

Diachronic NLP Analysis of “Therapy-Speak” in General Web Discourse: ADHD and Autism in Common Crawl

Jakob Lütchemeier

MSc Dissertation submitted in partial fulfilment of the requirements
for the degree of

MSc in Applied Social Data Science

School of Social Sciences and Philosophy

Department of Political Science

Supervisor: Dr Tom Paskhalis

Course Directors: Dr Jeffrey Ziegler; Dr Tom Paskhalis

Trinity College Dublin

10 August 2026

Note: Student ID is intentionally not printed (university guidance).

Declaration

I hereby declare that this MSc Dissertation is entirely my own work and that it has not been submitted as an exercise for a degree at this or any other university.

I have read and I understand the plagiarism provisions in the General Regulations of the University Calendar for the current year, found at <http://www.tcd.ie/calendar>.

I have also completed the Online Tutorial on avoiding plagiarism "Ready Steady Write", located at <http://tcd-ie.libguides.com/plagiarism/ready-steady-write>.

Signed: _____

Date: _____

Abstract

This dissertation investigates whether “therapy-speak” related to ADHD and autism has increased in general web discourse over the last decade, and how the surrounding language and framing may have shifted over time. Using Common Crawl as a large-scale web archive, I implement a reproducible pipeline that samples one crawl per year and extracts plaintext (WET) documents for trend estimation and a smaller, deeper corpus for downstream NLP analysis. The study emphasises methodological feasibility: principled sampling, robust filtering, contextual disambiguation, and domain caps to prevent large publishers from dominating the dataset.

Acknowledgements

I would like to thank my supervisor, Dr Tom Paskhalis, for guidance throughout this dissertation, and the MSc in Applied Social Data Science teaching team for methodological training and support.

Contents

Abstract	ii
Acknowledgements	iii
1 Introduction	1
2 Literature Review	2
3 Method	3
4 Results	4
5 Discussion	6
6 Conclusion	7
Appendices	8

List of Figures

4.1	Pilot hits by term label in the sampled WET documents.	5
-----	--	---

List of Tables

1 | Introduction

Research Questions

1. How has the prevalence of ADHD- and autism-related terminology in general web discourse changed over approximately the last 10 years?
2. How has the surrounding context and framing of these terms shifted over time (e.g., clinical vs. everyday usage; self-identification; moral/emotional framing)?

2 | Literature Review

3 | Method

No-copy-paste outputs

Figures are generated by code into `../reports/figures/` and included via `\includegraphics`. Tables are generated by code into `../reports/tables/` and included via `\input`.

Pilot validation outputs. Stage 1 pilot validation outputs are provided in the Appendix (see Appendix A), documenting throughput, hit-rate, domain concentration, and the observed boilerplate issue that Stage 1b will mitigate.

4 | Results

5 | Discussion

6 | Conclusion

Appendices

Stage 1 Pilot: Pipeline validation and data-quality diagnostics

Metric	Value
Docs scanned	157,461
Docs \geq min chars	146,995
Hits total	728
Hit rate (hits / docs \geq min chars)	0.005
Unique domains (all)	73,793
Unique domains (hits)	384
Capped removed	0

Domain	Hits	Share
blogspot.com	29	0.0398
medicalxpress.com	21	0.0288
twinkl.co.uk	21	0.0288
sciencedaily.com	16	0.0220
evitamins.com	14	0.0192
mycity4kids.com	10	0.0137
wordpress.com	10	0.0137
myaspergerschild.com	8	0.0110
holidays.net	7	0.0096
healthline.com	7	0.0096

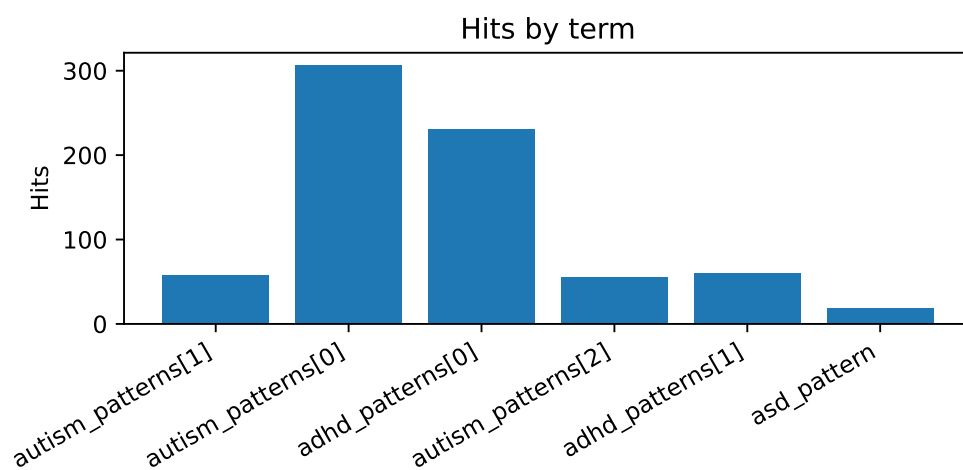


Figure 1: Pilot hits by term label in the sampled WET documents.