# Eliciting Topic Hierarchies from Large Language Models

Grace Li
gl2676@columbia.edu
Columbia University
New York, New York, USA

Tao Long
long@cs.columbia.edu
Columbia University
New York, New York, USA

Lydia B. Chilton
chilton@cs.columbia.edu
Columbia University
New York, New York, USA

## Abstract

Current research has explored how Generative AI can support the brainstorming process for content creators, but a gap remains in exploring support-tools for the pre-writing process. Specifically, our research is focused on supporting users in finding topics at the right level of specificity for their audience. This process is called topic scoping. Topic scoping is a cognitively demanding task, requiring users to actively recall subtopics in a given domain. This manual approach also reduces the diversity of subtopics that a user is able to explore. We propose using Large Language Models (LLMs) to support the process of topic scoping by iteratively generating subtopics at increasing levels of specificity: dynamically creating topic hierarchies. We tested three different prompting strategies and found that increasing the amount of context included in the prompt improves subtopic generation by 20 percentage points. Finally, we discuss applications of this research in education, content creation, and product management.

## CCS Concepts

• **Human-centered computing** → **Collaborative content creation**.

## Keywords

Content Creation, Generative AI, Topic Scoping, Writing-Support Tools, Human-AI Collaboration

## 1 Introduction

Generative AI has changed the way that content creators brainstorm ideas and structure content [22], [13]. But using models like ChatGPT to directly generate content often results in bland and inauthentic results [4]. As a result, many have instead focused on leveraging the ability of Generative AI in the formative stages of content creation [18]. Recent work has explored the effectiveness of AI in the brainstorming process, but not much research has been done to explore the prewriting process [13]. Before a writer begins the process of outlining, they must determine a topic of the right scope to fit the venue and medium that they will publishing their work on. They must consider the constraints such as character-count when publishing on Twitter and the length of their video when publishing on TikTok or Instagram reels. This process of finding a topic at the appropriate level of specificity is called topic scoping. Topic scoping is used by educators when creating lesson plans, journalists when finding angles on current events, and researchers when determining specific projects to pursue as part of a grant.

The main challenge of topic scoping is iteratively breaking a broad domain, like Computer Science, into smaller and smaller subtopics to help the user pick a more specific subject area. For example, when picking a topic to teach within User Interface Design, a teacher might decompose UI into Usability Heuristics, and then break it down further into Visual Information Design. Through the process of iterative topic scoping, we can generate topic hierarchies–information trees where each node is a subtopic under the root and each level of the tree represents an increasing level of subtopic specificity. The Dewey Decimal System for book classication in libraries and the Taxonomy of Life that classify organisms by their Kingdom, Phyllum, Class etc. are examples of topic hierarchies that have been established and exist in the world. But managing the creation of topic hierarchies remains a non-trivial task. Currently, the number established topic hierarchies are limited and static–the content often becoming outdated when they are not maintained. The maintenance of topic hierarchies is time and labor intensive task.

We explore how Large Language Models (LLMs) are able to dynamically elicit topic hierarchies to support the process of topic scoping. On their own, LLMs struggle with generating fine grain subtopics in more specific domains. Like users, as topics get more niche, machines struggle to generate subtopics that are unique, related, and specific to the given topic. We test three different prompting strategies on five different levels of topic specificity– with an emphasis on generating subtopics at the most specific level. The three conditions that we explored were
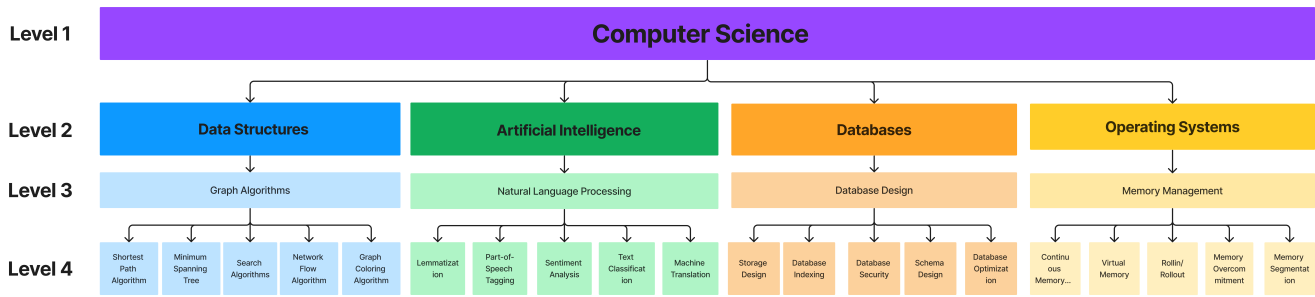
(1) Current Topic: "List 5 subtopics in *Natural Language Processing*."
(2) Root + Current Topic: "In *Computer Science*, list 5 subtopics in *Natural Language Processing*."
(3) Full Path + Current Topic: "In *Computer Science* and *Artificial Intelligence*, list 5 subtopics in *Natural Language Processing* "

We used two annotators to evaluate the subtopics generated from the three different prompting techniques and measure the appropriateness of the generated subtopics on the relatedness, repetitiveness, and specificity.

The Full Path + Current Topic prompting technique improved subtopic generation by 21 percentage points. This finding is in-line with previous findings that show that in-context learning improves language model generations [6]. This paper demonstrates that LLMs can assist in the topic scoping process, helping create a structured exploration of possible subtopics within a given domain and proposes future work around designing an interactive system to support users in the process of topic scoping.

## 2 Related Work

LLM-support in topic scoping and the creation of topic hierarchies extends research in the following domains: Approaches in Extracting and Structuring Information and LLMs on Various Knowledge Tasks.

**Figure 1: Topic hierarchy from Wikipedia's "Category: Computer science" with certain missing subtopics being filled in from college level curriculum. This Topic hierarchy is used as a test suite to evaluate the three different prompting strategies to generate subtopics in Level 2, Level 3, Level 4, and Level 5.**

## 2.1 Approaches in Extracting and Structuring Information

Topic hierarchies help facilitate the process of topic scoping by structuring knowledge into a tree-like structure. Topic hierarchies fit into a broader domain of knowledge structuring. Like knowledge graphs that seek to unify information across disparate sources [8], to streamline the process of knowledge retrieval [25], topic scoping seeks to add structure and organization to existing information. Furthermore, knowledge graphs have been used to support knowledge retrieval and brainstorming for different user centered systems [9]. Sun et. al. uses information propagation to add multimodal components to knowledge graphs for recommendation systems to support user-item interaction [20]. Ait-Mlouk et. al. used a support vector machine (SVM) for intent classification to query a knowledge base to support domain-specific question and answering [1].

Broadly, the idea of topic scoping is closely related to text mining–transforming unstructured knowledge into a structured format to find patterns and to make new discoveries [21]. Text mining has been explored within the context of education to support online learning [7]. Mansur et. al. used natural language processing techniques like SVM and K-means for text classification to group educational resources together [14]. Similarly, Crossley et. al. incorporated student demographic information as features in their automatic essay evaluation system to train a regression model [5]. Jin et. al. used the Blei's latent Dirichlet Allocation (LDA) to classify financial news articles and to obtain each article's topic distribution to train a linear regression model to make movement forecasts on the foreign currency market [10] [3].

Finally, the process of topic scoping is related to research in topic discovery and topic retrieval. Work in topic discovery includes characterizing the different topics within Twitter threads using a partially supervised learning model [17], revealing the implicit knowledge present in news streams using a multilayer clustering system to support similar topic exploration [16], and using fuzzy latent semantic analysis (FLSA) to eliminate topic redundancies in a medical corpora [11].

Previous research in constructing knowledge graphs, text mining, and topic retrieval have all used a combination natural language processing techniques and classical machine learning algorithms.

Our research focuses on using language models to explore their capabilities in generating topic hierarchies and supporting the process of topic scoping.

## 2.2 LLMs on Various Knowledge-Based Tasks

While many previous approachs leverage classical machine learning and natural language processing techniques to extract information from existing corpora, current research has also focused on leveraging LLMs to support various knowledge-based tasks. Wang et. al. has explored how LLMs can support finding conceptual relations between topics and connecting tangible scenes and experiences with abstract words [23]. Additionally, current work has also been done to support the process of using LLMs to support the clarification of an abstract concept into a semantically-related object [12], in sensemaking for complex topics by leveraging LLMs to support multilevel abstractions [19], and in retrieval-based knowledge tasks [24]. These research areas demonstrate the integration of LLMs to support various knowledge-based tasks, leveraging the ability of LLMs to generate diverse and creative connections between abstract and concrete topics. The topic of topic scoping is related to these works as it uses an LLM to support the systematic retrieval of information from the model itself.

## 3 Methodology

We explore the capabilities of LLMs to incrementally generate topic hierarchies through three different prompting strategies. We generate subtopics for up to 5 different levels of specificity, using a subset of Wikipedia's "Category of Computer Science" page as a test suite of topics. We used human annotators to evaluate the appropriateness of the generated subtopics for each of the prompting strategies.

## 3.1 5-Level Topic Classification System

To standardize the categorization of generated subtopics, we created a 5-level topic hierarchy to classify the level of specificity for a generated topic. Table 1 illustrates the 5-Level Topic Hierarchy with corresponding descriptions and examples. The table shows that Level 1 is the broadest level and contains topics related to broad domains of study, like Computer Science. The next

level, Level 2, contains more specific subtopics that explore general concepts within Computer Science, like Data Structures. Each level's specificity incrementally increases with Level 5 topics being the most specific and focused on specific implementations–like Dijkstra's algorithm as a specific implementation of Shortest Path Algorithms (Level 4 topic). We choose to set the depth of the table to be 5 because preliminary findings have shown that users struggled the most with manually brainstorming topics at Level 5.

| Level | Definition | Example in Computer Science |
|-------|------------|------------------------------|
| Level 1 | Topics related to domains areas of study | Computer Science |
| Level 2 | Subtopics that explore general topics | Data Structures |
| Level 3 | Subtopics that are general concepts | Algorithms |
| Level 4 | Subtopics exploring different use cases of general concepts | Shortest Path Algorithms |
| Level 5 | Subtopics that focus on specific implementations | Dijkstra's algorithm |

**Table 1: The corresponding level descriptions in the 5-Level topic classification system.**

## 3.2 Wikipedia's Topic Hierarchy

After defining the 5-Level classification system, we used Wikipedia's "Category of Computer Science" page as a reference to create test suite of topics to standardize the evaluation process for the different prompting strategies. Wikipedia's "Category of Computer Science" page is structured as a nested series of expandable subtopic lists where users are able to incrementally traverse through the different levels. To address any holes in Wikipedia's "Category of Computer Science" page, we also referenced online computer science syllabi to supplement any missing pieces of information.

To generate the test suite, we focused on topics in Computer Science (Level 1) as proof of concept. We choose to focus Data Structures, Artificial Intelligence, Databases, and Operating Systems as the four main Level 2 areas. We chose to study these areas because these are common courses in a computer science curriculum. Because we were interested in improving generations at Level 5 topic area, we choose to a total of 20 different subtopics. Figure 1 illustrates the complete test suite that we used to evaluate each prompting strategy. In total we tested 29 different topics in Computer Science. For each topic, we had an LLM generate 5 different subtopics because having multiple options to review is an important step of divergent brainstorming. As a result, a total of 145 generated subtopics were evaluated for each prompting strategy.

## 3.3 Prompting Strategies

We tested three different prompting techniques to help incrementally elicit topic hierarchies following the 5-Levels Topic Classification System. We specifically used OpenAI's GPT-4 API [15] as the LLM for this task and tested each prompting strategy on the Wikipedia test suite of Computer Science concepts. For all prompting strategies, we explicitly asked GPT-4 to generate 5 subtopics. The three prompting strategies are illustrated in Table 1. The Current Topic prompting strategy only contains the current topic when generating subtopics. The Root + Current Topic prompting strategy contains both the Level 1 topic, or root, and the current topic. Finally, the Full Path + Current Topic prompting strategy includes the entire chain of parent topics, leading up and including the current topic in the prompt.

## 3.4 Evaluation Criteria

The first author and an independent expert were tasked with annotating the generated subtopics for each prompting strategy. Both annotators were experts in computer science. We provided the evaluation rubric below to each annotator, along with detailed directions with how to annotate the generated subtopics. Annotations were done separately. Each annotator labeled 145 generated subtopics for each of the 3 different prompting strategies. We paid $16 an hour for 4 hours of work.

The evaluation criteria given to annotators cover issues of repetitiveness, specificity, and relatedness to the inputted topic:

(1) **Repetitive**: the generated subtopic repeats the same input topic;
(2) **Too specific**: the generated subtopic is too specific for the desired level;
(3) **Too general**: the generated subtopic is broader than the desired level;
(4) **Tangential**: the generated subtopic is at the correct level of specificity but is not directly related to the input topic;
(5) **Unrelated**: the generated subtopic is unrelated to the root level topic.

## 4 Results

The two annotators had a substantial inter-rater agreement on their assessment over three strategies, an average Cohen-Kappa of 0.61 across all annotation assignments. As a result we averaged the accuracy across the two annotators because of the high agreement. The Full Path + Current Topic yielded the highest average accuracy of **77%**. Followed by Root + Current Topic with an accuracy of 70%, and then Current Topic with an accuracy of 58%. These results are demonstrated in Figure 2. These results show that by including the full path of parent topics helped GPT-4 generate more concrete and specific subtopics. Including Root + Current Topic helped more than providing no additional information to the base prompt, Current Topic. This demonstrates that by providing more context in the prompt helps GPT-4 with generating specifically scoped and concrete subtopics.

We performed an analysis on the generated subtopics and found that the biggest problem was Too General errors, followed by Too Specific errors. Rarely were generated subtopics Tangential or Repetitive. For the Current Topic prompting strategy, 27% of the errors were due to the generated subtopics being Too General, while the Root + Current Topic and the Full Path + Current Topic yielded 14% and 10% improperly scoped topics due to Too General errors. While Too General errors were seen as the largest source of error across all three prompting strategies, Too Specific errors also accounted for the second largest error category, averaging about 9%

| Prompting Strategy | Base Prompt | Level 4 Sample Prompt |
|---|---|---|
| Current Topic | "List 5 subtopics of {curr_topic}" | "List 5 subtopics of *shortest path algorithms*." |
| Root + Current Topic | "In {level_1}, list 5 subtopics of {curr_topic}" | "In *computer science*, list 5 subtopics of *shortest path algorithms*." |
| Full Path + Current Topic | "In {level_1, ..., level_n-1}, list 5 subtopics of {curr_topic}" | "In *computer science, data structures, and graph algorithms*, list 5 subtopics of *Shortest path algorithms*." |

**Table 2: Three different prompting strategies that were used to elicit topic hierarchies from LLMs. The Table illustrates the name of the prompting strategy, the base prompt, and an example prompt using a Level 4 topic of shortest path algorithms.**

| | *Improperly Scoped Topics* | | | | |
|---|---|---|---|---|---|
| | **Too General** | Too specific | Unrelated | Tangential | Repetitive |
| **Current Topic** | **27%** | 8% | 4% | 2% | 1% |
| **Root + Current Topic** | **14%** | 9% | 3% | 3% | 0% |
| **Full Path + Current Topic** | **9%** | 9% | 0% | 2% | 0% |

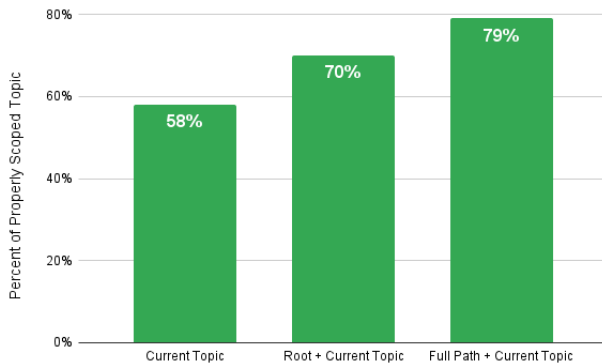**Table 3: Distribution of error category across each prompting strategy for all generated subtopics.**



**Figure 2: Bar chart demonstrating the percentage of properly scoped subtopics across all levels for each prompting strategy.**

| | Level 2 | Level 3 | Level 4 | Level 5 |
|---|---|---|---|---|
| **Current Topic** | 1% | 4% | 2% | 34% |
| **Root + Current Topic** | 1% | 3% | 4% | 21% |
| **Full Path + Current Topic** | 1% | 3% | 5% | 12% |

**Table 4: Distributions of errors across each level for each prompting strategy.**

strategy and remained largely similar across all prompting techniques. At Level 2 and Level 3, Too Specific was one of the most common errors generated. The reason that Too Specific errors occurred at the broader levels of specificity is probably due to LLMs struggling to incrementally generate subtopics. Specific examples of these errors are listed in the following sections.

### 4.1 Too General Error Examples

An example of a generated Level 5 subtopic being too General is for the topic of "Minimum Spanning Trees" and the generated subtopic "Randomized Algorithms" as a subtopic. "Minimum Spanning Trees" are graphs that connect all vertices with the minimum possible total edge weight while "Randomized Algorithms" is any algorithm that uses any degree of randomness in its logic. "Randomized Algorithms" are too general because they are not a specific implementation of minimum spanning trees, instead they are just a subset of types of algorithms. An appropriate subtopic of "Minimum Spanning Trees" would be "Kruskal's Algorithm" because that is a specific algorithm that is used for the purpose of finding minimum spanning trees.

Less frequently, Too General errors occur at the Level 3 area due to overlaps between the topic and the generated subtopic. For example, "Robotics" was generated as a subtopic of "Artificial Intelligence," a field that involves creating machines to emulate human intelligence. While there is an intersection between robots and AI techniques, the broad domain of robotics does not require robots to emulate human intelligence in the same way subtopics in "Artificial Intelligence" do. Wikipedia classifies Robotics as a subtopic

of the errors for all prompting approaches. While prompting strategy reduced many of the Too General errors, the number of Too Specific errors remained consistent across all prompting strategies. Table 3 contains the distribution of error type across all 3 different prompting strategies.

We also analyzed the specific levels that each type of error was occurring, as seen in Table 4. We found that that most errors occurred at Level 5. Too General errors were the most frequent error type at Level 5. This demonstrates that LLMs struggle to elicit specific information at that depth. Because LLMs operate on generating the next most probably word, it is possible that the more specific the subtopic, the less probable it is. As a result, LLMs will revert back to general topics when a certain level of specificity is reached because the generating a general topic might have a higher likelihood than generating a specific subtopic.

The Full Path + Current Topic prompting strategy reduced the errors in Level 5. But errors generated in Level 2, Level 3, and Level 4 were not reduced by the Full Path + Current Topic prompting

under Computer Engineering and not a subtopic under Artificial Intelligence. As a result, the generated subtopic of "Robotics" under "Artificial Intelligence" is too specific.

## 4.2    Too Specific Error Examples in Level 3

The Level 2 topic of "Artificial Intelligence" and the generated subtopic of "Neural Networks" is another example of the generated subtopics being Too Specific at the Level 3 area. Neural Networks are a specific structure that are used in machine learning and should be classified as a Level 4 topic which covers topics that explore different examples. To correct this example, the sequence of topics should go from "Artificial Intelligence" to "Machine Learning" and then to "Neural Networks." This sequence of topics first covers the general field of how machines can emulate human behaviors, before getting more specific into how machines can learn like humans, and finally how to use a technique to model how the human brain learns in computers.

## 5    Discussion

Building off previous research that use classical machine learning methods to extract topic hierarchies from unstructured data, we explored the effectiveness of LLMs in eliciting topic hierarchies and supporting the topic scoping processing: moving from abstract concepts to concrete examples. While LLMs are effective in generating a diverse set of subtopics, they still need assistance in structuring the topics into a topic hierarchy. We developed a 5-levels topic classification structure and used Wikipedia's "Category of Computer Science" page as a test suite of topics for each of the 5-levels. We explored 3 different prompting techniques to elicit topic hierarchies from LLMs and found that by including the entire sequence of parent topics helped reduce issues of improperly scoped topics.

### 5.1    Applications

There are many areas where dynamically generating topic hierarchies are beneficial. For example: educators can leverage the process of topic scoping to assist in curriculum development, content creators can use dynamically generated topic hierarchies to explore specific niches within larger themes, and product managers can use topic scoping to break down abstract goals into smaller, more concrete subtasks.

This work can be applied to the field of education to help educators with curriculum development. Traditionally, the curriculum design process requires the educator to conduct extensive planning and research based on the previous curricula and the specific needs of the students [2]. Through topic scoping, educators can better design a more audience-centric curriculum that fits their specific student demographics. The process of topic scoping and recreating knowledge hierarchies could help educators decide which topics to cover over the course of the semester that align with the grade level of their class.

Additionally, topic scoping can be used for content creators when brainstorming the types of topics they should cover. Topic scoping can help creators narrow down broad interests like reading into more specific ideas. For examples, Max is a content creator on Tiktok for book and reading content. He wants to explore what videos he should film for the upcoming month, he can use topic scoping to find different types of books to explore. Books can be

broken down by genres like Romance and Fiction, these genres can then be explored by setting, time period, and author. Max wants to explore Contemporary Fiction, from there he can further narrow down this topic into Translated, Contemporary Fiction, and even further into Japanese-Translated, Contemporary Fiction. Here Max is ready to start brainstorming his list of top 5 Japanese-Translated, Contemporary Fiction books to share with his Tiktok followers. Just like that, topic scoping helped Max find a niche within a larger reading community.

More broadly, the process of breaking an abstract goal into a series of concrete tasks and sub-goals can be applied to project and product managers as they track the development of a project by the smaller subgoals towards a more abstract goal. For example, Jenna is a product manager that is in charge of developing a new feature for a e-commerce website to improve traffic. The broad goal of feature development to improve traffic is abstract, so Jenna might break down that goal into smaller sub-goals like understanding current user traffic data on the site and even more specific goals like doing user studies and A/B tests. The applications of topic scoping are demonstrated in the process of iteratively, narrowing down a goal into specific and concrete tasks. By creating a topic hierarchy, Jenna is able to track the overall progress of the project and work in parallel with her teammates by each tackling one sub-goal category.

### 5.2    Limitations

While including the all parent topic yielded 78% of properly scoped subtopics, the main issue that caused improper subtopic generations were due to scope. Generated subtopics errors were mostly likely to be either Too General or Too Specific, demonstrating that more research can be done to improve the scoping issue. One approach is to formalize more fine-grained definitions for each level in the 5-levels classification system and including the definitions in in the prompt to improve Level 5 generations. Additionally, more work can be done to explore specific prompting techniques to reduce the amount of errors generated at Levels 2 and 3 since the current prompting strategies don't reduce the errors at the broader levels.

While there might be concerns around GPT-4 already being trained on Wikipedia's "Category of Computer Science" topics, our research isn't focused on novel topic generation instead we are exploring whether LLMs can incrementally generate topics at an increasing level of specificity. As a result, having a model trained on this data should not significantly impact performance. Additionally, our findings show that even if the model was trained on Wikipedia's "Category of Computer Science" topics, GPT-4 still struggles to generate specific and concrete subtopics at Level 5. This demonstrates that despite possibly training on this data, the model still struggles with replicating the hierarchical structure.

Finally, this work focuses on generating subtopics in Computer Science as a proof of concept. Future work can explore how these strategies can be generalized to other domains like Biology, Chemistry, and Physics. Additionally, an interactive, user-driven interface can be developed to support users in the process of topic scoping.

## 6    Conclusion

We found that it is possible to generate topic hierarchies from LLMs by incrementally generating subtopics at increasing levels

of specificity. This finding allows for future work to be done on dynamically generating, user-directed topic hierarchies to support a range of tasks like curriculum development for educators, content creation for creators, and product management tools.

# References

[1] Addi Ait-Mlouk and Lili Jiang. 2020. KBot: a Knowledge graph based chatBot for natural language understanding over linked data. *IEEE Access* 8 (2020), 149220–149230.

[2] Jill Anderson. 2022. The Challenges of Curriculum Design. https://www.gse.harvard.edu/news/22/06/challenges-curriculum-design

[3] David M Blei, Andrew Y Ng, and Michael I Jordan. 2003. Latent dirichlet allocation. *Journal of machine Learning research* 3, Jan (2003), 993–1022.

[4] Alex Calderwood, Vivian Qiu, Katy Ilonka Gero, and Lydia B Chilton. 2020. How novelists use generative language models: An exploratory user study.. In *HAI-GEN+ user2agent@ IUI*.

[5] Scott Crossley, Laura K Allen, Erica L Snow, and Danielle S McNamara. 2015. Pssst... textual features... there is more to automatic essay scoring than just you!. In *Proceedings of the fifth international conference on learning analytics and knowledge*. 203–207.

[6] Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and Zhifang Sui. 2022. A survey on in-context learning. *arXiv preprint arXiv:2301.00234* (2022).

[7] Rafael Ferreira-Mello, Máverick André, Anderson Pinheiro, Evandro Costa, and Cristobal Romero. 2019. Text mining in education. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* 9, 6 (2019), e1332.

[8] Claudio Gutierrez and Juan F Sequeda. 2021. Knowledge graphs. *Commun. ACM* 64, 3 (2021), 96–104.

[9] Shaoxiong Ji, Shirui Pan, Erik Cambria, Pekka Marttinen, and S Yu Philip. 2021. A survey on knowledge graphs: Representation, acquisition, and applications. *IEEE transactions on neural networks and learning systems* 33, 2 (2021), 494–514.

[10] Fang Jin, Nathan Self, Parang Saraf, Patrick Butler, Wei Wang, and Naren Ramakrishnan. 2013. Forex-foreteller: Currency trend modeling using news articles. In *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*. 1470–1473.

[11] Amir Karami, Aryya Gangopadhyay, Bin Zhou, and Hadi Kharrazi. 2018. Fuzzy approach topic discovery in health and medical corpora. *International Journal of Fuzzy Systems* 20 (2018), 1334–1345.

[12] Jiayi Liao, Xu Chen, Qiang Fu, Lun Du, Xiangnan He, Xiang Wang, Shi Han, and Dongmei Zhang. 2024. Text-to-image generation for abstract concepts. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 3360–3368.

[13] Tao Long, Dorothy Zhang, Grace Li, Batool Taraif, Samia Menon, Kynnedy Simone Smith, Sitong Wang, Katy Ilonka Gero, and Lydia B. Chilton. 2023. Tweetorial Hooks: Generative AI Tools to Motivate Science on Social Media. In *Proceedings of the 14th Conference on Computational Creativity (ICCC '23)*. Association for Computational Creativity.

[14] Andi Besse Firdausiah Mansur and Norazah Yusof. 2013. Social learning network analysis model to identify learning patterns using ontology clustering techniques and meaningful learning. *Computers & Education* 63 (2013), 73–86.

[15] OpenAI. 2023. GPT-4 Technical Report. arXiv:2303.08774 [cs.CL]

[16] Aurora Pons-Porrata, Rafael Berlanga-Llavori, and José Ruiz-Shulcloper. 2007. Topic discovery based on text mining techniques. *Information processing & management* 43, 3 (2007), 752–768.

[17] Daniel Ramage, Susan Dumais, and Dan Liebling. 2010. Characterizing microblogs with topic models. In *Proceedings of the international AAAI conference on web and social media*, Vol. 4. 130–137.

[18] Orit Shaer, Angelora Cooper, Osnat Mokryn, Andrew L Kun, and Hagit Ben Shoshan. 2024. AI-Augmented Brainwriting: Investigating the use of LLMs in group ideation. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. 1–17.

[19] Sangho Suh, Bryan Min, Srishti Palani, and Haijun Xia. 2023. Sensecape: Enabling multilevel exploration and sensemaking with large language models. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*. 1–18.

[20] Rui Sun, Xuezhi Cao, Yan Zhao, Junchen Wan, Kun Zhou, Fuzheng Zhang, Zhongyuan Wang, and Kai Zheng. 2020. Multi-modal knowledge graphs for recommender systems. In *Proceedings of the 29th ACM international conference on information & knowledge management*. 1405–1414.

[21] Ah-Hwee Tan et al. 1999. Text mining: The state of the art and the challenges. In *Proceedings of the pakdd 1999 workshop on knowledge disocovery from advanced databases*, Vol. 8. 65–70.

[22] Sitong Wang, Samia Menon, Tao Long, Keren Henderson, Dingzeyu Li, Kevin Crowston, Mark Hansen, Jeffrey V Nickerson, and Lydia B Chilton. 2024. ReelFramer: Human-AI co-creation for news-to-video translation. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. 1–20.

[23] Sitong Wang, Savvas Petridis, Taeahn Kwon, Xiaojuan Ma, and Lydia B Chilton. 2023. PopBlends: Strategies for Conceptual Blending with Large Language Models. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–19.

[24] Shirley Wu, Shiyu Zhao, Michihiro Yasunaga, Kexin Huang, Kaidi Cao, Qian Huang, Vassilis N Ioannidis, Karthik Subbian, James Zou, and Jure Leskovec. 2024. STaRK: Benchmarking LLM Retrieval on Textual and Relational Knowledge Bases. *arXiv preprint arXiv:2404.13207* (2024).

[25] Jihong Yan, Chengyu Wang, Wenliang Cheng, Ming Gao, and Aoying Zhou. 2018. A retrospective of knowledge graphs. *Frontiers of Computer Science* 12 (2018), 55–74.