- [Archive](#) 

Editor: [Open Journals](#) 

Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

License

Authors of papers retain copyright,
and release the work under a
Creative Commons Attribution 4.0
International License ([CC BY 4.0](#)).

Summary

Lextract is a Python pipeline that automatically locates, downloads, and extracts relevant market definitions from the European Commission's merger and antitrust decision PDFs. Relevant market definitions establish the specific scope of competition legislation and identify the specific set of product in an area, making them indispensable for economists, lawyers, and regulators when determining the effects of mergers and evaluating anticompetitive behavior. This pipeline has been designed for researchers and competition law experts who require a quick and scalable way to extract relevant market definitions from many cases at once. Considering that market definitions are highly sensitive and, to some extent, arbitrary, Lextract has been designed and implemented to extract definitions are accurately as possible, as a slight change in the language of a definition can drastically change its meaning. The process is split into five main steps: (1) fetch links, (2) exclude irrelevant decisions PDFs and scrape the text of the rest, (3) extract the market definitions section, (4) extract individual market definitions, and (5) clean and combine JSON files. This modular design allows for the scalable, reproducible extraction of market definitions with vast applications.

Statement of need

Competition authorities routinely delineate the relevant market as a first step in merger and antitrust assessments. Scholars analyze this language to track precedent, analyze trends in the scope of definitions, and identify the evolution of market definitions ([Robertson, 2019](#)). Despite its significance, only one commercial product addressing this need exists: LexisNexis's [Caselex Market Definitions Module](#) which suffers from being proprietary, immutable, and inaccessible to many academics.

The Commission has published over 6,000 merger and antitrust decisions ([Bernhardt & Dewenter, 2024](#)) and continues to add 280 annually ([Affeldt et al., 2021](#)), each formatted idiosyncratically and lengthy. Headings vary in language and definition placement within decisions is inconsistent. As a result, deterministic approaches such as regex are brittle and ineffective, while manual extraction is slow and irreproducible at scale. This pipeline rectifies this issue by providing a simple, open source way to extract market definitions that does not require manual guidance nor relies upon inaccurate patter matching techniques.

34 General workflow

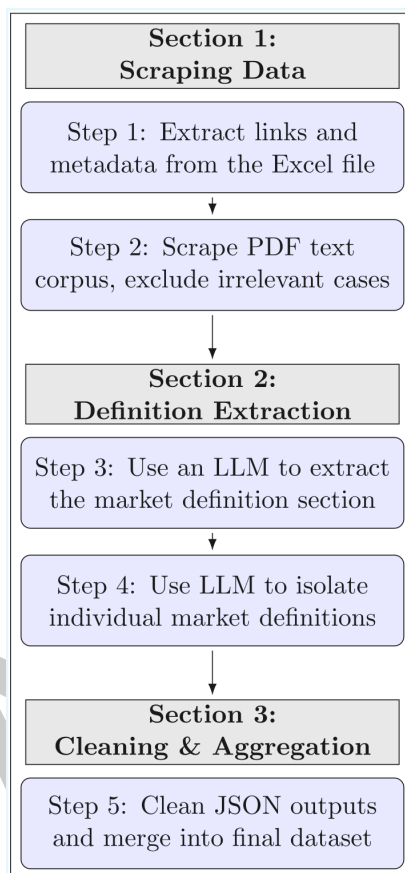


Figure 1: Workflow Diagram of Pipeline

35 The general workflow for extracting market definitions is split into 3 sections, and 5 steps.
 36 The first section involves the scraping of data, making use of regex: 1. A script processes
 37 an Excel file downloaded from the Commission's case search portal and extracts the links of
 38 decision documents and corresponding metadata (i.e. case number, year, policy area), saved in
 39 a plain text file. 2. Another script processes this file, and, using the decision document links,
 40 scrapes the decision text and convert them into a text corpus with the metadata, repeating
 41 this step for each link, while also sorting the corpus based on its length, with a break point a
 42 80,000 characters. It is during this process that decision documents without market definitions,
 43 identified by certain phrases or a page length less than three, are excluded.

44 The second section is responsible for the semantic extraction of market definitions: 3. Google
 45 Gemini is used to identify and extract only the section of the text corpus that contains the
 46 market definition section. 4. Afterwards, the process becomes more granular, with Gemini
 47 again being used, only this time to identify and isolate each individual market definition within
 48 those sections. Each definition is then tagged with a topic and saved in a structured JSON
 49 file, where each object contains all elements of the aforementioned metadata, a topic, and the
 50 market definition.

51 The third and final section improves the presentation of the data: 5. Each separate JSON file
 52 is cleaned to remove extraneous characters and then aggregated into a single file, which can
 53 then be used for research and analysis. By structuring the workflow this way, each processed
 54 case is consistently analyzed, reducing variability and improving accuracy. Lextract's code also

55 maintains a high level of accuracy, substantiated by a comprehensive test suite with 94% code
56 coverage.

57 While this pipeline makes use of Google Gemini, it is model-agnostic and, if properly refactored,
58 could utilize any LLM, including commercially hosted models like OpenAI's or locally deployed
59 ones such as LLaMA, Mistral, or DeepSeek. However, accuracy and consistency will vary
60 significantly deepening on model size and capabilities. Generally, smaller models, especially
61 local ones without a sufficient context length or reasoning ability, will tend to hallucinate
62 outputs, misidentify sections, or produce partial definitions.

63 Research applications

64 Lextract powers the database of [JurisMercatus](#), an open source search intercaze that allows
65 users to search for market definitions, leveraging the power of semantic search. The metadata
66 provided by Lextract enables filtering by year, policy area, and case number. Further, this
67 resource has the capability to support greater academic research and improve the accessibility
68 of market definitions.

69 Limitations

70 It should be noted that this system, as with all systems, is not perfect and contains inaccuracies.
71 First, with regards to step three, as a result of the previously mentioned fact that the heading
72 used to identify the market definition section is inconsistently phrased, what constitutes the
73 market definition is only heuristically defined, potentially leading to inaccuracies, especially
74 when the language of decision texts deviate significantly from expected pattern. Secondly, the
75 quality and reliability of the extraction is limited to that of the input. In other words, should
76 the input consist of missing pages or unconventional language, the model may be confused,
77 leading to partial or inaccurate results. Furthermore, due to the nature of LLMs, it is possible,
78 though unlikely, for output to be hallucinated or metadata to be misattributed. Lastly, though
79 it is understood that decisions are adjudicated in many different languages, with the European
80 Commission using multiple itself, the pipeline assumes that all decisions are provided in English
81 and excludes all others, thereby limiting its application to other languages, without at least a
82 moderate amount of modification.

83 Acknowledgements

84 I thank Professor Thibault Schrepel for his advisement and guidance throughout the project.

85 References

- 86 Affeldt, P., Duso, T., & Szücs, F. (2021). 25 years of European merger control. *International*
87 *Journal of Industrial Organization*, 76. <https://doi.org/10.1016/j.ijindorg.2021.102720>
- 88 Bernhardt, L., & Dewenter, R. (2024). The Impact of the More Economic Approach on
89 EU Merger Decisions. *Stanford Computational Antitrust*. [https://law.stanford.edu/](https://law.stanford.edu/wp-content/uploads/2024/12/Bernhardt-Dewenter.pdf)
90 [wp-content/uploads/2024/12/Bernhardt-Dewenter.pdf](https://law.stanford.edu/wp-content/uploads/2024/12/Bernhardt-Dewenter.pdf)
- 91 Robertson, V. (2019). The relevant market in competition law: a legal concept. *Journal of*
92 *Antitrust Enforcement*, 8(2). <https://doi.org/10.1093/jaenfo/jnz008>