

Finnish Media Scrapers

Eetu Mäkelä^{*1} and Pihla Toivanen^{†1}

¹ University of Helsinki

DOI: [10.21105/joss.03504](https://doi.org/10.21105/joss.03504)

Software

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

Editor: [Andrew Stewart](#) ↗

Reviewers:

- [@sara-shiho](#)
- [@pmyteh](#)

Submitted: 21 June 2021

Published: 16 July 2021

License

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

Summary

Finnish Media Scrapers is a package for extracting articles from Finnish journalistic media websites by the [University of Helsinki Human Sciences – Computing Interaction research group](#). Included are scrapers for the four biggest Finnish journalistic media: [YLE](#), [Helsingin Sanomat](#), [Iltalehti](#) and [Iltasanomat](#).

Statement of need

There is an increasing need for user-friendly computational tools in the humanities and social sciences. For example, a common workflow in media research is to collect a large amount of data and combine quantitative and qualitative methods in the analysis phase ([Koivunen et al. \(2021\)](#), [Weber & Monge \(2011\)](#)). This package responds to the research needs by providing easy-to-use tools for scraping Finnish media articles and extracting the article texts from the scraped HTML files. At the same time, the functionality has also been packaged as a Python module for the benefit of more computationally-savvy users.

The scripts were originally developed for a data journalism article ([Suomen Kuvalehti et al. \(2021\)](#)) analyzing how Finnish members of parliament were represented in the media in 2020. Further developing and packaging the scripts into a reusable package was based on an expressed interest from the Finnish computational science community. Since initial beta release a couple of months ago, the package is now known to be already used in at least two research projects targeting Finnish media analysis.

General workflow

The general workflow for using the scrapers is as follows: 1. The scrapers support specifying a keyword as well as a timespan for extraction, and output a CSV of all matching articles with links. 2. A second set of scripts then allows downloading the matched articles in HTML format. 3. Third, there are further scripts for extracting plain text versions of the article texts out of the HTML. 4. Finally, a script exists to post-filter the resulting plain texts again with keywords.

Important to know when applying the workflow is that due to the fact that all the sources use some kind of stemming for their search, they can often return also spurious hits. Further, if searching for multiple words, the engines often perform a search for either word instead of the complete phrase. The post-filtering script above exists to counteract this by allowing the refiltering of the results more rigorously and uniformly locally.

^{*}corresponding author

[†]co-first author

35 At the same time and equally importantly, the stemming for a particular media may not cover
36 e.g. all inflectional forms of words. Thus, it often makes sense to query for at least all common
37 inflected variants and merge the results. For a complete worked up example of this kind of
38 use, see the [members_of_parliament](#) folder, which demonstrates how one can collect and
39 count how many articles in each media mention the members of the Finnish Parliament.

40 Acknowledgements

41 We acknowledge contributions from the Suomen Kuvalehti team (Samuel Nyroos, Salla
42 Vuorikoski and Leena Sharma) during the testing phase of the scrapers.

43 References

- 44 Koivunen, A., Kanner, A., Janicki, M., Harju, A., Hokkanen, J., & Mäkelä, E. (2021). Emo-
45 tive, evaluative, epistemic: A linguistic analysis of affectivity in news journalism. *Journal-*
46 *ism*, 22(5), 1190–1206. <https://doi.org/10.1177/1464884920985724>
- 47 Suomen Kuvalehti, Mäkelä, E., & Toivanen, P. (2021). Vuosi valokeilassa: Kuka sai me-
48 dialta huomiota? Kuka jäi varjoon? Suomen kuvalehti selvitti tutkijoiden kanssa, miten
49 kansanedustajat näkyivät neljässä suuressa uutismediassa vuonna 2020. *Suomen Kuvale-*
50 *hti*, 24–33.
- 51 Weber, M. S., & Monge, P. (2011). The flow of digital news in a network of sources,
52 authorities, and hubs. *Journal of Communication*, 61(6), 1062–1081. <https://doi.org/10.1111/j.1460-2466.2011.01596.x>
- 53