

Commentary

Moving From Keywords to Contextual Meaning: A Commentary on Hybrid Bibliometric Synthesis in Health Research

Dimitrios Zikos, BSN, MSc, PhD

Department of Healthcare Management and Leadership, College of Health Professions, Texas Tech University Health Sciences Center, Lubbock, TX, United States

Corresponding Author:

Dimitrios Zikos, BSN, MSc, PhD
Department of Healthcare Management and Leadership
College of Health Professions, Texas Tech University Health Sciences Center
3601 4th Street
Lubbock, TX 79430
United States
Phone: 1 9894301787
Email: dzikos@ttuhsc.edu

Related Article:

Comment on: <https://www.jmir.org/2026/1/e86200>

Abstract

The fast growth of social media mining in health research has contributed to an invaluable but quite fragmented body of literature. As the amount of unstructured patient-reported data grows, traditional bibliometric analyses face methodological limitations, particularly regarding synonym fragmentation and arbitrary parameter selection. In their recent publication, “Thematic Mapping and Evolution of Social Media Mining in Health Research: Hybrid Bibliometric Synthesis,” Yang and Bohnet-Joschko attempt to address these flaws by introducing a semantic-structural (hybrid) bibliometric framework. This commentary evaluates the methodological innovations of their study and its departure from traditional syntactic keyword-matching tools. By combining citation-informed transformers (SPECTER2) and biomedical language models (PubMedBERT) and dimensionality reduction and density-based clustering, the authors created a reproducible pipeline. In their architecture, they start with foundational machine learning (statistical validity) before transitioning into large language models for qualitative synthesis. I will attempt to explain how this transition from syntactic mapping to semantic vector representation solves known challenges in evidence synthesis, naturally grouping conceptual synonyms without artificially forcing boundaries on the literature. Furthermore, I examine the practical implications of their temporal findings. Such real-time social media mining applications can be very useful for retrospective reporting and evaluating targeted public health interventions. While this pipeline offers high generalizability across disciplines, it also introduces a computational literacy barrier to some, and this re-emphasizes the need for data literacy for health professions. Ultimately, the study provides a transparent approach to informatics because mathematically validated frameworks are foundational for the future of evidence-driven public health policy and clinical decision-making.

J Med Internet Res 2026;28:e102159; doi: [10.2196/102159](https://doi.org/10.2196/102159)

Keywords: bibliometrics; machine learning; semantic mapping; public health surveillance; social media

Introduction

The fast growth of social media mining (SMM) in health research has provided new opportunities to study public health. However, the volume of this literature has created a fragmented knowledge base that resists traditional evidence

synthesis. In their recent article, “Thematic Mapping and Evolution of Social Media Mining in Health Research: Hybrid Bibliometric Synthesis,” Yang and Bohnet-Joschko [1] go beyond descriptive field mapping and introduce a more semantically focused framework for evidence synthesis by extending traditional bibliometric methods with a

hybrid semantic-structural pipeline. Since the health care data resources and knowledge base have already become massive and complex, synthesizing literature needs to go beyond just “cataloging” publications. This commentary explores how their methodology attempts to overcome statistical limitations, optimizes analytical architecture, and supports contextually relevant public health applications.

The Limitations of Traditional Keyword Matching

For over a decade, bibliometric reviews have mostly relied on out-of-the-box platforms like VOSviewer, but in some ways, those may reduce transparency regarding parameter sensitivity and clustering behavior [2]. They are created around prebuilt clustering algorithms (such as the Louvain method) that do not provide the researcher control over cluster thresholding, network normalizations, and parameter sensitivity. Additionally, these traditional methods are constrained by syntactic mapping: they treat keywords, titles, and abstracts of articles as isolated strings of text.

For example, traditional tools group papers by vocabulary (linking a Twitter flu paper with a bird flu paper), while the hybrid model groups them by scientific intent. Unless a researcher manually creates a detailed thesaurus, the algorithm fragments the literature based on vocabulary rather than meaning. Similarly, legacy topic modeling methods like latent Dirichlet allocation (LDA) rely on a rigid “bag-of-words” assumption that ignores contextual syntax and forces researchers to predefine the number of topics, artificially “bounding” the data [3].

Yang and Bohnet-Joschko [1] attempted to solve this by moving from syntactic to semantic mapping (Table 1). By choosing to use SPECTER2 and PubMedBERT embeddings, they ingest the entire context of a text and convert documents into high-dimensional vectors [4,5]. Because these models understand semantic meaning, they group conceptual synonyms together, making fewer errors than the traditional string-matching tools. Furthermore, because SPECTER2 is pretrained on very large amounts of citation data, it makes the approach an interesting hybrid of text mining and citation analysis.

Table 1. Comparison of traditional bibliometric methods to Yang and Bohnet-Joschko’s [1] pipeline.

Dimension	Traditional bibliometric methods	Yang and Bohnet-Joschko’s [1] hybrid pipeline
Data processing	Syntactic mapping of isolated text strings	Semantic vector representation using context
Core algorithms	Out-of-the-box platforms (eg, VOSviewer) using Louvain or latent Dirichlet allocation	Citation-informed transformers (SPECTER2), biomedical models (PubMedBERT), UMAP ^a , and HDBSCAN ^b
Clustering	Relies on forced boundaries, often requiring predefined topic counts	Density-based clustering that mathematically isolates unassigned noise
Vocabulary handling	Fragments literature based on terminology unless a manual thesaurus is built	Naturally groups conceptual synonyms together based on scientific intent

^aUMAP: uniform manifold approximation and projection.

^bHDBSCAN: hierarchical density-based spatial clustering of applications with noise.

Optimization of the Architecture

The strength of Yang and Bohnet-Joschko’s [1] methodology is in the sequencing of its analytical architecture. In the current era of artificial intelligence, it is tempting to feed raw data directly into a large language model. The study argues for an emerging informatics principle: structurally validated machine learning pipelines should precede large language model-assisted interpretation.

Yang and Bohnet-Joschko [1] validate their data structure using uniform manifold approximation and projection (UMAP) for dimensionality reduction [6] and hierarchical density-based spatial clustering of applications with noise (HDBSCAN) for structural clustering [7]. Unlike LDA, HDBSCAN does not require a predefined cluster count. It identifies naturally occurring dense regions of literature and mathematically isolates outlier papers rather than forcing them into irrelevant groups. This way, the subsequent thematic synthesis is grounded in reproducible data science.

From Methodology to Public Health Practice

Beyond the methodological advantages, the study’s [1] temporal slicing is, for SMM, an example of moving from computational experimentation to real-world public health engagement. The prominence of application-driven clusters, such as infodemiology and sociopsychological determinants, shows well how SMM can be a useful tool for community health surveillance [8]. Accurate SMM is essential for evaluating localized, real-world health initiatives. For instance, evaluating colorectal cancer screening prevention campaigns in Texas requires an understanding of localized, real-time community sentiment, barriers to access, and sociopsychological hesitation. SMM captures this narrative, augmenting structured retrospective federal tools that often lag by months or years. Such real-time sentiment surveillance may also support decision-making for health systems, including targeted outreach, misinformation monitoring, and fast evaluation of intervention uptake.

Furthermore, the Yang and Bohnet-Joschko [1] treatment of HDBSCAN’s unassigned noise (cluster 1) is a very

reasonable analytical choice. Rather than dismissing this noise, they identify it as a candidate “incubator pool.” I believe that embracing these topics can help researchers who want to predict the next wave of sociopsychological determinants before they become new established literature.

Limitations and Future Work

Despite its advantages over legacy systems, this pipeline introduces some challenges. The transition from graphical user interface-based tools to code-based learning models introduces a computational barrier to entry. There is a risk that evidence synthesis becomes gated behind advanced data science skills. Furthermore, while Yang and Bohnet-Joschko’s [1] dual-level validation ensures semantic coherence, the study lacks an empirical comparison against older baselines (such as LDA or Louvain clustering) on the identical dataset, which would be necessary to quantify the reduction in fragmentation. Additionally, future validation across multiple bibliographic databases (eg, PubMed, Scopus, and Web of Science) would help determine the stability of the pipeline.

Researchers must be careful not to treat these new machine learning algorithms as new inherently opaque systems. While

HDBSCAN eliminates the need to guess cluster counts, it introduces hypersensitivity to the minimum cluster size and the minimum sample size. Similarly, UMAP is dependent on neighbor and distance parameters. Future researchers who adopt this methodology will need to optimize and report their parameter selections to ensure their underlying models match their specific discipline. As a researcher and instructor, I would like to stress that addressing this requires shifts in health informatics education, promoting learning that relies on algorithmic logic and mechanics, and not just software operation.

Conclusion

Yang and Bohnet-Joschko’s [1] study moved successfully from syntactic keyword-matching to semantic vector representation, in an attempt to avoid the limitations of traditional bibliometrics. More importantly, they provided a transparent, reproducible blueprint for future studies. As the volume of medical literature and patient-reported data continues to increase, adopting machine learning pipelines is an imperative for any research designed to extract evidence-driven insights to guide patient care and public health policy.

Acknowledgments

Generative artificial intelligence was not used in any capacity during the preparation, writing, or editing of this manuscript.

Funding

The author declared no financial support was received for this work.

Conflicts of Interest

None declared.

References

1. Yang MJ, Bohnet-Joschko S. Thematic mapping and evolution of social media mining in health research: hybrid bibliometric synthesis. *J Med Internet Res*. May 8, 2026;28:e86200. [doi: [10.2196/86200](https://doi.org/10.2196/86200)] [Medline: [42115141](https://pubmed.ncbi.nlm.nih.gov/42115141/)]
2. van Eck NJ, Waltman L. Software survey: VOSviewer, a computer program for bibliometric mapping. *Scientometrics*. Aug 2010;84(2):523-538. [doi: [10.1007/s11192-009-0146-3](https://doi.org/10.1007/s11192-009-0146-3)] [Medline: [20585380](https://pubmed.ncbi.nlm.nih.gov/20585380/)]
3. Blei DM, Ng AY, Jordan MI. Latent Dirichlet allocation. *J Machine Learning Res*. 2003;3:993-1022. URL: <https://www.jmlr.org/papers/volume3/blei03a/blei03a.pdf> [Accessed 2026-06-01]
4. Cohan A, Feldman S, Beltagy I, Downey D, Weld DS. SPECTER: document-level representation learning using citation-informed transformers. *arXiv*. Preprint posted online on Apr 15, 2022. [doi: [10.48550/arXiv.2004.07180](https://doi.org/10.48550/arXiv.2004.07180)]
5. Gu Y, Tinn R, Cheng H, et al. Domain-specific language model pretraining for biomedical natural language processing. *ACM Trans Comput Healthcare*. Jan 31, 2022;3(1):1-23. [doi: [10.1145/3458754](https://doi.org/10.1145/3458754)]
6. McInnes L, Healy J, Melville J. UMAP: uniform manifold approximation and projection for dimension reduction. *arXiv*. Preprint posted online on Feb 9, 2018. [doi: [10.48550/arXiv.1802.03426](https://doi.org/10.48550/arXiv.1802.03426)]
7. Campello R, Moulavi D, Sander J. Density-based clustering based on hierarchical density estimates. In: Pei J, Tseng VS, Cao L, Motoda H, Xu G, editors. *Advances in Knowledge Discovery and Data Mining: 17th Pacific-Asia Conference, PAKDD 2013, Gold Coast, Australia, April 14-17, 2013, Proceedings, Part II*. 2013:160-172. [doi: [10.1007/978-3-642-37456-2_14](https://doi.org/10.1007/978-3-642-37456-2_14)]
8. Eysenbach G. Infodemiology and infoveillance tracking online health information and cyberbehavior for public health. *Am J Prev Med*. May 2011;40(5 Suppl 2):S154-S158. [doi: [10.1016/j.amepre.2011.02.006](https://doi.org/10.1016/j.amepre.2011.02.006)] [Medline: [21521589](https://pubmed.ncbi.nlm.nih.gov/21521589/)]

Abbreviations

HDBSCAN: hierarchical density-based spatial clustering of applications with noise
LDA: latent Dirichlet allocation
SMM: social media mining

UMAP: uniform manifold approximation and projection

Edited by Stephanie Law, Tiffany Leung; This is a non-peer-reviewed article; submitted 22.May.2026; accepted 22.May.2026; published 03.Jun.2026

Please cite as:

Zikos D

Moving From Keywords to Contextual Meaning: A Commentary on Hybrid Bibliometric Synthesis in Health Research
J Med Internet Res 2026;28:e102159

URL: <https://www.jmir.org/2026/1/e102159>

doi: [10.2196/102159](https://doi.org/10.2196/102159)

© Dimitrios Zikos. Originally published in the Journal of Medical Internet Research (<https://www.jmir.org>), 03.Jun.2026. This is an open-access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in the Journal of Medical Internet Research (ISSN 1438-8871), is properly cited. The complete bibliographic information, a link to the original publication on <https://www.jmir.org/>, as well as this copyright and license information must be included.