

Can Instructed Retrieval Models Really Support Exploration?

Piyush Maheshwari
University of Massachusetts
Amherst, USA
psmaheshwari@umass.edu

Sheshera Mysore
Microsoft
Seattle, USA
smysore@iesl.cs.umass.edu

Hamed Zamani
University of Massachusetts
Amherst, USA
zamani@cs.umass.edu

Abstract

Exploratory searches are characterized by under-specified goals and evolving query intents. In such scenarios, retrieval models that can capture user-specified nuances in query intent and adapt results accordingly are desirable — instruction-following retrieval models promise such a capability. In this work, we evaluate instructed retrievers for the prevalent yet under-explored application of aspect-conditional seed-guided exploration using an expert-annotated test collection. We evaluate both recent LLMs fine-tuned for instructed retrieval and general-purpose LLMs prompted for ranking with the highly performant Pairwise Ranking Prompting. We find that the best instructed retrievers improve on ranking relevance compared to instruction-agnostic approaches. However, we also find that instruction following performance, crucial to the user experience of interacting with models, does not mirror ranking relevance improvements and displays insensitivity or counter-intuitive behavior to instructions. Our results indicate that while users may benefit from using current instructed retrievers over instruction-agnostic models, they may not benefit from using them for long-running exploratory sessions requiring greater sensitivity to instructions.

CCS Concepts

• **Information systems** → **Users and interactive retrieval**; *Language models*; • **Human-centered computing** → **Natural language interfaces**.

Keywords

instructed retrievers; exploratory search; complex retrieval

ACM Reference Format:

Piyush Maheshwari, Sheshera Mysore, and Hamed Zamani. 2026. Can Instructed Retrieval Models Really Support Exploration?. In *2026 ACM SIGIR Conference on Human Information Interaction and Retrieval (CHIIR '26)*, March 22–26, 2026, Seattle, WA, USA. ACM, New York, NY, USA, 6 pages. <https://doi.org/10.1145/3786304.3787888>

1 Introduction

The rise of instruction following capabilities has enabled users to interact with powerful LLMs through natural language instructions [23]. Following their emergence in generative models, retrieval models that follow user instructions were developed [2, 28]. These models promise users the ability to specify nuanced relevance criteria, narratives, or aspects with their query. This has extended

the frontier of retrieval models from being optimized for semantic similarity to targeting more complex requests.

To enable the development of instructed retrievers, recent work has constructed large-scale benchmarks spanning standard keyword search as well as more complex requests ranging from code or theorem retrieval, tip of tongue retrieval, and literature search, among others [16, 29, 33]. However, little work has investigated the use of instructed retrievers in the exploratory search scenarios where users start at a seed-document and explore a corpus of documents. Our work bridges this gap through a focused evaluation of instructed retrievers for seed-guided exploration.

Seed-guided exploration has been the focus of significant research in the design community [6, 25] and remains common in exploring document/item collections such as research papers, e-commerce products, and books [19, 36]. Despite this prevalence, seed-guided exploration remains under-supported in search systems [36]. The development of highly performing instructed retrievers promises to benefit this prevalent mode of exploration. Specifically, seed-guided exploration is characterized by long document queries where users desire retrieval only based on a few aspects of the document. This calls for retrievers that can perform retrieval based on nuanced aspects of long seed documents [20, 36] – instructed retrievers offer a promising solution to this challenge. Similarly, the dynamic nature of intents and goals in exploration calls for easily steerable retrieval models [35] – natural language instructions offer a promising interaction mode for such steerability.

In this paper, we choose the case of scientific document exploration and conduct a careful evaluation of instructed retrievers for their ranking relevances as well as instruction following ability. Retrieval in science offers compelling use cases for accelerating scientific discovery [11, 13] and aiding researchers in coping with an exploding literature [1, 7]. It also enables us to leverage a high-quality, aspect conditional seed-guided exploration test collection [20, CSFCube] for our experiments. CSFCube consists of expert annotations for paragraph-length queries (paper abstracts) annotated for relevance w.r.t multiple aspects (problem, solution, result) per query. This enables evaluation of the important scenario in seed-guided exploration, where users desire retrieval only based on some aspects of the seed document. Further, CSFCubes’ multiple instructions per query (Figure 1) also enables the evaluation of *instruction following* in instructed retrievers. Instruction following forms a core part of users’ experience interacting with instructed retrievers and remains a notably missing aspect in multiple prior benchmarking efforts for instructed retrievers [34]. To the best of our knowledge, no other datasets enable both evaluations for exploration applications [9]. In our experiments, we evaluate recent LLMs fine-tuned for instructed retrieval, general-purpose LLMs used with Pairwise Ranking Prompting [26], and small parameter LMs finetuned for aspect conditional retrieval.



This work is licensed under a Creative Commons Attribution 4.0 International License. CHIIR '26, Seattle, WA, USA

© 2026 Copyright held by the owner/author(s).

ACM ISBN 979-8-4007-2414-5/2026/03

<https://doi.org/10.1145/3786304.3787888>

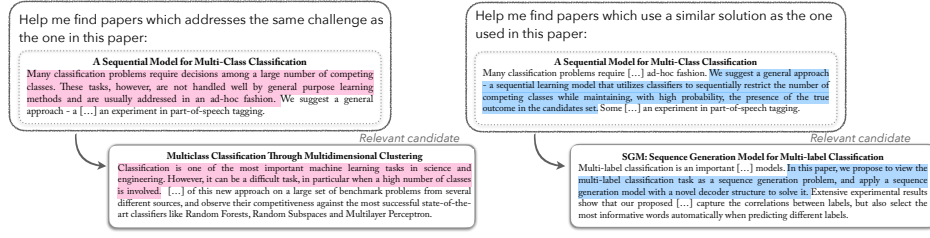


Figure 1: CSFCube contains the same query document annotated by experts w.r.t candidate documents for relevance with different instructions. This enables evaluation of both ranking relevance and instruction following for instructed retrievers.

We find that the best commercial and open-LLM-based instructed retrievers (gpt-4o and GritLM-7B [17]) improve significantly over instruction agnostic approaches in ranking relevance. However, ranking relevance trends are not reflected in instruction following performance. Models with poor ranking relevance (Llama-3-70B) tended to show better instruction following performance. In experiments examining the sensitivity of the best models to instruction variants, we find that while nuanced instructions improve over generic instructions, extensive tuning of the instructions doesn’t deliver improvements. Finally, across experiments, models’ sensitivity to instructions remains close to that of instruction-agnostic models, indicating that current models may lack the sensitivity for long-running recall-oriented exploratory sessions. In sum, we contribute the first experimental analysis of instructed retrievers in terms of ranking relevance *and* instruction following for the under-explored application of seed-guided exploratory search. Our work forms an important pre-requisite for human-centered research studying instructed retrievers for exploratory search applications.

2 Related Work

Instruction following in retrieval models. The emergence of conversational interfaces for information access has led to an increased focus on building models and systems to support natural language queries [24]. Early work in the IR community explored these in the guise of “verbose” queries or “narrative” queries for applications in web and e-commerce applications [4, 12]. In recent work, instructed retrievers have unified several forms of complex retrieval interactions spanning constraint-based retrieval, seed-guided retrieval, tip-of-tongue retrieval, and others [16]. Our work extends the understanding of this emerging form of interaction and provides an analysis of instructed retrieval models for aspect conditional seed-guided exploration.

Benchmarking complex retrieval. The emergence of complex natural language queries has been accompanied by important benchmarking efforts for such tasks [31]. Sun et al. [30, MAIR] collate 126 retrieval tasks from multiple domains, which have instructions alongside queries, and benchmark several retrieval models. Our experiments (Section 3) use the best performing instructed retriever from their analysis. Going beyond merely including instructions with queries, recent work has focused on curating benchmarks for complex and reasoning-intensive retrieval – a new frontier unlocked by instructed retrieval. Su et al. [29, BRIGHT] benchmark models on complex requests needing multiple reasoning steps. Closer to our work, Killingback and Zamani [16, CRUMB] and Wang

et al. [33, BRICO] introduce benchmarks for complex retrieval with multi-aspect queries. Oh et al. [21, INSTRUTIR] and Weller et al. [34, FollowIR] go beyond benchmarking ranking relevance and highlight the importance of evaluating *instruction following*, a core part of user experience with instructed retrievers. Weller et al. advocate for having multiple instructions and accompanying relevance judgments per query to evaluate instruction following and introduce the p-MRR metric for instruction following. The multi-aspect relevance judgments of CSFCube (Section 3) enable evaluation of this important aspect of instructed retrievers.

Retrieval for Science Applications. A large body of work has focused on developing retrieval and exploration systems aimed at scientific texts. Early work related to ours focused on citation recommendation [10] and analogical search engines to aid scientific brainstorming [15]. More recent work has benchmarked LLM-based retrievers for complex literature search queries [1] or developed retrieval approaches for RAG systems used for automated scientific discovery [11]. Our work extends these efforts by analyzing instructed retrievers for aspect conditional seed-guided exploration. The closest work to ours is presented by Wang et al. [32, DORIS-MAE], who introduce LLM-annotated multi-aspect queries for scientific retrieval. In contrast to Wang et al. [32], we leverage an expert-annotated test collection enabling greater reliability [8]. Further, our multi-aspect queries *and* relevances enable evaluation of instruction following (elaborated in Section 3). Finally, we contribute an evaluation of multiple families of LLM-based instructed retrievers, going beyond dense retrievers alone.

3 Experimental Setup

Test collection. We use the CSFCube [20] test collection, annotated for aspect conditional seed-guided retrieval by experts in our experiments. It consists of 50 queries and a document collection of 800k documents for retrieval. All queries and documents in CSFCube are drawn from computer science and engineering domains (see examples in Figure 1). Specifically, CSFCube contains 50 pairs of query aspects and documents. The query aspect indicates which aspect of the query/seed document a retriever should use for retrieval. We format these aspects as instructions for instructed retrievers in our experiments. The aspects focus on common aspects in most scientific documents [5, 15] and are diverse enough to generalize to a broader set of user-defined aspects: “*Background*”, “*Method*”, “*Result*”. Of special note, 32 of the query-aspect pairs in CSFCube annotate the same query document with relevance along two different aspects, which enables CSFCube to be used for instruction

following evaluation [34]. Each query-aspect pair is annotated by at least two experts for relevance against 200 documents on average. This provides a reliable and high-quality test collection to evaluate both ranking relevance and instruction following.

Retrieval Models. We include a diverse set of models in our evaluation. We evaluate instruction agnostic dense retrievers (SPECTER2, SciNCL) specialized for scientific document retrieval [22, 27]. We also include an aspect conditional multi-vector retrieval model (OTASPIRE) to evaluate a prior generation of exploratory search model [18]. Both these model families are built on small parameter models (110M parameters) and form our baseline models. For an embedding-based instructed retriever, we use GritLM-7B, the best performing instructed retriever in prior benchmarking efforts [30]. Finally, we use a range of open weight (Mistral-7B, Llama-3-8B, Llama-3-70B) and commercial (gpt-3.5-turbo, gpt-4o) LLMs for ranking with Pairwise Ranking Prompting (PRP) [26]. PRP ranks documents by their win-rate in pairwise relevance comparisons of a pair of candidate documents for a query. It strongly outperforms other prompting methods, such as pointwise scoring, side-steps the long-context limitations of listwise LLM re-ranking, and approaches the performance of trained ranking models. We use an efficient heapsort-based implementation of PRP in our experiments [26, Sec 3.3]. We also include a point-wise prompting approach due to its simplicity and wide use. Given the size of the instructed retrievers, we use them as re-rankers to re-rank the first 200 documents from the performant SciNCL first-stage-ranker. We use the expert annotations in CSFCube to report “Human” performance for ranking relevance and instruction following.

Evaluation Metrics. We evaluate both ranking relevance and instruction following of models. We use NDCG@20 to evaluate ranking relevance and report statistical significance with a paired t-test at $p < 0.05$. We use p-MRR to evaluate instruction following [34]. The metric ranges from -1 to $+1$ and indicates worst to best instruction following. p-MRR compares changes in ranking for a single query (q) w.r.t two instructions i and i' , each with different relevance judgements. The metric is formulated below, where R indicates a document’s rank (1 is the top rank), and is computed for every document that is relevant for (i, q) but not (i', q):

$$p - MRR = \begin{cases} \frac{R_{i'}}{R_i} - 1 & \text{if } R_i > R_{i'} \\ 1 - \frac{R_{i'}}{R_i} & \text{otherwise} \end{cases} \quad (1)$$

The metric per document is first averaged for a query and then averaged across queries to result in a dataset-level metric. The metric evaluates how well a model follows instruction i , with positive scores indicating that the model results in relevant documents ranked at higher positions for (i, q). On the other hand, negative scores indicate that (i', q) results in documents relevant to (i, q) at better ranks, i.e., counterintuitive instruction following behavior. A p-MRR of 0 indicates that the ranked lists for (i, q) and (i', q) are identical, as expected for an instruction agnostic model. We use the relevance annotations for multiple aspects per query in CSFCube to report p-MRR. For example, we report instruction following for “Background” treating it as i and {“Method”, “Result”} as i' , etc. Our code details our experiments further: <https://github.com/MSheshera/exploration-instructionfollowing>

4 Results

We report ranking relevance and instruction following results in Tables 1 and 2. In Section 4.2 we probe the sensitivity of instructed retrievers to instruction variants, an important aspect of the user experience interacting with instructed retrievers.

4.1 Relevance and Instruction Following

Ranking relevance. We begin by noting that pointwise (denoted pw) prompting approaches underperform PRP prompting. Similarly, a general-purpose LLM (gpt-4o_{prp}), prompted with PRP, approaches a trained instructed retrieval model (GritLM-7B). These results mirror prior work [26]. The best performing instructed retrievers (GritLM-7B and gpt-4o_{prp}) improve upon the ranking quality of baseline models. We also see that several open-weight (Mistral, Llama-3) and commercial (gpt-3.5-turbo) models underperform baseline first-stage rankers. In sum, **our results indicate that the best instructed retrievers meaningfully improve upon baseline models for seed-guided exploration.** However, we also see that all models lag behind *Human* performance by a large margin, indicating significant room for improvement.

Instruction following. We see that trends in ranking relevance (Table 1) are not mirrored in instruction following (Table 2). While Llama-3-70B shows poor ranking relevance, it shows strong instruction following. Similarly, while gpt-4o_{prp} performs well for ranking relevance, it shows less consistent instruction following performance, varying between different aspects. On the other hand, GritLM-7B balances between ranking relevance and instruction following, displaying more consistent performance on both metrics. We also see OTASPIRE, an aspect conditional retrieval model that uses aspect-specific query sentences to perform retrieval, shows stronger instruction following performance. In Section 4.2, we see that this strategy also results in improved instruction following for GritLM-7B. As with ranking relevance, we see that *Human* performance leaves a significant margin of improvement for all retrieval models. Further, **p-MRR scores nearing 0 or negative values indicate that the best instructed retrievers display instruction-agnostic or counter-intuitive instruction following behavior.** This indicates that, despite strong ranking relevance, there remains significant room to improve models for exploratory search applications, which require nuanced forms of relevance and instruction following behavior that is intuitive to users.

Aspect-specific variation. Finally, we also examine differences across different aspects in Tables 1 and 2, given that users tend to examine different aspects depending on their application and expertise [14]. *Background* instructions see both the best relevance and instruction following, and *Method* sees the worst performance. These results are reflected in multiple studies using the CSFCube test collection [9, 20]. *Background* performance is most similar to topical/keyword similarity that models are optimized for and consequently displays strong performance; on the other hand, strong performance on *Method* aspects requires abstract forms of similarity beyond term overlap, a capability most models don’t yet possess.

4.2 Instruction Sensitivity

Given that instructions serve as the primary mode of user interaction with instructed retrievers, we examine the impact of common

Table 1: Ranking relevance on CSFCube. [×] indicates lack of a significant difference from GritLM-7B.

Aspects →	Agg.	Background	Method	Result
Models	NDCG@20	NDCG@20	NDCG@20	NDCG@20
Human	77.91	82.75	70.38	81.00
SciNCL	36.43	45.22	26.19 [×]	38.24 [×]
SPECTER2	36.42	43.44	26.34 [×]	40.43 [×]
OTASPIRE	36.02	45.33	23.89 [×]	39.35
GritLM-7B	42.38	53.15	28.82	45.90
gpt-3.5-turbo _{pw}	20.40	32.16	09.73	20.20
gpt-4o _{pw}	23.71	41.79	21.46	36.87
Mistral-7B _{prp}	27.34	36.51	15.58	30.35
Llama-3-8B _{prp}	24.40	18.51	20.94	33.48
Llama-3-70B _{prp}	34.79	43.50	24.68 [×]	36.71
gpt-3.5-turbo _{prp}	30.82	37.33	19.11	36.40
gpt-4o _{prp}	<u>41.07[×]</u>	<u>52.40[×]</u>	<u>30.04[×]</u>	<u>41.43[×]</u>

Table 2: Instruction following results on CSFCube using p-MRR [34]. We report scores between [−1, 1] scaled by 100.

Aspects →	Agg.	Background	Method	Result
Models	p-MRR	p-MRR	p-MRR	p-MRR
Human	+25.3	+29.8	+22.2	+24.0
SciNCL	0.0	0.0	0.0	0.0
SPECTER2	0.0	0.0	0.0	0.0
OTASPIRE	+3.88	+4.53	+3.13	+3.98
GritLM-7B	+2.02	+2.99	+1.14	+1.92
gpt-3.5-turbo _{pw}	+0.80	+2.60	−3.74	+3.55
gpt-4o _{pw}	+0.85	+1.17	+1.13	+0.24
Mistral-7B _{prp}	−1.10	+3.29	−8.67	+2.07
Llama-3-8B _{prp}	−4.00	+0.68	−3.06	−9.59
Llama-3-70B _{prp}	+4.44	+8.23	+0.37	+4.73
gpt-3.5-turbo _{prp}	+0.86	+2.03	−4.20	+4.74
gpt-4o _{prp}	+0.97	+3.79	+1.50	−2.39

instruction/input variants on both NDCG@20 and p-MRR in Table 3. We select GritLM-7B given that it balances between strong ranking relevance and instruction following for this analysis. We avoid experimenting with gpt-4o for instruction variants, given its significant expense. We compare the Base instruction, i.e., short aspect conditional instructions (see Figure 1) to the following instruction variants: (1) Generic: Uses a generic instruction rather than an aspect conditional one, representing a instruction-agnostic approach with instructed retrievers. (2) Definition: Uses a detailed aspect definition resulting in long form instructions. Our definitions are drawn from Mysore et al. [20]. (3) Paraphrases: Reports the average metric over three separate paraphrases of the instruction. (4) Aspect subset: Selects the subset of sentences in the query document that correspond to an aspect (available in CSFCube) rather than using the entire aspect. We also use a short aspect-specific instruction.

We note several trends in Table 3. We see that Generic instructions result in similar NDCG@20 metrics to the base aspect-specific instruction; however, we also see Generic instructions resulting in p-MRR nearing 0. Indicating that while GritLM-7B may not display

Table 3: Examining the sensitivity of the best performing instructed retriever (GritLM-7B) to instruction variants.

Aspects →	Agg.	Background	Method	Result
Models	NDCG@20	NDCG@20	NDCG@20	NDCG@20
GritLM-7B (Base)	42.38	53.15	28.82	45.90
↪ Generic	42.08	52.46	28.94	45.53
↪ Definition	41.43	52.54	28.18	44.25
↪ Paraphrases	42.03	52.24	28.87	45.38
↪ Aspect subset	38.21	47.70	25.48	42.07
	p-MRR	p-MRR	p-MRR	p-MRR
GritLM-7B (Base)	+2.02	+2.99	+1.14	+1.92
↪ Generic	+0.18	+0.26	+0.20	+0.09
↪ Definition	+2.01	+2.38	+1.22	+2.44
↪ Paraphrases	+1.53	+2.70	+0.95	+0.93
↪ Aspect subset	+5.89	+7.60	+5.96	+4.10

strong sensitivity to nuanced instructions, they provide greater steerability than generic and unnuanced instructions. On the other hand, we see that both aspect specific instructions (Definition) or Paraphrases impact NDCG@20 and p-MRR by smaller amounts. This indicates that GritLM-7B remains insensitive to the precise wording of an instruction, likely easing the burden on users to craft the best possible instruction, diverging from interactions with general-purpose LLMs [3]. Finally, we see that the Aspect subset strategy delivers the best instruction following performance. However, we see this results in worse ranking relevance, likely due to lost query document context. Taken together, these results indicate that the best current instructed retrievers offer meaningful gains in ranking relevance and offer some sensitivity to nuanced instructions, but users may not benefit from extensive instruction tuning and may require more sensitive models to aid long-running recall-oriented exploratory sessions.

5 Conclusions

In this paper, we present a focused evaluation of instructed retrievers for the under-explored application of aspect-conditional exploratory search. We make use of an expert-annotated dataset, CSFCube to evaluate both ranking relevance and instruction following in instructed retrievers. Our results indicate that while instructed retrievers improve ranking relevance over instruction-agnostic models, current models show counterintuitive or insensitive behavior to instructions. Instructed retrievers promise an advanced retrieval model for building and studying future human-centered interactive exploratory search systems. However, little work has evaluated human centered aspects of instructed retrievers such as instruction following and sensitivity to instruction variants for exploratory search applications – our work fills this gap and forms an important empirical basis for future human-centered research.

Acknowledgments

We thank anonymous reviewers for their feedback. This work was supported in part by the Center for Intelligent Information Retrieval. Any opinions, findings and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect those of the sponsor.

References

- [1] Anirudh Ajith, Mengzhou Xia, Alexis Chevalier, Tanya Goyal, Danqi Chen, and Tianyu Gao. 2024. LitSearch: A Retrieval Benchmark for Scientific Literature Search. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (Eds.). Association for Computational Linguistics, Miami, Florida, USA, 15068–15083. <https://doi.org/10.18653/v1/2024.emnlp-main.840>
- [2] Akari Asai, Timo Schick, Patrick Lewis, Xilun Chen, Gautier Izacard, Sebastian Riedel, Hannaneh Hajishirzi, and Wen-tau Yih. 2023. Task-aware Retrieval with Instructions. In *Findings of the Association for Computational Linguistics: ACL 2023*, Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (Eds.). Association for Computational Linguistics, Toronto, Canada, 3650–3675. <https://doi.org/10.18653/v1/2023.findings-acl.225>
- [3] Shubham Atreja, Joshua Ashkinaze, Lingyao Li, Julia Mendelsohn, and Libby Hemphill. 2025. What's in a Prompt?: A Large-Scale Experiment to Assess the Impact of Prompt Design on the Compliance and Accuracy of LLM-Generated Text Annotations. *Proceedings of the International AAAI Conference on Web and Social Media* 19 (June 2025), 122–145. <https://doi.org/10.1609/icwsm.v19i1.35807>
- [4] Toine Bogers and Marijn Koolen. 2017. Defining and Supporting Narrative-driven Recommendation. In *Proceedings of the Eleventh ACM Conference on Recommender Systems (Como, Italy) (RecSys '17)*. Association for Computing Machinery, New York, NY, USA, 238–242. <https://doi.org/10.1145/3109859.3109893>
- [5] Joel Chan, Joseph Chee Chang, Tom Hope, Dafna Shahaf, and Aniket Kittur. 2018. Solvent: A mixed initiative system for finding analogies between research papers. *Proceedings of the ACM on Human-Computer Interaction* 2, CSCW (2018), 1–21.
- [6] Duen Horng Chau, Aniket Kittur, Jason I. Hong, and Christos Faloutsos. 2011. Apolo: making sense of large network data by combining rich user interaction and machine learning. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (Vancouver, BC, Canada) (CHI '11). Association for Computing Machinery, New York, NY, USA, 167–176. <https://doi.org/10.1145/1978942.1978967>
- [7] Johan S. G. Chu and James A. Evans. 2021. Slowed canonical progress in large fields of science. *Proceedings of the National Academy of Sciences* 118, 41 (2021), e2021636118. <https://doi.org/10.1073/pnas.2021636118> arXiv:<https://www.pnas.org/doi/pdf/10.1073/pnas.2021636118>
- [8] Laura Dietz, Oleg Zendel, Peter Bailey, Charles L. A. Clarke, Ellese Cotterill, Jeff Dalton, Faegheh Hasibi, Mark Sanderson, and Nick Craswell. 2025. Principles and Guidelines for the Use of LLM Judges. In *Proceedings of the 2025 International ACM SIGIR Conference on Innovative Concepts and Theories in Information Retrieval (ICTIR) (Padua, Italy) (ICTIR '25)*. Association for Computing Machinery, New York, NY, USA, 218–229. <https://doi.org/10.1145/3731120.3744588>
- [9] Heejin Do, Sangwon Ryu, Jonghwi Kim, and Gary Geunbae Lee. 2024. Multi-Facet Blending for Faceted Query-by-Example Retrieval. *arXiv preprint arXiv:2412.01443* (2024).
- [10] Michael Färber and Adam Jatowt. 2020. Citation recommendation: approaches and datasets. *International Journal on Digital Libraries* 21, 4 (2020), 375–405.
- [11] Aniketh Garikaparthi, Manasi Patwardhan, Aditya Sanjiv Kanade, Aman Hassan, Lovekesh Vig, and Arman Cohan. 2025. MIR: Methodology Inspiration Retrieval for Scientific Research Problems. *arXiv preprint arXiv:2506.00249* (2025).
- [12] Manish Gupta and Michael Bendersky. 2015. Information Retrieval with Verbose Queries. In *Proceedings of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval (Santiago, Chile) (SIGIR '15)*. Association for Computing Machinery, New York, NY, USA, 1121–1124. <https://doi.org/10.1145/2766462.2767877>
- [13] Kexin Huang, Serena Zhang, Hanchen Wang, Yuanhao Qu, Yingzhou Lu, Yusuf Roohani, Ryan Li, Lin Qiu, Gavin Li, Junze Zhang, Di Yin, Shruti Marwaha, Jennefer N. Carter, Xin Zhou, Matthew Wheeler, Jonathan A. Bernstein, Mengdi Wang, Peng He, Jingtian Zhou, Michael Snyder, Le Cong, Aviv Regev, and Jure Leskovec. 2025. Biomni: A General-Purpose Biomedical AI Agent. *bioRxiv* (2025). <https://doi.org/10.1101/2025.05.30.656746> arXiv:<https://www.biorxiv.org/content/early/2025/06/02/2025.05.30.656746.full.pdf>
- [14] Emi Ishita, Yasuko Hagiwara, Yukiko Watanabe, and Yoichi Tomiura. 2018. Which Parts of Search Results do Researchers Check when Selecting Academic Documents?. In *Proceedings of the 18th ACM/IEEE on Joint Conference on Digital Libraries (Fort Worth, Texas, USA) (JCDL '18)*. Association for Computing Machinery, New York, NY, USA, 345–346. <https://doi.org/10.1145/3197026.3203867>
- [15] Hyeonsu B. Kang, Xin Qian, Tom Hope, Dafna Shahaf, Joel Chan, and Aniket Kittur. 2022. Augmenting Scientific Creativity with an Analogical Search Engine. *ACM Trans. Comput.-Hum. Interact.* 29, 6, Article 57 (Nov. 2022), 36 pages. <https://doi.org/10.1145/3530013>
- [16] Julian Killingback and Hamed Zamani. 2025. Benchmarking Information Retrieval Models on Complex Retrieval Tasks. *arXiv preprint arXiv:2509.07253* (2025).
- [17] Niklas Muennighoff, Chongjin SU, Liang Wang, Nan Yang, Furu Wei, Tao Yu, Amanpreet Singh, and Douwe Kiela. 2025. Generative Representational Instruction Tuning. In *The Thirteenth International Conference on Learning Representations*. <https://openreview.net/forum?id=BC4llvfSzv>
- [18] Sheshera Mysore, Arman Cohan, and Tom Hope. 2022. Multi-Vector Models with Textual Guidance for Fine-Grained Scientific Document Similarity. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Marine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz (Eds.). Association for Computational Linguistics, Seattle, United States, 4453–4470. <https://doi.org/10.18653/v1/2022.naacl-main.331>
- [19] Sheshera Mysore, Mahmood Jasim, Haoru Song, Sarah Akbar, Andre Kenneth Chera Randall, and Narges Mahyar. 2023. How Data Scientists Review the Scholarly Literature. In *Proceedings of the 2023 Conference on Human Information Interaction and Retrieval (Austin, TX, USA) (CHIIR '23)*. Association for Computing Machinery, New York, NY, USA, 137–152. <https://doi.org/10.1145/3576840.3578309>
- [20] Sheshera Mysore, Tim O’Gorman, Andrew McCallum, and Hamed Zamani. 2021. CSFCube - A Test Collection of Computer Science Research Articles for Faceted Query by Example. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*. <https://openreview.net/forum?id=8Y50dBmGU>
- [21] Hanseok Oh, Hyunji Lee, Seonghyeon Ye, Haebin Shin, Hansol Jang, Changwook Jun, and Minjoon Seo. 2024. Instructir: A benchmark for instruction following of information retrieval models. *arXiv preprint arXiv:2402.14334* (2024).
- [22] Malte Ostendorff, Nils Rethmeier, Isabelle Augenstein, Bela Gipp, and Georg Rehm. 2022. Neighborhood Contrastive Learning for Scientific Document Representations with Citation Embeddings. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (Eds.). Association for Computational Linguistics, Abu Dhabi, United Arab Emirates, 11670–11688. <https://doi.org/10.18653/v1/2022.emnlp-main.802>
- [23] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems*, S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (Eds.), Vol. 35. Curran Associates, Inc., 27730–27744. https://proceedings.neurips.cc/paper_files/paper/2022/file/b1efde53be364a73914f58805a001731-Paper-Conference.pdf
- [24] Andrea Papenmeier, Dagmar Kern, Daniel Hienert, Alfred Sliwa, Ahmet Aker, and Norbert Fuhr. 2021. Starting Conversations with Search Engines - Interfaces that Elicit Natural Language Queries. In *Proceedings of the 2021 Conference on Human Information Interaction and Retrieval (Canberra ACT, Australia) (CHIIR '21)*. Association for Computing Machinery, New York, NY, USA, 261–265. <https://doi.org/10.1145/3406522.3446035>
- [25] Antoine Ponsard, Francisco Escalona, and Tamara Munzner. 2016. PaperQuest: A Visualization Tool to Support Literature Review. In *Proceedings of the 2016 CHI Conference Extended Abstracts on Human Factors in Computing Systems* (San Jose, California, USA) (CHI EA '16). Association for Computing Machinery, New York, NY, USA, 2264–2271. <https://doi.org/10.1145/2851581.2892334>
- [26] Zhen Qin, Rolf Jagerman, Kai Hui, Honglei Zhuang, Junru Wu, Le Yan, Jiaming Shen, Tianqi Liu, Jialu Liu, Donald Metzler, Xuanhui Wang, and Michael Bendersky. 2024. Large Language Models are Effective Text Rankers with Pair-wise Ranking Prompting. In *Findings of the Association for Computational Linguistics: NAACL 2024*, Kevin Duh, Helena Gomez, and Steven Bethard (Eds.). Association for Computational Linguistics, Mexico City, Mexico, 1504–1518. <https://doi.org/10.18653/v1/2024.findings-naacl.97>
- [27] Amanpreet Singh, Mike D’Arcy, Arman Cohan, Doug Downey, and Sergey Feldman. 2022. SciRepEval: A Multi-Format Benchmark for Scientific Document Representations. In *Conference on Empirical Methods in Natural Language Processing*. <https://api.semanticscholar.org/CorpusID:254018137>
- [28] Hongjin Su, Weijia Shi, Jungo Kasai, Yizhong Wang, Yushi Hu, Mari Ostendorff, Wen-tau Yih, Noah A. Smith, Luke Zettlemoyer, and Tao Yu. 2023. One Embedder, Any Task: Instruction-Finetuned Text Embeddings. In *Findings of the Association for Computational Linguistics: ACL 2023*, Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (Eds.). Association for Computational Linguistics, Toronto, Canada, 1102–1121. <https://doi.org/10.18653/v1/2023.findings-acl.71>
- [29] Hongjin Su, Howard Yen, Mengzhou Xia, Weijia Shi, Niklas Muennighoff, Han yu Wang, Liu Haisu, Quan Shi, Zachary S Siegel, Michael Tang, Ruoxi Sun, Jinsung Yoon, Sercan O Arik, Danqi Chen, and Tao Yu. 2025. BRIGHT: A Realistic and Challenging Benchmark for Reasoning-Intensive Retrieval. In *The Thirteenth International Conference on Learning Representations*. <https://openreview.net/forum?id=ykucs5q381b>
- [30] Weiwei Sun, Zhengliang Shi, Wu Jiu Long, Lingyong Yan, Xinyu Ma, Yiding Liu, Min Cao, Dawei Yin, and Zhaochun Ren. 2024. MAIR: A Massive Benchmark for Evaluating Instructed Retrieval. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (Eds.). Association for Computational Linguistics, Miami, Florida, USA, 14044–14067. <https://doi.org/10.18653/v1/2024.emnlp-main.778>

- [31] Johanne R. Trippas, Sara Fahad Dawood Al Lawati, Joel Mackenzie, and Luke Gallagher. 2024. What do Users Really Ask Large Language Models? An Initial Log Analysis of Google Bard Interactions in the Wild. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval* (Washington DC, USA) (*SIGIR '24*). Association for Computing Machinery, New York, NY, USA, 2703–2707. <https://doi.org/10.1145/3626772.3657914>
- [32] Jianyou Wang, Kaicheng Wang, Xiaoyue Wang, Prudhvira Naidu, Leon Bergen, and Ramamohan Paturi. 2023. Scientific Document Retrieval using Multi-level Aspect-based Queries. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*. <https://openreview.net/forum?id=XjaWEAyToL>
- [33] Xiaoyue Wang, Jianyou Wang, Weili Cao, Kaicheng Wang, Ramamohan Paturi, and Leon Bergen. 2024. Birco: A benchmark of information retrieval tasks with complex objectives. *arXiv preprint arXiv:2402.14151* (2024).
- [34] Orion Weller, Benjamin Chang, Sean MacAvaney, Kyle Lo, Arman Cohan, Benjamin Van Durme, Dawn Lawrie, and Luca Soldaini. 2025. FollowIR: Evaluating and Teaching Information Retrieval Models to Follow Instructions. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, Luis Chiruzzo, Alan Ritter, and Lu Wang (Eds.). Association for Computational Linguistics, Albuquerque, New Mexico, 11926–11942. <https://doi.org/10.18653/v1/2025.naacl-long.597>
- [35] Ryen W White and Resa A Roth. 2009. *Exploratory search: Beyond the query-response paradigm*. Number 3. Morgan & Claypool Publishers.
- [36] Huiwen Zhang, Dana Mckay, Michael Twidale, and George Buchanan. 2024. Something Just Like This: A Secret History of the Role of Analogues in Information Seeking. In *Proceedings of the 2024 Conference on Human Information Interaction and Retrieval* (Sheffield, United Kingdom) (*CHIIR '24*). Association for Computing Machinery, New York, NY, USA, 189–198. <https://doi.org/10.1145/3627508.3638317>