



HAL
open science

Mining tortured acronyms from the scientific literature

Alexandre Clause, Guillaume Cabanac, Pascal Cuxac, Cyril Labbé

► **To cite this version:**

Alexandre Clause, Guillaume Cabanac, Pascal Cuxac, Cyril Labbé. Mining tortured acronyms from the scientific literature. 2023. hal-04311600

HAL Id: hal-04311600

<https://hal.science/hal-04311600>

Preprint submitted on 28 Nov 2023

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 Attribution 4.0 International License

Mining tortured acronyms from the scientific literature

[Alexandre CLAUSSE](#)¹, [Guillaume CABANAC](#)^{1, 2}, [Pascal CUXAC](#)³, and [Cyril LABBÉ](#)⁴

¹ Université Toulouse 3 – Paul Sabatier, IRIT UMR 5505 CNRS, Toulouse, France

² Institut universitaire de France (IUF), Paris, France

³ INIST – CNRS, UAR76, Vandœuvre-lès-Nancy, France

⁴ Univ. Grenoble Alpes, CNRS, Grenoble INP, LIG, 38000 Grenoble, France

Objective. The ‘Problematic Paper Screener’ (PPS, WCRI’22, <https://doi.org/10.48550/arXiv.2210.04895>) supports the human re-assessment of scientific articles flagged as suspicious. The ‘tortured detector’ tabulates 12k papers containing tortured phrases: established scientific concepts paraphrased with synonyms, such as ‘butt-centric waterway’ for ‘anal canal.’ Some acronyms are even tortured, such as ‘Convolutional Brain Organisation (CNN)’ for ‘Convolutional Neural Network (CNN).’ This abstract tackles the following task: discover and classify all acronyms from any given article: tortured or genuine.

Method. First, we built a test collection by sampling 75 tortured articles in open access in engineering. The first author visually studied each article to highlight and extract all acronyms, each one was labelled as tortured or not. Second, we designed heuristics to classify acronyms probing initials (mis)match. Third, we benchmarked the algorithm’s output to compute recall and precision scores.

Results. The publicly released test collection (<https://doi.org/10.5281/zenodo.10014634>) contains 975 acronyms with 355 tortured and 53 suspects ones. Most tortured acronyms mismatch their initials (98%), genuine acronyms match their initials (97%), both may contain compounds and stop words. These statistics informed the designing of the classifier that yields 76% precision and 89% recall. Incremental failure analysis identifies false positives (e.g., examples given in brackets) and false negatives (i.e., matching acronym–initials yet irrelevant). As a positive consequence, we extended the list of tortured phrases used by the PPS with 190 new tortured phrases (4%). This extended list now informs the detection of papers featuring tortured acronyms.

Conclusion. The ‘acronym classifier’ allows the screening of incoming manuscripts or peer-reviewed articles and spot questionable passages. It is as agnostic as it requires no prior knowledge of tortured acronyms already found. Publishers may add this screener in their editorial workflow, scientific sleuths can use it to screen articles for post-publication review. We welcome developers to use the test collection (and its foreseeable extended versions) to benchmark their own implementations and compare to the baseline performance reported in this WCRI abstract.

Acknowledgements. GC and CL acknowledge the [NanoBubbles](#) project that has received Synergy grant funding from the European Research Council (ERC), within the European Union’s Horizon 2020 program, grant agreement no. [951393](#).