

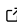


# Lextract: A Python Pipeline for the Automated Extraction of European Commission Market Definitions

Shriyan S. Yamali <sup>1</sup>

<sup>1</sup> Independent Researcher, 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Published: unpublished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 extracts relevant market definitions from the European Commission's merger and antitrust decision PDFs. Relevant market definitions establish the scope of competition legislation and identify the specific set of products in an area ([Tsangaris, 2017](#)), which make them indispensable for economists, lawyers, and regulators when determining the effects of mergers and evaluating anticompetitive behavior ([Patakyová, 2020](#)). This pipeline has been designed for researchers and competition law experts who require a quick and accurate way to extract relevant market definitions from many cases at once. This level of accuracy is accomplished by using strict natural language processing and rule-based pattern recognition ([Braun et al., 2025](#); [Breton et al., 2025](#)) to identify market definitions while excluding all irrelevant information. By automating this process, Lextract enables merger and antitrust research at scale and contributes to more efficient competition policy analysis.

## Statement of need

Competition authorities routinely delineate relevant markets as a first step in merger and antitrust assessments. The definition of the relevant market establishes the market position of firms operating within it, helping regulators and courts control mergers and evaluate potential abuses of dominant positions; this makes defining the relevant market a predominant step in competition law analysis. For instance, in the 2025 case *United States v. Google LLC*, the outcome of the decision was impacted by how the relevant market was defined and whether or not Google and its services were found to hold a dominant position within that market ("[Antitrust — Sherman Act, Section 2](#)," 2025).

Despite its significance, only one commercial product addressing the need to quickly access relevant market definitions exists: LexisNexis's [Caselex Market Definitions Module](#) which suffers from being proprietary, immutable, and inaccessible to many academics.

Furthermore, the Commission has published over 6,000 merger and antitrust decisions ([Bernhardt & Dewenter, 2024](#)) and continues to add 280 more annually ([Affeldt et al., 2021](#)), each structured and formatted idiosyncratically, with inconsistent placement of definitions and headings that vary in language. As a result, deterministic approaches such as regex are brittle and ineffective ([Wang et al., 2020](#)) for extracting market definitions, 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 rely upon inaccurate pattern-matching techniques.

## General workflow

The general workflow for extracting market definitions is split into three sections and five steps. The first section involves the scraping of data, which makes use of regex: 1. A script processes an Excel file downloaded from the [Commission's case search portal](#) and extracts the links of decision documents and corresponding metadata (i.e., case number, year, policy area), saved in a plain text file. 2. Another script processes this file, and, using the decision document links, scrapes the decision text and converts it into a text corpus with the metadata, repeating this step for each link, while also sorting the corpus based on its length, with a breakpoint at 80,000 characters. It is during this process that decision documents without market definitions, identified by certain phrases or a page length less than three, are excluded.

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

The third and final section improves the presentation of the data: 5. Each separate JSON file is cleaned to remove extraneous characters and then aggregated into a single file, which can then be used for research and analysis. By structuring the workflow this way, each processed case is consistently analyzed, reducing variability and improving accuracy ([Pastrana et al., 2025](#)). It also allows Lextract's code to maintain a high level of accuracy, substantiated by its comprehensive test suite with 94% code coverage.

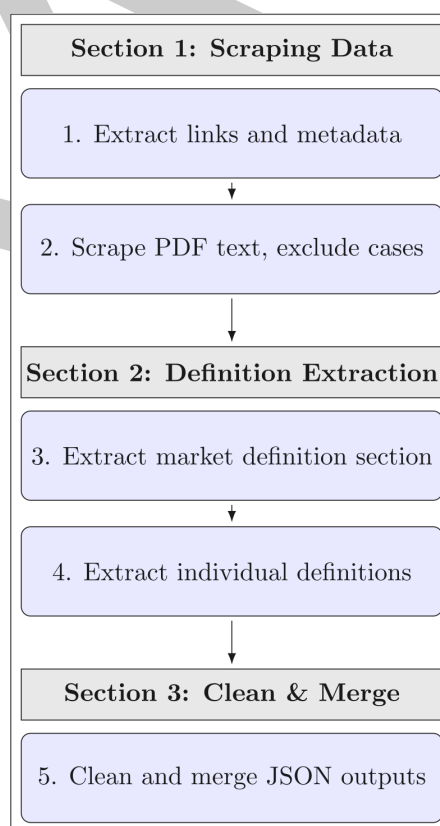


Figure 1: Workflow Diagram of Pipeline

## Research applications

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

## Limitations

It should be noted that this system, as with all systems, is not perfect and contains inaccuracies. First, with regards to step three, as a result of the previously mentioned fact that the heading used to identify the market definition section is inconsistently phrased, what constitutes a market definition heuristically and arbitrarily defined, potentially leading to inaccuracies, especially when the language of decision texts deviates significantly from the expected pattern. Secondly, the quality and reliability of the extraction are limited to that of the input. In other words, should the input consist of missing pages or unconventional language, the model may be confused, leading to partial, hallucinated, or inaccurate results ([Valentin et al., 2024](#)). Additionally, though it is understood that decisions are adjudicated in many different languages, with the European Commission using multiple itself, to maintain accuracy, the pipeline excludes all decisions that are not provided in English, thereby limiting its application to other languages without at least a moderate amount of modification.

Lastly, while this pipeline makes use of Google Gemini, it is model-agnostic and, if properly refactored, could utilize any LLM. This includes commercially hosted models like OpenAI's or locally deployed ones such as LLaMA, Mistral, or DeepSeek. However, accuracy and consistency will vary significantly depending on model size and capabilities. Generally, smaller models, especially local ones without a sufficient context length or reasoning ability, will tend to hallucinate outputs, misidentify sections, or produce partial definitions ([Sun et al., 2025](#)).

Model Type	Accuracy	Context Length	Speed	Cost
Hosted L (eg. GPT-4o)	High	Very High	Moderate	High
Hosted S (eg. Gemini Flash)	Moderate	High	Fast	Moderate
Local L (eg. DeepSeek 67B)	Moderate	Medium	Slow	Low
Local S (eg. LLaMA 3-8B)	Low	Low	Moderate	Low

Table I compares the relative capabilities (accuracy, context length, speed, cost) of different LLMs when applied to the task of extract relevant market definitions. "L" = Large models (>30B parameters); "S" = Small models (<30B parameters).

## Acknowledgements

I am grateful to Professor Thibault Schrepel of Stanford Law School for his invaluable advisement and guidance throughout the course of this project. This research received no funding from any government agency, university, company, or non-profit organization.

## Conflict of Interest

The author declares no conflict of interest.

## References

- Affeldt, P., Duso, T., & Szücs, F. (2021). 25 years of european merger control. *International Journal of Industrial Organization*, 76. <https://doi.org/10.1016/j.ijindorg.2021.102720>
- Antitrust — sherman act, section 2. (2025). *Harvard Law Review*, 138(3). <https://harvardlawreview.org/print/vol-138/united-states-v-google-llc>
- Bernhardt, L., & Dewenter, R. (2024). The impact of the more economic approach on EU merger decisions. *Stanford Computational Antitrust*. <https://law.stanford.edu/wp-content/uploads/2024/12/Bernhardt-Dewenter.pdf>
- Braun, C., Lilienbeck, A., & Mentjukov, D. (2025). *The hidden structure – improving legal document understanding through explicit text formatting*. <https://doi.org/10.48550/arXiv.2505.12837>
- Breton, J., Billami, M. M., Chevalier, M., Nguyen, H. T., Satoh, K., Trojahn, C., & Zin, M. M. (2025). Leveraging LLMs for legal terms extraction with limited annotated data. *Artificial Intelligence and Law*. <https://doi.org/10.1007/s10506-025-09448-8>
- Pastrana, M., Ordoñez, H., Cobos-Lozada, C. A., & Muñoz, M. (2025). Best practices evidenced for software development based on DevOps and scrum: A literature review. *Applied Sciences*, 15(10). <https://doi.org/10.3390/app15105421>
- Patakyová, M. T. (2020). Competition law in digital era – how to define the relevant market? *4th International Scientific Conference on Economics & Management*. <https://doi.org/10.31410/EMAN.2020.171>
- Sun, C., Li, Y., Wu, D., & Boule, B. (2025). *OnionEval: An unified evaluation of fact-conflicting hallucination for small-large language models*. <https://doi.org/10.48550/arXiv.2501.12975>
- Tsangaris, P. (2017). Competition law enforcement. In *Capacity withdrawals in the electricity wholesale market: Between competition law and regulation* (pp. 37–103). Springer Nature Link. [https://doi.org/10.1007/978-3-662-55513-2\\_3](https://doi.org/10.1007/978-3-662-55513-2_3)
- Valentin, S., Fu, J., Detommaso, G., Xu, S., Zappella, G., & Wang, B. (2024). *Cost-effective hallucination detection for LLMs*. <https://www.amazon.science/publications/cost-effective-hallucination-detection-for-llms>
- Wang, P., Brown, C., Jennings, J. A., & Stolee, K. T. (2020). An empirical study on regular expression bugs. *Proceedings of the 17th International Conference on Mining Software Repositories*, 103–113. <https://doi.org/10.1145/3379597.3387464>