



**HAL**  
open science

## Screening articles for tortured phrases with a regular expressions-based detector

Alexandre Clause, Guillaume Cabanac, Pascal Cuxac, Ophélie Fraisier-Vannier,  
Cyril Labbé

### ► To cite this version:

Alexandre Clause, Guillaume Cabanac, Pascal Cuxac, Ophélie Fraisier-Vannier, Cyril Labbé. Screening articles for tortured phrases with a regular expressions-based detector. PRC'25: 10th International Congress on Peer Review and Scientific Publication, Sep 2025, Chicago, United States. . <hal-05273043>

**HAL Id: hal-05273043**

**<https://hal.science/hal-05273043v1>**

Submitted on 22 Sep 2025

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons CC BY 4.0 - Attribution - International License

# Screening articles for tortured phrases with a regular expressions based detector

Alexandre Clause<sup>1</sup>, Guillaume Cabanac<sup>1,2</sup>, Pascal Cuxac<sup>3</sup>, Ophélie Fraisier-Vannier<sup>1</sup>, Cyril Labbé<sup>4</sup>

<sup>1</sup> IRIT UMR 5505, University of Toulouse, Toulouse, France

<sup>2</sup> Institut Universitaire de France, Paris, France


<sup>3</sup> INIST-CNRS UAR 67, Vandoeuvre-lès-Nancy, France

<sup>4</sup> LIG UMR 5217, Univ. Grenoble Alpes, Grenoble, France


<b>Objective</b>	Provide a free-to-use tortured phrases detector to identify problematic articles
<b>Contribution</b>	A standalone algorithm using regular expressions on a list of identified fingerprints
<b>Results</b>	Comparable performance to the current PPS-Dimensions screening process

7k+ fingerprints	100k+ Hindawi articles
<b>Experiment</b>	
Between 2020 and 2022	18 disciplines (STEM)

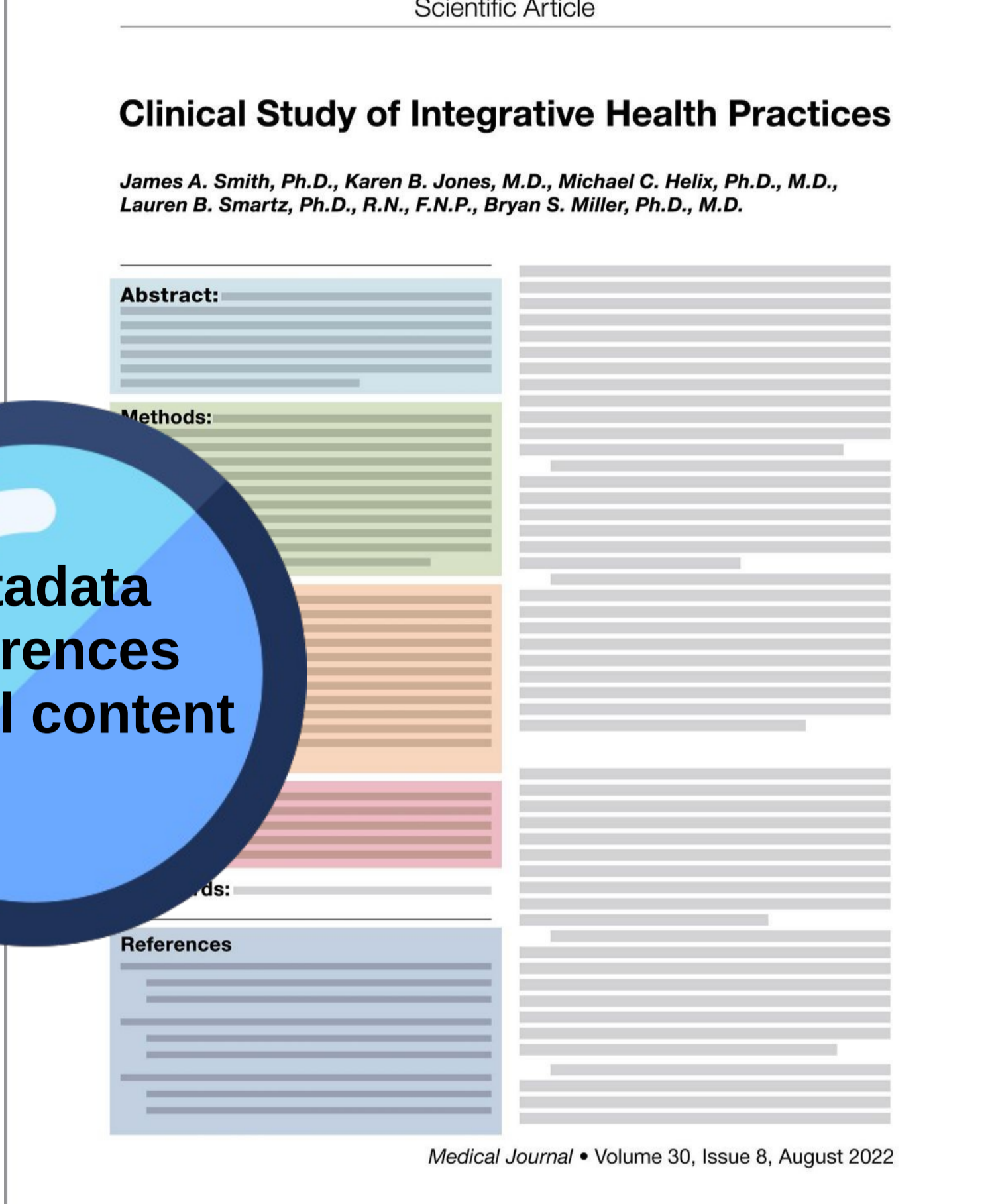
## Detection of tortured phrases in the PPS



**Problematic Paper Screener**  
Est. February 27<sup>th</sup>, 2021

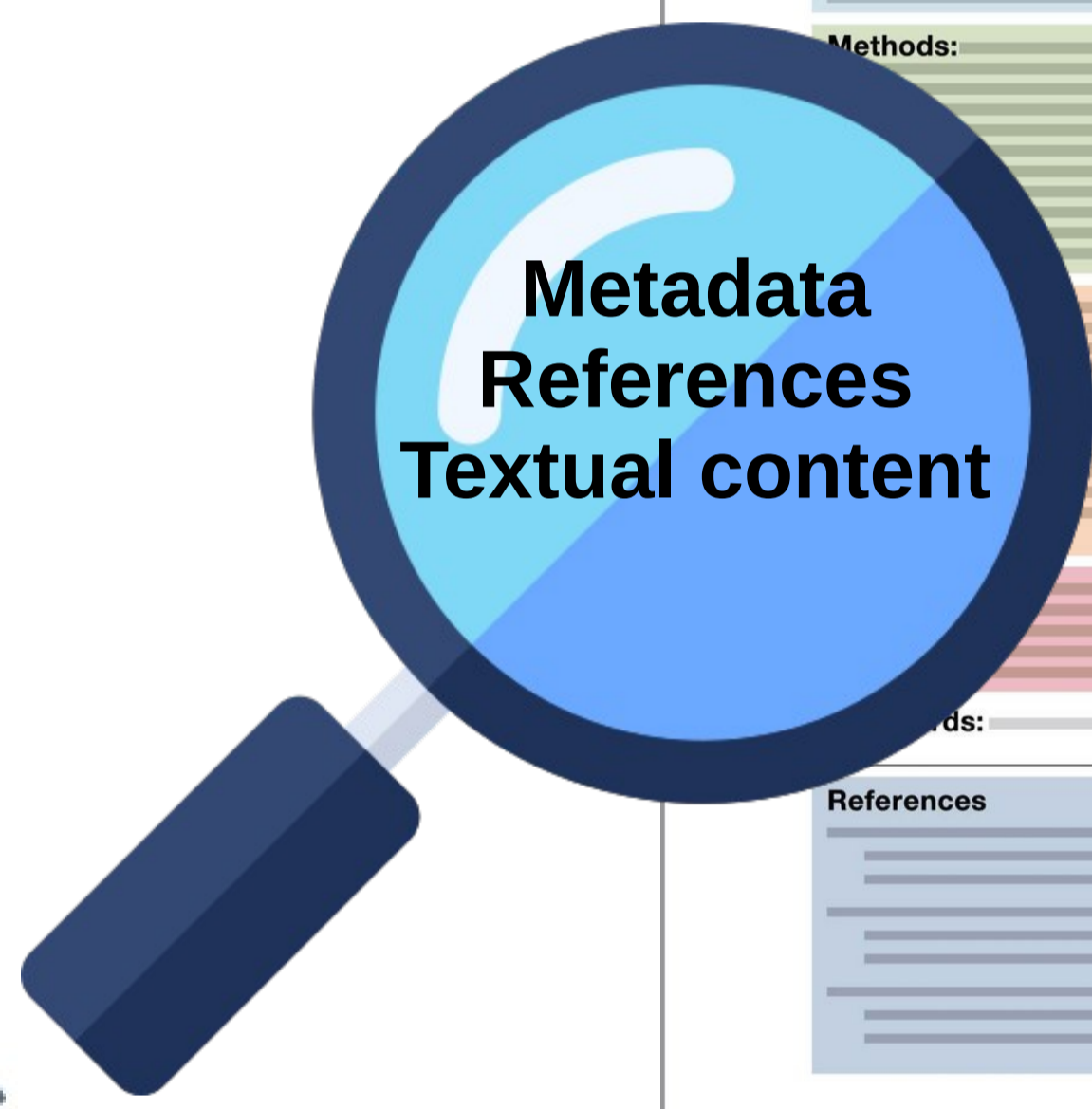


**Dimensions**  
A Digital Science Solution

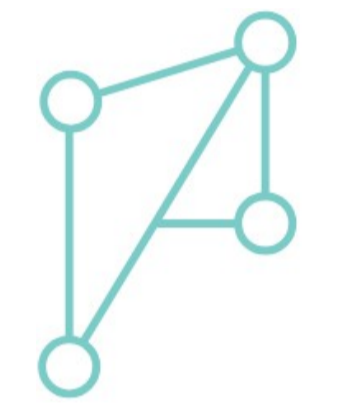


Detectors:	Fingerprint Tortured phrase	Expected text
1 Annulled	"Information mining"	Data mining
1 Tortured	"Grouping methods"	Classification methods
2 SCigen	"Huge information"	Big data
3 Mathgen	"Crown epidemic"	Corona epidemic
4 SBIR	"Data technology (IT)"	Information technology
5 Suspect		
6 Seek&Blastn		
7 Problematic Cell Lines		
8 Citejacked		
9 Feet of Clay		

7k+ fingerprints over 130M articles



**Metadata  
References  
Textual content**



**PUBPEER**  
The online Journal club

This Wiley article cited 7 times contains several **tortured phrases** that make some passages hard to parse. These typically result from an attempt to avoid plagiarism detection using a paraphrasing software. So far, the following have been spotted:

Tortured Phrases (found)	Established Phrases (expected)
Fourier change	Fourier transform
malignant growth AND cancer	cancer
receptive oxygen species	reactive oxygen species
sub-atomic weight	molecular weight
warmth exchange	heat transfer
statically significant	A typo: statistically significant

Can the author explain why they departed from the established phrases?

Screening process

## Results on 100k+ Hindawi articles

401 articles found in Dimensions but missing from the XML corpus metadata

Hindawi XML corpus: 100k+ articles

Flagged using regular expressions: 2,455 articles

1,948 mutual results from 139 journals: 58% overlap

Flagged by the PPS: 3,400 articles

Top-200 results: 48 false positives & 52 false negatives due to indexing issues

## Conclusion

- We obtained comparable results to the current screening process
- We invite publishers to include our method in their evaluation process
- We provided feedback for the PPS and Dimensions
- Future work will focus on screening figures, tables, and references

## References

1. Cabanac G, Labbé C, Magazinov A. The 'Problematic Paper Screener' automatically selects suspect publications for post-publication (re)assessment. Preprint. 2022. doi:10.48550/arXiv.2210.04895
2. Cabanac G, Labbé C, Magazinov A. Tortured phrases: A dubious writing style emerging in science. Evidence of critical issues affecting established journals. Preprint. 2021. doi:10.48550/arXiv.2107.06751
3. Herzog C, Hook D, Konkiel S. Bringing down barriers between scientometricians and data. Quantitative Science Studies. 2020;1:387–395. doi:10.1162/qss\_a\_00020

✉ alexandre.clause@irit.fr    🔗 irit.fr/iris



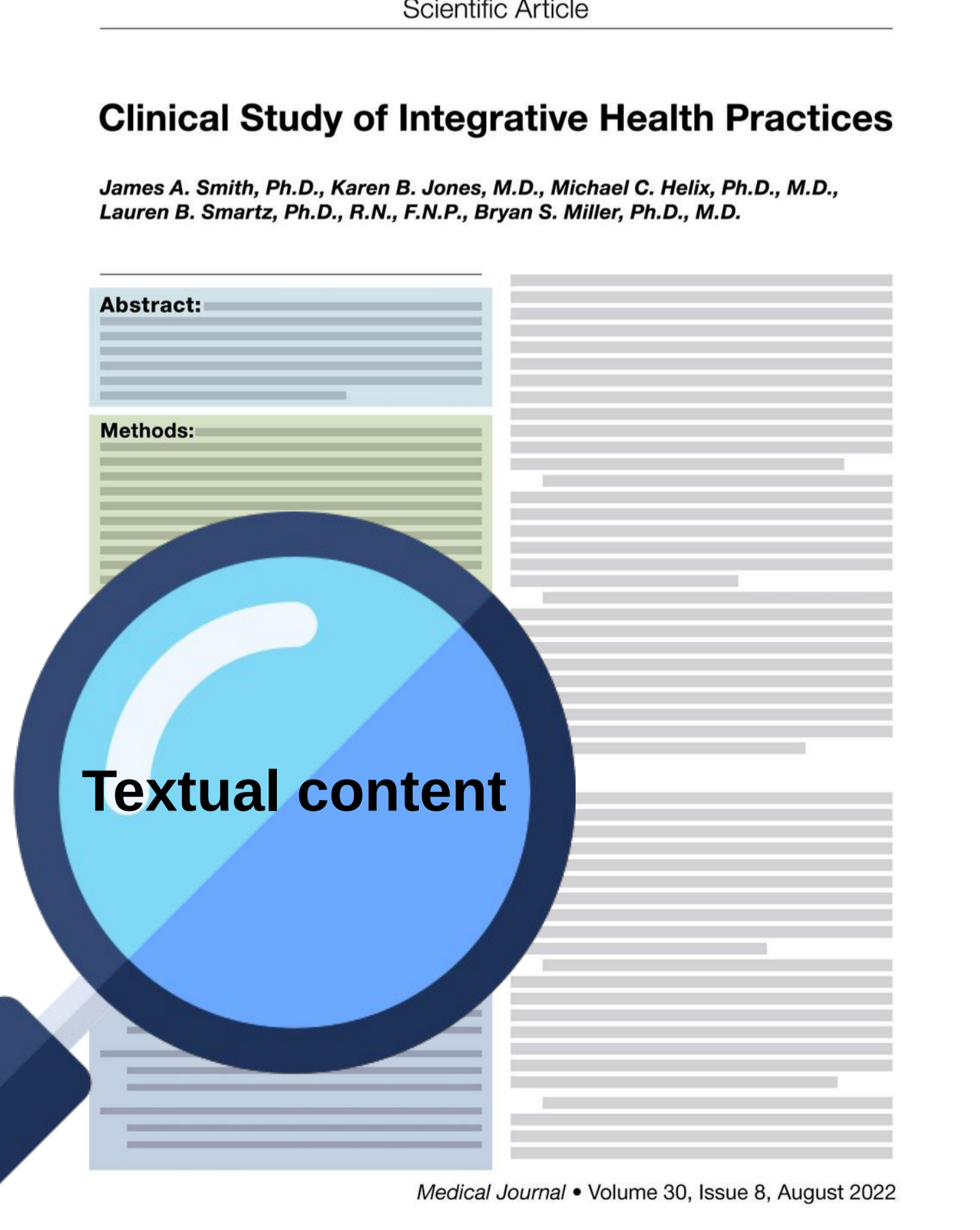
The NanoBubbles project has received Synergy grant funding from the European Research Council (ERC), within the European Union's Horizon 2020 programme, grant agreement no. 951393.

## Detection of tortured phrases in our approach

Fingerprint Tortured phrase	Expected text	Fingerprint type
"Man-made consciousness"	Artificial intelligence	A
"Profound learning" AND "deep learning"	Deep learning	A AND B
"128 pieces"~5	128 bits	A B~X
"Concentrate highlight" NOT specular	Extract feature	A NOT B
"Affirmed recuperated demise" OR "affirmed recuperated passing"	Confirmed, recovered and death (cases)	A OR B
"Partial administrator" AND ("Baleanu" OR "Riemann")	Fractional operator	A AND (B OR C)

**Regular expression detection**

```
n_matches(\bAs?\b.?\bAs?\b) > 0
min(n_matches(\bAs?\b.?\bAs?\b), n_matches(\bBs?\b.?\bBs?\b)) > 0
n_matches(\bAs?\b.?\bAs?\b (\b\w+\b ){0,X} \bBs?\b.?\bBs?\b) > 0
n_matches(\bAs?\b.?\bAs?\b) > 0
and
n_matches(\bBs?\b.?\bBs?\b) = 0
max(n_matches(\bAs?\b.?\bAs?\b), n_matches(\bBs?\b.?\bBs?\b)) > 0
max(
  min(n_matches(\bAs?\b.?\bAs?\b), n_matches(\bBs?\b.?\bBs?\b)),
  min(n_matches(\bAs?\b.?\bAs?\b), n_matches(\bCs?\b.?\bCs?\b))
) > 0
```



**Textual content**

DOI	10.1155/2020/1280632
Tortured phrases	Brilliant AND keen AND savvy AND srewd Conveyance framework Nucleic corrosive Number juggling
Tips	4

Screening process