

SciNetBENCH: A RELATION-AWARE BENCHMARK FOR SCIENTIFIC LITERATURE RETRIEVAL AGENTS

Anonymous authors

Paper under double-blind review

ABSTRACT

The rapid development of AI agent has spurred the development of advanced research tools, such as Deep Research. Achieving this require a nuanced understanding of the relations within scientific literature, surpasses the scope of keyword-based or embedding-based retrieval. Existing retrieval agents mainly focus on the content-level similarities and are unable to decode critical relational dynamics, such as identifying corroborating or conflicting studies or tracing technological lineages, all of which are essential for a comprehensive literature review. Consequently, this fundamental limitation often results in a fragmented knowledge structure, misleading sentiment interpretation, and inadequate modeling of collective scientific progress. To investigate relation-aware retrieval more deeply, we propose **SciNetBench**, the first **Scientific Network Relation-aware Benchmark** for literature retrieval agents. Constructed from a corpus of over 18 million AI papers, our benchmark systematically evaluates three levels of relations: ego-centric retrieval of papers with novel knowledge structures, pair-wise identification of scholarly relationships, and path-wise reconstruction of scientific evolutionary trajectories. Through extensive evaluation of three categories of retrieval agents, we find that their accuracy on relation-aware retrieval tasks often falls below 20%, revealing a core shortcoming of current retrieval paradigms. Notably, further experiments on the literature review tasks demonstrate that providing agents with relational ground truth leads to a substantial 23.4% performance improvement in the review quality, validating the critical importance of relation-aware retrieval. We publicly release our benchmark at <https://anonymous.4open.science/r/SciNetBench/> to support future research on advanced retrieval systems.

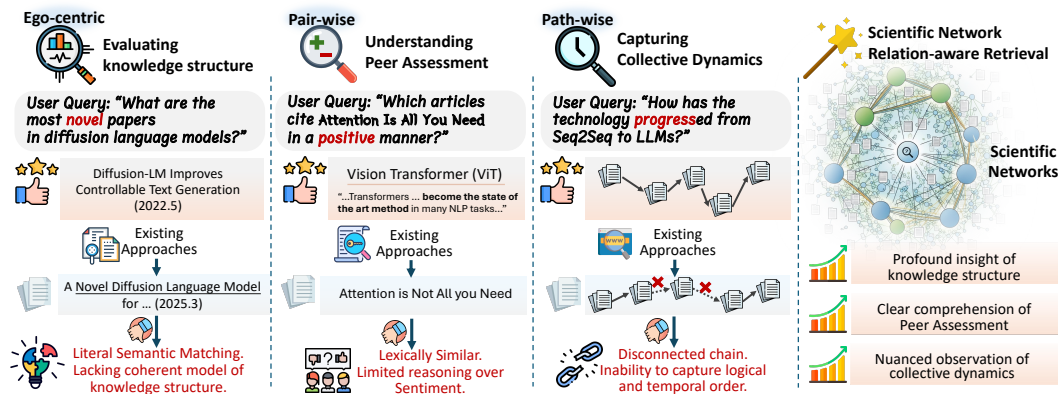


Figure 1: Importance of Scientific Networks in Literature Retrieval Scenarios

1 INTRODUCTION

The rapid development of AI agent has given rise to advanced research tools, such as *Deep Research* (OpenAI, 2025a), fostering the progress of automated scientific systems which are often referred to as AI Scientists (Yamada et al., 2025; Lu et al., 2024). The effective operation of such systems relies on a core capability: high-quality literature retrieval, which is essential to accurately

054 identify related work and position research projects within the existing literature. However, litera-
 055 ture retrieval in scientific research is non-trivial, as it requires accurately understanding the deeper
 056 relations in the scientific network. Figure 1 illustrates three representative retrieval cases across
 057 different scholarly scenarios: *evaluating knowledge structure*, *understanding peer assessment*, and
 058 *capturing collective dynamics*. All these retrieval heavily rely on the scientific network, highlighting
 059 the critical importance of relation-aware retrieval.

060 However, existing retrieval agents generally struggle to achieve relation-aware retrieval, as also il-
 061 lustrated in Figure 1. Specifically, embedding-based agents (Beltagy et al., 2019; Cohan et al.,
 062 2020; Huang et al., 2020) rely solely on static representations of the literature, limiting them to shal-
 063 low semantic matching and preventing multi-step, agentic reasoning. Meanwhile, Deep Research
 064 agents (He et al., 2025; Lála et al., 2023; Skarliniski et al., 2024), despite their iterative pipelines,
 065 cannot explicitly model and fully leverage the relations encoded in the scientific network. Further-
 066 more, existing literature retrieval benchmarks also overlooked the deep relations in the scientific
 067 networks. Most benchmarks focus primarily on semantic precision, evaluating whether retrieved re-
 068 sults are relevant in terms of domain or topic (Ajith et al., 2024; He et al., 2025). Although STARK
 069 introduces the notion of network (Wu et al., 2024), it remains limited to structured hops between en-
 070 tities such as authors, sections, and papers, without addressing the deeper relations that exist directly
 071 between the publications themselves.

072 Motivated by this, we propose **SciNetBench**, the first **Scientific Network Relation-aware**
 073 **Benchmark** to systematically evaluate how well literature retrieval agents capture these relations
 074 within the scientific network. Owing to the need of open-access publications and the fact that agent-
 075 based tools are most prevalent in the discipline of AI, we concentrate our benchmark on AI. Con-
 076 structed from a corpus of 18,639,141 AI papers, SciNetBench cover all major AI subfields. To
 077 capture different levels of relational structures, we design three categories of tasks:

- 078 • **Ego-centric**: Retrieving papers based on individual intrinsic scientific properties, such as iden-
 079 tifying the most novel or disruptive work within a research field.
- 080 • **Pair-wise**: Identifying the specific relational context between two papers, such as whether one
 081 supports, contradicts, or builds upon the other.
- 082 • **Path-wise**: Retrieving citation paths that capture the evolution of scientific ideas, such as
 083 reconstructing the development trajectory from a foundational concept to a SOTA method.

084 We conduct extensive evaluations of eight retrieval methods across three categories on our bench-
 085 mark: (1) retrieval via embedding models, (2) retrieval via agentic models, (3) retrieval via deep
 086 research pipelines. Experimental results demonstrate that current approaches consistently fail across
 087 all three proposed tasks. We further evaluate the impact on downstream applications through *litera-*
 088 *ture review*, demonstrating that the lack of relation-aware retrieval results in incomplete or inaccurate
 089 literature summaries and highlighting the crucial importance of efficient exploitation of the literature
 090 network. Our contributions can be summarized as follows:

- 091 • We propose the **first benchmark** for relational retrieval in scientific literature, constructed
 092 from a corpus of over 18 million AI papers. It provides standardized tasks, queries, and
 093 evaluation metrics to assess capabilities beyond semantic matching.
- 094 • We introduce a **novel taxonomy** of scientific relational retrieval through three granularities:
 095 ego-centric, pair-wise, and path-wise. This taxonomy provides a conceptual framework for
 096 understanding scholarly relationships.
- 097 • Through extensive evaluation of eight retrieval methods, we **quantify the limitations** of these
 098 methods and demonstrate the critical importance of relational retrieval. Additional experi-
 099 ments further validate the benefits of literature network for downstream applications.

100 2 BENCHMARK OVERVIEW AND EVALUATED MODELS

101 2.1 BENCHMARK OVERVIEW

102 *Scientific network* refers to a semantically enriched graph of scholarly publications, constructed from
 103 citation relations among papers and citation contexts within documents, so as to capture not only the
 104 presence of citation links but also their underlying semantic nature. To build such a network, we
 105 first leverage OpenAlex (RELEASE 2025-07-07), a comprehensive open scholarly dataset, to obtain
 106
 107

108 citation relations among scientific papers. We downloaded the complete data snapshot directly from
 109 its official Amazon S3 bucket* , which contains metadata for 269,091,010 papers, including titles,
 110 abstracts, authorship, citations, and publication information. To construct an AI-specific corpus,
 111 we identified 177 AI-related subtopics (e.g., three-dimensional reconstruction, tool-augmented rea-
 112 soning) using OpenAlex’s topic classification system through manual review and aggregation (see
 113 A.2 for the full list). This filtering resulted in a candidate pool of 18,639,140 AI papers, which
 114 to our knowledge represents the largest curated literature base for retrieval benchmarking to date.
 115 Additionally, we incorporated the complete arXiv PDF corpus (as of July 7, 2025) for auxiliary text
 116 extraction and validation.

117 From this corpus, we constructed 1,087 high-quality queries across three relational tasks: 354
 118 (32.6%) for ego-centric retrieval, 600 (55.2%) for pair-wise relation identification, and 133 (12.2%)
 119 for path-wise evolutionary analysis. Queries were first generated through structured rules (e.g.,
 120 leveraging citation and topics) and subsequently validated and refined through manual expert re-
 121 view, ensuring both coverage and reliability. The distribution across tasks reflects their relative
 122 prevalence and importance in real-world scientific reasoning, with pair-wise relations being most
 123 common in scholarly discourse and path-wise trajectories representing more complex, higher-level
 124 reasoning. Our benchmark contains no private or non-public information.

125 2.2 EVALUATED MODELS

126 We evaluate a total of eight retrieval models across three categories:

127
 128 **Category I: Retrieval via Embedding Models: SciBERT** (Beltagy et al., 2019) is the first embed-
 129 ding model specifically trained on scientific literature. It was pre-trained on a corpus of 3.17 billion
 130 tokens, predominantly from the biomedical domain. It also constructed a new WordPiece vocabu-
 131 lary directly from the scientific corpus, which significantly enhanced the model’s representational
 132 capacity. More recently, with the rapid advances in large language models (LLMs), their embedding
 133 layers have also been regarded as reliable embedding models. In our benchmark, we include the
 134 newly released and powerful **Qwen3-8B-Embedding** (Zhang et al., 2025) model in evaluation.

135 **Category II: Retrieval via Agentic Models:** This category encompasses frameworks that employ
 136 agentic workflows for information retrieval and synthesis. We include **paperQA2** (Skarlinski et al.,
 137 2024), a recently released system by FutureHouse designed for high-accuracy, retrieval-augmented
 138 QA over scientific documents. Its agent-driven framework integrates vector retrieval with LLM-
 139 based comprehension, first segmenting the corpus into discrete text chunks and indexing them in-
 140 dividually, then executing a pipeline that includes evidence gathering and answer generation with
 141 explicit citation support. Also selected is **PaSa** (He et al., 2025), which operates through a Crawler
 142 and a Selector. The Crawler autonomously generates search queries from user input, retrieves rel-
 143 evant papers, and iteratively expands the search scope through citation tracking. The Selector then
 144 evaluates the relevance of the retrieved papers to the query. Further included are **gpt-4o-mini-**
 145 **search** and **gpt-4o-search** (OpenAI, 2025b), which leverage LLMs to perform multi-step reason-
 146 ing. These models are equipped with powerful web search tools, enabling them to autonomously
 147 generate queries, retrieve information from diverse online sources, and synthesize responses.

148 **Category III: Retrieval via Deep Research Agents:** We selected **o4-mini-deep-research** and **o3-**
 149 **deep-research** (OpenAI, 2025a). Deep Research agents operate through a deeply iterative pipeline,
 150 in which they autonomously decompose complex queries into sub-tasks and dynamically adapt re-
 151 trieval strategies. They are capable of executing dozens of iterative search-and-reasoning cycles
 152 before the final answer.

153 3 EGO-CENTRIC RELATION: EVALUATING KNOWLEDGE STRUCTURES 154 THROUGH SCIENTOMETRIC INSIGHTS

155 3.1 EVALUATION PROTOCOL

156
 157
 158 **Construction:** Beyond merely finding related papers, deep scientific inquiry often requires eval-
 159 uating the intrinsic scholarly value of individual publications, such as identifying truly novel or
 160 disruptive work. To address this need, we introduce Ego-centric Retrieval, a category focused on
 161

*<https://docs.openalex.org/download-all-data/download-to-your-machine>

162 assessing papers based on their knowledge structure. Literature attributes like novelty often arise
 163 from distinctive configurations in a paper’s knowledge structure, such as the pioneering combina-
 164 tion of concepts from previously disconnected fields. By leveraging scientometric indicators, we
 165 can quantify such structural characteristics to answer queries like, “Which is the most novel paper
 166 in diffusion language models?”, a task beyond semantic matching, as it requires comprehension of
 167 abstract, structure-derived properties.

168 To perform this task, we make use of the collective knowledge embedded in citation networks. For
 169 example, a paper’s citation patterns, how it is cited by later work, can serve as a reliable indicator
 170 of its novelty, and disruptiveness. Building on this idea, we convert two established scientometric
 171 indicators, the *novelty* and the *disruption index*, into concrete retrieval queries. This allows us to
 172 evaluate whether current retrieval models can correctly interpret and respond to queries aimed at
 173 capturing the deeper scientific value of papers.

174 Specifically, we draw on the method proposed
 175 by Uzzi et al. (Uzzi et al., 2013) for measuring
 176 **novelty**: by analyzing co-citation pairs within
 177 a paper’s references, they quantify the extent
 178 to which the work combines rare or “atypical”
 179 knowledge components. This is calculated by
 180 converting each co-citation pair’s frequency into
 181 a Z-score relative to the disciplinary norm, with
 182 the paper’s final novelty score being the 10th per-
 183 centile ($p_{10,z}$) of these scores. A lower score thus
 184 signifies a more novel combination of knowl-
 185 edge. In parallel, we adopt the **disruption in-
 186 dex** introduced by Funk and Owen-Smith (Funk
 187 & Owen-Smith, 2017). This metric evaluates
 188 whether subsequent publications continue to cite
 189 both the focal paper and its predecessors, or in-
 190 stead shift to citing only the focal paper. The in-
 191 dex is calculated as $(N_i - N_j)/(N_i + N_j)$, where N_i is the count of papers citing only the focal work
 192 and N_j is the count citing both the focal work and its references, thereby characterizing whether the
 work extends existing trajectories or fundamentally disrupts prior research.

193 **Evaluation:** For evaluation, each query is used to retrieve a ranked list of 50 candidate papers. We
 194 employ a multi-faceted assessment strategy covering three distinct aspects. First, for our sciento-
 195 metrics indicators, we calculate the raw novelty and disruption scores for each retrieved paper. The
 196 raw novelty score is a Z-score that is theoretically unbounded, where more negative values indicate
 197 higher novelty, while the disruption index ranges from -1 (consolidating) to +1 (disruptive).
 198 Given the distinct and unintuitive scales of these raw scores, we convert each into a percentile rank
 199 against a global reference corpus of millions of papers. The average of these percentile ranks for the
 200 candidates constitutes our final *Novelty-SoS* and *Disruption-SoS* metrics.

201 Second, we incorporate LLMs (GPT-5) to provide complementary semantic judgments, yielding the
 202 *Novelty-LLM* and *Disruption-LLM* metrics. These models assign a score from 1 to 10 for each
 203 concept based solely on the paper’s title and abstract. Third, we establish ground truth labels to
 204 measure retrieval performance. For each of the 177 AI subfields, we identify the top 50 most novel
 205 and top 50 most disruptive papers based on their scientometric ranks, creating two distinct ground
 206 truth sets. Performance is then measured using *Novelty-Recall@50* and *Disruption-Recall@50*,
 207 which evaluate the system’s ability to include these key papers within its top 50 results.

209 3.2 EXPERIMENTAL RESULTS

210
 211 As shown in Table 1, a clear performance hierarchy is evident across all evaluation metrics for
 212 ego-centric retrieval. Deep Research systems (e.g., o3-deep-research) consistently achieved the best
 213 results, followed by web search-based agentic models (e.g., gpt-4o-search), while other models
 214 demonstrated substantially weaker performance. This demonstrates that agentic workflows and flex-
 215 ible use of web tools can effectively enhance performance in literature retrieval. Nevertheless, even
 the top-performing systems struggled significantly according to recall-based evaluation: the best

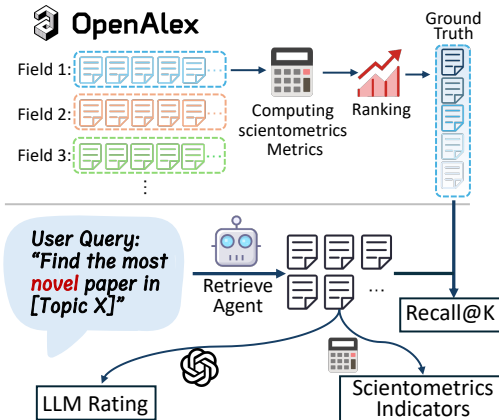


Figure 2: Evaluation Protocol of Ego-Centric Retrieval.

Category	Models	Novelty -SoS	Novelty -LLM	Novelty -Recall@50	Disruption -SoS	Disruption -LLM	Disruption -Recall@50
Embedding	SciBERT	3.686	2.382	0.00%	4.299	2.201	0.00%
	Qwen3-Embed	4.115	2.836	0.20%	3.658	3.045	2.25%
Agentic	PaSa	4.646	5.059	0.10%	6.050	3.831	2.68%
	PaperQA	5.512	5.491	0.17%	5.675	3.970	2.44%
	gpt-4o-mini-search	6.352	6.196	0.60%	6.656	5.874	3.12%
	gpt-4o-search	6.447	6.275	1.10%	6.711	5.968	4.35%
DeepResearch	o4-mini-deep-research	6.815	6.490	1.30%	6.910	7.060	4.60%
	o3-deep-research	6.951	6.491	1.42%	6.956	6.970	3.92%

Table 1: Performance for Frontier Models and Agents on Ego-Centric Tasks

recall@50 for novelty was only 1.42%, and for disruption only 4.60%, indicating that over 95% of truly groundbreaking papers were missed by all systems.

This pattern was consistent across both scientometrics scores and LLM-based assessment, revealing that current retrieval approaches are misaligned with the demands of scientific practice, where accurate assessment of papers based on their intrinsic scientific properties is essential. The observed failures underscore the necessity of developing relation-aware retrieval models capable of understanding scholarly networks and capturing the deeper scientific value of papers.

4 PAIR-WISE RELATION: UNDERSTANDING PEER ASSESSMENT THROUGH CITATION CONTEXTS

4.1 EVALUATION PROTOCOL

Building upon the scientometrics indicators discussed previously, which primarily focus on the statistical properties of individual nodes within the scholarly network, this section extends the analysis to pairwise relations between papers, with particular emphasis on fine-grained semantic associations derived from citation contexts (peer assessment). Accurately identifying such relational information is critical for multiple downstream scientific applications: it enables high-precision literature recommendation by moving beyond topical similarity to capture nuanced scholarly dialogues; it significantly improves the quality of retrieval-augmented generation (RAG) systems by providing evidence chains with explicit sentiment and contextual labels; and it supports the construction of richly structured knowledge graphs that reflect the true discursive landscape of a field, facilitating advanced analyses such as trend detection, controversy mapping, and knowledge gap identification.

To operationalize this focus, we designed a suite of pairwise retrieval tasks encompassing two critical types of scholarly relationships. The first type involves **sentiment-oriented queries**, such as “Which papers cite Paper XX positively?”, requiring systems to distinguish between critical, supportive, or neutral citations based on contextual sentiment. The second type targets **context-based co-mention queries**, exemplified by “Which papers are frequently mentioned together with Paper XX within the same paragraph?”. This task demands the identification of papers jointly referenced within a coherent narrative segment (e.g., a paragraph in the related work section), thereby capturing methodological comparisons, or thematic contrasts within the scientific literature.

Evaluation: Retrieval quality is assessed through four complementary metrics. *Cite-Acc* measures whether a retrieved paper is actually cited by the query paper. *Cite-Sentiment* extends the evaluation by analyzing the sentiment of the citation: each retrieved paper’s PDF is obtained from arXiv and parsed with *GROBID* Lopez & Romary (2013), which provides both the reference list and in-text citation links; verified citations are then traced to their surrounding paragraphs, where GPT-5 determines whether the citation is positive, negative, or neutral. *CoMention-Acc* captures contextual co-citation by checking in the citation network whether a source article cites both the retrieved paper and the query paper. Building on this, *CoMention-Paragraph* requires stronger evidence by parsing the co-citing article with *GROBID* to confirm that both citations not only appear but also co-occur within the same paragraph, thereby ensuring that co-citation evidence is grounded in explicit textual context rather than inferred solely from the network.

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

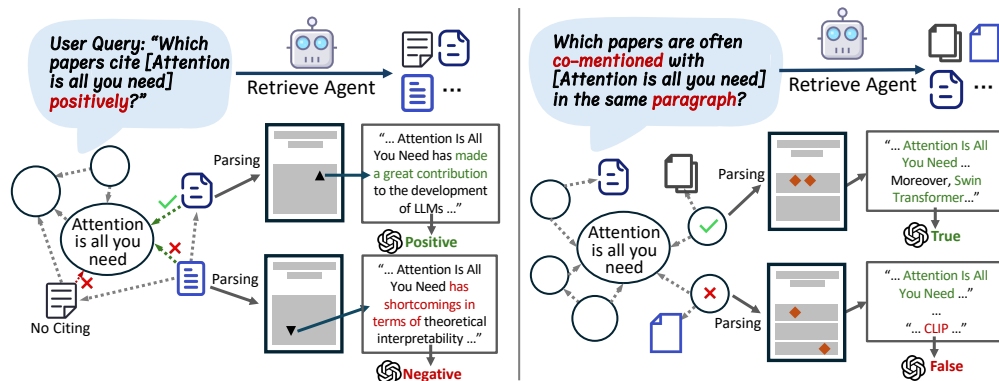


Figure 3: Evaluation Protocol of Pair-Wise Retrieval.

Category	Models	Cite-Acc	Cite-Sentiment	CoMention-Acc	CoMention-Paragraph
Embedding	SciBERT	1.80%	0.00%	0.00%	0.00%
	Qwen3-Embed	10.65%	3.07%	23.95%	5.32%
Agentic	PaSa	18.31%	4.89%	26.31%	7.63%
	PaperQA	8.20%	3.45%	8.25%	7.11%
	gpt-4o-mini-search	46.17%	12.28%	66.04%	6.89%
	gpt-4o-search	46.36%	9.28%	68.82%	7.32%
DeepResearch	o4-mini-deep-research	61.89%	16.77%	68.60%	10.10%
	o3-deep-research	62.84%	13.80%	76.88%	16.30%

Table 2: Performance for Frontier Models and Agents on Pair-Wise Tasks

4.2 EXPERIMENTAL RESULTS

Table 2 presents the performance of different systems on the pair-wise tasks. The results indicate that current agentic retrieval approaches, including both web search tools and reasoning-based agents, provide improvements in retrieval quality, with deep research platforms yielding even more substantial gains. For example, *Cite-Acc* reaches around 46% for agentic models and exceeds 60% for deep research systems, while *CoMention-Acc* can be as high as 77%. These findings suggest that leveraging external search capabilities or reasoning mechanisms enables models to more effectively identify citation links and co-citation patterns compared to purely embedding-based methods.

However, substantial challenges remain. Metrics such as *Cite-Sentiment* and *CoMention-Paragraph* continue to show low performance across all approaches, indicating that capturing citation sentiment and grounding co-mentioned papers within the same paragraph remain difficult. Overall, while advanced retrieval methods enhance citation detection, relational and context-aware reasoning is still far from solved, clearly highlighting the necessity of leveraging the literature network for document retrieval.

5 PATH-WISE RELATION: CAPTURING COLLECTIVE DYNAMICS OF SCIENTIFIC EVOLUTION

5.1 EVALUATION PROTOCOL

Construction: The previously introduced ego-centric and pair-wise tasks assess models’ ability to capture intrinsic properties and binary relations. However, scientific progress typically unfolds as an evolving narrative, where new ideas build upon prior work in multi-step trajectories, forming the collective dynamics of scientific knowledge. To capture this essential aspect, we propose our third category: *Path-Wise Retrieval*, which evaluates whether a system can reconstruct the evolutionary path connecting a sequence of papers. For example, a researcher might ask: “What are the key milestones linking the seminal Transformer paper to today’s large language models?” Answering such queries requires understanding not just paper relevance, but also the logical and citational dependencies that form a coherent developmental chain.

The importance of this task lies in its centrality to literature reviews and research trend analysis. A system that merely outputs unordered related papers cannot reveal the intellectual structure of a field. Yet, existing retrieval methods are almost entirely incapable of constructing such paths, as they lack mechanisms for modeling temporal progression or causal reasoning in scholarly lineage. By introducing the path-wise task, our benchmark offers the first rigorous testbed for evaluating retrieval systems on their ability to reconstruct scientific evolution, pushing them beyond shallow retrieval toward genuine knowledge synthesis.

To construct meaningful technological evolution queries, we leveraged the aforementioned 177 AI subfields. For each subfield, we retrieved the top 50 most-cited papers of all time (treated as classical papers) and the top 10 most-cited papers since 2024 (treated as emerging papers). By randomly pairing classical and emerging papers, we generated a large set of candidate pairs. We then applied the OpenAlex citation network to filter out paper pairs that are topologically connected, followed by manual inspection to ensure that the paired papers remain thematically coherent. This yielded a total of 133 high-quality queries, such as: “What is the most influential citation path from *Attention Is All You Need* to *An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale*?”

Evaluation: For each query, we constructed the ground truth path using a breadth-first search (BFS) algorithm. Specifically, we enumerated all connecting paths between the two endpoint papers and computed the cumulative citation count of all papers along each path. The path with the highest total citation count was selected as the ground truth. This approach is well justified, as it closely parallels the notion of collective credit allocation proposed by Hua-Wei Shen et al. (Shen & Barabási, 2014), where

citations are interpreted as community votes that represent collective recognition of a research trajectory. We evaluated retrieval results along three complementary dimensions. First, **Consistency** measures the degree of overlap between the retrieved path and the ground-truth path. Second, **Connectivity** evaluates whether the retrieved papers form a connected citation path linking the query endpoints within the citation network. Third, **Rationality** measures the plausibility of the retrieved path. We first evaluate this using an LLM, which is prompted with the titles and abstracts of the retrieved papers to judge whether the sequence forms a coherent and reasonable evolutionary narrative. To further enhance reliability, we complement this with a human evaluation: for a subset of 50 representative queries, three AI-expert annotators independently scored the retrieved paths from each model on a 1–10 scale based on rationality. The scores are then averaged to form a **Rationality-Human** metric, which is included as an additional column in the table.

5.2 EXPERIMENTAL RESULTS

Results shown in Table 3 also reveal a pronounced performance gap between traditional embedding-based methods and more advanced retrieval paradigms on the path-wise task. Embedding models such as SciBERT and Qwen3-Embed essentially fail, achieving near-zero **Consistency** and **Connectivity**, and very low scores in both LLM-judged (**Rationality-LLM**) and human-judged (**Rationality-Human**) evaluations. This indicates that these models cannot capture sequential dependencies or reconstruct coherent scientific trajectories beyond surface-level semantic similarity.

In contrast, web search and deep research systems demonstrate substantially stronger performance. Models such as *gpt-4o-search* and *o3-deep-research* achieve over 60% **Consistency** and significantly higher **Connectivity**, successfully retrieving papers that lie along true evolutionary paths. Deep re-

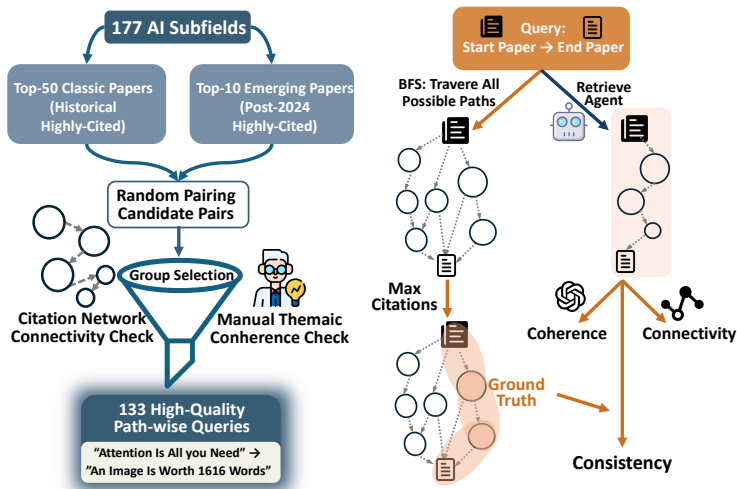


Figure 4: Evaluation Protocol of Path-Wise Retrieval.

Category	Models	Consistency	Connectivity	Rationality-LLM	Rationality-Human
Embedding	SciBERT	0%	0%	1.195	1.60
	Qwen3-Embed	2.11%	3.00%	3.024	2.25
Agentic	PaSa	4.30%	2.1%	2.405	2.50
	PaperQA	5.56%	2.50%	1.988	2.14
	gpt-4o-mini-search	51.76%	7.23%	4.451	4.98
	gpt-4o-search	65.45%	10.24%	4.651	5.25
DeepResearch	o4-mini-deep-research	61.33%	12.00%	6.700	6.65
	o3-deep-research	63.10%	14.00%	6.840	7.04

Table 3: Performance for Frontier Models and Agents on Path-Wise Tasks

search models also outperform others in *Rationality*, with both LLM and human evaluations confirming that the retrieved sequences form more coherent and plausible narratives of scientific progress. These results highlight that reconstructing intellectual lineages requires relation-aware retrieval, and they validate the path-wise task as a rigorous benchmark for assessing advanced retrieval capabilities that go beyond surface-level semantic matching.

6 DEMONSTRATING REAL-WORLD VALUE OF SCINETBENCH VIA DOWNSTREAM APPLICATIONS

To demonstrate the practical value of our benchmark, we conducted a case study on **automatic literature review**. Through this experiment, we aim to illustrate, from an application perspective, the critical importance of relation-aware retrieval.

Experimental Setup: We selected a set of 20 representative queries from our path-wise benchmark, each corresponding to a research path with ground-truth papers and abstracts. For each query, the ground-truth sequence provides a reference trajectory of the evolution of ideas, allowing for systematic evaluation of survey generation methods.

Four different approaches were evaluated:

Base LLM: A LLM that generates surveys solely from the input query and paper abstracts, without any external retrieval or ground-truth information. *Search-enabled LLM:* LLM leverages web-based literature retrieval tools to generate survey reports without access to ground-truth paper sequences. *Deep Research System:* Our proposed system, which explicitly models literature evolution and performs multi-hop retrieval over relational structures. *Base LLM with Ground Truth:*

The same model as above, but provided with the ground-truth paper sequences for each query to assess the upper-bound performance achievable when the full evolution path is known. All methods used identical prompt templates, emphasizing structured survey writing in academic Markdown style, highlighting the progression and connections between papers, and focusing on relevance, completeness, depth, logical flow, and overall usefulness.

Evaluation Metrics: To quantify survey quality, we adopted two complementary protocols: (1) *LLM Automatic scoring:* Each generated report was evaluated along five dimension: *Relevance*, *Completeness*, *Depth*, *Logical Consistency*, and *Usefulness*. Scores were assigned on a scale of 1–10, and aggregated averages across queries were computed for each method (Figure 5). (2) *Human preference ranking:* Three domain experts were asked to comparatively rank the four systems’ outputs for each query (from most to least useful). Table 4 summarizes the distribution of ranks and average scores.

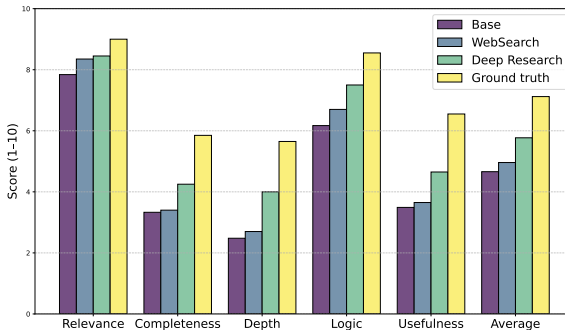


Figure 5: LLM evaluation of survey generation quality.

Results and Analysis: Across both automatic evaluation metrics (Figure 5) and human preference rankings (Table 4), the ground-truth-assisted model attains the highest scores on every evaluated dimension, confirming the clear upper bound when accurate research trajectories are available. Deep-Research system is the strongest baseline: it substantially outperforms the rest of the models, most notably in Completeness (4.25) and Depth (4.00). By contrast, the Base LLM and Search-enabled LLM lag behind; the search-enabled model attains comparatively high relevance but exhibits low completeness and depth, while the base model shows only modest gains in logical coherence. The concordance between automatic evaluation and human rankings indicates that relation-aware retrieval materially improves survey quality, yet the remaining gap in completeness and depth between Deep Research and Ground Truth highlights the need for better modeling of literature relations.

Taken together, these findings underscore that **capturing literature relations is critical for practical downstream applications**. Systems with access to richer relational context produce more coherent, informative, and useful survey reports. Beyond literature review, literature relations can also provide significant benefits

in other applications such as automated experiment design, innovation ideation, and scientific knowledge discovery. This highlights the practical value of our benchmark in supporting the development of systems capable of nuanced scientific reasoning that directly benefits real-world applications.

Method	Avg. Rank ↓	#Rank=1	#Rank=2	#Rank=3	#Rank=4
Ground Truth	1.25	15	4	1	0
Deep Research System	2.00	3	12	4	1
Base LLM	3.10	2	3	10	5
Search-enabled LLM	3.65	0	1	5	14

Table 4: Human preference rankings across 20 queries. Lower average rank indicates better overall preference.

7 RELATED WORKS

Several benchmarks have been proposed to advance scientific literature retrieval. Here, we review representative efforts and highlight how our work emphasizes relational understanding beyond the automatic or entity-centric retrieval.

LitSearch (Ajith et al., 2024) constructs queries using two complementary strategies. *Inline-citation questions* sample paragraphs with citations from research papers, then GPT-4 rewrites them into literature search questions answerable by the cited works. *Author-written questions* are crafted by paper authors, guided by realism, specificity, and resistance to trivial keyword-based resolution.

The PASA (He et al., 2025) benchmark similarly generates queries from *related work* sections using LLMs (e.g., GPT-4o) and expands candidate sets via conventional and academic search engines, search-augmented ChatGPT, and LLM rewriting, with manual expert filtering. Both PASA and LitSearch primarily focus on topical localization, identifying papers in specific domains. In contrast, our benchmark emphasizes reasoning over scholarly relationships, such as methodological influence, disruptive contributions, and conceptual development. STARK (Wu et al., 2024) builds a semi-structured database from the Microsoft Academic Graph, supporting structured knowledge queries that require multi-hop reasoning over predefined entities. Unlike STARK, our benchmark evaluates the discovery of implicit, semantically rich connections among scientific works, providing a more natural testbed for deep scientific reasoning.

8 CONCLUSION

In this paper, we propose a benchmark, **SciNetBench**, for relation-aware retrieval in scientific literature. SciNetBench evaluates retrieval systems across three granularities: **ego-centric** tasks that focus on individual papers’ intrinsic scientific properties, **pair-wise** tasks that assess the relationships between two papers, and **path-wise** tasks that reconstruct citation paths to capture the evolution of scientific ideas. Our experiments demonstrate that current retrieval methods struggle to capture these relational structures, and that this deficiency can substantially degrade downstream applications such as literature review. By emphasizing relational understanding over isolated texts or semantic similarity alone, SciNetBench highlights the practical necessity of integrating scholarly relations, providing a foundation for developing retrieval systems capable of nuanced scientific reasoning and more reliable knowledge synthesis.

486 REPRODUCIBILITY STATEMENT
487

488 We provide all resources and code necessary to use our SciNetBench. All datasets and models
489 employed are fully open-source or publicly accessible, and no privacy or copyright concerns are in-
490 volved. All datasets and models are cited in Section 2.2. Our project code and dataset, including the
491 implementation of evaluation scripts, queries, and metrics, are publicly available via the following
492 anonymous link: <https://anonymous.4open.science/r/SciNetBench/>.
493

494 REFERENCES
495

496 Anirudh Ajith, Mengzhou Xia, Alexis Chevalier, Tanya Goyal, Danqi Chen, and Tianyu Gao. Lit-
497 search: A retrieval benchmark for scientific literature search. *arXiv preprint arXiv:2407.18940*,
498 2024.

499 Iz Beltagy, Kyle Lo, and Arman Cohan. SciBERT: A pretrained language model for scientific text.
500 In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (eds.), *Proceedings of the 2019 Con-
501 ference on Empirical Methods in Natural Language Processing and the 9th International Joint
502 Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 3615–3620, Hong Kong,
503 China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1371.
504 URL <https://aclanthology.org/D19-1371/>.
505

506 Arman Cohan, Sergey Feldman, Iz Beltagy, Doug Downey, and Daniel Weld. SPECTER:
507 Document-level representation learning using citation-informed transformers. In Dan Jurafsky,
508 Joyce Chai, Natalie Schluter, and Joel Tetreault (eds.), *Proceedings of the 58th Annual Meet-
509 ing of the Association for Computational Linguistics*, pp. 2270–2282, Online, July 2020. As-
510 sociation for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.207. URL <https://aclanthology.org/2020.acl-main.207/>.
511

512 Russell J Funk and Jason Owen-Smith. A dynamic network measure of technological change. *Man-
513 agement science*, 63(3):791–817, 2017.

514 Yichen He, Guanhua Huang, Peiyuan Feng, Yuan Lin, Yuchen Zhang, Hang Li, et al. Pasa: An llm
515 agent for comprehensive academic paper search. *arXiv preprint arXiv:2501.10120*, 2025.
516

517 Jui-Ting Huang, Ashish Sharma, Shuying Sun, Li Xia, David Zhang, Philip Pronin, Janani Pad-
518 manabhan, Giuseppe Ottaviano, and Linjun Yang. Embedding-based retrieval in facebook
519 search. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge
520 Discovery & Data Mining, KDD '20*, pp. 2553–2561, New York, NY, USA, 2020. Associa-
521 tion for Computing Machinery. ISBN 9781450379984. doi: 10.1145/3394486.3403305. URL
522 <https://doi.org/10.1145/3394486.3403305>.

523 Pierre Lopez and Laurent Romary. Grobid: Combining automatic bibliographic data recognition
524 and term extraction for scholarship publications. In *Proceedings of the 13th ACM/IEEE-CS joint
525 conference on Digital libraries*, pp. 343–344. ACM, 2013.
526

527 Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The ai scien-
528 tist: Towards fully automated open-ended scientific discovery. *arXiv preprint arXiv:2408.06292*,
529 2024.

530 Jakub Lála, Odhran O’Donoghue, Aleksandar Shtedritski, Sam Cox, Samuel G. Rodrigues, and
531 Andrew D. White. Paperqa: Retrieval-augmented generative agent for scientific research. *arXiv
532 preprint arXiv:2312.07559*, 2023. URL <https://doi.org/10.48550/arXiv.2312.07559>.
533

534 OpenAI. Deepresearch models: o4-mini-deep-research and o3-deep-research, 2025a.
535 <https://platform.openai.com/docs/models/deep-research>.

536 OpenAI. Web-enhanced large language model (llm) search systems: gpt-4o-mini-search and gpt-
537 4o-search, 2025b. <https://platform.openai.com/docs/models/gpt-4o>.
538

539 Hua-Wei Shen and Albert-László Barabási. Collective credit allocation in science. *Proceedings of
the National Academy of Sciences*, 111(34):12325–12330, 2014.

540 Michael D. Skarlinski, Sam Cox, Jon M. Laurent, James D. Braza, Michaela Hinks, Michael J.
541 Hammerling, Manvitha Ponnampati, Samuel G. Rodrigues, and Andrew D. White. Language agents
542 achieve superhuman synthesis of scientific knowledge. *arXiv preprint arXiv:2409.13740*, 2024.
543 URL <https://doi.org/10.48550/arXiv.2409.13740>.
544

545 Brian Uzzi, Satyam Mukherjee, Michael Stringer, and Ben Jones. Atypical combinations and scien-
546 tific impact. *Science*, 342(6157):468–472, 2013.

547 Shirley Wu, Shiyu Zhao, Michihiro Yasunaga, Kexin Huang, Kaidi Cao, Qian Huang, Vassilis N
548 Ioannidis, Karthik Subbian, James Zou, and Jure Leskovec. Stark: Benchmarking llm retrieval
549 on textual and relational knowledge bases. *Advances in Neural Information Processing Systems*,
550 37:127129–127153, 2024.

551 Yutaro Yamada, Robert Tjarko Lange, Cong Lu, Shengran Hu, Chris Lu, Jakob Foerster, Jeff Clune,
552 and David Ha. The ai scientist-v2: Workshop-level automated scientific discovery via agentic tree
553 search. *arXiv preprint arXiv:2504.08066*, 2025.
554

555 Yanzhao Zhang, Mingxin Li, Dingkun Long, Xin Zhang, Huan Lin, Baosong Yang, Pengjun Xie,
556 An Yang, Dayiheng Liu, Junyang Lin, et al. Qwen3 embedding: Advancing text embedding and
557 reranking through foundation models. *arXiv preprint arXiv:2506.05176*, 2025.
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593

A APPENDIX

A.1 EXPERIMENTAL DETAILS

For both SciBERT and Qwen3-Embedding Model, we processed a total of 18,639,140 papers individually, concatenating each paper’s title and abstract into a single string in the format “title: <title>, abstract: <abstract>”. We then obtained embeddings for each concatenated string using the respective embedding model. The embedding dimensionality is 768 for SciBERT and 4096 for Qwen3-Embed.

To enable efficient large-scale similarity search over these high-dimensional embeddings, we constructed a dedicated index using the Faiss library. The preprocessing pipeline first performs L2 normalization on all embeddings, which ensures that maximum inner-product search is equivalent to finding the highest cosine similarity, a standard metric for semantic relevance. We employ an IndexIVFPQ structure, which combines an inverted file system (IVF) for coarse partitioning of the search space and product quantization (PQ) for compact vector representation. Specifically, the algorithm first partitions the entire vector space into 16,384 cells using k-means clustering, where each cell is represented by a centroid vector. This IVF structure enables a substantial pruning of the search space by only examining a small subset of cells closest to the query vector. Within each cell, vectors are further compressed using product quantization: each vector is split into 32 sub-vectors, and each sub-vector is quantized into an 8-bit code pointing to the nearest centroid in a learned codebook.

This two-level scheme, coarse quantization via IVF followed by fine-grained PQ, significantly reduces memory footprint while accelerating retrieval, achieving a favorable trade-off between search accuracy and efficiency. All embeddings and queries are computed on a single NVIDIA A100 GPU.

For PASA and PaperQA2, we strictly followed the implementations provided in their respective GitHub repositories. PASA was deployed on a local A100 GPU using its pretrained pasa-7b-crawler and pasa-7b-selector models, with API keys configured for Google Search and other relevant tools. For PaperQA2, we performed segmentation and embedding on all papers to fully leverage the model’s capabilities. For all OpenAI models, we accessed them using official API keys from the OpenAI platform.

A.2 AI RESEARCH FIELDS AND SUBFIELDS

- **Reinforcement Learning**

- Model-Free Reinforcement Learning
- Model-Based Reinforcement Learning
- Offline Reinforcement Learning
- Hierarchical Reinforcement Learning
- Inverse Reinforcement Learning
- Safe Reinforcement Learning
- Multi-Agent Reinforcement Learning
- Reinforcement Learning with Function Approximation
- Reinforcement Learning for Control and Robotics

- **Bandits and Multi-Armed Bandits**

- Stochastic Bandits
- Contextual Bandits
- Linear Bandits
- Combinatorial Bandits
- Bandits with Knapsacks
- Pure Exploration and Best-Arm Identification

- **Imitation Learning**

- Behavior Cloning
- Inverse Reinforcement Learning
- Generative Adversarial Imitation Learning
- Dataset Aggregation
- Offline Imitation Learning

- **Self-Supervised Learning**

- Contrastive Learning
- Masked Autoencoders
- Predictive Coding
- Bootstrap Your Own Latents
- Clustering-based Self-Supervised Learning

- 648 • **Contrastive Learning**
- 649 – Information Noise Contrastive Estimation Loss and Mutual Information Estimation
- 650 – Supervised Contrastive Learning
- 651 – Hard Negative Mining
- 652 – Cross-Modal Contrastive Learning
- 653 • **Federated Learning**
- 654 – Federated Averaging
- 655 – Personalized Federated Learning
- 656 – Privacy-Preserving Federated Learning
- 657 – Communication-Efficient Federated Learning
- 658 – Federated Learning with Heterogeneous or Non-IID Data
- 659 • **Continual Learning and Lifelong Learning**
- 660 – Regularization-based Methods
- 661 – Replay-based Methods
- 662 – Parameter Isolation Methods
- 663 – Task-Free and Online Continual Learning
- 664 – Concept Drift Adaptation
- 665 • **Online Learning**
- 666 – Online Gradient Descent
- 667 – Online Convex Optimization
- 668 – Online-to-Batch Conversion
- 669 – Adaptive Learning Rates
- 670 • **Multi-Task Learning**
- 671 – Hard Parameter Sharing
- 672 – Soft Parameter Sharing
- 673 – Task Relation Learning
- 674 – Multi-Objective Optimization
- 675 • **Active Learning**
- 676 – Uncertainty Sampling
- 677 – Query by Committee
- 678 – Expected Model Change
- 679 – Bayesian Active Learning
- 680 • **Transfer Learning**
- 681 – Feature-based Transfer
- 682 – Fine-Tuning Pretrained Models
- 683 – Domain Adaptation
- 684 – Multi-Domain Training
- 685 • **Domain Adaptation and Generalization**
- 686 – Unsupervised Domain Adaptation
- 687 – Semi-Supervised Domain Adaptation
- 688 – Domain Generalization
- 689 – Adversarial Domain Adaptation
- 690 • **Meta-Learning**
- 691 – Gradient-based Meta-Learning
- 692 – Metric-based Meta-Learning
- 693 – Black-box Meta-Learning
- 694 – Meta-Optimization
- 695 • **Human-in-the-Loop and Human Feedback**
- 696 – Interactive Labeling
- 697 – Preference Learning
- 698 – Human-Guided Model Editing
- 699 • **Reinforcement Learning from Human Feedback**
- 700 – Reward Model Training
- 701 – Preference Elicitation
- 702 – Policy Optimization with Human Feedback
- 703 – Direct Preference Optimization
- 704 • **Generative Models**
- 705 – Autoregressive Models
- 706 – Normalizing Flows
- 707 – Energy-Based Models
- 708 – Score-Based Models
- 709 • **Diffusion Models**
- 710 – Denoising Diffusion Probabilistic Models
- 711 – Latent Diffusion Models
- 712 – Score Matching and Stochastic Differential Equation Approaches
- 713 – Guided Diffusion for Text-to-Image Generation
- 714 • **Graph Neural Networks**
- 715 – Graph Convolutional Networks
- 716 – Graph Attention Networks
- 717 – Graph Isomorphism Networks
- 718 – Spatio-Temporal Graph Neural Networks
- 719 – Heterogeneous Graph Learning
- 720 • **Transformers**

- 702 – Encoder-Only Transformers
- 703 – Decoder-Only Transformers
- 704 – Encoder-Decoder Transformers
- 705 • **Deep Learning Theory**
- 706 – Generalization Bounds
- 707 – Double Descent Phenomenon
- 708 – Neural Tangent Kernel Theory
- 709 – Implicit Bias in Gradient Descent
- 710 • **Natural Language Processing**
- 711 – Language Modeling
- 712 – Text Classification
- 713 – Named Entity Recognition
- 714 – Question Answering
- 715 – Information Extraction
- 716 – Machine Translation
- 717 – Text Summarization
- 718 – Dialogue Systems
- 719 • **Large Language Models**
- 720 – Pretraining Objectives
- 721 – Scaling Laws
- 722 – Fine-Tuning Methods
- 723 – Alignment Methods
- 724 – Evaluation and Benchmarking of Language Models
- 725 • **Large Language Model Agents**
- 726 – Tool-Augmented Reasoning
- 727 – Planning and Task Decomposition
- 728 – Memory-Augmented Agents
- 729 – Function Calling and API Orchestration
- 730 • **Prompt Engineering and In-Context Learning**
- 731 – Zero-Shot Prompting
- 732 – Few-Shot Prompting
- 733 – Chain-of-Thought Reasoning
- 734 – Automatic Prompt Search
- 735 • **Computer Vision**
- 736 – Image Classification
- 737 – Object Detection
- 738 – Semantic and Instance Segmentation
- 739 – Pose Estimation
- 740 – Three-Dimensional Reconstruction
- 741 • **Vision-Language Models**
- 742 – Contrastive Pretraining
- 743 – Image Captioning
- 744 – Visual Question Answering
- 745 – Grounded Language Understanding
- 746 • **Time Series Analysis and Forecasting**
- 747 – Autoregressive Models for Time Series
- 748 – Neural Forecasting Models
- 749 – Multivariate Time Series Analysis
- 750 – Anomaly Detection in Time Series
- 751 • **Autonomous Driving and Robotics**
- 752 – Perception for Autonomous Vehicles
- 753 – Path Planning and Control
- 754 – Simulation-to-Real Transfer
- 755 – Multi-Sensor Fusion
- 756 – Learning from Demonstration
- 757 • **Explainable Artificial Intelligence and Interpretability**
- 758 – Feature Attribution
- 759 – Counterfactual Explanations
- 760 – Concept-based Explanations
- 761 – Causal Interpretability
- 762 • **Adversarial Machine Learning**
- 763 – White-Box and Black-Box Attacks
- 764 – Adversarial Training
- 765 – Certified Robustness
- 766 – Physical Adversarial Attacks
- 767 • **Fairness, Accountability, and Transparency**
- 768 – Fairness Metrics
- 769 – Bias Mitigation Techniques
- 770 – Transparent Reporting
- 771 – Algorithmic Auditing
- 772 • **Privacy-Preserving Machine Learning**
- 773 – Differential Privacy
- 774 – Federated Privacy
- 775 – Homomorphic Encryption
- 776 – Secure Multi-Party Computation
- 777 • **Artificial Intelligence Safety and Reliability**

- 756 – Robustness Evaluation
- 757 – Out-of-Distribution Detection
- 758 • **Optimization for Machine Learning**
- 759 – Stochastic Gradient Descent and Variants
- 760 – Adaptive Optimizers
- 761 • **Model and System Efficiency**
- 762 – Model Pruning
- 763 – Quantization
- 764 – Knowledge Distillation
- 765 • **Causal Inference and Discovery**
- 766 – Structural Causal Models
- 767 – Causal Representation Learning
- 768 • **Artificial Intelligence for Science**
- 769 – Molecular Simulation
- 770 – Protein Folding
- 771 • **Artificial Intelligence for Drug Discovery and Healthcare**
- 772 – Drug-Target Interaction Prediction
- 773 – Molecule Generation
- 774 • **Physics-Informed Machine Learning**
- 775 – Physics-Informed Neural Networks
- 776 – Differentiable Simulators
- 777 • **Artificial Intelligence for Social Good**
- 778 – Disaster Response
- 779 – Public Health Modeling
- 780 – Alignment Research
- 781 – Verification and Formal Methods
- 782 – Second-Order Optimization Methods
- 783 – Non-Convex Optimization
- 784 – Neural Architecture Search
- 785 – Low-Rank Approximations
- 786 – Invariant Risk Minimization
- 787 – Counterfactual Reasoning
- 788 – Climate Modeling
- 789 – Material Science
- 790 – Medical Image Analysis
- 791 – Clinical Decision Support
- 792 – Scientific Machine Learning
- 793 – Education and Accessibility
- 794 – Environmental Monitoring

782 A.3 PROMPTS

783 The LLM prompts for novelty evaluation in the Ego-Centric task are as follows:

```

787 You are an expert academic reviewer. Your task is to evaluate a
788 scientific paper on its Novelty.
789
790 ### Definition of Novelty:
791 Definition: Novelty refers to the uniqueness and originality of the
792 research question, methodology, data, or conclusions relative
793 to existing research.
794 Focus: Does the paper introduce new ideas, perspectives, or methods
795 within the existing body of knowledge? For example, applying a
796 method from Field A to Field B for the first time.
797 Scoring Criteria: A score of 0 represents completely derivative
798 work, while a score of 10 represents a highly original and
799 groundbreaking idea.
800
801 ---
802
803 Please evaluate the Novelty of the following paper based on the
804 definition provided.
805
806 **Title**: [Paper Title]
807
808 **Abstract**: [Paper Abstract]
809 (or: [No abstract provided. Please evaluate based on the title
alone.])

```

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

```
Your response MUST be a single JSON object with 'score' (an integer from 0 to 10) and 'reasoning' (a brief explanation).
```

The LLM prompts for disruptiveness evaluation in the Ego-Centric task are as follows:

```
You are an expert academic reviewer. Your task is to evaluate a scientific paper on its Disruptiveness.

### Definition of Disruptiveness:
Definition: Disruptiveness refers to the way a paper influences subsequent research-does it cause future work to cite the paper itself, rather than the previous works it was built upon?
Focus: Does the paper change the direction of a research field or its methodologies, causing prior work to be marginalized? For example, the foundational papers on CRISPR gene-editing technology opened new research avenues and made previous editing methods obsolete.
Scoring Criteria: A score of 0 represents no disruptive potential (e.g., a review paper), while a score of 10 represents the potential to highly transform a field.

---

Please evaluate the Disruptiveness of the following paper based on the definition provided.

Title: [Paper Title]

Abstract: [Paper Abstract]
(or: [No abstract provided. Please evaluate based on the title alone.])

---

Your response MUST be a single JSON object with:
- "score": an integer from 0 to 10
- "reasoning": a brief explanation supporting the score
```

In the Pair-Wise task, the code for classifying the sentiment tendency of citation context using LLM is as follows:

```
You are an expert in academic literature analysis. Your task is to classify the sentiment of a citation context.

Please analyze the following citation context from a research paper, which mentions the target paper titled "[Target Paper Title]".

Classify the context into one of three categories:
- Positive: The citing paper praises, builds upon, or confirms the findings of the target paper.
- Negative: The citing paper criticizes, questions, or points out limitations of the target paper.
- Neutral: The citing paper simply mentions or describes the target paper as background information without expressing a strong opinion.
```

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

```
Your response MUST BE only ONE of the three category names:  
Positive, Negative, or Neutral.  
  
Context to analyze:  
---  
[Insert citation context here]  
---
```

The prompts for LLM evaluation in the Path-Wise task are as follows:

```
You are an expert in scientometrics and academic research. Your  
task is to evaluate the quality of a proposed citation path  
based on its technical evolution.  
  
### Core Task:  
Assess if the provided sequence of papers represents a logical and  
meaningful technological or conceptual evolution from the start  
paper to the end paper.  
  
### What constitutes a good technical evolution path? (Key  
Principles)  
1. Thematic Consistency: All papers must strictly revolve around  
the same core research topic defined by the query. Deviations  
into unrelated subjects indicate a poor path.  
2. Content Cohesion & Logical Flow: The content of adjacent papers  
must be closely related. Each paper should logically follow  
from the previous one, building upon its ideas, refining its  
methods, or addressing its limitations.  
3. Progressive Development: The path must demonstrate clear  
progress. Later papers should represent advancements,  
extensions, or significant new applications of the concepts  
introduced in earlier papers. The path should tell a story of  
innovation.  
4. Represents a Main Line of Inquiry: The path should follow a  
significant and recognized line of development within the  
research field, not an obscure or tangential branch.  
  
### Scoring Criteria (0-10):  
- Score 9-10 (Excellent): A perfect or near-perfect path. It is  
thematically consistent, shows clear progressive development,  
and represents a major line of inquiry. The logical flow is  
impeccable.  
- Score 7-8 (Good): A strong, coherent path. Most papers are  
relevant and show progression, but there might be a minor  
logical gap or a less influential paper included.  
- Score 4-6 (Mediocre): The path has some relevance but lacks  
strong cohesion. It may include several tangential papers, the  
logical progression is weak, or it fails to capture the main  
developmental thread.  
- Score 1-3 (Poor): The path is largely incoherent. Papers are  
thematically disconnected, show no clear progress, or are  
mostly irrelevant to the query.  
- Score 0 (Failure): A completely random collection of papers with  
no logical or thematic connection.  
  
---  
  
Please evaluate the following citation path based on the detailed  
criteria provided.
```

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

```
Original Request: [original_query]

--- Proposed Citation Path ---
[Paper list will be inserted here]
-----

Your response MUST be a single JSON object with 'score' (an integer
  from 0 to 10) and 'reasoning' (a detailed explanation for your
  score, critiquing the path based on the four key principles).
```

In the survey generation experiment 6, the prompts for generating the survey are as follows:

```
You are a helpful academic assistant that generates surveys using
retrieved literature.

Please generate a survey report in Markdown format based on the
following information:

Domain: [Domain Name]

Start paper:
Title: "[Start Paper Title]"
Abstract: [Start Paper Abstract]

End paper:
Title: "[End Paper Title]"
Abstract: [End Paper Abstract]

Task:
Generate a survey report describing the technological evolution
  from the start paper to the end paper.
Include key technical developments, major milestones, and method
  evolution.
Organize the report in a clear Markdown format.

Important:
- Do not use any ground-truth paths.
- Rely only on information retrieved via search capabilities.
```