
Guest editorial: Extraction and evaluation of knowledge entities in the age of artificial intelligence

Guest editorial

433

1. Introduction

Scientific documents serve as an essential mediator for research achievements and scientific knowledge. With the increasing availability of full-text data and advanced data analytical techniques, bibliometric researchers have started to focus on granular content within scientific documents, transitioning their foci from the external metadata of these documents to the internal knowledge they encompass. In scientific documents, knowledge consists of many interconnected units, known as knowledge entities (Ding *et al.*, 2013).

Ding *et al.* (2013) classified knowledge entities into macro-level (e.g. author, journal, references), meso-level (e.g. keywords) and micro-level (e.g. dataset, method, biomedical entities) categories. These categories refer to traditional bibliographic entities for evaluation purposes and knowledge entities found within the full-text content, respectively. The extraction and evaluation of knowledge entities can improve existing knowledge services and meet the needs of researchers for rapid and accurate access to scientific knowledge (Ma *et al.*, 2023; Mayr *et al.*, 2014). Scholars can evaluate the scientific, social and economic attributes of knowledge and assess its role in science, technology, innovation and policy. Analyzing dynamic changes over time and identifying geographical differences can provide insights to understand the patterns of knowledge use and dissemination, which can promote knowledge discovery and generate new scientific activities.

The 2nd Workshop on the Extraction and Evaluation of Knowledge Entities from Scientific Documents (EEKE2021) was held online, co-located with the ACM/IEEE Joint Conference on Digital Libraries 2021 (JCDL2021) on September 30, 2021. The EEKE workshop series (<https://eekeworkshop.github.io/>) aims to engage relevant communities in addressing open problems in the extraction and evaluation of knowledge entities from scientific documents. Participants have contributed tremendous accomplishments in identifying knowledge entities, exploring entity features, analyzing entity relationships, and developing extraction platforms or knowledge bases. The workshop has provided scholars, especially early career researchers, with knowledge recommendations and other knowledge entity-based services (Zhang *et al.*, 2021).

2. Topics in this special issue

This special issue collected articles presented at the EEKE2021 workshop and external submissions. A total of 18 papers were submitted, of which 9 were accepted after a rigorous review process. These 9 articles, contributed by 36 authors from 4 countries (China, USA, South Korea and Thailand), are included in this collection. In this editorial, we provide an overview of these papers and introduce the themes and contributions of this Special Issue. The papers are grouped into three topics: entity extraction and entity relations extraction



Aslib Journal of Information
Management
Vol. 75 No. 3, 2023
pp. 433-437
© Emerald Publishing Limited
2050-3806
DOI 10.1108/AJIM-05-2023-507

Chengzhi Zhang acknowledges the National Natural Science Foundation of China (Grant No. 72074113), and Yi Zhang acknowledges the Discovery Early Career Researcher Award granted by the Australian Research Council (Grant No. DE190100994).

(3 articles); annotation tools and knowledge entity graph construction (2 articles) and applications of knowledge entities (4 articles).

2.1 Entity extraction and entity relations extraction

Knowledge entities such as concepts, algorithms and methods are included in scientific literature (Ding *et al.*, 2013; Wang *et al.*, 2022). Extracting knowledge entities and entity relationships from scientific literature plays a crucial role in scientific information retrieval, fine-grained scientific information recommendation and scientific evaluation.

To extract “band gap information” from academic papers, Ghosh and Lu (2023) collected 1.44 million titles and abstracts of scholarly articles related to materials science. They filtered the collection to 11,939 articles that may contain core information on materials and their band gap values. The results revealed that the current system can accurately extract information from 51.32% of the articles, partially extract from 36.62% of the articles and incorrectly extract from 12.04% of the articles.

To tackle the issues of insufficient training corpus and the quality control of annotations, Yan *et al.* (2023) proposed a novel integrated solution for Chinese historical Named Entity Recognition (NER), including automatic entity extraction and man-machine cooperative annotation. This solution is valuable for enhancing the effectiveness of Chinese historical NER and promoting the development of low-resource information extraction.

“Problem-solving” stands out as one of the most fundamental and critical insights of scientific research. Using computer vision as an example, Chen *et al.* (2023) built a “problem-solving” knowledge graph of scientific domains by extracting four entity relation types, namely problem-solving, problem hierarchy, solution hierarchy and association. They illustrated the utility of the extracted relations in constructing domain knowledge graphs and uncovering historical research trends.

2.2 Annotation tools and the construction of knowledge entity graphs

Knowledge entity annotation and knowledge graph construction are essential tasks for knowledge entity applications, including the development of annotation tools, corpora construction and the creation of knowledge graphs.

To explore the use of entity descriptions and network structures in enhancing knowledge graph completion with a high generalization ability across datasets, Yu *et al.* (2023) proposed an entity-description augmented knowledge graph completion model (EDA-KGC). The authors conducted extensive experiments on the FB15K, WN18, FB15K-237 and WN18RR datasets to validate the effectiveness of the model.

To address challenges posed by growing volumes of pre-annotated literature and diverse annotations, such as teamwork, quality control and time management, Wang *et al.* (2023) developed the Bureau for Rapid Annotation Tool (Brat), an annotation collaboration workbench that includes an enhanced semantic constraint system, Vim-like shortcut keys, an annotation filter and a graph-visualizing annotation browser.

2.3 Applications of knowledge entities

Knowledge entities could seed not only granular information retrieval and recommendation (Mayr *et al.*, 2014), intelligent bibliometrics (Zhang *et al.*, 2020), but also new cognitive and practical applications such as the innovative evaluation of scientific document (Liu *et al.*, 2022) and predicting future research directions (Zhang *et al.*, 2023a, 2023b, 2023c).

To explore mental health information entities and the connections between the biomedical, psychological and social domains of bipolar disorder (BD), Timakum *et al.* (2023) used Reddit posts and full-text papers from PubMed Central to extract BD entities and their relationships in the datasets using a dictionary-based and rule-based approach. The findings indicate that

the drug side effects entity was frequently identified in both datasets as a mental health information entity.

Zhang *et al.* (2023a, 2023b, 2023c) introduced a crucial factor in topic selection, i.e. topic popularity, to investigate its relationship with team performance. The authors used gene/protein entities as proxies for topics and extracted them to monitor the development of topic popularity. By comparing various dimensions of team performance, the study explored the relationship between the phase of selected topic popularity and the academic performance of research teams.

Zeng *et al.* (2023) utilized the knowledge elements extracted through the Lexicon-LSTM model to measure the interdisciplinary characteristics of Chinese research in library and information science (LIS). They constructed a subject knowledge graph to support the searching and classification of knowledge elements. The results showed that in LIS, the interdisciplinary diversity indicator exhibited an upward trend from 2011 to 2021, while the disciplinary balance and difference indicators showed a downward trend.

Focusing on discovery of topic evolution path and semantic relationship, Zhang *et al.* (2023a, 2023b, 2023c) identified entities that have the same semantics but different expressions for accurate topic evolution path discovery. They also revealed semantic relationships of topic evolution to better understand what leads to topic evolution. This work provided a new perspective for topic evolution analysis by considering the semantic representation of patent entities.

3. Conclusions

The EEKE workshop series and this Special Issue reinforced the increasing interest of the community in the extraction and evaluation of entities from scientific literature. Endeavors in this field are advancing rapidly, from developing new methods and techniques to their innovative applications in broad practical domains. In today's AI age, three key factors that are driving AI development – namely, data, algorithms and computing power – must work corporately to promote and support one another to achieve success and create value.

Currently, Large-scale pre-trained Language Models (LLMs), such as ChatGPT and GPT-4, have gained widespread popularity due to their broad applications across various industries and fields. There is no doubt that LLMs play an exceptionally important role in scientific documents mining and analysis. Some scholars have tested and evaluated the performance of ChatGPT and GPT-4 in entity extraction (Hu *et al.*, 2023; González-Gallardo *et al.*, 2023) and entity relationship recognition (Rehana *et al.*, 2023) in special domains. With AI's continued advancements and the increasing promotion of open access movements, we can anticipate a more extensive utilization of LLMs, such as Generative Pretrained Transformer (GPT), for extraction and evaluation of knowledge entities from scientific documents in the future.

Chengzhi Zhang

*Department of Information Management, Nanjing University of Science and Technology,
Nanjing, China*

Philipp Mayr

GESIS – Leibniz Institute for the Social Sciences, Cologne, Germany

Wei Lu

School of Information Management, Wuhan University, Wuhan, China, and

Yi Zhang

*Faculty of Engineering and Information Technology, Australian Artificial Intelligence
Institute, University of Technology Sydney, Sydney, Australia*

References

- Chen, G., Peng, J., Xu, T. and Xiao, L. (2023), "Extracting entity relations for 'problem-solving' knowledge graph of scientific domains using word analogy", *Aslib Journal of Information Management*, Vol. 75 No. 3, pp. 481-499, doi: [10.1108/AJIM-03-2022-0129](https://doi.org/10.1108/AJIM-03-2022-0129).
- Ding, Y., Song, M., Han, J., Yu, Q., Yan, E., Lin, L. and Chambers, T. (2013), "Entitymetrics: measuring the impact of entities", *PLoS One*, Vol. 8 No. 8, e71416, doi: [10.1371/journal.pone.0071416](https://doi.org/10.1371/journal.pone.0071416).
- Ghosh, S. and Lu, K. (2023), "Band gap information extraction from materials science literature – a pilot study", *Aslib Journal of Information Management*, Vol. 75 No. 3, pp. 438-454, doi: [10.1108/AJIM-03-2022-0141](https://doi.org/10.1108/AJIM-03-2022-0141).
- González-Gallardo, C.E., Boros, E., Girdhar, N., Hamdi, A., Moreno, J.G. and Doucet, A. (2023), "Yes but.. Can ChatGPT identify entities in historical documents?", arXiv preprint arXiv:2303.17322 doi: [10.48550/arXiv.2303.17322](https://doi.org/10.48550/arXiv.2303.17322).
- Hu, Y., Ameer, I., Zuo, X., Peng, X., Zhou, Y., Li, Z., Li, Y., Li, J., Jiang, X. and Xu, H. (2023), "Zero-shot clinical entity recognition using ChatGPT", arXiv preprint arXiv:2303.16416 doi: [10.48550/arXiv.2303.16416](https://doi.org/10.48550/arXiv.2303.16416).
- Liu, M., Bu, Y., Chen, C., Xu, J., Li, D., Leng, Y., Freeman, R.B., Meyer, E.T., Yoon, W., Sung, M., Jeong, M., Lee, J., Kang, J., Min, C., Song, M., Zhai, Y. and Ding, Y. (2022), "Pandemics are catalysts of scientific novelty: evidence from COVID-19", *Journal of the Association for Information Science and Technology*, Vol. 73 No. 8, pp. 1065-1078, doi: [10.1002/asi.24612](https://doi.org/10.1002/asi.24612).
- Ma, Y., Liu, J., Lu, W. and Cheng, Q. (2023), "From 'what' to 'how': extracting the procedural scientific information toward the metric-optimization in AI", *Information Processing and Management*, Vol. 60 No. 3, 103315, doi: [10.1016/j.ipm.2023.103315](https://doi.org/10.1016/j.ipm.2023.103315).
- Mayr, P., Scharnhorst, A., Larsen, B., Schaer, P. and Mutschke, P. (2014), "Bibliometric-enhanced information retrieval", in de Rijke, M., Kenter, T., de Vries, A.P., Zhai, C., de Jong, F., Radinsky, K. and Hofmann, K. (Eds), *Advances in information retrieval 36th European conference on information retrieval*, Amsterdam, the Netherlands, doi: [10.1007/978-3-319-06028-6_99](https://doi.org/10.1007/978-3-319-06028-6_99).
- Rehana, H., Çam, N.B., Basmaci, M., He, Y., Özgür, A. and Hur, J. (2023), "Evaluation of GPT and BERT-based models on identifying protein-protein interactions in biomedical text", arXiv preprint arXiv:2303.17728, doi: [10.48550/arXiv.2303.17728](https://doi.org/10.48550/arXiv.2303.17728).
- Timakum, T., Song, M. and Kim, G. (2023), "Integrated entitymetrics analysis for health information on bipolar disorder using social media data and scientific literature", *Aslib Journal of Information Management*, Vol. 75 No. 3, pp. 535-560, doi: [10.1108/AJIM-02-2022-0090](https://doi.org/10.1108/AJIM-02-2022-0090).
- Wang, Y., Zhang, C. and Li, K. (2022), "A review on method entities in the academic literature: extraction, evaluation, and application", *Scientometrics*, Vol. 127 No. 5, pp. 2479-2520, doi: [10.1007/s11192-022-04332-7](https://doi.org/10.1007/s11192-022-04332-7).
- Wang, Z., Xu, S., Wang, Y., Chai, X. and Chen, L. (2023), "Bureau for Rapid Annotation Tool: collaboration can do more among variance annotations", *Aslib Journal of Information Management*, Vol. 75 No. 3, pp. 523-534, doi: [10.1108/AJIM-01-2022-0046](https://doi.org/10.1108/AJIM-01-2022-0046).
- Yan, C., Tang, X., Yang, H. and Wang, J. (2023), "A deep active learning-based and crowdsourcing-assisted solution for named entity recognition in Chinese historical corpora", *Aslib Journal of Information Management*, Vol. 75 No. 3, pp. 455-480, doi: [10.1108/AJIM-03-2022-0107](https://doi.org/10.1108/AJIM-03-2022-0107).
- Yu, C., Zhang, Z., An, L. and Li, G. (2023), "A knowledge graph completion model integrating entity description and network structure", *Aslib Journal of Information Management*, Vol. 75 No. 3, pp. 500-522, doi: [10.1108/AJIM-01-2022-0031](https://doi.org/10.1108/AJIM-01-2022-0031).
- Zeng, J., Cao, S., Chen, Y., Pan, P. and Cai, Y. (2023), "Measuring the interdisciplinary characteristics of Chinese research in library and information science based on knowledge elements", *Aslib Journal of Information Management*, Vol. 75 No. 3, pp. 589-617, doi: [10.1108/AJIM-03-2022-0130](https://doi.org/10.1108/AJIM-03-2022-0130).
- Zhang, Y., Porter, A., Cunningham, S.W., Chiavetta, D. and Newman, N. (2020), "Parallel or intersecting lines? Intelligent bibliometrics for investigating the involvement of data science in

policy analysis”, *IEEE Transactions on Engineering Management*, Vol. 68 No. 5, pp. 1259-1271, doi: [10.1109/TEM.2020.2974761](https://doi.org/10.1109/TEM.2020.2974761).

Guest editorial

Zhang, C., Mayr, P., Lu, W. and Zhang, Y. (2021), “Preface to the 2nd workshop on extraction and evaluation of knowledge entities from scientific documents at JCDL 2021”, *Proceedings of the 2nd Workshop on Extraction and Evaluation of Knowledge Entities from Scientific Documents (EEKE 2021) co-located with JCDL 2021*, pp. 1-4, available at: <https://ceur-ws.org/Vol-3004/preface.pdf>.

Zhang, C., Xiang, Y., Hao, W., Li, Z., Qian, Y. and Wang, Y. (2023a), “Automatic recognition and classification of future work sentences from academic articles in a specific domain”, *Journal of Informetrics*, Vol. 17 No. 1, 101373, doi: [10.1016/j.joi.2022.101373](https://doi.org/10.1016/j.joi.2022.101373).

Zhang, J., Liu, Y., Jiang, L. and Shi, J. (2023b), “Discovery of topic evolution path and semantic relationship based on patent entity representation”, *Aslib Journal of Information Management*, Vol. 75 No. 3, pp. 618-642, doi: [10.1108/AJIM-03-2022-0124](https://doi.org/10.1108/AJIM-03-2022-0124).

Zhang, T., Tan, F., Yu, C., Wu, J. and Xu, J. (2023c), “Understanding relationship between topic selection and academic performance of scientific teams based on entity popularity trend”, *Aslib Journal of Information Management*, Vol. 75 No. 3, pp. 561-588, doi: [10.1108/AJIM-03-2022-0135](https://doi.org/10.1108/AJIM-03-2022-0135).