

# How has IxD&A Evolved in the Last Decade? A Topic Modeling Approach to Identify Salient Themes and Cultural Shifts.

Razvan Paroiu<sup>1\*</sup>, Mihai Dascalu<sup>1,2,4</sup> and Carlo Giovannella<sup>3,4</sup>

<sup>1</sup>Computer Science and Engineering Department, National University of Science and Technology Politehnica of Bucharest, 313 Splaiul Independentei, Bucharest, 060042, Romania.

<sup>2</sup>Academy of Romanian Scientists, Nr. 3 Ilfov, Bucharest, 050044, Romania.

<sup>3</sup>University of Rome Tor Vergata Rome, Italy and <sup>4</sup>ASLERD,  
razvan.paroiu@upb.ro (\*corresponding author)

**Abstract.** The volume of published articles has grown exponentially, making research documentation considerably harder. To address this, we have developed and present in this paper a toolkit that can effectively process full-text manuscripts starting from their PDFs, extract topics, and provide correlations between manuscripts based on their common topics. The effectiveness of the toolkit is demonstrated by applying it to the "Interaction Design and Architecture(s)" (IxD&A) journal's full-text articles published between 2013 and 2024 (N = 450). Topic modeling was performed using BERTopic and Llama as LLM to generate coherent topics. As a result, we extracted 246 topics from the entire corpus, which were automatically filtered for specificity using Llama-3.1. In-depth visualizations and analyses with a focus on a human-centered, smart learning perspective, are presented. We release our toolkit as open-source on GitHub to enable users to easily apply our method in other relevant contexts.

**Keywords:** topic modelling, BERTopic, text analysis, Llama 3.1, scientific journal archive

## 1 Introduction

The rapid growth of digital content and academic publications has led to an overwhelming volume of unstructured text across various domains. Scientific journals, educational repositories, and digital libraries now host thousands of articles, making it increasingly difficult for researchers to navigate and extract meaningful insights from these large collections. This overload of information highlights the need for automated methods that can help organize and summarize content on a scale.

One such solution is topic modeling, a statistical or machine learning technique used to discover underlying topics within a collection of documents. Existing approaches automatically group words that frequently appear together across different documents, aiming to discover clusters of words that signify distinct topics. Some of

the first popular algorithms for topic modeling include Latent Dirichlet Allocation (LDA) [1, 2] and Non-Negative Matrix Factorization (NMF) [3, 4]. These models uncover main themes, trends, or subjects in unstructured text, making them essential to modern applications like text mining [5–7], summarization [8–10], and sentiment analysis [11–13].

An alternative to LDA is BERTopic [14], which leverages state-of-the-art language models to extract topics in large text corpora. Developed as an extension of traditional topic modeling approaches, BERTopic utilizes Transformer-based contextualized vector representations of words (embeddings), particularly from BERT (Bidirectional Encoder Representations from Transformers) [15–17], to capture rich semantic information from the text. By combining these embeddings with dimensionality reduction algorithms like UMAP (Uniform Manifold Approximation and Projection for Dimension Reduction) [18] and clustering algorithms like HDBSCAN (Hierarchical Density-Based Spatial Clustering of Applications with Noise) [19], BERTopic can identify coherent and contextually relevant topics with improved accuracy when compared with LDA [20]. Due to its modularity, BERTopic is also highly recommended for large datasets, as the algorithms within the pipeline can be easily modified according to specific requirements.

The objective of this research is to employ topic modeling, more specifically a custom version of BERTopic, to investigate the major changes in a scientific community around a specific journal across the years - in our case, the IxD&A multidisciplinary journal, whose focus of interest shifted given the technological development and the resulting cultural paradigm shifts. These past few years have marked significant turning points in science since they have seen a shift in research interests in tandem with our society's need to adapt to changing circumstances (i.e., COVID pandemic, rise of generative AI). As such, we present our analysis performed on 431 full-text articles published between 2013 and 2024.

The relevance of our goal is demonstrated by Griffiths and Steyvers [21] that applied LDA and argued for its utility in organizing and categorizing large collections of text. The authors used LDA to discover underlying topics in a corpus of scientific articles, which were then used to classify documents based on their topic distribution. Topic modeling was also employed to analyze trends and dynamics in scientific research across time. Moreover, by applying LDA to bibliometric data from scientific papers, Suominen and Toivanen [22] identified emerging research areas and track shifts in scientific focus. Their study showcased how topic modeling can be used to map the evolution of scientific disciplines, informed strategic decisions for research funding and collaboration, and also opened the path for future advancements in topic modeling applied to research papers [23–26].

Our main contributions to this research topic are the following:

- The introduction of an automated toolkit to analyze the evolving topics within IxD&A, aiming to understand trends, shifts in focus, and emerging areas of interest over time;
- An analysis of the topic evolution over the years that we realized by prompting a large language model and by selecting the most representative topics from the BERTopic's generated list based on the journal description;
- The public release of our source code on <https://github.com/upb-nlp/journaltopicfinder>, to provide a resource to others wishing to build upon,

refine, and expand the code. By making our code open-source, we also facilitate a rapid adaptation of our method to other journal websites.

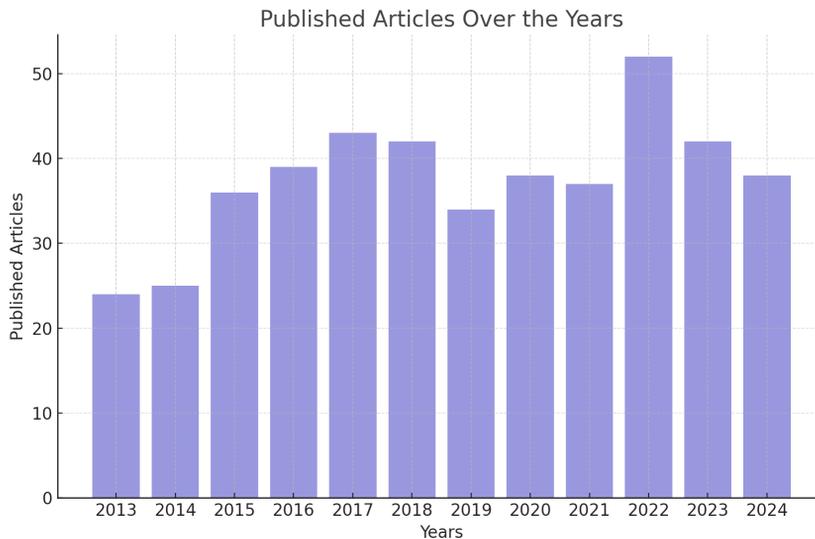
Next section provides a detailed description of the method and data used to extract the evolution of topics. In the third section, we present the results obtained using our method, and in the fourth and last sections, we discuss the limitations of our method and comment our findings.

## 2 Method

### 2.1 Data

We first extracted all the articles from IxD&A in PDF format and parsed them into text. The article's content was extracted by the parser, which employed heuristics based on the distribution of font and size usage to separate the text into sections and paragraphs. The reference section was then eliminated, and the section headers, image captions, and footnotes were also removed. The same parser was also utilized in other state-of-the-art research papers [27]. Ultimately, 15 out of 465 articles, ranging from 2013 to 2024, could not be processed due to various parsing errors and were disregarded. The distribution of the remaining articles per year is shown in Figure 1. A total of 27,314 paragraphs were extracted and analyzed from the entire corpus.

Among these, the longest extracted paragraph had 512 words, while the mean length of all paragraphs across the corpus was 464 words.



**Fig. 1.** IxD&A journal article distribution

## 2.2 Topic Modelling

BERTopic comprises five key steps (see Figure 2) contributing to the overall topic modeling process. Each step has a corresponding configurable sub-model that plays a specific role in handling different aspects of the process, from embedding generation to clustering and dimensionality reduction.



**Fig. 2.** Representation of the five key sub-models of the BERTopic that contribute to the overall process of topic modeling

**Text Encoding.** The first step in BERTopic consists of generating text embeddings, which capture the semantic meaning of words and documents. The embedding model in BERTopic can be customized; however, by default, it utilizes pre-trained Transformer-based models, such as BERT or its variants from HuggingFace’s sentence-transformers library. These models can produce high-dimensional embeddings representing the contextual relationships between words and phrases within the text. The quality of these embeddings is crucial, as they form the basis for clustering similar documents based on their content. For this reason, the research presented in this paper used all-mpnet-base-v2 to embed the text into arrays of size 512, a model with the best results in the current rankings of sentence Transformers [28].

**Embedding Dimensionality Reduction.** Since Transformer-based embeddings are high-dimensional, dimensionality reduction is applied to make the data more manageable for clustering algorithms. In our research, the UMAP algorithm is used to reduce the dimensionality of the embeddings from 512 to 100, but other alternatives such as PCA or Truncated SVD are also available. Besides being predominantly used by BERTopic, UMAP is also used in other state-of-the-art research papers [29]. Reducing dimensionality helps speed up computations and makes the clustering more efficient, while still retaining the critical semantic information that distinguishes topics.

**Clustering of Similar Documents.** Once the dimensionality is reduced, the next step involves clustering the embeddings to identify different topics. BERTopic uses HDBSCAN as its default clustering algorithm, but K-means is also available as an alternative. HDBSCAN is well-suited for this task because it can handle clusters of varying sizes and densities, and it does not require users to specify the number of clusters beforehand. The algorithm also effectively separates noise (outliers) from meaningful clusters, making the topics more robust and coherent.

As part of the postprocessing pipeline, CountVectorizer is used to convert raw text into a document-term matrix. This matrix counts the occurrences of each word in the documents, which helps construct a term-document matrix that serves as input for the

clustering algorithms used in BERTopic. Using CountVectorizer, BERTopic can also consider restrictions by removing stop words and imposing an ngram range (i.e., sequence length of adjacent words considered in subsequent steps). With stop words, all the English words such as “the”, “or” from the topics have been removed, whereas the ngram range was further restricted to a maximum topic length of 2.

To further refine the topic representations, BERTopic uses a custom version of the Term Frequency-Inverse Document Frequency (TF-IDF) approach, called class-based TF-IDF (c-TF-IDF). Unlike traditional TF-IDF, which measures word importance across documents, c-TF-IDF focuses on distinguishing words between different topic clusters. It calculates word importance specifically for each topic, enabling an improved interpretability of topics, as it highlights the words most uniquely associated with each one.

As for analyzing the evolution of topics over time in the academic papers, we used the Dynamic Topic Modeling feature from BERTopic. This feature enables us to observe how topics emerge, grow, or decline across different periods. By incorporating a temporal dimension into the analysis, BERTopic dynamically adjusts the topic distributions as new data becomes available, providing insights into shifts of thematic focuses. This is done by segmenting the data into different time intervals (all the articles from each year, ranging from 2013 to 2024) and applying our analysis to each segment, allowing for a comparison of topic evolution across these intervals.

**LLM Post-processing.** As a final refinement, LLMs can generate better topic labels or improve coherence by reinterpreting topic-word distributions in a broader contextual format.

Fine-tuning with LLMs in BERTopic results in more meaningful and human-readable topics, especially in datasets where the context or language usage is highly specialised or diverse, such as in scientific articles. For this reason, we used

Llama 3.1 8B Instruct [30] model to reinterpret the word distributions generated by HDBScan. The system prompt we used is:

```
< |start header id| > system < |end header id| >  
You are a helpful, respectful, and honest assistant for  
labeling  
topics.< |eot id| >  
Topics that are too general should be excluded.
```

To improve the accuracy of the model in topic generation, both the keywords produced by BERTopic and the corresponding source documents were provided as input. The main prompt, where BERTopic is replacing [DOCUMENTS] and [KEYWORDS] with the raw topics generated by HDBSCAN, is:

```
< |start header id| > user < |end header id| >  
I have a topic that contains the following documents:  
[DOCUMENTS] The topic is described by the following  
keywords: '[KEYWORDS]'.
```



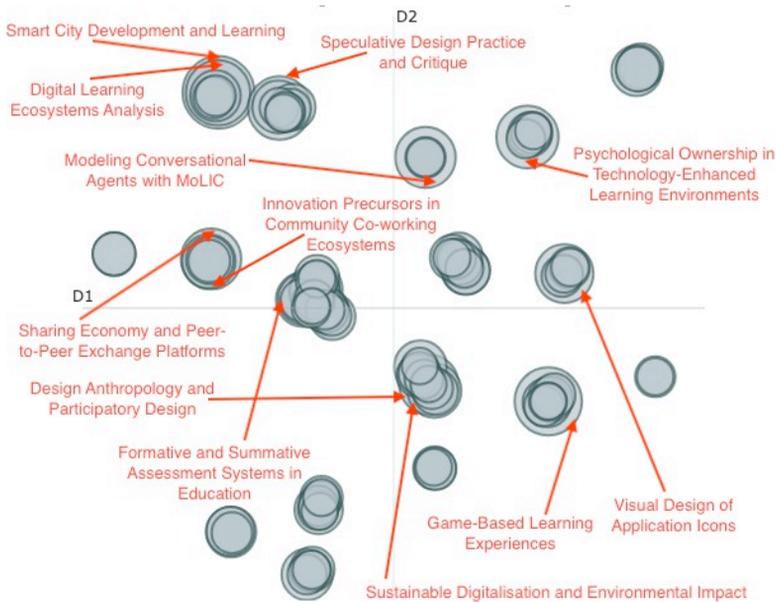
Based on the information above, please create a short label for this topic. Make sure you only return the label and nothing else. < |eot id| >

< |start header id| > assistant < |end header id| >

### 3 Results

Running the entire BERTopic pipeline on the entire corpus, which consists of 450 articles, generated 246 topics. Because of the large number of generated topics, we reduced the results to the top 100. Appendix A displays the final findings, while Figure 3 displays the top 12 topics visualized by the frequency of the top 10 words.

The trained model was then used to extract all the articles relevant to a given topic. This script, together with the findtopics function provided by BERTopic, which identifies similar topics to a given topic, can be easily used to generate correlations between documents. As the codebase is open-source, we actively support initiatives aimed at extending this functionality to other journals.



**Fig. 4.** IxD&A's topics map.

Figure 4 displays the topic map of IxD&A journal. For plotting the embeddings in a 2D graph, we used UMAP [31], which is effective at maintaining the overall arrangement of clusters, ensuring their relative positions make sense globally. Distance is relative, but similar topics are grouped together, while semantically different ones are further away [32].

We then used BERTopic's Dynamic Topic Modeling to track how topics evolve across time, a feature which also helps us adjust distributions as new data is added. By segmenting the data by year (2013–2024), we compared topic trends, revealing how themes emerged, grew, or declined across different periods.

Since displaying 100 topics on a line chart would be overwhelming, we reduced them to 10 by using the Llama 3.1 70b model, which we prompted with the following text:

```
Description:  
[journal description extracted from Scimago Journal &  
Country  
Rank website [33]]  
Given the previous journal description, specify the  
topics from the list below that best define the journal:  
List:  
[enumerating 100 topics]
```

The model then extracted 10 topics considered the most representative according to the journal description. The evolution across time of these topics is depicted in Figure 5.

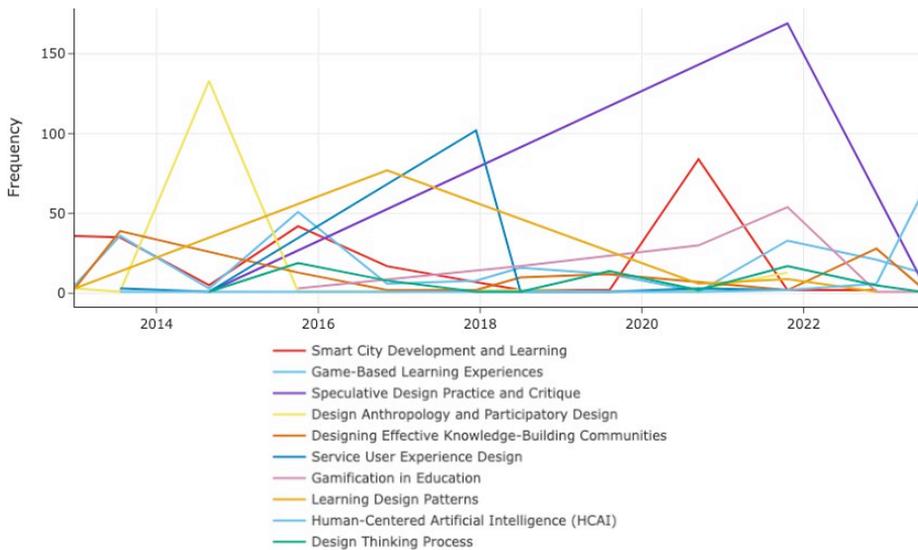


Fig. 5. The evolution of IxD&A's top 10 most representative topics

## 4 Discussion

The main advantages of the proposed topic modeling process are the integration of Transformer-based models capable of efficiently clustering topics (BERTopic) with

Large Language Models, recognized for their great performance in analyzing text data.

Overall, the topics generated by BERTopic are reasonable since they largely refer to design thinking, technology augmented experiences, technology enhanced learning, and the social implications of such topics that can be considered the major pillars of the multidisciplinary cultural framework of reference for the IxD&A journal (see Appendix A). The results were validated in collaboration with a member of the IxD&A editorial board, and the following limitations were identified.

#### 4.1 Limitations

A few of the identified topics are very specific and/or unusual. Examples of problems encountered include:

1. The need to group some topics under a more general, common label: for example, both "Service Design Workshop Process" and "Qualitative Research Methods - Thematic Analysis" could be categorized together as "Methodologies", to refer to the understanding of the research methods involved.
2. The need to shorten some topics: for instance, "Game-Based Learning Experiences" could be easily shortened to "Game-Based Learning", without losing the essence of the subject matter.
3. The need to understand the reasons why topics like "Penn State's Transition to Remote Education During the Pandemic" that appear somewhat peculiar and very specific are identified; indeed no more than a couple of papers refer to such topic, which raises questions about its significance within the broader academic landscape of IxD&A journal.

Additional experiments were conducted using few-shot learning to reduce the occurrence of such examples. However, even though the generated topics were shorter, more instances of unrelated topics, such as "introduction" and "references", appeared. From such experiments, we can also conclude that an important limiting factor of our method is Llama's capability to generalize the raw outputs extracted by BERTopic. As BERTopic's raw results are not in an easily human-readable natural language format, as LDA results are, the model used for naming the final topic represents an important step in BERTopic's pipeline. The following examples serve as evidence to support this claim:

1. The raw topic "pwads, room, pwad, rooms, daily, daily activities, ot, disorientation, healthcare, self orientation" is described as "Assistive Technology and Architectural Barriers in Long-Term Healthcare Centers for People with Alzheimer's Disease, which could be a hallucination of Llama.
2. The raw topic "display, swimming, displays, appropriation, swimming center, ubi, center, unanticipated, public display, display network" is described as "Public Display Appropriation in Oulu", where Oulu does not appear anywhere in the raw topic.

The way in which the ten topics - whose time trends were studied and shown in Figure 5 - is an example of how it is necessary to provide guidance elements to Llama to improve its outputs. Almost all of the topics that were removed from the list, in

fact, are topics with too specific descriptions, like "Psychological Ownership in Technology-Enhanced Learning Environments", "Modeling Conversational Agents with MoLIC" or "Sharing Economy and Peer-to-Peer Exchange Platforms".

On the other hand, through human inspection, such a filtering operation can be used to highlight potential shifts in the journal's positioning within its cultural context of reference. The fact that rather well-defined topics, such as "Digital Learning Ecosystems Analysis" or "Augmented Reality in Education" were not identified as relevant during the comparison operated against the journal description might indicate that over time more attention has been shaping up and more space has been given to articles devoted to technologically augmented environments.

Another aspect to consider when interpreting Figure 5 is that extremely pronounced peaks in the time trend may be induced by the publication of special issues on specific topics, rather than by trends traceable in international research developments. However, this is an effect due to the journal's editorial choices, rather than problems related to the topic modeling process, and, if anything, it confirms its reliability.

## 5 Conclusions and Future Work

We have introduced a novel approach to topic modeling based on the integration of BERTopic and Llama-3.1 in handling large unstructured text datasets and argued for the effectiveness of the overall model by applying it to the corpus of the articles published by a multidisciplinary scientific journal - IxD&A - in the last decade, which is a respectable test case due to the variety of topics covered by this multidisciplinary journal.

In general, the topics and themes identified by BERTopic align well with the expected outcomes of a literature review based on full-text research papers from the IxD&A journal. Notably, it is particularly encouraging to observe that BERTopic was able to generate topics that focus on integrating design thinking and the most advanced information technologies applied to educational contexts and environments, and the implications this may have for social innovation, reflecting the main scopes of the journal. However, the Llama's capability to generalize the raw outputs extracted by BERTopic presents several criticalities and represents the main limiting factor of the overall process.

The toolkit's effectiveness was assessed based on the full-text articles, which may provide too many insights into the information. To ensure a more focused evaluation, we will focus solely on abstracts as they already provide a human summary of the article. Future work will also consider improving the toolkit's integration within other journal websites and optimizing its interoperability with existing platforms. This would contribute to a more efficient and user-friendly experience for researchers and publishers, facilitating wider adoption and practical application.

Future improvements will also account for the variability introduced by differing journal issues. Currently, the method does not explicitly incorporate contextual distinctions, which may lead to biased results. Furthermore, since the performance of the LLaMA 3.1 model plays a significant role in the current findings, a comparative analysis involving additional large language models should be conducted. This will

provide a clearer understanding of the strengths and weaknesses of our approach across different models.

In conclusion, our work offers a new perspective and a new tool for topic modeling of large unstructured text datasets which, since our toolkit has been released as open-source on GitHub (<https://github.com/upb-nlp/journaltopicfinder>), can be applied to analyse archives of other scientific journals as well as other collection of data - e.g. repositories of educational materials, medical data. facilitating easier recognition by the end user of the links and similarities between different documents.

**Acknowledgement.** This research was supported by the project “Romanian Hub for Artificial Intelligence - HRIA”, Smart Growth, Digitization and Financial Instruments Program, 2021-2027, MySMIS no. 334906.

### **CRedit author statement**

**Razvan Paroiu:** Methodology, Software, Validation, Formal analysis, Investigation, Data curation, Resources, Writing – original draft, Visualization; **Mihai Dascalu:** Conceptualization, Writing – review and editing, Supervision, Funding acquisition; **Carlo Giovanella:** Conceptualization, Writing – review and editing.

## **References**

1. Blei, D.M., Ng, A.Y., Jordan, M.I.: Latent dirichlet allocation. In: Proceedings of the 2003 Conference on Neural Information Processing Systems (NIPS), pp. 601–608 (2003). <http://www.jmlr.org/papers/volume3/blei03a/blei03a.pdf>
2. Blei, D.M., Griffiths, J.B., Blei, D.M., Jordan, M.I.: Probabilistic topic mod-els. Communications of the ACM 55(4), 77–84 (2007) <https://doi.org/10.1145/2133806.2133826>
3. Lee, D.D., Seung, H.S.: Learning the parts of objects by non-negative matrix factorization. Nature 401(6755), 788–791 (1999) <https://doi.org/10.1038/44565>
4. Cichocki, A., Zdunek, R., Amari, S.-i.: Nonnegative matrix and tensor factor-izations: Applications to exploratory multi-way data analysis and blind source separation. In: Wiley Encyclopedia of Electrical and Electronics Engineering (2009). <https://doi.org/10.1002/9780470747278>
5. Jansson, P., Liu, S.: Distributed representation, lda topic modelling and deep learning for emerging named entity recognition from social media. In: NUT@EMNLP (2017). <https://doi.org/10.18653/v1/W17-4420>
6. Khan, S.K., Ahmed, F., Mubeen, M.: A text-mining research based on lda topic modelling: A corpus-based analysis of pakistan’s un assembly speeches (1970-2018). Int. J. Humanit. Arts Comput. 16, 214–229 (2022) <https://doi.org/10.3366/ijhac.2022.0291>
7. Uthirapathy, S.E., Sandanam, D.: Topic modelling and opinion analysis on cli-mate change twitter data using lda and bert model. Procedia Computer Science (2023) <https://doi.org/10.1016/j.procs.2023.01.071>
8. Rani, R., Lobiyal, D.K.: An extractive text summarization approach using tagged-lda based topic modeling. Multimedia Tools and Applications 80, 3275–3305 (2020) <https://doi.org/10.1007/s11042-020-09549-3>
9. Onah, D.F.O.: A data-driven latent semantic analysis for automatic text summa-rization using lda topic modelling. 2022 IEEE International Conference on Big Data (Big Data), 2771–2780 (2022) <https://doi.org/10.1109/BigData55660.2022.10020259>

10. Sawahata, H., Nishino, T.: Automatic extractive summarization for japanese academic papers by lda. *Information Engineering Express* (2023) <https://doi.org/10.52731/iee.v9.i2.759>
11. Wang, H.-i., Sun, K., Wang, Y.: Exploring the chinese public's perception of omicron variants on social media: Lda-based topic modeling and sentiment analysis. *International Journal of Environmental Research and Public Health* 19 (2022) <https://doi.org/10.3390/ijerph19148377>
12. Mei, Y., A.Hernandez, A.: Sentiment analysis of lijiang ancient town attraction reviews based on lda. *Proceedings of the 2023 3rd International Conference on Big Data, Artificial Intelligence and Risk Management* (2023) <https://doi.org/10.1145/3656766.3656774>
13. Erniyati, E., Harsani, P., Mulyati, M., Fahriza, L.D.: Topic modeling lda and svm in sentiment analysis of hotel reviews. *Komputasi: Jurnal Ilmiah Ilmu Komputer dan Matematika* (2023) <https://doi.org/10.33751/komputasi.v20i2.7604>
14. Grootendorst, M.: Bertopic: Neural topic modeling with bert. arXiv preprint arXiv:2009.04822 (2020)
15. Devlin, J., Chang, M.-W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805 (2018)
16. Sanh, H., Debut, L., Chaumond, J., Wolf, T.: Distilbert, a distilled version of bert: Smaller, faster, cheaper and lighter. arXiv preprint arXiv:1910.01108 (2019)
17. Liu, Y., Ott, M., Goyal, N., Du, J., Cohn, T., Chen, D., Manning, C.D.: Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692 (2019)
18. McInnes, L., Healy, J., Melville, J.: UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction (2020). <https://arxiv.org/abs/1802.03426>
19. McInnes, L., Healy, J., Melville, J.: Hdbscan: Hierarchical density-based spatial clustering of applications with noise. In: *Proceedings of the 2017 International Conference on Data Mining (ICDM)*, pp. 544–553 (2017). <https://doi.org/10.1109/ICDM.2017.79>
20. Liu, Y.: Comparison of lda and bertopic in news topic modeling: A case study of the new york times' reports on china. *Pacific International Journal* 7, 47–51 (2024) <https://doi.org/10.55014/pij.v7i3.616>
21. Griffiths, T.L., Steyvers, M.: Finding scientific topics. In: *Proceedings of the National Academy of Sciences (PNAS)*, vol. 101, pp. 5228–5235 (2004). <https://doi.org/10.1073/pnas.0307752101>
22. Suominen, A., Toivanen, H.: Map of science with topic modeling: Comparison of unsupervised learning and human-assigned subject classification. *Journal of the Association for Information Science and Technology* 67 (2016) <https://doi.org/10.1002/asi.23596>
23. Trevisani, M., Tuzzi, A.: Learning the evolution of disciplines from scientific literature: A functional clustering approach to normalized keyword count trajectories. *Knowl. Based Syst.* 146, 129–141 (2018) <https://doi.org/10.1016/j.knosys.2018.01.035>
24. Uban, A.S., Caragea, C., Dinu, L.P.: Studying the evolution of scientific topics and their relationships. In: *Findings* (2021). <https://doi.org/10.18653/v1/2021.findings-acl.167>
25. Gerasimenko, N., Chernyavskiy, A., Nikiforova, M., Ianina, A., Vorontsov, K.: Incremental topic modeling for scientific trend topics extraction. *COMPUTATIONAL LINGUISTICS AND INTELLECTUAL TECHNOLOGIES* (2023)
26. Ionita, R.F., Corlatescu, D.G., Ruseti, S., Dascalu M., T.-M.S., N., T., Banica, C.K.: Comprehensive sociograms of the scientific bulletin community. *Scientific Bulletin* 81, 3–12 (2019)
27. Paroiu, R., Ruseti, S., Dascalu, M., Trausan-Matu, S., McNamara, D.S.: Asking questions about scientific articles—identifying large n studies with llms. *Electronics* 12(19) (2023) <https://doi.org/10.3390/electronics12193996>
28. developers: Pretrained Models 2014; Sentence Transformers documentation — sbert.net. <https://www.sbert.net/docs/sentence-transformer/pretrained-models.html>. [Accessed 19-09-2024]

29. Izzidien, A.: Using the interest theory of rights and hohfeldian taxonomy to address a gap in machine learning methods for legal document analysis. *Humanit Soc Sci Commun* 10(251) (2023) <https://doi.org/10.1057/s41599-023-01693-z>
30. Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., et al.: The Llama 3 Herd of Models (2024). <https://arxiv.org/abs/2407.21783>
31. McInnes, L., Healy, J., Melville, J.: UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction (2018). <https://arxiv.org/abs/1802.03426>
32. Andy, C., Adam, P.: Understanding UMAP (n.d.). <https://pair-code.github.io/understanding-umap/>
33. developers: Interaction Design and Architecture(s) — scimagojr.com. <https://www.scimagojr.com/journalsearch.php?q=21100440523&tip=sid&exact=no>. [Accessed 24-09-2024]

## Appendix A

**Table 1.** The top 100 topics identified by Llama 3.1 on the base of the top 10 raw keywords extracted by BERTopic.

Nr.	Raw topics extracted by Bertopic	Topics processed by Llama 3.1
1	smart, smart city, city, cities, smart cities, urban, city learning, ethical smart, city projects, citizens	Smart City Development and Learning
2	games, game, learning, players, educational, game based, learning goals, students, playing, educational games	Game-Based Learning Experiences
3	ecosystem, ecosystems, learning ecosystem, learning ecosystems, digital learning, smart learning, learning, smart, digital, smartness	Digital Learning Ecosystems Analysis
4	speculative, speculative design, design, critical, critical design, fiction, design fiction, speculative critical, discursive, critical speculative	Speculative Design Practice and Critique
5	ownership, psychological ownership, psychological, self, control, selfregulated, regulated, regulated learning, srl, agency	Psychological Ownership in Technology-Enhanced Learning Environments
6	molic, conversational, conversational agents, agents, bixby, modeling, ana, user, conversation, scene	Modeling Conversational Agents with MoLIC
7	sharing, exchange, peer, services, service, peer peer, sharing economy, te, economy, sharing services	Sharing Economy and Peer-to-Peer Exchange Platforms
8	anthropology, design anthropology, da, participatory design, participatory, disruption, pd, design, future orientation, future	Design Anthropology and Participatory Design
9	icon, visual, icons, application, first, visual design, first exposure, exposure, application icon, brand	Visual Design of Application Icons
10	sustainable, digitalisation, sustainability, sustainable development, sdgs, fairphone, planetary, planetary boundaries, waste, phones	Sustainable Digitalisation and Environmental Impact
11	assessment, quiz, formative, stumbling, assessment systems, formative assessment, stumbling blocks, juxtalearn, learning, students	Formative and Summative Assessment Systems in Education
12	cces, innovation, cce, social, social innovation, interviewee, social entrepreneurs, entrepreneurs, tribe, skunkworks	Innovation Precursors in Com-munity Co-working Ecosystems (CCEs)
13	validity, reliability, loadings, variance, discriminant, ave, table, items,	Measurement Model Validation

	discriminant validity, index	
14	data, data exploration, exploration, design inquiry, method, inquiry, exploration design, dataset, data techniques, design method	Design Method for Data Exploration
15	workshop, workshops, participants, esc, proto, proto card, prototyping, toolkit, prototypes, materials	Service Design Workshop Process
16	reflection, reflective, children, collaborative, reflective practices, reflection action, reviewed, et al, et, al	Designing Technology-Mediated Reflection in Children's Collaborative Interactions
17	ar, ar application, augmented reality, reality, augmented, location based, motivation, location, application, use ar	Augmented Reality in Education
18	feedback, eat, writing, essays, comma, students, automated, mette, automated feedback, students writing	Automated Feedback in Writing Instruction
19	hackathon, hackathons, civic, civic hackathons, open data, prototypes, civic hackathon, data, open, hackmytown	Civic Hackathons
20	knowledge, discourse, knowledge building, community, formal, communities, kbcs, informal, building, ideas	Designing Effective Knowledge-Building Communities
21	book, rare, books, rare historic, historic, rare book, rare books, reading, printed books, printed	Human Interaction with Rare Historic Books
22	service, ux, service design, persona, touchpoints, journey, user journey, touchpoint, user, attributes	Service User Experience Design
23	gamification, gamification teaching, university students, teaching, teachers, gamification activity, university, gamification activities, school students, students	Gamification in Education
24	youth, gis, places, garden, map, power places, community, neighborhood, maps, favorite	Engaging Youth in Urban Planning through GIS-Based Community Mapping
25	heritage, cultural, cultural heritage, sites, intangible, visitors, cultural wiki, assets, semantic, information	Digital Cultural Heritage Preservation and Education
26	children, design fiction, fiction, diversity, inclusion, cci, diversity inclusion, researchers, design, role process	Design Fiction for Children's Empowerment and Diversity
27	maintenance, technicians, technician, guidance, maintenance work, industrial, ar, maintenance technicians, reporting, demonstration	Industrial Maintenance Technology and Technician Support
28	robotics, robot, robots, children, educational robotics, ev3, distance, robot tutee, math, project	Designing Educational Robotics Projects for Children
29	libraries, incubator, library, social innovation, librarians, program, innovation, social innovations, design	Social Innovation in Public Libraries Incubators

	thinking, innovations	
30	remote, faculty, remote education, remote teaching, teaching, remote learning, transition, emergency remote, transition remote, online	Penn State's Transition to Remote Education During the Pandemic
31	pattern, patterns, design patterns, pattern language, design pattern, language, ldshe, design, learning design, lds pattern	Learning Design Patterns
32	lds, teachers, til, tpd, program, creation, rm, practices, tpd program, school university	Teacher Professional Development Program
33	simulation, land, market, energy, gdp, ecosystem, energy choices, areas, ses, defined	Agent-Based Modeling of Ecosystem Services (SES) and Land Use
34	math, preschool, children, early math, play learn, learn game, preschool teachers, early, play, magical garden	Effectiveness of Digital Play & Learn Games in Preschool Math Education
35	thermal, temperature, auditory, vibration, sd, baseline, skin, modality, stimuli, temperature change	Assistance Systems for Hazard Detection Tasks
36	programming, teaching programming, teachers, pupils, teaching, teach, teacher, teach programming, challenges, teaching plan	Teachers' Experiences with Teaching Programming
37	uss, specification, ime, agile, vagueness, ms, mvms, ucd, epics, epic	User Story Specification Practice
38	hrc, human factors, factors, human, cobot, cobots, factors hrc, enrich, factors enrich, safety	Human Factors in Human-Robot Collaboration (HRC)
39	mobile social, social media, media, elective, twitter, tweets, mobile, course, autmsm2014, social	Designing a Mobile Social Media Elective Course
40	usability, ux, experience, evaluation, user experience, user, ux evaluation, triangulation, product, methods	Usability and User Experience Evaluation Methods
41	visual impairment, impairment, people visual, shopping, visual, sighted, grocery, visually impaired, impaired, people	Accessibility and Inclusive Design for People with Visual Impairment
42	minecraft, carbon, materials, soil, earthen, embodied car-bon, building, earthen materials, embodied, low carbon	Designing with Earthen Materials in Virtual Environments
43	feeler, monitoring, self monitoring, techno monitoring, self, techno, using feeler, learners, prototype, participants	Self-Monitoring Technologies for Learning and Education
44	parents, parent, messages, parents teachers, child, family ties, parent engagement, family, app, teachers	Parental Involvement in Early Childhood Education
45	parents, parent, messages, parents teachers, child, family ties, parent engagement, family, app, teachers	Indigenous Education and Technology Integration
46	bedouins, bedouin, ich, egypt, community, heritage, bedouin ich, poems, cultural, egyptian	Digitally Mediated Documentation of Bedouin Intangible Cultural Heritage
47	pedestrian, walking, walkway, speed, campus, pedestrians, pedestrian speed,	Factors Affecting Pedestrian Speed in

	type, master plan, university	University Campuses
48	focus group, qualitative, themes, analysis, thematic, focus, thematic analysis, coding, interviews, data	Qualitative Research Methods - Thematic Analysis
49	tas, tas cas, cas, orchestration, tool, orchestration tool, groups, ta, expert teachers, classroom	Classroom Orchestration Tool Development
50	drama, information literacy, literacy, condition, vaccine, experimental condition, play, control condition, theater, experimental	Information Literacy through Drama-Based Education
51	food, recaa, blog, ageing, elders, food blog, mediatized, food culture, healthy, mandate	Food Blogging as Activism for Successful Ageing
52	ninja, moral, road safety, conversational, road, disengagement, rule, moral disengagement, safety, decision	Moral Disengagement in Road Safety Decision Making
53	duration, estimated, estimated duration, time, subjective duration, perception, subjective, mean estimated, physical spatial, mean	Time Perception in Physical-Spatial Environments
54	entrepreneurs, transistórias, sustainability, design students, futures, tapada, mercês, tapada das, das mercês, das	Designing Sustainability through Co-Creation with Local Entrepreneurs
55	business, business models, models, business model, com-munity based, community, based business, firm, building blocks, interviewees	Business Model Innovation in Community-Based Ventures
56	engagement, mathematics, behavioral, behavioral engagement, initial, items, level, cognitive, subscale, emotional	Mathematics Student Engagement Analysis
57	targets, physical activity, active participants, active, physically active, mcr, prompts, activity, physically, performance	Physical Activity Promotion through Mobile Notifications
58	creation, design, users, creation design, remote, process, design process, participatory, designers, user	Co-Creation and Co-Design Processes
59	competencies, competence, skills, competency, knowledge skills, life skills, knowledge, competences, skill, assessing	Competency Frameworks and Assessment
60	trackaware, photoframe, peripheral, displays, photoframe users, peripheral displays, focused, researchers, peripheral perception, focused attention	Evaluation of Trackaware for Situated Displays
61	memory, collective memory, collective, city, memories, clio, collective city, city memory, virtual city, records	Urban Memory and Design Learning
62	appropriation, technology, tactics, meaning potential, user tactics, users, design strategies, technology	Technology Appropriation Research

	appropriation, meaning, appropriations	
63	activity, activity theory, activity systems, theory, motive, elicitation, ev, leontev, leont, object	Activity Theory and Knowledge Elicitation
64	site, local data, data, local, augmented reality, architects, augmented, architect, architectural, architectural design	Utilizing Location-Based Augmented Reality in Architectural Design
65	pwdads, room, pwad, rooms, daily, daily activities, ot, disorientation, healthcare, self orientation	Assistive Technology and Architectural Barriers in Long-Term Healthcare Centers for People with Alzheimer's Disease
66	ai, humans, fairness, ai tools, human, intelligence, artificial, ai humans, ia, ethical	Human-Centered Artificial Intelligence (HCAI)
67	learning design, ou, ld, designs, learning, rienties, learning designs, toetenel, activities, design activities	Learning Design in Online Education
68	citizen, citizen science, science, observatory, citizen observatory, cobweb, observatories, sensor, citizen observatories, data	Citizen Science Education
69	app, feedback reports, web app, feedback, reports, participation, pedestrians, web, mobile application, mobile	Field Trial of Mobile Participation Application
70	digital competence, competence, digital, tdc, gender, teachers, org 10, doi org, https doi, doi	Teacher Digital Competence and Gender Differences
71	summer school, summer, urban sensing, air, sensing, urban, urban data, air quality, pollution, school	Urban Sensing Education
72	games, video games, mar, gameplay, space, video, ar, game, gameworld, mar video	Hybrid Video Games and Spatial Interaction
73	tests, groupware, test, task, user tests, usability, tycho, testers, protocol, volunteers	User Test Facilitation and Protocol Design
74	design, design education, teaching practices, assignment, studio, students, teaching, course, courses, practices	Design Education and Experimental Teaching Practices
75	augmented reality, augmented, reality, urban, ar, citizens, councils, urban transformation, space, cities	Augmented Reality in Urban Transformation
76	glasses, picking, ar, smart glasses, order picking, shift, shelf, scenario, hmd, scenario ar	Smart Glasses in Order Picking Processes
77	oer, ownership, adoption, adoption oer, institutional, educational resources, emotional ownership, authorship, oer adoption, agency	Open Educational Resources Adoption in Higher Education Institutions
78	thinking, design thinking, problem, design, problems, designers, process, innovation, complex, design phase	Design Thinking Process
79	perceived, perceived usefulness, usefulness, intention, intention use,	Mobile Learning and Self-Assessment Acceptance

	mobile, continuance, mobile self, acceptance, self assessment	
80	insider, fieldwork, ethnographic, performative, ethnographer, performative knowledge, field, researcher, analyst, actor	Reflexive Insider Research Methods in Ethnography
81	od, skills, od skills, learning approaches, data, skills learning, learning outcome, data skills, using od, context skills	Open Data Education and Literacy
82	wireless, signal, wi fi, wi, fi, buildings, signal strength, propagation, connectivity, space	Wireless Network Signal Propagation in Buildings
83	creativity, muse, coursework, instructors, creative, csds, creativity assessment, assessment, cat, instructor	Assessing Creativity in Coursework
84	avs, av, vehicle, vehicles, driving, persona, mobility, future, cars, transportation	Social Relationships with Autonomous Vehicles (AVs)
85	elderly, children elderly, children, user groups, contextmapping, vulnerable, vulnerable generations, generations, needs, user	Designing for Vulnerable User Groups: Intergenerational Design
86	display, swimming, displays, appropriation, swimming center, ubi, center, unanticipated, public display, display net-work	Public Display Appropriation in Oulu
87	illich, conviviality, media, convivial, platforms, platform, power, tools, commons, media scholars	Rethinking Platform Capital-ism through Conviviality and Participatory Design
88	palermo, game, games, mobility, urban, urban games, mov, citizens, city, common goods	Urban Games in Palermo
89	values, values interaction, interaction design, teaching values, values based, values design, teaching, assessment, based design, interaction	Teaching Values-Based Interaction Design in Higher Education
90	odours, reading, olfactory, odour, testing sessions, multisensory, subjects, reading experience, olfactory display, testing	Multisensory Learning and Reading Experience with Olfactory Stimuli
91	creation, tel, students, levels creation, creation design, design, creation methods, levels, modes, better	TEL Co-Creation Methods and Student Participation
92	corridor, adolescents, corridors, employees, activities, residents, interactions, activities corridor, observed adolescents, sleco	Adolescent behavior in institutional corridors
93	geocaching, cache, geocaching com, caches, geocaches, com, geocache, com website, relatedness, geocachers	Geocaching and its impact on senior citizens' well-being
94	science, science learning, science center, science education, outside classroom, center, learning outside, non formal, outside, design workshop	Science Learning Outside the Classroom
95	strength, ontostrength, selfit,	OntoStrength: A Personalized Strength

	movement, strength development, trainee, ontology, strength skills, signature, psychomotor	Development Ontology
96	light, coloured, coloured light, light scenarios, 000, 500, orange, self perceived, scenarios, green	Self-Perceived Satisfaction and Energy Levels under Different Colored Light Scenarios
97	jade, board, game, ergonomics, board game, version jade, version, criteria, software, ergonomic	Ergonomic Education through Board Games
98	maker, diy, makerspaces, making, expertise, maker culture, culture, open, production, maker cultures	Maker Culture and Expertise
99	cbt, satisfaction, cbt centre, quality, service, centre, overall satisfaction, service quality, service satisfaction, overall	Mediating role of CBT centre service satisfaction in predicting overall CBT satisfaction
100	observations, mmla, lesson, observation, la, data, observata, multimodal, lesson observation, datasets	Lesson Observation Data Integration