

Chain of Retrieval: Multi-Aspect Iterative Search Expansion and Post-Order Search Aggregation for Full Paper Retrieval

Anonymous ACL submission

Abstract

Scientific paper retrieval, particularly framed as document-to-document retrieval, aims to identify relevant papers in response to a long-form query paper, rather than a short query string. Previous approaches to this task have focused exclusively on abstracts, embedding them into dense vectors as surrogates for full documents and calculating similarity between them. Yet, abstracts offer only sparse and high-level summaries, and such methods primarily optimize one-to-one similarity, overlooking the dynamic relations that emerge across relevant papers during the retrieval process. To address this, we propose Chain of Retrieval (CoR), a novel iterative framework for full-paper retrieval. Specifically, CoR decomposes each query paper into multiple aspect-specific views, matches them against segmented candidate papers, and iteratively expands the search by promoting top-ranked results as new queries, thereby forming a tree-structured retrieval process. The resulting retrieval tree is then aggregated in a post-order manner: descendants are first combined at the query level, then recursively merged with their parent nodes, to capture hierarchical relations across iterations. To validate this, we present SCIFULLBENCH, a large-scale benchmark providing both complete and segmented contexts of full papers for queries and candidates, and results show that CoR significantly outperforms existing retrieval baselines.

1 Introduction

Information Retrieval (IR) is the task of searching for query-relevant documents from a large external corpus, evolving from sparse keyword matching (Sparck Jones, 1972; Robertson et al., 1995) to dense representation-based similarity (Karpukhin et al., 2020; Izacard et al., 2021). Notably, in the era of Large Language Models (LLMs) (Achiam et al., 2023; Team et al., 2023; Dubey et al., 2024; DeepSeek-AI et al., 2025), IR has become increas-

ingly important, which allows LLMs to utilize up-to-date external information (Lewis et al., 2020).

In contrast to conventional retrieval tasks, whose queries are short (such as questions or keywords), scientific paper retrieval poses unique challenges. Specifically, queries are structured, long-form documents that encapsulate diverse aspects, ranging from research motivation and proposed methodology to experimental design and empirical findings. Also, relevance in this setting is inherently multi-faceted, as a paper may be considered relevant for various reasons, such as pursuing similar research objectives or employing comparable methods.

However, previous work has focused on abstract-level representations as a proxy for full papers (Cohan et al., 2020; Yasunaga et al., 2022; Ostendorff et al., 2022; Singh et al., 2023; Zhang et al., 2023a), fine-tuning models on citation-linked pairs of papers, each represented by its abstract and title. While they are effective to some extent, abstracts are inherently sparse and offer only high-level summaries, making them insufficient to capture the nuanced, multi-faceted relationships between scientific works. As a result, abstract-based approaches are limited to shallow similarity and struggle with tasks requiring deeper contextual understanding of full papers beyond surface-level abstracts, such as identifying complementary works, synthesizing literature reviews, or generating novel research directions (Asai et al., 2024; Baek et al., 2024; Chamoun et al., 2024; Jin et al., 2024; D’Arcy et al., 2024).

Nevertheless, it is non-trivial to represent and retrieve full scientific papers. On the one hand, full papers often exceed 100K tokens in length, far surpassing the context limits of most embedding models. Even if encoded into a single vector, representing its diverse aspects (e.g., motivation, methods, and experiments) within a single embedding may lead to oversimplification and blur fine-grained distinctions, particularly when relevance is tied to a specific aspect. Moreover, existing approaches typi-

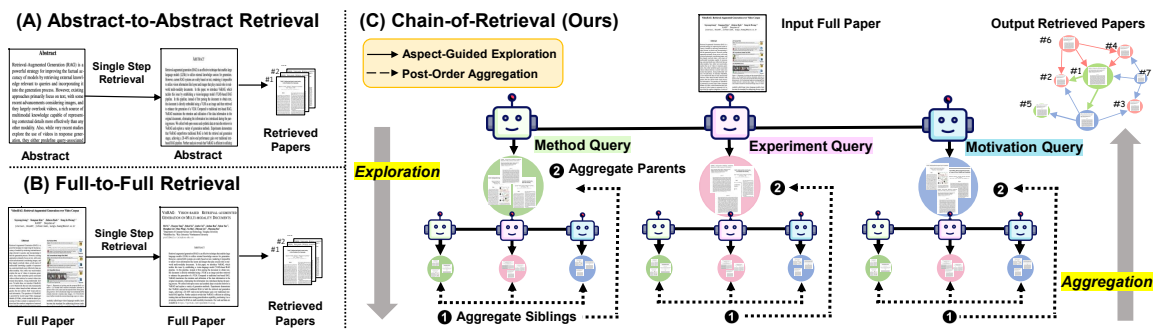


Figure 1: Conceptual illustration of our retrieval method (C) compared to prior retrieval approaches (A and B).

084 cally optimize semantic similarity between isolated
 085 pairs of papers, overlooking the broader, collec-
 086 tive relations that emerge across relevant works
 087 during retrieval. This design stands in contrast to
 088 human information-seeking behavior in the infor-
 089 mation foraging theory (Pirolli and Card, 1995,
 090 1999), which is dynamic and iterative: humans con-
 091 tinuously refine their understanding by comparing
 092 newly encountered works with prior knowledge.

093 To address this, we propose Chain of Retrieval
 094 (COR), a novel framework designed to capture dy-
 095 namic inter-paper relevance while leveraging the
 096 full paper context, operated through a two-phase
 097 process: *Exploration* and *Aggregation*, illustrated
 098 in Figure 1. In the *Exploration* phase, the query
 099 paper is decomposed into multiple aspect-specific
 100 views, and each aspect is handled by a specialized
 101 query optimizer, which is further trained via Direct
 102 Preference Optimization (DPO) (Rafailov et al.,
 103 2023) on a self-generated preference to mitigate the
 104 shortage of realistic supervision for multi-retrieval
 105 systems. After retrieving with these aspect-centric
 106 queries, the retrieved (non-duplicate) papers most
 107 semantically aligned with the original query are
 108 promoted for the next iteration, recursively expand-
 109 ing the search with a tree structure.

110 In the *Aggregation* phase, results across different
 111 depths are merged to capture higher-order relations
 112 among retrieved works. Specifically, inspired by
 113 the notion of triadic closure in the network theo-
 114 ry (Granovetter, 1973), we design an algorithm
 115 that first merges top- k results among sibling nodes
 116 (papers retrieved from the same parent), then com-
 117 bines these aggregated results with those of their
 118 parent node, and continues this process bottom-up
 119 until convergence at the root. This hierarchical ag-
 120 gregation reinforces strong multi-hop connections
 121 while naturally attenuating weaker, indirect signals,
 122 yielding a robust ranking of relevant papers.

123 We evaluate COR by extending existing paper re-

124 trieval benchmarks, as they are originally designed
 125 for abstract-only retrieval, adapting them to include
 126 full papers and more recent publications (to miti-
 127 gate the potential model contamination) in the ML
 128 and NLP domains, which we refer to as SCIFULL-
 129 BENCH. Across experiments on SCIFULLBENCH,
 130 COR consistently outperforms retrieval baselines
 131 (that either rely on abstracts or naively encode full
 132 papers), while remaining compatible with diverse,
 133 domain-agnostic embedding models, highlighting
 134 both the robustness and generality of our COR.

2 Related Work 135

136 **Scientific Paper Retrieval** Scientific paper re-
 137 trieval aims to identify relevant works given a query,
 138 which may take the form of short queries (Ajith
 139 et al., 2024), research proposals (Garikaparthi et al.,
 140 2025), or full papers, and the latter being the focus
 141 of this work. Pioneering work in scientific paper-
 142 to-paper retrieval used bibliometric statistics, such
 143 as influence scores (Zhou et al., 2012; Mohapatra
 144 et al., 2019), co-citations (Small, 1973; Haruna
 145 et al., 2018), and bibliographic coupling (Kessler,
 146 1963). Recently, thanks to the capability of neural
 147 models, many studies have focused on calculating
 148 semantic similarities between abstracts of respec-
 149 tive documents (Bhagavatula et al., 2018; Osten-
 150 dorff, 2020), with the embedding models optimized
 151 to this domain. For example, Cohan et al. (2020)
 152 and Ostendorff et al. (2022) fine-tune BERT-based
 153 models (Devlin et al., 2019) using abstract pairs
 154 extracted from the citation graph, and Mysore et al.
 155 (2022) further consider the Wasserstein distance
 156 between sentence segments within abstract pairs.
 157 Moreover, recent studies target multiple tasks (such
 158 as paper classification and citation prediction in ad-
 159 dition to paper retrieval) in a unified framework by
 160 considering their respective representations (Singh
 161 et al., 2023; Zhang et al., 2023a). Complementary
 162 to these, hybrid approaches integrate bibliomet-

ric signals with vector-based retrieval, reweighting ranks via post-hoc adjustments (Hamedani et al., 2016; Guo et al., 2022). However, such statistical indicators (e.g., citations) primarily reflect long-term scholarly accumulation and tend to capture retrospective importance rather than the contextual information expressed in the documents themselves. In contrast, we model paper relations dynamically during inference using aspect-driven queries generated from the full content of scientific papers.

Query Optimization with LLMs Effective formulation of queries is central to retrieval performance, and has long been studied from classical relevance feedback approaches (Rocchio Jr, 1971; Salton and Buckley, 1990) to query expansion techniques (Kuzi et al., 2016; Nogueira et al., 2019). More recent methods either leverage LLMs themselves to reformulate queries (Yu et al., 2023a; Wang et al., 2023; Gao et al., 2023), or further augment them with external query-relevant information retrieved in an auxiliary step (Yu et al., 2023b; Shen et al., 2024; Park and Lee, 2024; Lei et al., 2024). Another line of work seeks to reduce query ambiguity by decomposing a complex query into smaller queries (Zheng et al., 2024; Korikov et al., 2024). We note that our approach aligns with this decomposition paradigm but extends it to a more challenging setting: instead of handling short, user-issued queries, we decompose long-form scientific documents into multiple aspect-specific subqueries.

3 Methodology

We define the paper retrieval problem and present Chain of Retrieval (COR), illustrated in Figure 2.

3.1 Preliminary

Paper-to-Paper Retrieval Given a query paper D (with its abstract as D_{abstract}), the goal of paper retrieval is to return a ranked list of relevant papers from the corpus \mathcal{C} . In contrast to existing studies that use D_{abstract} , we utilize its complete version D for query formulation and corpus construction.

3.2 Aspect-Aware Multi-Vector Retrieval

Aspect-Aware Query Optimization Scientific papers encapsulate multiple facets, such as research motivation, method design, and experimental validation following Baek et al. (2024) and Moussa et al. (2025), which might not be represented by a single query. To capture this diversity, we formulate the query optimization process that transforms

the full paper D into a set of aspect-specific queries. Formally, we define a family of query optimization functions: $\mathcal{F} = \{f_R, f_M, f_E\}$, where each function f maps D to a query $q = f(D)$ that targets a specific aspect. Notably, each function is instantiated with an LLM guided by an aspect-specific template (see Appendix I). In addition to these fine-grained queries, we consider the abstract D_{abstract} , since it reflects a broad view of the overall content (that can complement specialized views), yielding the final query set: $\mathcal{Q} = \{f_i(D) \mid f_i \in \mathcal{F}\} \cup \{D_{\text{abstract}}\}$.

Retrieval with Multi-View Corpora Once the optimized query set \mathcal{Q} is formulated, we perform retrieval individually for each query $q \in \mathcal{Q}$. Specifically, for the abstract-based query D_{abstract} , we use a corpus $\mathcal{C}_{\text{abstract}}$ that contains candidate abstracts (as they are comparable in length and structure). For the other aspect-specific queries derived from the full paper, we use a segmented version of the full corpus: $\mathcal{C}_{\text{chunked}}$, where each paper is split into fixed-length chunks (likely to capture its specific aspect) and indexed using multi-vector representations (Khattab and Zaharia, 2020; Santhanam et al., 2021). Each query then retrieves its top- k relevant segments, which are mapped back to their source papers via the mapping function $h(x)$, as follows: $\mathcal{R}_q = \{h(x) \mid x \in \text{Retrieve}(q, \mathcal{C}, k)\}$. Finally, the resulting ranked lists of candidate papers from all queries are aggregated to form the multi-view retrieval pool, denoted as follows: $\mathcal{R} = \bigcup_{q \in \mathcal{Q}} \mathcal{R}_q$.

3.3 Iterative Chain-of-Retrieval

Iterative Aspect-Aware Search Expansion The exploration phase aims to iteratively expand the search space by promoting promising retrieved results as new queries. Specifically, starting from the root query paper, we branch into aspect-specific queries (namely, motivation, methods, and experiments), using the process described in Section 3.2, while excluding the abstract view. Then, each query retrieves a set of top- k candidate papers from the corpus, which serve as the first layer of expansion. At each subsequent depth r , every aspect-specific branch selects a single representative work from its top- k results and promotes it as the next query paper, ensuring continuity of semantic relevance while avoiding divergence from the original query intent, and this query paper is further decomposed into aspect-specific queries for the next iteration. At the end, this process grows a tree-structured retrieval space: after r iterations, the search expands

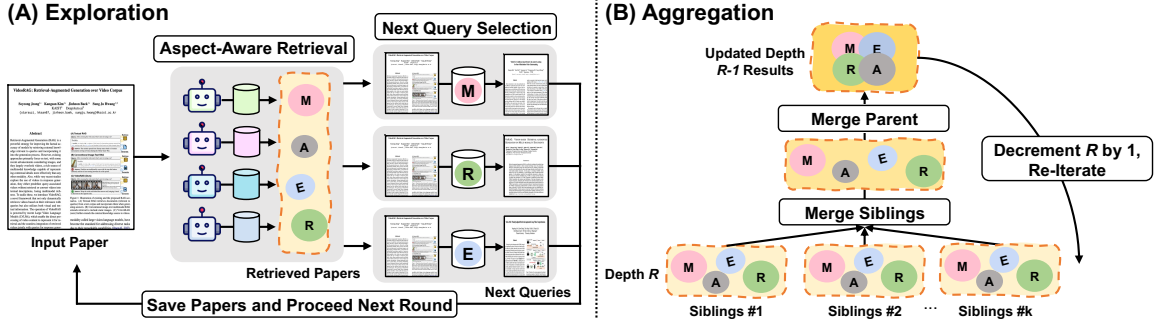


Figure 2: Overview of Chain-of-Retrieval (CoR), consisting of (A) Exploration and (B) Aggregation (formalized in Algorithm 1).

into $4(3)^{r-1}$ distinct exploration paths, each corresponding to an aspect-driven exploration chain.

To avoid redundancy and ensure efficiency, we store previously selected papers in aspect-specific caches: $\text{CACHE}[a]$ for $a \in \{R, M, E\}$, where each cache records documents with respect to the initial branching aspects of the root query paper, so that any paper stored in the corresponding cache is excluded from future expansions within that branch (to avoid repeated visits to the same works). Further, every retrieved top- k set across depths is preserved in a memory M , which latter supports aggregation across multiple levels of the search tree.

Recursive Post-Order Aggregation After R rounds of expansion, we aggregate results in a bottom-up, post-order fashion. Starting from the leaf nodes (i.e., the final retrieved sets at depth R), we first merge sibling results that share the same parent using Reciprocal Rank Fusion (RRF) (Cormack et al., 2009), defined as follows: $\text{RRF}(P) = \sum_{q \in \mathcal{Q}} \frac{1}{k + \text{rank}_q(P)}$, where $\text{rank}_q(P)$ denotes the position of paper P in the ranked list retrieved by query q , and k is a smoothing constant. This step (MERGESIBLINGS in Algorithm 1) emphasizes documents consistently ranked highly across different aspect queries. The merged sibling sets are then recursively combined upward with their parent results (MERGEEDGES), gradually consolidating evidence across the retrieval tree until the root.

It is worth noting that, through this recursive process, rank signals from deeper layers are naturally decayed: RRF suppresses the influence of lower-ranked items, and thus papers retrieved only at later depths contribute less to the final ranking. Formally, the effective weight of ranks from depth R decreases at a rate proportional to $\sigma(2\sigma)^{R-1}$, where σ is the number of query sets per retrieval. As a result, papers strongly and repeatedly supported across multiple aspects and iterations are reinforced, while weaker or spurious signals are

Algorithm 1 Chain-of-Retrieval (CoR).

Please see Appendix F for the detailed algorithms.

Require: Root paper D ; Max depth R ; Top- k per query K ;
Corpora $\mathcal{C} = \{\mathcal{C}_{\text{abstract}}, \mathcal{C}_{\text{chunked}}\}$

Ensure: Final top- k candidates for D

```

1:  $Q \leftarrow [(D, \text{ROOT})]$ 
2:  $\mathcal{M} \leftarrow []$ 
3:  $\text{CACHE}[a] \leftarrow \emptyset \quad \forall a \in \{R, M, E\}$ 
4: for  $r = 0$  to  $R - 1$  do ▷ Exploration
5:    $\mathcal{M}[r] \leftarrow [], Q_{\text{next}} \leftarrow []$ 
6:   for all  $(d, \text{PARENT}) \in Q$  do
7:      $P \leftarrow \text{ONEHOPRETRIEVAL}(d, \text{PARENT}, \mathcal{C}, K)$ 
8:      $\mathcal{M}[r].\text{EXTEND}(P)$ 
9:      $Q_{\text{next}} \leftarrow \text{SELECTNEXTQUERY}(P, \text{CACHE}, \gamma)$ 
10:   $Q \leftarrow Q_{\text{next}}$ 
11: for  $r = R - 1$  downto  $1$  do ▷ Aggregation
12:   $S \leftarrow \text{MERGESIBLINGS}(\mathcal{M}[r])$ 
13:   $\mathcal{M}[r - 1] \leftarrow \text{MERGEEDGES}(\mathcal{M}[r - 1], S)$ 
14: return  $\text{MERGESIBLINGS}(\mathcal{M}[0])$ 

```

attenuated, yielding a robust ranking for the query. 301

3.4 Multi-Aspect Preference Optimization 302

Offline Policy Exploration To improve the quality of aspect-specific queries, we further train each query-generation agent f with the Direct Preference Optimization (DPO) algorithm (Rafailov et al., 2023). A key challenge, however, lies in constructing reliable agent-specific preference datasets, as collecting human-labeled preferences for multi-agent systems is both costly and impractical. Moreover, recent studies emphasize that initializing DPO with the Supervised Fine-Tuned (SFT) model is crucial to mitigate distributional shift issues (Feng et al., 2024; Xu et al., 2024; Meng et al., 2024). 303-314

To address these challenges, we follow the synthetic preference construction procedure of Meng et al. (2024), treating an off-the-shelf instruction-tuned LLM (\mathcal{P}_θ) as the SFT policy (\mathcal{P}_{SFT}). Specifically, for each input paper D , we generate a set of k rolled-out candidate queries from every functional agent: $Q_d^{\text{Exp}} = \bigcup_{f \in \{f_R, f_M, f_E\}} \{q_{(D,f)}^{(1)}, \dots, q_{(D,f)}^{(k)}\}$. After that, among these candidates, we identify 315-322

the best and worst queries according to a reward function: $q_{(D,f)}^{\text{Best}} = \arg \max_j \text{Reward}(q_{(D,f)}^{(j)})$ and $q_{(D,f)}^{\text{Worst}} = \arg \min_j \text{Reward}(q_{(D,f)}^{(j)})$. In addition, a preference pair is retained only if the reward gap between those two exceeds a margin τ , ensuring that $\mathcal{D}_{\text{pref}}$ contains discriminative signals for training.

Reward Formulation Designing an effective reward function is crucial for preference optimization, as it determines how query-generation agents are guided to improve. Our principle is to provide reward signals that are reliable and reproducible, cost-efficient to compute, and directly aligned with retrieval quality; thus, instead of relying on qualitative judgments from LLMs (LLMs-as-judge), which are costly and potentially inconsistent, we leverage citation links as environmental feedback, providing direct, objective reward signals (where we use Recall@ k in practice). Also, for each agent $f \in \{f_R, f_M, f_E\}$, we define an independent reward function Reward_f , such that the reward for one agent does not depend on the outputs of others, under the assumption that reinforcing the quality of aspect-specific queries individually is sufficient, since the final retrieval results are aggregated across agents via a linear combination of their outputs.

4 Experiments

4.1 Experimental Setup

We describe here the setup used to validate our method. Further implementation details and training procedures are in Appendix B and Appendix C.

Datasets To evaluate full paper-to-paper retrieval, we build upon and extend existing abstract-only benchmarks, which we refer to as SCIFULLBENCH. Specifically, it preserves the original retrieval formulation and relevance definition, while augmenting each query-candidate pair with full-text content and refreshing the corpus with more recent publications. In addition, we also evaluate COR under abstract-only configurations on SCIFULLBENCH (Table 17) and on an existing benchmark (Table 19), as well as in a full-paper-to-full-paper variant of the same benchmark (Table 20), to ensure that our gains are not specific to the full-paper scenario and generalize across retrieval regimes, on which COR consistently outperforms the baselines.

To build SCIFULLBENCH, we follow prior work on scientific paper retrieval (Singh et al., 2023; Medić and Snajder, 2022; Cohan et al., 2020) and define relevance using neighboring relationships in

the academic citation graph, including both incoming citations and outgoing references. Query papers are collected from major ML (NeurIPS, ICLR) and NLP (ACL, EMNLP) venues using OpenReview API¹, ACL-Anthology², and SEA (Yu et al., 2024). We then retain papers with at least ten relevant neighbors and randomly sample 400 query papers per venue. For the retrieval corpus, we collect CS-related papers from arXiv (2020-2025)³, where full text is publicly available, resulting in approximately 40K papers. Lastly, we post-process papers by removing reference sections, filtering citation markers, and eliminating their in-text mentions, to avoid hindsight leakage. Further dataset construction details are provided in Appendix A.

Retrieval Models We compare our method with existing retrievers developed for scientific paper retrieval, as follows: **SPECTER2 Base** (Singh et al., 2023); **SciNCL** (Ostendorff et al., 2022); **SciMultiMHAExpert** (Zhang et al., 2023a); **SPECTER2 Adapters + MTL CTRL** (Singh et al., 2023). We also consider off-the-shelf general-purpose embedding models, such as **Jina-Embeddings-V2-Base-EN** (Günther et al., 2023), **BGE-M3** (Chen et al., 2024a), and **Inf-Retriever-v1-1.5B** (Yang et al., 2025), leveraged further as backbones of COR.

Retrieval Units We consider various retrieval units, where we denote **A** for *Abstract*, **F** for *Full paper*⁴, and **C** for *chunked context* (i.e., segmented units of the full paper, each capped at 3K tokens). Using them, we consider four retrieval setups: **A2A** (Abstract-to-Abstract), **F2F** (Full-to-Full), **A2C** (Abstract-to-Chunk), and **F2C** (Full-to-Chunk).

Query Optimizers We instantiate query optimizers using LLMs: the open-weight instruct models Llama-3.2-3B-Instruct (Meta, 2024) and QWEN-2.5-3B-Instruct (Yang et al., 2024; Team, 2024). For reproducibility, we set the temperature to 0 and the maximum generation tokens to 2,000.

Training and Inference Configuration For each query, we retrieve the top-300 candidate papers. To prevent repeatedly selecting the same candidate as the next query, we always skip the highest-ranked results (selected previously) and instead promote the next-best. During preference pair construction,

¹<https://openreview.net/>

²<https://aclanthology.org/>

³<https://arxiv.org/>

⁴Since full papers often exceed the context length of embedding models, we truncate them up to the maximum length.

Table 1: Main results of various domain-specific and domain-agnostic retrievers on SCIFULLBENCH.

Methods	ICLR-NeurIPS				ACL-EMNLP			
	References		Citations		References		Citations	
	nDCG	Recall	nDCG	Recall	nDCG	Recall	nDCG	Recall
Lexical-Based Retrievers								
BM-25 (A2A)	28.29	38.21	27.52	37.41	21.02	30.79	20.86	29.37
BM-25 (F2F)	26.34	37.02	36.48	47.47	23.38	35.10	28.08	38.58
Domain-Specific Retrievers (A2A)								
SciNCL	36.24	51.80	33.07	47.90	26.12	40.71	25.58	38.91
SPECTER2-Base	35.24	50.41	34.48	49.09	25.07	39.02	26.85	40.21
SPECTER2-Adapter-MTL CTRL	36.22	51.07	33.12	47.38	25.32	38.86	25.48	38.54
SciMult-MHAExpert	30.95	45.04	28.32	41.07	23.11	36.12	22.57	33.79
Jina-Embeddings-v2-BASE-EN								
Abstract-to-Abstract (A2A)	36.78	49.92	33.85	47.89	24.96	37.90	26.45	38.92
Full-to-Full (F2F)	37.05	50.60	35.80	49.63	27.55	41.69	28.34	40.65
COR w/ Llama-3.2-3B-Instruct	38.88	55.50	38.58	56.21	28.34	44.44	31.80	47.03
COR w/ QWEN-2.5-3B-Instruct	38.91	55.78	39.36	56.97	28.21	44.83	32.54	48.30
BGE-M3								
Abstract-to-Abstract (A2A)	32.08	44.09	30.31	42.85	22.95	34.55	24.41	36.28
Full-to-Full (F2F)	33.71	45.84	32.22	44.01	24.13	35.88	24.85	35.87
COR w/ Llama-3.2-3B-Instruct	33.81	49.39	34.54	49.27	24.62	39.22	28.87	42.67
COR w/ QWEN-2.5-3B-Instruct	34.84	50.43	35.72	50.77	25.22	39.84	29.33	43.72
Inf-Retriever-v1-1.5B								
Abstract-to-Abstract (A2A)	45.84	61.72	38.61	54.63	33.43	50.10	30.94	45.53
Full-to-Full (F2F)	31.09	42.27	35.58	48.79	32.86	48.00	30.31	43.64
COR w/ Llama-3.2-3B-Instruct	47.31	65.31	43.66	61.54	34.48	53.63	36.18	52.95
COR w/ QWEN-2.5-3B-Instruct	47.21	64.78	44.21	62.06	34.41	53.30	36.87	53.59

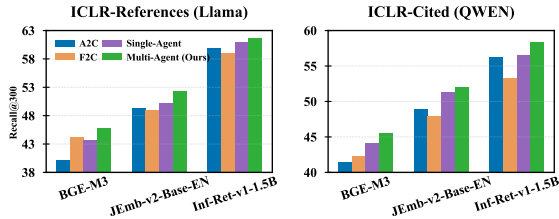


Figure 3: Retrieval performance comparing different query formulations against our multi-agent COR framework (after a single depth of retrieval and with untrained query optimizers).

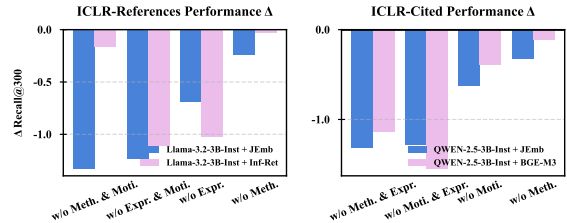


Figure 4: Change (Δ) in retrieval performance (relative to using all aspects) when excluding individual scientific aspects (Moti = Motivation, Meth = Method, Expr = Experiment).

we generate 16 candidate queries per aspect and evaluate them based on Recall@30, while excluding pairs whose reward margin is smaller than 3%.

4.2 Experimental Results and Analysis

In this section, we provide an in-depth analysis of the experimental results of our method. Please refer to Appendix E for more experiments and analyses.

Main Results Table 1 presents the main results, where COR outperforms all baselines across various settings, validating the effectiveness of our proposed framework for full paper-to-paper retrieval⁵. Notably, when using the same domain-agnostic retriever, COR surpasses abstract-to-abstract (A2A) baselines by an average of 6.37% in Recall, demonstrating that simply relying on abstracts is suboptimal compared to our aspect-driven approach. Also, interestingly, general-purpose embedding models such as Jina-Embeddings-V2-Base-EN and BGE-

⁵We report results using the top-300 documents, and additionally provide results for top-100 and top-200 documents and additional metrics in the Appendix E.

M3 (whose A2A performance lags behind representative domain-specific retrievers) achieve competitive or superior performance once integrated into COR, emphasizing that COR delivers strong performance regardless of the retrievers by serving as a plug-and-play framework. Lastly, when compared against F2F retrieval with long-context embedding models, COR achieves average gains of 7.82% in Recall, confirming that simply embedding entire papers is insufficient, and that structured, iterative search expansion with post-order aggregation is crucial for effective full paper-to-paper retrieval.

Effectiveness of Specialized Query Optimizers

Our results in Figure 3 confirm the effectiveness of using multiple specialized LLM agents, each dedicated to a single aspect. More specifically, using multiple specialized agents outperforms a single base agent tasked with generating multiple queries (without aspect separation). Also, when contrasted with setups that use either full papers or abstracts as queries (matched with the same chunked corpus) and that retrieve the same number of total candi-

Table 2: Performance comparison with and without the Aspect-Aware Cache in the selection algorithm, reported after retrieval depths of three using DPO-trained query optimizers. We report the mean and standard deviation over three independent runs, and underline statistically significant improvements.

{Optimizer} + {Retriever}	ICLR-References	ICLR-Citations
	Recall@300	Recall@200
Llama-3.2-3B-Inst + JEmb-v2		
w/o Aspect-Aware Cache	53.63 ± 0.11	46.80 ± 0.03
w/ Aspect-Aware Cache	<u>54.15 ± 0.39</u>	<u>47.08 ± 0.04</u>
Llama-3.2-3B-Inst + BGE-M3		
w/o Aspect-Aware Cache	46.88 ± 0.16	40.82 ± 0.18
w/ Aspect-Aware Cache	<u>47.35 ± 0.13</u>	<u>41.16 ± 0.27</u>
Llama-3.2-3B-Inst + Inf-Ret-v1		
w/o Aspect-Aware Cache	62.65 ± 0.10	51.83 ± 0.05
w/ Aspect-Aware Cache	<u>62.98 ± 0.19</u>	<u>52.31 ± 0.11</u>

dates as COR of a single retrieval round, our multi-aspect queries generated with specialized agents achieve gains of 3.36% and 3.33%, respectively, highlighting the advantage of decomposing noisy, monolithic queries into aspect-aware subqueries.

Importance of Multi-Aspect Coverage To examine the role of aspect diversity in query formulation, we analyze how the retrieval performance changes as we vary the set of participating query optimizer agents. As shown in Figure 4, excluding any single aspect leads to measurable performance degradation, which supports the necessity of capturing the multifaceted nature of scientific papers. Not all aspects, however, contribute equally: excluding experimental or research-motivation queries causes larger drops in performance than omitting methodological queries. We attribute this to the inherent characteristics of each aspect. Specifically, experimental and research-motivation views often yield stronger semantic overlap across related papers, grounded in shared entities such as tasks, baselines, and datasets. In contrast, methodological details are more heterogeneous and less entity-driven, demanding deeper semantic abstraction and reasoning to identify commonalities.

Role of Aspect-Aware Cache Table 2 shows the benefits of incorporating the aspect-aware cache into our *Next Query Selection* strategy (described in Section 3.3), which enables preventing redundant exploration within each aspect branch and thus encouraging more diverse search trajectories. From this, COR with the aspect-aware cache consistently outperforms those without it, demonstrating its effectiveness in steering the retrieval process toward broader coverage and stronger overall performance.

Efficacy of Query Preference Optimization To assess the contribution of reinforcing aspect-aware

Table 3: Performance comparison of specialized agents trained with DPO versus untrained query optimizers, evaluated after a single round of retrieval (i.e., one depth). We report the mean and standard deviation across three independent trials, and underline statistically significant improvements.

{Optimizer} + {Retriever}	ICLR-References		NeurIPS-Citations	
	Recall@100	Recall@300	Recall@100	Recall@300
Llama-3.2-3B-Inst + JEmb-v2				
w/o DPO	36.20 ± 0.00	52.14 ± 0.00	41.70 ± 0.00	57.36 ± 0.00
w/ DPO	<u>37.10 ± 0.12</u>	<u>52.60 ± 0.05</u>	41.33 ± 0.06	<u>57.95 ± 0.10</u>
Llama-3.2-3B-Inst + BGE-M3				
w/o DPO	31.32 ± 0.00	45.61 ± 0.00	36.24 ± 0.00	50.85 ± 0.00
w/ DPO	<u>31.59 ± 0.02</u>	<u>45.91 ± 0.02</u>	<u>36.55 ± 0.07</u>	50.49 ± 0.09
Llama-3.2-3B-Inst + Inf-Ret-v1				
w/o DPO	45.15 ± 0.05	61.75 ± 0.05	47.91 ± 0.00	63.04 ± 0.00
w/ DPO	<u>45.92 ± 0.06</u>	<u>62.15 ± 0.06</u>	<u>48.57 ± 0.05</u>	<u>63.61 ± 0.07</u>
QWEN-2.5-3B-Inst + JEmb-v2				
w/o DPO	36.72 ± 0.00	52.43 ± 0.00	41.82 ± 0.00	58.20 ± 0.00
w/ DPO	<u>37.41 ± 0.17</u>	<u>53.00 ± 0.10</u>	<u>42.21 ± 0.15</u>	<u>58.64 ± 0.12</u>
QWEN-2.5-3B-Inst + BGE-M3				
w/o DPO	30.94 ± 0.00	45.68 ± 0.00	37.09 ± 0.00	51.23 ± 0.00
w/ DPO	<u>31.52 ± 0.14</u>	<u>46.60 ± 0.04</u>	<u>37.75 ± 0.15</u>	<u>52.08 ± 0.21</u>

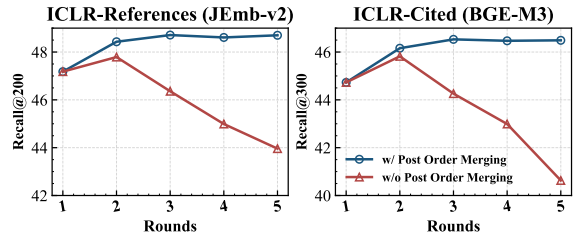


Figure 5: Retrieval results as a function of retrieval depths, with DPO-trained Llama-3.2-3B-Instruct as optimizers.

query optimizers with DPO, we conduct an ablation study isolating the DPO training. As shown in Table 3, integrating DPO consistently improves performance across diverse LLM-retriever configurations in both the reference and citation setups. Notably, although the preference pairs are constructed in the Reference (outgoing-links) setting using Recall as the reward signal, the gains transfer to other evaluation setups of Citation (incoming-links).

Analysis on Post-Order Aggregation To verify whether our aggregation mechanism is responsible for performance improvements beyond those obtained from multi-aspect expansion, we conduct an analysis. As shown in Figure 5, COR yields performance gains across successive retrieval rounds, with the gap over baselines widening as depth increases. This confirms the effectiveness of our post-order aggregation, which progressively attenuates weak signals from distant connections and thereby enhances robustness against noise in deeper rounds. By contrast, the baseline strategy (merging results from all branches at once) fails to adaptively penalize such signals, leading to steadily diminishing performance. However, improvements plateau beyond depth 3, which suggests diminishing returns, likely due to added noise with deeper expansions.

Diversity Analysis Recall that COR is designed to expand the search space iteratively across multiple aspects of the query paper, and we further

Table 4: Analysis on the semantic coverage of the Top-300 retrieved documents over ground truth documents with Convex Hull Volume. The results of CoR are reported after 3 rounds (i.e., depth) with Llama-3.2-3B-Inst as query optimizers.

{Type} w/ {Retriever}	ACL (Citations)	NeurIPS (Citations)
	CHV Ratio (Δ)	CHV Ratio (Δ)
Relative Coverage over A2A		
F2F w/ JEmb-v2	1.126 (+12.6%)	1.135 (+13.5%)
F2F w/ BGE-M3	1.088 (+8.8%)	1.069 (+6.9%)
F2F w/ Inf-Ret-v1	1.128 (+12.8%)	1.097 (+9.7%)
Relative Coverage over A2A		
CoR w/ JEmb-v2	1.190 (+19.0%)	1.198 (+19.8%)
CoR w/ BGE-M3	1.107 (+10.7%)	1.097 (+9.7%)
CoR w/ Inf-Ret-v1	1.167 (+16.7%)	1.140 (+14.0%)
Relative Coverage over F2F		
CoR w/ JEmb-v2	1.095 (+9.5%)	1.078 (+7.8%)
CoR w/ BGE-M3	1.044 (+4.4%)	1.042 (+4.2%)
CoR w/ Inf-Ret-v1	1.060 (+6.0%)	1.072 (+7.2%)

analyze this by measuring the semantic coverage of retrieved candidates over ground-truth in the embedding space using the Convex Hull Volume Ratio (CHV Ratio), which compares coverage between two configurations, with the same set (see Appendix C.3 for details). As shown in Table 4, CoR consistently improves the diversity of semantic coverage among the retrieved documents: relative to the baseline A2A setup, it achieves a 14.98% gain in CHV Ratio, surpassing the 10.72% gain of the F2F setup. Moreover, when directly compared against F2F, CoR still provides an additional 6.52% improvement. These results highlight that beyond improving retrieval accuracy, the proposed CoR enables broader semantic exploration, which allows discovering relevant papers that would likely be missed under the vanilla retrieval strategies.

Generalization to Other Domains To examine the generality of CoR beyond scientific paper retrieval, we further evaluate its performance on a different but challenging document-to-document retrieval task: patent-to-patent retrieval, which has substantially long queries (please see Appendix D for details on task design and implementation). As shown in Table 5, CoR delivers robust gains over both the domain-agnostic and domain-specific baselines, surpassing F2F retrieval by 2.01% and A2A retrieval by 6.96%. Importantly, these improvements are achieved without any training (i.e., using the off-the-shelf GPT-4o for query optimization), highlighting the versatility and broad applicability of CoR for long document-to-document retrieval.

Human Evaluation To further assess the quality of retrieved results beyond the citation-based ground truth in SCIFULLBENCH, we conduct a human evaluation. A total of 15 participants evaluate the relevance of the Top-5 retrieved candidates for three methods (A2A, F2F, CoR) across 15

Table 5: Experimental results on patent-to-patent retrieval using PATENTFULLBENCH (self-constructed; See Appendix D). Results for CoR are reported after two rounds of retrieval.

Methods	References	Citations
	Recall@100	Recall@100
Baseline		
PAT-SPECTER (A2A) (Ghosh et al., 2024)	47.36	52.16
PaECTER (A2A) (Ghosh et al., 2024)	52.77	56.25
BGE-M3 (A2A)	46.42	49.57
BGE-M3 (F2F)	52.30	55.81
Inf-Ret-v1 (A2A)	52.78	57.94
Inf-Ret-v1 (F2F)	56.23	58.59
JEmb-v2-Base-EN (A2A)	49.39	53.27
JEmb-v2-Base-EN (F2F)	56.57	59.61
CoR (Ours): {Optimizer} + {Retriever}		
GPT-4o-Mini-2024-0718 + BGE-M3	54.71	56.78
GPT-4o-Mini-2024-0718 + Inf-Ret-v1	58.62	63.06
GPT-4o-Mini-2024-0718 + JEmb-v2	57.23	60.74

Table 6: Human study on the relevance of ground truth candidates. The results are reported after 3 rounds (i.e., depth) using DPO-trained QWEN-2.5-3B-Instruct models for query optimizer, and jina-embeddings-v2-base-en model for retriever.

Retrieval Methods	Precision@5
Abstract-to-Abstract (A2A)	60.83
Full-to-Full (F2F)	53.33
Chain-of-Retrieval (CoR; Ours)	64.17

query papers, assigning one of four labels to each candidate (motivation-relevant, method-relevant, experiment-relevant, or irrelevant). We then compute per-query precision by treating any of the three aspect labels as relevant and the remaining label as irrelevant. As shown in Table 6, CoR achieves the highest precision, surpassing A2A and F2F by 3.34% and 10.84%, respectively, indicating that its improvements extend beyond citation-based metrics to human-perceived relevance. To ensure the reliability of human judgments, we measure inter-annotator agreement on it with Cohen’s Kappa coefficient (Cohen, 1960), obtaining a score of 0.43, which corresponds to moderate agreement and supports the reliability of our human evaluation. For further details, please refer to Appendix C.4.

5 Conclusion

In this work, we introduced Chain of Retrieval (in short, CoR), a novel framework for paper-to-paper retrieval that iteratively expands the search space via aspect-aware query optimization (reinforced with DPO training) and recursively merges results via post-order aggregation. For evaluation, we presented SCIFULLBENCH, a large-scale dataset extended from existing benchmarks for full-context scientific retrieval, and the results show that CoR consistently outperforms the abstract-level and full-context baselines, even with off-the-shelf embedding models and for a different patent retrieval task, highlighting its robustness and generality.

591 **Limitations**

592 While our work introduces a novel retrieval ap-
593 proach that iteratively expands the search space
594 via multi-aspect query optimization over the full
595 context of scientific papers, it still has room for
596 future work. First, an adaptive branch pruning
597 mechanism could be further considered during the
598 exploration phase to more dynamically guide the
599 search results with the semantic alignment of indi-
600 vidual papers with the root query, which would be
601 an interesting direction for future work. Second,
602 while our framework retrieves diverse and relevant
603 papers, deeper exploration and repeated LLM infer-
604 ence introduce additional computational overhead.
605 While we partially mitigate this with a lightweight
606 selection mechanism (in place of expensive LLM
607 verifiers), further optimizations on it would be a
608 valuable direction, which we leave as future work.

609 **Ethics Statement**

610 Although our Chain of Retrieval (COR) frame-
611 work improves retrieval performance over prior
612 approaches, it can still retrieve irrelevant or mis-
613 leading papers, which may propagate incorrect or
614 harmful information to both human users and down-
615 stream AI models. To ensure the development of
616 trustworthy automated systems, future work should
617 incorporate robust verification or re-ranking mech-
618 anisms that can effectively filter irrelevant or poten-
619 tially harmful documents from the retrieved set.

620 **References**

621 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama
622 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
623 Diogo Almeida, Janko Altenschmidt, Sam Altman,
624 Shyamal Anadkat, and 1 others. 2023. Gpt-4 techni-
625 cal report. *arXiv preprint arXiv:2303.08774*.

626 Anirudh Ajith, Mengzhou Xia, Alexis Chevalier, Tanya
627 Goyal, Danqi Chen, and Tianyu Gao. 2024. *Lit-
628 Search: A retrieval benchmark for scientific literature
629 search*. In *Proceedings of the 2024 Conference on
630 Empirical Methods in Natural Language Processing*,
631 pages 15068–15083, Miami, Florida, USA. Associa-
632 tion for Computational Linguistics.

633 Akari Asai, Jacqueline He, Rulin Shao, Weijia Shi,
634 Amanpreet Singh, Joseph Chee Chang, Kyle Lo,
635 Luca Soldaini, Sergey Feldman, Mike D’arcy, and
636 1 others. 2024. Openscholar: Synthesizing scien-
637 tific literature with retrieval-augmented lms. *arXiv
638 preprint arXiv:2411.14199*.

639 Parul Awasthy, Aashka Trivedi, Yulong Li, Mihaela
640 Bornea, David Cox, Abraham Daniels, Martin Franz,

Gabe Goodhart, Bhavani Iyer, Vishwajeet Kumar, 641
Luis Lastras, Scott McCarley, Rudra Murthy, Vignesh P, Sara Rosenthal, Salim Roukos, Jaydeep Sen, 642
Sukriti Sharma, Avirup Sil, and 3 others. 2025. *Gran- 643
ite embedding models*. *Preprint*, arXiv:2502.20204. 644

Jinheon Baek, Sujay Kumar Jauhar, Silviu Cucerzan, 645
and Sung Ju Hwang. 2024. Researchagent: Iter- 646
ative research idea generation over scientific liter- 647
ature with large language models. *arXiv preprint 648
arXiv:2404.07738*. 649

Chandra Bhagavatula, Sergey Feldman, Russell Power, 650
and Waleed Ammar. 2018. *Content-based citation 651
recommendation*. In *Proceedings of the 2018 Con- 652
ference of the North American Chapter of the Asso- 653
ciation for Computational Linguistics: Human Lan- 654
guage Technologies, Volume 1 (Long Papers)*, pages 655
238–251, New Orleans, Louisiana. Association for 656
Computational Linguistics. 657

Peter Brown and Yaoqi Zhou. 2019. Large expert- 658
curated database for benchmarking document simi- 659
larity detection in biomedical literature search. 660
Database, 2019:baz085. 661

Eric Chamoun, Michael Schlichtkrull, and Andreas Vla- 662
chos. 2024. *Automated focused feedback generation 663
for scientific writing assistance*. In *Findings of the As- 664
sociation for Computational Linguistics: ACL 2024*, 665
pages 9742–9763, Bangkok, Thailand. Association 666
for Computational Linguistics. 667

Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu 668
Lian, and Zheng Liu. 2024a. *Bge m3-embedding: 669
Multi-lingual, multi-functionality, multi-granularity 670
text embeddings through self-knowledge distillation*. 671
Preprint, arXiv:2402.03216. 672

Jianlyu Chen, Nan Wang, Chaofan Li, Bo Wang, Shi- 673
tao Xiao, Han Xiao, Hao Liao, Defu Lian, and 674
Zheng Liu. 2024b. *Air-bench: Automated hetero- 675
geneous information retrieval benchmark*. *Preprint*, 676
arXiv:2412.13102. 677

Arman Cohan, Sergey Feldman, Iz Beltagy, Doug 678
Downey, and Daniel Weld. 2020. *SPECTER: 679
Document-level representation learning using 680
citation-informed transformers*. In *Proceedings 681
of the 58th Annual Meeting of the Association 682
for Computational Linguistics*, pages 2270–2282, 683
Online. Association for Computational Linguistics. 684

Jacob Cohen. 1960. A coefficient of agreement for 685
nominal scales. *Educational and psychological mea- 686
surement*, 20(1):37–46. 687

Gordon V. Cormack, Charles L. A. Clarke, and Stefan 688
Büttcher. 2009. *Reciprocal rank fusion outperforms 689
condorcet and individual rank learning methods*. *Pro- 690
ceedings of the 32nd international ACM SIGIR con- 691
ference on Research and development in information 692
retrieval*. 693

Michael Han Daniel Han and Unsloth team. 2023. *Un- 694
sloth*. 695

697	Mike D’Arcy, Tom Hope, Larry Birnbaum, and Doug Downey. 2024. Marg: Multi-agent review generation for scientific papers. <i>arXiv preprint arXiv:2401.04259</i> .	<i>for Computational Linguistics (Volume 1: Long Papers)</i> , pages 28614–28659, Vienna, Austria. Association for Computational Linguistics.	754 755 756
701	DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Jun-Mei Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiaoling Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, and 179 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. <i>ArXiv</i> , abs/2501.12948.	Mainak Ghosh, Sebastian Erhardt, Michael E. Rose, Erik Buunk, and Dietmar Harhoff. 2024. Paecter: Patent-level representation learning using citation-informed transformers. <i>Preprint</i> , arXiv:2402.19411.	757 758 759 760
709	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.	Ronald L. Graham. 1972. An efficient algorithm for determining the convex hull of a finite planar set. <i>Info. Proc. Lett.</i> , 1:132–133.	761 762 763
718	Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. 2024. The faiss library.	Mark S Granovetter. 1973. The strength of weak ties. <i>American journal of sociology</i> , 78(6):1360–1380.	764 765
722	Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony S. Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, and 510 others. 2024. The llama 3 herd of models. <i>ArXiv</i> , abs/2407.21783.	Michael Günther, Jackmin Ong, Isabelle Mohr, Alaeddine Abdessalem, Tanguy Abel, Mohammad Kalim Akram, Susana Guzman, Georgios Mastrapas, Saba Sturua, Bo Wang, and 1 others. 2023. Jina embeddings 2: 8192-token general-purpose text embeddings for long documents. <i>arXiv preprint arXiv:2310.19923</i> .	766 767 768 769 770 771 772
730	William Fedus, Barret Zoph, and Noam Shazeer. 2022. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. <i>Journal of Machine Learning Research</i> , 23(120):1–39.	Hongmei Guo, Zhesi Shen, Jianxun Zeng, and Na Hong. 2022. Hybrid methods of bibliographic coupling and text similarity measurement for biomedical paper recommendation. In <i>MEDINFO 2021: One World, One Health—Global Partnership for Digital Innovation</i> , pages 287–291. IOS Press.	773 774 775 776 777
734	Duanyu Feng, Bowen Qin, Chen Huang, Zheng Zhang, and Wenqiang Lei. 2024. Towards analyzing and understanding the limitations of dpo: A theoretical perspective. <i>arXiv preprint arXiv:2404.04626</i> .	Masoud Reyhani Hamedani, Sang-Wook Kim, and Dong-Jin Kim. 2016. Simcc: A novel method to consider both content and citations for computing similarity of scientific papers. <i>Information Sciences</i> , 334:273–292.	778 779 780 781 782 783
738	Luyu Gao, Xueguang Ma, Jimmy Lin, and Jamie Callan. 2023. Precise zero-shot dense retrieval without relevance labels. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1762–1777, Toronto, Canada. Association for Computational Linguistics.	Khalid Haruna, Maizatul Akmar Ismail, Abdullahi Baffa Bichi, Victor Chang, Sutrisna Wibawa, and Tutut Herawan. 2018. A citation-based recommender system for scholarly paper recommendation. In <i>Computational Science and Its Applications—ICCSA 2018: 18th International Conference, Melbourne, VIC, Australia, July 2-5, 2018, Proceedings, Part I 18</i> , pages 514–525. Springer.	784 785 786 787 788 789 790
745	Zhaolin Gao, Kianté Brantley, and Thorsten Joachims. 2024. Reviewer2: Optimizing review generation through prompt generation. <i>arXiv preprint arXiv:2402.10886</i> .	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, and 1 others. 2022. Lora: Low-rank adaptation of large language models. <i>ICLR</i> , 1(2):3.	791 792 793 794 795
749	Aniketh Garikiparthi, Manasi Patwardhan, Aditya Sanjiv Kanade, Aman Hassan, Lovekesh Vig, and Arman Cohan. 2025. MIR: Methodology inspiration retrieval for scientific research problems. In <i>Proceedings of the 63rd Annual Meeting of the Association</i>	Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2021. Unsupervised dense information retrieval with contrastive learning. <i>arXiv preprint arXiv:2112.09118</i> .	796 797 798 799 800
753		Yiqiao Jin, Qinlin Zhao, Yiyang Wang, Hao Chen, Kaijie Zhu, Yijia Xiao, and Jindong Wang. 2024. Agentreview: Exploring peer review dynamics with llm agents. <i>arXiv preprint arXiv:2406.12708</i> .	801 802 803 804
		Anshul Kanakia, Zhihong Shen, Darrin Eide, and Kuansan Wang. 2019. A scalable hybrid research paper recommender system for microsoft academic. In <i>The world wide web conference</i> , pages 2893–2899.	805 806 807 808

809	Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 6769–6781, Online. Association for Computational Linguistics.	867
810		868
811		869
812		870
813		871
814		872
815		
816	Maxwell Mirton Kessler. 1963. Bibliographic coupling between scientific papers. <i>American documentation</i> , 14(1):10–25.	
817		
818		
819	Omar Khattab and Matei Zaharia. 2020. Colbert: Efficient and effective passage search via contextualized late interaction over bert. In <i>Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval</i> , pages 39–48.	873
820		874
821		875
822		876
823		
824		
825	Anton Korikov, George Saad, Ethan Baron, Mustafa Khan, Manav Shah, and Scott Sanner. 2024. Multi-aspect reviewed-item retrieval via LLM query decomposition and aspect fusion . In <i>Proceedings of the Workshop Information Retrieval’s Role in RAG Systems (IR-RAG 2024) co-located with the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2024), Washington DC, USA, 07 18, 2024</i> , volume 3784 of <i>CEUR Workshop Proceedings</i> , pages 23–33. CEUR-WS.org.	877
826		878
827		
828		
829		
830		
831		
832		
833		
834		
835		
836	Saar Kuzi, Anna Shtok, and Oren Kurland. 2016. Query expansion using word embeddings. In <i>Proceedings of the 25th ACM international on conference on information and knowledge management</i> , pages 1929–1932.	879
837		880
838		881
839		882
840		883
841	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In <i>Proceedings of the 29th Symposium on Operating Systems Principles</i> , pages 611–626.	884
842		885
843		886
844		887
845		888
846		889
847		890
848	Anne Lauscher, Brandon Ko, Bailey Kuehl, Sophie Johnson, David Jurgens, Arman Cohan, and Kyle Lo. 2021. MultiCite: Modeling realistic citations requires moving beyond the single-sentence single-label setting . <i>Preprint</i> , arXiv:2107.00414.	891
849		892
850		893
851		
852		
853	Yibin Lei, Yu Cao, Tianyi Zhou, Tao Shen, and Andrew Yates. 2024. Corpus-steered query expansion with large language models . In <i>Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 393–401, St. Julian’s, Malta. Association for Computational Linguistics.	894
854		895
855		896
856		897
857		898
858		899
859		900
860	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, and 1 others. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. <i>Advances in neural information processing systems</i> , 33:9459–9474.	901
861		902
862		903
863		904
864		905
865		906
866		907
		908
		909
		910
		911
		912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922

923	Bowen Peng, Jeffrey Quesnelle, Honglu Fan, and Enrico Shippole. 2023. Yarn: Efficient context window extension of large language models. <i>arXiv preprint arXiv:2309.00071</i> .	Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, and 1 others. 2023. Gemini: a family of highly capable multimodal models. <i>arXiv preprint arXiv:2312.11805</i> .	976
924			977
925			978
926			979
927	Peter Pirolli and Stuart Card. 1995. Information foraging in information access environments. In <i>Proceedings of the SIGCHI conference on Human factors in computing systems</i> , pages 51–58.	Qwen Team. 2024. Qwen2.5: A party of foundation models .	980
928			981
929			
930			
931	Peter Pirolli and Stuart Card. 1999. Information foraging. <i>Psychological review</i> , 106(4):643.	Laurens Van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-sne. <i>Journal of machine learning research</i> , 9(11).	982
932			983
933	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. <i>Advances in neural information processing systems</i> , 36:53728–53741.	Liang Wang, Nan Yang, and Furu Wei. 2023. Query2doc: Query expansion with large language models . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 9414–9423, Singapore. Association for Computational Linguistics.	984
934			985
935			986
936			987
937			988
938	Stephen E Robertson, Steve Walker, Susan Jones, Micheline M Hancock-Beaulieu, Mike Gatford, and 1 others. 1995. Okapi at trec-3. <i>Nist Special Publication Sp</i> , 109:109.	Qingyun Wang, Doug Downey, Heng Ji, and Tom Hope. 2024. SciMON: Scientific inspiration machines optimized for novelty . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 279–299, Bangkok, Thailand. Association for Computational Linguistics.	989
939			990
940			991
941			992
942	Joseph John Rocchio Jr. 1971. Relevance feedback in information retrieval. <i>The SMART retrieval system: experiments in automatic document processing</i> .	Siwei Wu, Yizhi Li, Kang Zhu, Ge Zhang, Yiming Liang, Kaijing Ma, Chenghao Xiao, Haoran Zhang, Bohao Yang, Wenhua Chen, Wenhao Huang, Noura Al Moubayed, Jie Fu, and Chenghua Lin. 2024. SciMMIR: Benchmarking scientific multi-modal information retrieval . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 12560–12574, Bangkok, Thailand. Association for Computational Linguistics.	993
943			994
944			995
945	Gerard Salton and Chris Buckley. 1990. Improving retrieval performance by relevance feedback. <i>Journal of the American society for information science</i> , 41(4):288–297.	Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan, Lingfeng Shen, Benjamin Van Durme, Kenton Murray, and Young Jin Kim. 2024. Contrastive preference optimization: Pushing the boundaries of llm performance in machine translation. <i>arXiv preprint arXiv:2401.08417</i> .	996
946			997
947			1000
948			1001
949	Keshav Santhanam, Omar Khattab, Jon Saad-Falcon, Christopher Potts, and Matei Zaharia. 2021. Colbertv2: Effective and efficient retrieval via lightweight late interaction. <i>arXiv preprint arXiv:2112.01488</i> .	An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, and 40 others. 2024. Qwen2 technical report. <i>arXiv preprint arXiv:2407.10671</i> .	1002
950			1003
951			1004
952			1005
953			1006
954	Tao Shen, Guodong Long, Xiubo Geng, Chongyang Tao, Yibin Lei, Tianyi Zhou, Michael Blumenstein, and Daxin Jiang. 2024. Retrieval-augmented retrieval: Large language models are strong zero-shot retriever . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 15933–15946, Bangkok, Thailand. Association for Computational Linguistics.	Junhan Yang, Jiahe Wan, Yichen Yao, Wei Chu, Yinghui Xu, and Yuan Qi. 2025. inf-retriever-v1 (revision 5f469d7) .	1007
955			1008
956			1009
957			1010
958			1011
959			1012
960			
961	Amanpreet Singh, Mike D’Arcy, Arman Cohan, Doug Downey, and Sergey Feldman. 2023. SciRepEval: A multi-format benchmark for scientific document representations . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 5548–5566, Singapore. Association for Computational Linguistics.	Michihiro Yasunaga, Jure Leskovec, and Percy Liang. 2022. LinkBERT: Pretraining language models with document links . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8003–8016, Dublin, Ireland. Association for Computational Linguistics.	1013
962			1014
963			1015
964			1016
965			1017
966			1018
967			1019
968	Henry Small. 1973. Co-citation in the scientific literature: a new measure of the relationship between documents. <i>J. Am. Soc. Inf. Sci.</i> , 42:676–684.		1020
969			1021
970			1022
971	Karen Sparck Jones. 1972. A statistical interpretation of term specificity and its application in retrieval. <i>Journal of documentation</i> , 28(1):11–21.		1023
972			1024
973			1025
974	Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan		1026
975			1027
			1028
			1029

1030 Jianxiang Yu, Zichen Ding, Jiaqi Tan, Kangyang Luo,
1031 Zhenmin Weng, Chenghua Gong, Long Zeng, Ren-
1032 Jing Cui, Chengcheng Han, Qiushi Sun, Zhiyong
1033 Wu, Yunshi Lan, and Xiang Li. 2024. [Automated
1034 peer reviewing in paper SEA: Standardization, eval-
1035 uation, and analysis](#). In *Findings of the Association
1036 for Computational Linguistics: EMNLP 2024*, pages
1037 10164–10184, Miami, Florida, USA. Association for
1038 Computational Linguistics.

1039 Wenhao Yu, Dan Iter, Shuohang Wang, Yichong Xu,
1040 Mingxuan Ju, Soumya Sanyal, Chenguang Zhu,
1041 Michael Zeng, and Meng Jiang. 2023a. [Generate
1042 rather than retrieve: Large language models are
1043 strong context generators](#). In *The Eleventh Inter-
1044 national Conference on Learning Representations*.

1045 Wenhao Yu, Zhihan Zhang, Zhenwen Liang, Meng
1046 Jiang, and Ashish Sabharwal. 2023b. Improving lan-
1047 guage models via plug-and-play retrieval feedback.
1048 *arXiv preprint arXiv:2305.14002*.

1049 Dun Zhang, Panxiang Zou, and Yudong Zhou. 2025a.
1050 [Dewey long context embedding model: A technical
1051 report](#). *Preprint*, arXiv:2503.20376.

1052 Yanzhao Zhang, Mingxin Li, Dingkun Long, Xin Zhang,
1053 Huan Lin, Baosong Yang, Pengjun Xie, An Yang,
1054 Dayiheng Liu, Junyang Lin, Fei Huang, and Jingren
1055 Zhou. 2025b. Qwen3 embedding: Advancing text
1056 embedding and reranking through foundation models.
1057 *arXiv preprint arXiv:2506.05176*.

1058 Yu Zhang, Hao Cheng, Zhihong Shen, Xiaodong Liu,
1059 Ye-Yi Wang, and Jianfeng Gao. 2023a. [Pre-training
1060 multi-task contrastive learning models for scientific
1061 literature understanding](#). In *Findings of the As-
1062 sociation for Computational Linguistics: EMNLP
1063 2023*, pages 12259–12275, Singapore. Association
1064 for Computational Linguistics.

1065 Yu Zhang, Bowen Jin, Xiusi Chen, Yanzhen Shen, Yunyi
1066 Zhang, Yu Meng, and Jiawei Han. 2023b. Weakly
1067 supervised multi-label classification of full-text sci-
1068 entific papers. In *Proceedings of the 29th ACM
1069 SIGKDD Conference on Knowledge Discovery and
1070 Data Mining*, pages 3458–3469.

1071 Huaixiu Steven Zheng, Swaroop Mishra, Xinyun Chen,
1072 Heng-Tze Cheng, Ed H. Chi, Quoc V Le, and Denny
1073 Zhou. 2024. [Take a step back: Evoking reasoning via
1074 abstraction in large language models](#). In *The Twelfth
1075 International Conference on Learning Representa-
1076 tions*.

1077 Yan-Bo Zhou, Linyuan Lü, and Menghui Li. 2012.
1078 Quantifying the influence of scientists and their pub-
1079 lications: distinguishing between prestige and popu-
1080 larity. *New Journal of Physics*, 14(3):033033.

A SciFullBench

In this section, we provide more details on the construction process of SciFullBench, with its statistics and comparison against existing benchmarks.

A.1 Construction Procedure

Step 1 We first crawled research papers from top-tier machine learning venues including ICLR 2024, 2025, NeurIPS 2024, 2023, ACL 2024, 2023, and EMNLP 2024, 2023 for our query documents. We scraped existing publications from the main conference for ACL and EMNLP, including both long- and short- papers from the ACL-Anthology website⁶, and used the OpenReview⁷ API to obtain submitted research papers for ICLR and NeurIPS. Moreover, for ICLR 2024 and NeurIPS 2023, we included a subset of documents uploaded by authors from SEA (Yu et al., 2024). Afterwards, we obtained metadata from the arXiv database uploaded to the Kaggle⁸ website, with available data up to January 2025. Moreover, we constructed our initial temporary raw corpus where its id starts with 20 to 25 which refers to its uploaded date on arXiv database, limiting our target corpus to fairly recent papers. Also, we only filtered papers that are in the machine learning domain, to those belong in cs.AI, cs.LG, cs.CL, cs.CV, cs.NE, cs.IR, cs.DS, cs.CC, cs.DL, cs.HC, cs.RO, cs.MM, cs.CG, cs.SY since our queries are sampled from ML venues.

Step 2 After devising raw data for both queries and candidates, we subsequently formulate gold candidates using citations to represent contextual proximity between scientific documents. We collect the title of papers that cite our potential query documents using the Semantic Scholar API⁹ and check the existence of a paper with a matched title (case-insensitive) in our arXiv metadata, and ruled out any potential query documents with fewer than ten gold candidates that met such conditions. As for benchmark split with references, we use Allen AI Science Parse¹⁰ released under Apache 2.0 license to obtain reference information for respective documents, and filtered documents based on the implemented criteria from above. Using such data, we formulate our raw (query document, gold candidate) pairs where its abstract informa-

tion is intact and where more than 10 deduplicated candidates exist, organized by splits in respective years belonging to four venues (ICLR, NeurIPS, ACL and EMNLP). For each split, we aggregate all the filtered (query, candidate) set from each venue without regarding its published year. For example, we merge all the candidates in the citation split for ICLR 2024 and ICLR 2025 into a single set of ICLR-citations split, not ICLR-2024-citations or ICLR-2025-citations. From this filtered pool, we randomly sampled 500 potential query candidates and formulated a corpus containing gold candidates for each query. This process enables our target corpus to be kept within manageable size. Next, we assembled the original pdf files of the papers in our temporarily formulated target corpus. Since arXiv APIs do not support large-scale requests, we used the Google Cloud API¹¹ to access the full paper dump directly and collect the latest updated versions for each paper by matching its unique ids within the arXiv database.

Step 3 Subsequently, using the Allen AI Science Parse tool, we parsed the query and candidate pdf files collected in Step 2 into a json format file containing title, abstract, main content list containing dictionaries with header and contents as keys, and list of reference paper information along with its mentions for both queries and candidates. Since we use citation information as the main signal for measuring document similarities, excluding the citations and the entire reference section was crucial to validate the fairness of our benchmark. Although Allen AI Science Parse generally separates reference from main content, several issues persist, where direct reference information is intermixed within the main content. Thus, we applied an auxiliary filtering algorithm to remove such issues. Since our tool parses a pdf file into a structured json file, we were able to obtain text information segmented into multiple passages. Hence, prior to reformatting it back to complete text, we eliminated sections or headers that contain at least one reference title (case-insensitive) completely while reformatting into a full paper. However, for some cases the reference section was left vacant. Since our filtering algorithm depends on the set of titles within parsed reference section, it was unable to correctly exclude cases intermixed with reference information when such data are inaccessible. Hence, we ruled out any query documents or pa-

⁶<https://aclanthology.org/>

⁷<https://openreview.net/>

⁸<https://www.kaggle.com/>

⁹<https://www.semanticscholar.org>

¹⁰<https://github.com/allenai/science-parse>

¹¹<https://cloud.google.com/apis?hl=ko>

Table 7: Statistics of query papers in the SCIFULLBENCH benchmark.

	ICLR		NeurIPS		EMNLP		ACL	
	References	Citations	References	Citations	References	Citations	References	Citations
domain	ML		ML		ML/CL		ML/CL	
years included	2024,2025		2023,2024		2023,2024		2023,2024	
# of papers per year	(141,259)	(324,76)	(170, 230)	(353,47)	(226, 174)	(334, 66)	(184,216)	(256,144)
# of tokens per paper	9402.35	8183.37	11184.39	9035.43	5908.36	6226.45	6190.06	6614.6
average # of candidates per sample	21.74	42.25	24.58	40.37	19.64	39.24	18.05	33.89
minimum # of candidates per sample	10	10	10	10	10	10	10	10

Table 8: Statistics of the corpus in the SCIFULLBENCH.

	SCIFULLBENCH
Total # of papers	40,782
Avg. # of tokens per abstract	200.39
Avg. # of tokens per paper	6987.67
Total # of segmented corpus	115,044
Avg. # of tokens per segment	2479.90

pers within our target corpus that did not include more than four references in the respective parsed documents. In this way, it resolves problematic cases in which we give a pass on perturbed documents that still contained reference information. Moreover, due to our filtering heuristics, there remains a problem where numerous documents experience severe loss of their original content, since any passage or header that had at least 1 reference title was excluded. To mitigate such concerns, we excluded documents in which such problematic passages comprise more than half of the total main-content list. Furthermore, we ensured to eliminate any residing title (case-insensitive) included in the reference section.

Step 4 Based on the preprocessed full document, we formulate a pseudo-definitive benchmark consisting of (query, gold candidate), and target corpus set by random sampling 400 query documents per split in each venue that still 10 gold candidates residing subsequent to the filtering procedure. Ultimately, we formulate a total of 3200 query documents, along with large-scale target corpora consisting of all the gold candidates of such queries. Finally, we remove citations and interlinked mentions from the entire set of formatted papers and finalize our benchmark. Since this process does not lead to exclusion of query document from our benchmark, we have called the previous stage a pseudo-definitive benchmark. Citation mentions are removed using the information provided from our parsing tool, along with 10 different citation patterns using the Python regular expression. Con-

sequently, we devise a definitive, final benchmark in which all documents for both query and candidate corpus consist of title, abstracts, full-paper text, and list of segmented content of corresponding full-paper text. The overall statistics of our SciFullBench can be seen at Table 7 and Table 8.

A.2 Comparison with Prior Benchmarks

Previous benchmarks for scientific literature search, where the full context is available for both query and candidates, rarely exist. Kanakia et al. (2019) introduce expert-annotated paper recommendation benchmark using abstract and citations within Microsoft Academic Graph (MAG). SciDocs (Cohan et al., 2020) and SciRepEval (Singh et al., 2023) reveal an evaluation set to search relevant scientific literature within a pool of 30 candidates per query document with 5 gold labels in the computer science domain, while MDCR (Medić and Snajder, 2022) disclose a benchmark with 60 candidates per query from 19 scientific fields sourced from MAG. RELISH (Brown and Zhou, 2019) also provides expert-annotated gold candidates that are relevant to the respective input documents in the biomedical domain. In addition, Mysore et al. (2021) formulates a CFSCube benchmark to evaluate fine-grained sentence-wise alignment between abstract passages. SciMMIR (Wu et al., 2024) also presents a multimodal document retrieval evaluation set to evaluate the performance of figure-wise document retrieval frameworks. However, such benchmarks do not include candidates or queries in which their entire lexical content is fully disclosed. Although Zhang et al. (2023b) intend to evaluate their classifier pipeline using the complete scientific context, it contains five potential candidates per query, which is infeasible to evaluate large-scale literature retrieval frameworks. Meanwhile, MultiCite dataset (Lauscher et al., 2021) provides full context with fine-grained citation-intent annotation on the paper’s citation context (sentences), to evaluate the

model’s ability to accurately identify citation intents. However, its scope is inherently limited to a trivial classification task, and does not reveal meta-data information (especially full context) on the cited papers nor properly filtered and preprocessed, inappropriate for evaluating real-world paper-to-paper retrieval frameworks. Conversely, our proposed benchmark consists of fully-disclosed document context, along with massive number of target candidates, adhering to the realistic setting for evaluating scientific literature retrieval.

B Training Procedures

In this section, we provide details for the training data construction process, followed by the preference set construction and detailed settings, such as hyper-parameters and training environment.

B.1 Query and Corpus Construction

Composition Our input query document for offline rollout includes publications and submissions from leading venues in ML likewise SciFullBench, but with different published or submitted years. We incorporate publications and submissions from **ICLR** 2017 to 2023, **NeurIPS** 2016 to 2022 from the dataset revealed in REVIEWER2 (Gao et al., 2024), **ACL** 2017 to 2022, **EMNLP** 2017 to 2022, and **NAACL** 2017 to 2022. As for the target candidate corpus construction, we follow a similar approach to SciFullBench in which we crawled ML papers in the arXiv database, with dates ranging from 2020 to 2025 as mentioned in Appendix A. Ultimately, we formulate 15K input query documents with at least 10 ground-truth candidates, along with 97k target corpus consisting of all the ground-truth candidates of input query documents.

Formulation Workflow The formulation procedure for building the training set is similar to the process explained in Appendix A, with several key differences. **First**, to prepare query documents that do not overlap with our SCIFULLBENCH benchmark (to prevent data leakage), we collected publications and submissions from different venues (or years) from ones in our evaluation set. Yet, there still exists a slim possibility that duplicate papers may exist in both the benchmark and train set since research papers can be resubmitted to other venues if rejected. To preclude such corner cases, we excluded query documents from our train set where its abstracts match (using a custom matching function, including lowering case, removing empty

sequences, etc.) the ones in query document of SCIFULLBENCH. We particularly used abstracts to check duplicate papers because several queries in our SCIFULLBENCH may not contain title information. **Secondly**, compared to the benchmark construction process that filters query documents with at least 10 ground truth documents for each citation and reference split, we rule out queries with at least 10 ground truth documents including both citation and reference split. Moreover, since the number of our input query documents drastically exceeds those in SCIFULLBENCH (nearly 5 times), our target candidate corpus size inevitably becomes exponentially larger. Thus, we randomly sampled 140K samples from the initial filtered target corpus, and reformulated our queries accordingly with at least 10 candidates as well. Finally, another minor difference in train set construction compared to SCIFULLBENCH was its order in applying the filtering algorithm (**reference removing, citation removing**), where we executed both algorithms for query documents immediately after acquiring full content, followed by repeated citation removing algorithm to further enhance the extent of removal.

B.2 Preference Set Construction

Using the query documents and the formulated corpus from Appendix B.1, we exclusively used references as gold labels. Furthermore, we excluded query documents that exist in the SciFullBench corpus, since our trained query optimizer agents are provided with documents from our benchmark corpus during the multi-round retrieval process. To ensure fairness in evaluating the effectiveness of our algorithm, we excluded potential query documents from our training set that exist in the corpus used for the inference setting, as the documents in the corpus are used as query for later rounds of retrieval. Moreover, to ensure that data leakage does not occur, we additionally checked the SciFullBench queries for those whose title information is available, and additionally filtered queries in our train set with similar or same titles using the matching function introduced earlier.

During the offline policy exploration process, we set the maximum input tokens to 60k and the maximum completion tokens to 2000. Moreover, we set the temperature to 0.7 and the nucleus sampling ratio to 0.8 to promote exploration, with the branching factor of 16 per aspect-aware LLM agent, using a single 48GB NVIDIA A6000 GPU with VLLM gpu memory utilization rate of 0.8 for each respec-

1349 tive LLM agent + Retriever setup. Note that we ar- 1398
 1350 bitrarily stopped exploration under our notion that 1399
 1351 enough preference data was accumulated through 1400
 1352 self-exploration, indicating that we did not use the
 1353 entire dataset constructed in Appendix B.1. Since
 1354 our objective is to demonstrate the effectiveness of
 1355 DPO in our system compared to untrained query
 1356 optimizers instead of directly demonstrating the ef-
 1357 fectiveness of our data, such process was sufficient
 1358 to validate the effectiveness of our approach while
 1359 gaining reasonable efficiency during the training
 1360 process. Subsequent to offline exploration, we mea-
 1361 sured the recall@30 reward of generated queries
 1362 and paired the query response with the highest re-
 1363 ward and the lowest reward, sampling only those
 1364 with a recall@30 difference of at least 3%.

1365 B.3 Training Details and Hyper-Parameters

1366 We trained two instruct models, Llama-3.2-3B- 1401
 1367 Instruct and QWEN-2.5-3B-Instruct, using the Un- 1402
 1368 sloth (Daniel Han and team, 2023) library to reduce 1403
 1369 memory consumption. We trained the respective 1404
 1370 query optimizer agents with LoRA (Hu et al., 2022) 1405
 1371 adapters and 4-bit quantization on the default mod- 1406
 1372 els. In addition, we set the maximum total number 1407
 1373 of tokens to 4000, allocating 2000 tokens each for 1408
 1374 input and output. This ensures that the responses 1409
 1375 from our preference set match the output token lim- 1410
 1376 its in the inference setup, allowing us to exclusively 1411
 1377 assess the effect of training rather than other factors. 1412
 1378 Meanwhile, shortened input tokens significantly re- 1413
 1379 duce training time, mitigating overfitting to details 1414
 1380 within the query document, and encourage learning 1415
 1381 of favorable patterns from the preferred response 1416
 1382 queries. Furthermore, we trained each LLM agent 1417
 1383 with a linear scheduler, a learning rate of 1e-5, a 1418
 1384 weight decay rate of 0.01, and a maximum gradient 1419
 1385 norm of 0.6 for 3 epochs, with a global batch size 1420
 1386 of 32 on a single NVIDIA A5000 GPU with 24GB 1421
 1387 size VRAM, while occasionally using NVIDIA 1422
 1388 A6000 GPU with 48GB size VRAM to speedup 1423
 1389 training process. Under such setup, it required ap- 1424
 1390 proximately 14 hours to train a single agent. As 1425
 1391 for the LoRA configuration, we set both the rank 1426
 1392 and α scale to 64, and applied the adapters to both 1427
 1393 the attention layers and the feedforward network, 1428
 1394 without leveraging dropout or bias factors. 1429

1395 C Experiment Details

1396 In this section, we provide detailed setup for exper- 1430
 1397 iments, including prompts, neural retrievers, query 1431

1398 optimizers, experiment environment, metrics, and 1399
 1400 implementation details. We also provide a detailed 1401
 description for human evaluation results in Table 6.

1401 C.1 Inference Settings

1402 **Prompt Construction** We constructed prompts 1403
 1404 for scientific paper-to-paper retrieval based on the 1404
 1405 three major aspects that comprise a scientific paper: 1405
Methodology, Experiments, and Research Ques- 1406
tions, inheriting the approaches in conventional 1407
 1408 research for the AI-for-Scientific Discovery do- 1408
 1409 main (Baek et al., 2024; Wang et al., 2024; Moussa 1409
 1410 et al., 2025), widely regarded as the most common 1410
 1411 and foundational aspects that constitute scientific 1411
 1412 papers. The initial raw prompts were first created 1412
 1413 using ChatGPT, and human-in-the loop revision 1413
 1414 was conducted, to devise the optimal prompt style 1414
 1415 suitable for our framework. The prompts used for 1415
 CoR can be found in Appendix I.

1416 **Retrieval Setting** We primarily use the Eu- 1416
 1417 clidean distance to measure similarities between 1417
 1418 our adaptively generated queries and candidates, 1418
 1419 where we acquired candidates specifically with 1419
 1420 the minimum L2 distance. Since minimizing the 1420
 1421 L2 distance was objective for most of our base- 1421
 1422 line domain-specific embedding models such as 1422
 1423 SciNCL, SPECTER2-Base, SPECTER2-Adapter- 1423
 1424 MTL CTRL, we matched such settings when we 1424
 1425 used jina-embeddings-v2-base-en, BGE-M3, and 1425
 1426 Inf-Retriever-v1-1.5b. However, for SciMult, we 1426
 1427 measured MIPs (Maximum Inner Product), since 1427
 1428 it was trained to maximize MIPs likewise in DPR. 1428
 1429 In addition, we utilize the FAISS library (Douze 1429
 1430 et al., 2024), which enables efficient retrieval from 1430
 1431 a large-scale corpora. We implemented an efficient 1431
 1432 L2 search using FlatL2 and FlatIP for MIPs. 1432

1433 In addition, we concatenated the title and ab- 1433
 1434 stract information for our baseline A2A retrieval 1434
 1435 settings, following the conventional retrieval setup 1435
 1436 from Cohan et al. (2020) and Singh et al. (2023) 1436
 1437 for both query and target candidates. For our F2F 1437
 1438 setting, we utilize the entire query and candidate pa- 1438
 1439 pers and embed them into single vectors. As for our 1439
 1440 CoR framework, when using abstracts as query and 1440
 1441 candidates, we use the concatenated version of title 1441
 1442 and abstracts for both query and candidates in like- 1442
 1443 wise in baseline A2A setup, with the exception of 1443
 1444 root depth to balance the effect of title information 1444
 1445 of root query paper, where the query optimizers 1445
 1446 receive full paper that occasionally includes title 1446
 1447 information. This implementation allows objective 1447

comparison on the effect of aspect-aware decomposition of root paper on a single depth.

Moreover, when mapping the chunked content of different papers back to the paper representation, there were several minor overlaps of left over content (11 out of total segments) consisting of single "." character and parts of check list for NeurIPS papers, even though they were from different papers. When mapping back to the original paper content for next query selection and evaluation, we map those segments to an initial origin single paper that appeared the earliest. Since such contents do not reveal essential semantics within a paper, we map such segments to a single paper for convenience.

Query Optimizers We set the repetition penalty to 1.2. In addition, we used the VLLM (Kwon et al., 2023) library for faster inference of open source models. For Llama-3.2-3B-Instruct, we used the default context window size of 131072, while using YARN (Peng et al., 2023) to extrapolate the default context window size of QWEN-2.5-3B-Instruct from 32768 to 131072. Moreover, for inference, we used combinations of three NVIDIA A5000/RTX3090/4090 GPUs with 24GB VRAM size, along with vllm gpu memory utilization rate of 1.0, to deploy three trained aspect-aware query optimizer LLM agents simultaneously, using fp16 precision for faster inference. As for the inference setup for untrained query optimizer models, we experimented on a combination of single NVIDIA A6000/RTX 4090 GPU with 48GB VRAM using the same context window and fp16 precision, along with a fixed gpu memory utilization rate of 0.7. When it comes to experiments using GPT-4o-Mini-2024-0718 and GPT-4o-2024-11-20 as query optimizers, we used temperature 0 and default hyper parameters of the OpenAI client.chat.completions API¹², while completely using the default hyperparameter settings for gpt-4.1-2025-04-14 (temperature is also set to default). Moreover, 60 is used as the hyperparameter k for Reciprocal Rank Fusion. For base agent experiments in Table 3, we equalized the hyper parameters to those used in its respective counterpart comparison groups, and are prompted to generate three different queries in a structured manner, using the Python Pydantic BaseModel¹³ module.

¹²<https://platform.openai.com/docs/guides/>

¹³<https://docs.pydantic.dev/latest/api>

C.2 Models

Domain Specific Retrievers We compare our method with retrieval using domain-specific embedding models, as mentioned in the main section, to demonstrate robust improvement over approaches that seek to devise optimal representations of paper abstracts through extensive training. For fair assessment, we mainly compare our method with models that demonstrated SoTA performance in previous citation link prediction and paper retrieval benchmarks. SPECTER2 models with adapters and its multitask control codes have been reported to have achieved SoTA performance on the MDCR benchmark (Medić and Snajder, 2022). We used proximity control code adapter, as it was specialized in acquiring and ranking relevant papers, and also base models uploaded to huggingface¹⁴. In addition, we experiment with SPECTER2-base models, trained with triplet loss by inducing neural retrievers to favor positive abstract pairs over vice versa given abstract of query documents sampled from discrete citation graphs. SciNCL¹⁵ is also trained in a similar manner, where its abstract pairs are derived from the continuous citation embedding space, achieving the strongest performance on the SciDocs benchmark to date. In addition, we compare our method with the SciMult-MHAExpert model implemented with Mixture-of-Experts (Ferdus et al., 2022) architecture which compartmentalize internal transformer layers for different tasks in scientific literature tasks, and is reported to outperform previous models on the recommendation benchmark (Kanakia et al., 2019). For SciMult-MHAExpert, we experimented with the model released on the official SciMult github repository¹⁶, utilizing an expert model specially trained for link prediction tasks that predicts linked documents given the abstract of the source document.

Domain Agnostic Retrievers As for our domain-agnostic embedding models, we particularly chose off-the-shelf retrievers with long context window. This was to ensure that we provide fair experimental setting for our baselines, most notably full document-to-full document retrieval setting. Since general full-documents exceed the context window of most embedding models, we specifically chose long-context window embedding models to mitigate unfair penalization on our full

¹⁴<https://huggingface.co/allenai/specter2>

¹⁵<https://huggingface.co/malteos/scincl>

¹⁶<https://github.com/yuzhimanhua/SciMult>

document-to-document retrieval baseline. In particular, for our main experiment, we use three different neural retrievers, Jina-Embeddings-v2-Base-EN, and BGE-M3. both highly cited general-purpose embedding models with a context window of 8192 tokens, showing that our method is compatible with the most widely used retriever models. In addition, we used Inf-Retriever-v1-1.5B, a lightweight version of Inf-Retriever-v1, demonstrating SoTA performance among open source embedding models for long-context retrieval(LongEmbed) in MTEB (Muennighoff et al., 2022)¹⁷ as of September 12th, 2025, with context window size of 32768, along with SoTA performance on AIR-Bench (general QA tasks) benchmark (Chen et al., 2024b). This further shows that our method can perform robustly with the most recent general purpose embedding models with advanced long-context handling capabilities.

Furthermore, we provide additional experiments using other recent general-purpose embedding models varying in context window size, such as QWEN3-Embedding-0.6B (Zhang et al., 2025b), Dewey-EN-Beta (Zhang et al., 2025a), and Granite-Embeddings-English-R2 (Awasthy et al., 2025). QWEN3-Embedding-0.6B is a light weight version of QWEN3-Embedding model series, recently released model that displays SoTA performance on MTEB-multilingual evaluation split, while Granite-Embeddings-English-R2 is the most recently released model as of August 2025. In addition, Dewey-EN-Beta is another recent model that has a context window size of 131072, which is the longest among up-to-date existing models, demonstrating outstanding performance on long-context retrieval tasks, attaining SoTA performance on the MTEB LongEmbed benchmark split.

C.3 Metrics

In our experiments, to measure the precision and relevance of our retrieved candidates, we adopt three primary metrics for our main results. **Recall@K** is used as our main metric, which measures the ratio of correct candidates within the Top@k retrieved results. This aligns with our objective, where we seek to acquire a more relevant pool of papers using various aspect-aware queries. Moreover, we report **Normalized Discount Cumulative Gain**, **nDCG@K**, **Mean Average Precision**, **mAP@K**, and **MRR@K** to validate the robustness

of our method. Unlike the standard nDCG metric, which assigns higher weights to gold candidates with higher significance, we treat all gold candidates for a given query equally, assigning each a uniform weight of 1. Moreover, we report our results using relatively high K values (such as 100, 200, 300), to ensure the fairness of our evaluation suite. This is due to the large average number of annotated ground truth candidates per query in our self-constructed SCIFULLBENCH benchmark (see Table 7). Hence, to ensure that the retrieved candidates are capable of covering all the ground-truth candidates, we set K to a relatively high value.

In addition, we provide **Convex Hull Volume** (Graham, 1972) as a metric to quantitatively measure the dispersion of the embedding representations of retrieved candidates, which computes the volume of the smallest possible bounding convex closure of the data points after projecting the embedding vectors to the three-dimensional vector space, presenting the semantic coverage that each retrieval method presents. Notably, we employ convex hull volume in two different ways. As for the analysis provided in Table 4, we measured the relative convex hull volume ratio of ground-truth candidates present among the Top@300 retrieved results using the two equal samples (benchmark data with same query paper, ground truth candidates) used for binary comparison. This metric represents the relative semantic coverage of relevant documents from respective retrieval systems in a more fine grained manner. To ensure the fairness of our analysis, we excluded any sample that were unable to compute convex hull volume using retrieved ground truth candidates (due to its inability to retrieve sufficient number of ground-truth candidates) from either A2A, F2F, and COR on a given benchmark split, assuring that the comparisons were held using the same samples for respective three configurations. On the other hand, as for the experiments regarding Figure 9, we measured the absolute convex hull volume of total Top@k retrieved samples, to analyze the trend of semantic coverage with respect to retrieval depth progression.

C.4 Human Evaluation Setup

For sampling query papers for human study, we use NeurIPS, ACL, EMNLP split in SCIFULLBENCH where its query paper exists for both reference and citations split, and randomly sampled 15 different query papers. For A2A, F2F, and COR variant, we used result variants where jina-embeddings-v2-

¹⁷<https://huggingface.co/spaces/mteb/leaderboard>

base-en is utilized as backbone retriever, and DPO-trained QWEN-2.5-3B-Instruct model for query optimizers in CoR, with retrieval depth of 3, and utilized retrieved results from Citations (inward citation) split. When selecting the top-5 results for each variant, we excluded the 0th index result, since the same paper is mostly retrieved due to the structure of our target corpus. Moreover, we exclude ground truth candidates for reference and citations split, in order to measure the efficacy of CoR on candidates that are regarded as incorrect by our benchmark, and formulated top-k results for qualitative evaluation. For each paper we made two duplicate tasks for evaluation and originally hired 15 human evaluators studying the field of machine learning, who willingly consented to perform evaluation on two different query papers, for approximately 20 U.S. Dollars for respective evaluator. During the evaluation process, each participant is provided three retrieval variants with its order randomly shuffled, each variant consisting of 15 pairs of query, retrieved paper with title and abstracts. The participant is then instructed to choose one of the following (**Method**, **Experiment**, **Motivation/Research Question**, and **Irrelevant**) for every retrieved result, checking the relevance with the query paper. The details on human evaluation guideline can be found at Appendix H.

D Patent-to-Patent Retrieval

In this section, we provide detailed procedures for constructing the benchmark and experimental settings for the patent-to-patent retrieval task.

D.1 PatentFullBench Construction Workflow

Step 1) Collecting Raw Data: We initially obtained a full XML dump from the USPTO Open Data portal website¹⁸, and initially crawled US patents from 2020 to 2025 June 17th.

Step 2) Formatting Raw Data: We parsed XML dump from the USPTO Open Data portal website, obtaining publication number, publication date, application number, application date, title, abstract, inventors, priority claims, claims, descriptions, references (citations from given patent), and main classification for respective patents. As for citations and references, we initially filtered U.S patents exclusively, and ones where its date is over 2020, and obtained its classified country, date, document

Table 9: Statistics of query documents in PatentFullBench.

	References	Citations
Nation	United States	
# of documents per split	400	400
# of tokens per abstract	119.50	120.31
# of tokens per document	190335.46	186504.29
average # of candidates per sample	37.9	53.93

number (id), citation patterns (et al.), collecting a group of raw patents where each patent at least had 10 citations from the parsed XML dump.

Step 3) Initial Gold Label Annotation and Target Corpus Construction:

After forming a pool of raw extracted patent contents from XML dump data, we obtained patent invention title, abstracts, id, and date, using PatentsView API¹⁹ for respective cited patents from given pool of raw extracted content via querying its respective document ids. Since PatentsView API does not support full-content access, we made sure to exclude patents where its id did not belong in our pool of extracted content from **step 2**. Furthermore, we made sure to only include references where its returned year from PatentsView API was the same as the queried year, and where there exists a patent whose title matches the returned title from PatentsView API. After obtaining actual reference to full patents per input query document, we attempt to further label each query document with inward citations (incoming links) by labeling with gold documents that cite the query document, compared to the labels obtained from **step 3**, where the query document cites the gold labels. Using the reference information constructed from **step 3**, we further process citation relations through traversing all the patent files, and record the query patent B that cites A patent as B as ground-truth candidate for inward citation of A, by matching the exact patent title and document number. We ensured to check document number to strictly follow the information provided in the metadata. Since the ID format of the source in which we obtain relevance signals (inward, outward citation signal) and paper information (title, abstract, full paper) is homogeneous unlike SciFullBench construction, we directly use document numbers to label relevant documents, and use exact title-wise string-by-string match for finding relevant patent documents.

Subsequent to obtaining both inward citations and outward citations (references) for respective

¹⁸<https://data.uspto.gov/bulkdata/datasets>

¹⁹<https://patentsview.org/apis/purpose>

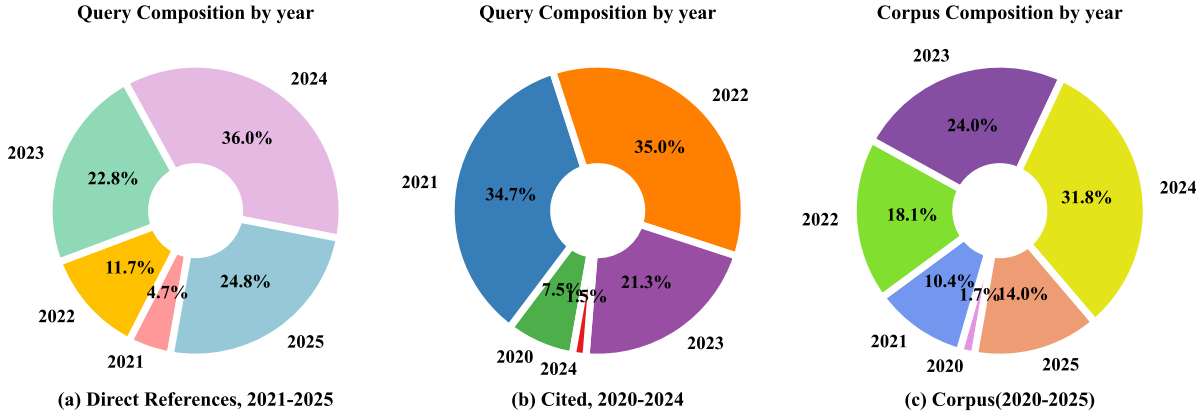


Figure 6: PatentFullBench composition by published year.

query patent documents, we removed duplicate labels that had either the same title or abstract content (considering the two labels identical if only one among title, abstract is the same, leaving only the first one as ground truth label), and ensure that the patent ids of ground-truth candidates match the ids of the ids of the raw patent information accumulated in **step 2**.

For respective splits, we excluded query documents from our benchmark set that had less than 10 ground truth candidates, and initially randomly sampled 800 query documents for respective splits. Successively, we formulated our target corpus via accumulating the labels that exist within our randomly sampled query document splits without duplicates, excluding documents that had either had the exact same title and abstracts (using the single, earliest found candidate with novel title, abstract information). Then, we made sure to filter out gold candidates within the initially sampled query documents that did not exist in the target corpus due to the exclusion of candidates within the target corpus. This process ensures that our potential gold candidates for queries in our benchmark always exist in the target corpus. Finally, we randomly sampled 400 queries that had at least 10 gold candidates after the above filtering step, and reformulated our target corpus making sure that duplicate candidates do not exist (considering the two candidates as identical if one of id, title, abstract are all same, using the) in our corpus. Finally, we obtain a corpus consisting a total of 5171 documents.

When comparing duplicates, we compare the content string-by-string without pre-processing, since the publication numbers (document id) from were the priority cue for finding matches, and exact string-by-string match is a conservative mech-

Table 10: Corpus statistics in PatentFullBench.

PATENTFULLBENCH	
Total # of documents	5150
Avg. # of tokens per abstract	120.97
Avg. # of tokens per document	29141.85
Total # of segmented corpus	52584
Avg. # of tokens per segment	2865.06

anism for filtering only the correct documents as relevant labels, and formulate ground-truth candidates. Since the ground truth candidates for each query is strictly different string-by-string, our formulated corpus also consists of candidates that are strictly different string-by-string.

Step 4) Formatting Full Context of Query and Candidate Documents and Final Filtering Successively, we further processed the full context of documents that are going to be either used as query or candidates for our evaluation. As for query documents, we concatenated title, abstract, claims, and description content, to formulate a full context for respective patents. Likewise the construction process of SciFullBench, we made sure to eliminate potential cues that reveal hindsight information on cited patents. First, using python regular expression, we eliminated information from the section within "**CROSS REFERENCE TO RELATED APPLICATIONS**", sentences that start with "**US**", ending with next line sign or phrase "**reference herein**". Moreover, using the citation information obtained from the XML dump in **step 2**, we eliminated citation patterns, and document numbers within the full context of patents. Finally, using the title and abstract information of constructed gold candidates from **Step 3**, we eliminated title and abstract information that exist in

1797 the full content of respective patents. During the
1798 processing of candidate documents, we followed
1799 the same procedure as query documents, where in-
1800 stead of using title and abstract information from
1801 the constructed gold labels, we removed every title
1802 and abstract information from the query documents
1803 that use the candidate document as its gold candi-
1804 dates, making sure that the candidate documents do
1805 not reveal direct hints on the relevant patents. After
1806 sanitizing the full documents, we segmented the
1807 full context of candidate documents for every 3000
1808 tokens (based on NLTK, similar to SciFullBench
1809 construction). Since each full-context of patents
1810 had partial overlap, we eliminated the patent candi-
1811 date from our corpus that had at least a single du-
1812 plicate partial chunk, and removed such patent candi-
1813 dates from the gold labels (candidates) of queries
1814 in our benchmark split. Finally, we constructed
1815 chunked corpus containing a total of 5150 different
1816 patent documents with both full context and title,
1817 abstracts, and 52584 chunks of patent documents
1818 along with 3000 token granularity. The statistics
1819 for PATENTFULLBENCH can be observed at Ta-
1820 ble 9 and Table 10.

1821 D.2 Experiment Details

1822 **Prompt Construction** We constructed prompts
1823 for patent-to-patent retrieval based on the three
1824 major aspects that comprise a patent: **Claim**, **Back-**
1825 **grounds**, and **Method**. **Method** aspect refers to
1826 technical details and implementation structure of
1827 the novel invention claimed by the given patent,
1828 while **Backgrounds** provide background informa-
1829 tion on the technical problems of prior works and
1830 its limited technical objectives. Moreover, **Claim**
1831 aspects comprise of the scope of its legal protection,
1832 and the main technical contributions of the inven-
1833 tion that respective patents intend to claim. The
1834 prompt construction process consumed roughly one
1835 hour to devise appropriate prompts for query op-
1836 timizers, with the help of ChatGPT to generate
1837 initial raw prompts and revising the content and
1838 structure with human feedback to ensure that the
1839 prompts are suitable for patent-to-patent retrieval
1840 featuring multiple aspects. For more details, please
1841 see Appendix I.

1842 **Inference Setting** We evaluate our patent to
1843 patent retrieval setting with the same domain-
1844 agnostic embedding models that we utilized for
1845 our evaluation, additionally reporting abstract-to-
1846 abstract retrieval results with domain-specific re-

1847 trievers, PAT-SPECTER and PaECTER (Ghosh
1848 et al., 2024), which is the most recent model for
1849 patent retrieval trained via minimizing the triplet
1850 loss using the binary patent abstract samples, simi-
1851 lar to the training methods employed in devising
1852 domain specific retrievers for paper retrieval in the
1853 academic domain, with its dataset sampled from the
1854 European Patent Office(EPO)²⁰ website, compris-
1855 ing of 30,000 English patents ranging from 1985
1856 to 2022. Moreover, we utilized proprietary LLM,
1857 GPT-4o-Mini-2024-0718 for query optimizer and
1858 truncated the input document context to 120000
1859 tokens using the tiktoken²¹ library since patent con-
1860 text length is significantly longer than scientific
1861 papers, occasionally exceeding the context window
1862 of GPT-4o-Mini-2024-0718 model. As for the hy-
1863 per parameters, we set the temperature to 0, same
1864 as mentioned in Appendix C.1.

1865 E Supplementary Experiments

1866 Here, we provide additional experiments to further
1867 support the results and analysis in the main paper.

1868 **Consistent Generalization Across Embedding**
1869 **Models** In our main results, we showed three
1870 variations of embedding models to demonstrate
1871 the effectiveness of our CoR framework. In this
1872 section, we further provide additional experiments
1873 on other existing long-context embedding models,
1874 namely **Granite-Embeddings-R2**, **QWEN3-0.6B-**
1875 **Embedding**, and **Dewey-en-Beta** in Table 11. The
1876 experimental results strongly support the robust-
1877 ness and general applicability of our research find-
1878 ings, as even without a trained model, by using pro-
1879 prietary **GPT-4o-2024-11-20** and **GPT-4.1-2025-**
1880 **04-14** with just a single round (depth) of retrieval,
1881 we were able to observe consistent performance
1882 improvements compared to naive F2F and A2A
1883 retrieval setups.

1884 The improvements achieved using Granite-
1885 Embeddings-R2 indicate that our framework can
1886 perform robustly under a context window of 8192
1887 as repeatedly demonstrated in our main results, and
1888 is compatible with even the most recently released
1889 model. Moreover, the results from the QWEN3-
1890 0.6B-Embedding model imply that our model can
1891 be applied to embedding models with longer con-
1892 text windows (32768) likewise Inf-Retriever-v1-
1893 1.5B in Table 1, as well as to popular state-of-
1894 the-art long-context embedding models for gen-

²⁰<https://www.epo.org/en>

²¹<https://github.com/openai/tiktoken>

Table 11: Performance of CoR with different LLM backbones and neural retrievers under single-round retrieval.

Model + Retriever	ICLR-Citations	NeurIPS-Citations
	Recall@300	Recall@100
Baseline		
Granite-Emb-Eng-R2 (A2A)	50.19	41.39
Granite-Emb-Eng-R2 (F2F)	54.00	44.37

QWEN3-0.6B-Emb (A2A)	51.47	41.76
QWEN3-0.6B-Emb (F2F)	54.03	44.81

Dewey-en-Beta (A2A)	47.76	40.54
Dewey-en-Beta (F2F)	53.15	42.06
Ours (CoR)		
GPT-4o + Granite-Emb-Eng-R2	55.64	47.22
GPT-4o + QWEN3-0.6B-Emb	56.94	47.15
GPT-4.1 + Dewey-en-Beta	56.73	47.99

Table 12: Performance analysis of ours when the original abstract to abstract retrieval is not utilized in our pipeline.

Retrieval Method	ACL-Citations	ICLR-Citations
	Recall@200	Recall@200
Baseline		
JEmb-v2 (A2A)	35.33	39.19
JEmb-v2 (F2F)	35.88	42.44

BGE-M3 (A2A)	33.76	35.16
BGE-M3 (F2F)	32.64	36.06

Inf-Ret-v1 (A2A)	41.77	45.42
Inf-Ret-v1 (F2F)	38.86	43.00
Ours w/o A2A		
QWEN-2.5-3B-Inst + JEmb-v2	41.13	44.41
Llama-3.2-3B-Inst + BGE-M3	35.83	37.49
Llama-3.2-3B-Inst + Inf-Ret-v1	47.05	50.65

eral domain-agnostic multilingual retrieval tasks in MTEB benchmark.

Furthermore, we provide evaluation results on Dewey-En-Beta, which is an embedding model with currently the longest context window among existing models, while demonstrating State-of-the-Art performance on LongEmbed benchmark. This suggests that our paper-retrieval framework is applicable to models with the longest context windows up-to-date, indicating the effectiveness and robustness of our suggested method to handle long context retrieval.

Robustness without Abstract-to-Abstract Retrieval Meanwhile, in Table 12, the results are reported when abstract-to-abstract retrieval (A2A) is not utilized within CoR. Despite A2A setup being omitted from its original multi-aspect queries, consistent improvement over baselines can be observed, although not to the extent when abstract-abstract retrieval is integrated into our pipeline, indicating the positive, balancing effect that A2A retrieval presents, along with robustness of our method. This in turn highlights the effectiveness of our pipeline in generating multiple aspect-aware queries in the absence of complementary abstract to abstract retrieval, strongly supporting our research contributions.

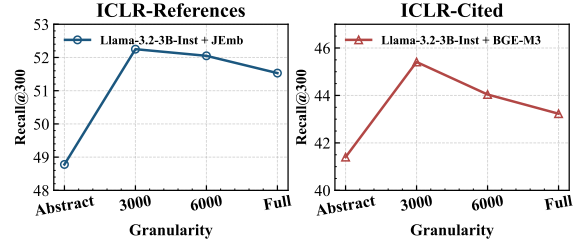


Figure 7: Results across candidate paper granularities.

Comparison with Naive Chaining Table 13 shows the benefits of employing multi-round retrieval with CoR compared to that of naive chain-of-retrieval methods using Full Paper-to-Full Paper and Abstract-to-Abstract retrieval setup. The results indicate that our aspect-guided chain-of-retrieval along with post-order recursive aggregation is more effective, and trivial iterative retrieval is not sufficient. Note that the baseline experiments were conducted with under the same hyper parameter setting, with a cache to avoid redundant sampling likewise our framework, with 300 candidates retrieved and aggregated via Reciprocal Rank Fusion. This shows that our approach in exploration and aggregation guided by aspect-centric queries is more effective compared to iterative top-k retrieval using raw content as queries, supporting the validity of our proposed retrieval system.

Starting Index Hyper Parameter Ablation In this section, we provide experiments on the hyper parameter selection starting index γ that we used for **Next Query Selection Algorithm**, to choose the starting point of choosing the representative paper with highest semantic similarity. In our case, value 1 is used to report our results in the main paper (always start from the second highest rank), to select the paper with highest similarity while avoiding the same papers as input query that may exist in our corpus (which inevitably appears at rank 0, since same papers display 100% semantic similarity). To further validate the benefits of selecting paper with highest-semantic similarity for subsequent retrieval, variations of CoR with different selection starting index γ is reported in Figure 8. The results demonstrate steady decrease in performance as selection starting index is increased, implying that chain-of-retrieval using the nearest-possible paper for next query selection is more desirable compared to vice versa, due to the semantic drift that occurs when relatively unrelated papers are selected as query for subsequent round of retrieval.

Table 13: Ablation study on the effect of our search framework compared to naive chain-of retrieval methods based on A2A and F2F retrieval setting.

Model + Retriever	ICLR-References	ICLR-Citations
	Recall@100	Recall@200
Baseline		
JEmb-v2 (A2A)	33.04	39.10
JEmb-v2 (F2F)	33.19	41.54

BGE-M3 (A2A)	28.71	34.72
BGE-M3 (F2F)	31.00	36.12

Inf-Ret-v1 (A2A)	42.50	44.67
Inf-Ret-v1 (F2F)	36.67	42.38

Ours (w/o DPO)		
Llama-3.2-3B-Inst + JEmb-v2	38.34	46.73
Llama-3.2-3B-Inst + BGE-M3	32.62	41.84
Llama-3.2-3B-Inst + Inf-Ret-v1	46.74	52.42

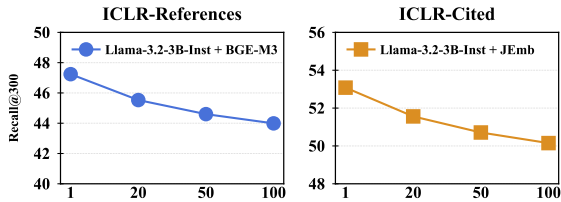


Figure 8: Hyperparameter Ablation Study on the effect of selection starting index γ on the overall performance. We use DPO-trained query optimizers on 3 rounds of retrieval.

Corpus Granularity Ablation Study Furthermore, in order to see how the granularity of candidate representations affects the performance on our multi-vector retrieval system, we conduct further analysis. The results in Figure 7 demonstrate the effectiveness of our multi-vector approach: segmenting each candidate document into 3K-token chunks yields the best performance, which outperforms not only abstract-only and full-document single-vector representations but also coarser chunking strategies (such as 6K-token segments), suggesting that finer-grained representations are effective in capturing diverse and localized signals within full-length papers, which enables our query to capture a more multi-faceted dimension of target papers, resulting in an enhanced performance.

Effect of Reciprocal Rank Fusion We conduct an additional ablation study on the effect of RRF (Reciprocal Rank Fusion) in our hybrid retrieval system (single round of COR), compared to embedding-level merging approach in prior multi-vector retrieval approaches. We primarily compare with two traditional approaches in embedding-level merging strategy, the naive aggregation strategy that forcefully computes similarities (in our case L2 distance) between all the sub-vectors of respective query and candidates and naively sums it up to acquire similarity score of original query and documents. In our experiments, to mitigate the unfair penalization due to document length, we normal-

Table 14: Analysis on the impact of incorporating rank fusion to attain unified results of our hybrid retrieval framework on the EMNLP-Citations split.

Retrieval Method	Recall@100	MRR@50
Ours w/o RRF & w/o A2A		
Naive Aggregation		
GPT-4o-Mini-2024-0718 + BGE-M3	21.61	32.61
Late Interaction (MaxSim)		
GPT-4o-Mini-2024-0718 + BGE-M3	26.00	31.53

Ours w/o A2A		
GPT-4o-Mini-2024-0718 + BGE-M3	26.88	34.54

Table 15: Performance comparison against baselines augmented with metadata (i.e., author information).

Retrieval Method	ACL-References	ACL-Citations
	Recall@200	Recall@200
Baseline w/ Author Information		
JEmb-v2 (A2A)	34.25	34.84
JEmb-v2 (F2F)	36.54	36.04

BGE-M3 (A2A)	32.16	33.73
BGE-M3 (F2F)	31.62	32.66

Inf-Ret-v1 (A2A)	46.66	42.36
Inf-Ret-v1 (F2F)	44.31	38.92

Ours (COR)		
Llama-3.2-3B-Inst + JEmb-v2	39.81	42.93
Llama-3.2-3B-Inst + BGE-M3	35.75	39.48
Llama-3.2-3B-Inst + Inf-Ret-v1	49.30	48.23

ize the aggregated results with the total number of subvectors used for similarity calculation.

Meanwhile, we also present comparison with late interaction strategy (Khattab and Zaharia, 2020; Santhanam et al., 2021), where only the maximum similarity for subquery and its corresponding sub-documents is aggregated to form original query, document similarity (in our case, minimum l2 distance). Table 14 illustrates a performance drop when queries optimized in various aspects are used as subvector representations of original documents and are used to compute a single similarity value when matched with the segmented corpus of target candidates. This supports the validity of our design choice, where our ranking-merging system is more robust to noise, while allowing candidates that are strongly aligned in one aspect but unaligned otherwise to be retrieved as shown in two metrics Recall@100 and MRR@50, demonstrating the effectiveness in surfacing the most relevant candidates among the top-k candidates while assuring that more relevant candidates are retrieved, supporting our motivation for incorporating RRF compared to vector-merging approaches.

Additional Diversity Analysis In this section, we provide further analysis on the improved diversity of semantics that our retrieval system presents. Figure 9 presents a linearly increasing trend where the Top@300 absolute Convex Hull Volume of 400 samples per benchmark split proportionally increases with respect to the retrieval depth, indi-

1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024

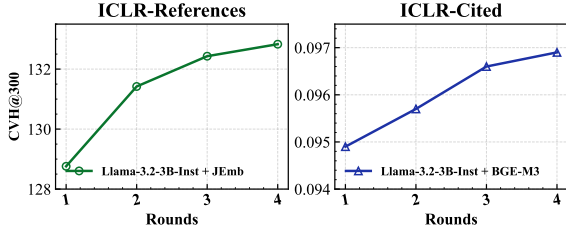


Figure 9: Diversity Analysis with retrieval depth progression. Note that we perform evaluation using DPO-trained Query Optimizers.

Table 16: Latency comparison across retrieval depths. The results are reported using jina-embeddings-v2-base-en, and untrained Llama-3.2-3B-Instruct query optimizer.

Depth	A2A	F2F	COR
1	0.38 (s)	1.35 (s)	35.29 (s)
2	0.67 (s)	2.28 (s)	189.99 (s)
3	0.54 (s)	3.12 (s)	493.45 (s)

cating that our Chain of Retrieval framework is able to diversify the semantics of search results as the rounds progress. Moreover, along with the quantitative metrics to represent the enhanced coverage that COR presents, we visualize the distribution of the embedding vectors of Top@300 retrieved documents for eight sampled queries on a two-dimensional plane using **t-distributed Stochastic Neighbor Embedding (t-SNE)** (Van der Maaten and Hinton, 2008) and its distribution boundaries via a two-dimensional convex hull area of the projected data points in Figure 10. We simultaneously visualize the distribution of t-SNE data points for both ground truth candidates, and baseline retrieved results (including A2A retrieved results and F2F results using the same embedding model), and retrieved candidates from COR. The t-SNE Convex Hull visualization results also advocate our claim, where our method typically displays a more dispersed embedding distribution compared to dispersion of retrieved document embeddings from baselines. In addition, the cases in Figure 10 demonstrate that through improved diversity within the retrieved pool of papers, COR is able to retrieve more relevant papers that exist in ground-truth candidates, as the coverage of the documents has increased compared to previous retrieval approaches.

Latency Analysis While the primary objective of our work is to improve retrieval effectiveness rather than efficiency, we also report the latency of COR in Table 16, measured on a single NVIDIA RTX A6000 GPU. Although latency reduction lies

Table 17: Results of applying COR to the A2A retrieval setting on SCIFULLBENCH with two retrieval rounds.

Retrieval Method	ICLR-References	ICLR-Citations
	Recall@300	Recall@300
Baseline		
JEmb-v2 (A2A)	48.38	44.49
BGE-M3 (A2A)	41.95	39.61
Inf-Ret-v1 (A2A)	59.66	51.56
Ours (CoR)		
GPT-4o-Mini + JEmb-v2	49.47	46.41
GPT-4o-Mini + BGE-M3	43.12	41.92
GPT-4.1 + Inf-Ret-v1	60.95	53.27

Table 18: Performance analysis when a single comprehensive query covering the three aspects is used.

Retrieval Method	ACL-Citations	ICLR-Citations
	Recall@200	Recall@200
Single Comprehensive Query w/o A2A		
QWEN-2.5-3B-Inst + JEmb-v2	38.66	41.68
Llama-3.2-3B-Inst + BGE-M3	32.45	34.03
Llama-3.2-3B-Inst + Inf-Ret-v1	45.36	48.55
Ours (CoR) w/o A2A		
QWEN-2.5-3B-Inst + JEmb-v2	41.13	44.41
Llama-3.2-3B-Inst + BGE-M3	35.83	37.49
Llama-3.2-3B-Inst + Inf-Ret-v1	47.05	50.65

Table 19: Performance of COR on the SciDocs benchmark under the A2A setup; statistically meaningful gains are underlined based on paired t-tests.

Retrieval Method	Recall@10	mAP@10	nDCG@10
A2A Retrieval			
JEmb-v2 (A2A)	92.46 ± 0.00	85.53 ± 0.00	91.66 ± 0.00
BGE-M3 (A2A)	90.36 ± 0.00	80.88 ± 0.00	88.70 ± 0.00
Inf-Ret-v1 (A2A)	95.69 ± 0.00	90.02 ± 0.00	94.44 ± 0.00
Ours (CoR) w/ A2A			
GPT-4.1 + JEmb-v2	93.20 ± 0.09	86.20 ± 0.08	92.18 ± 0.06
GPT-4.1 + BGE-M3	90.52 ± 0.07	81.15 ± 0.10	88.88 ± 0.04
GPT-4.1 + Inf-Ret-v1	96.16 ± 0.10	90.40 ± 0.03	94.80 ± 0.04

outside the scope of our work, exploring techniques such as intermediate-branch or search-space pruning would be an exciting direction for future work.

Generalization to the A2A Retrieval Setting

While our framework is primarily designed for the full-paper-to-full-paper retrieval setting, we additionally evaluate whether COR is applicable to the A2A setting. To this end, we report its performance on the A2A configuration of SCIFULLBENCH in Table 17, where title-abstract pairs are used as inputs and retrieval is performed over an abstract-only corpus, while keeping the same query-optimization prompts. The results demonstrate that COR is not limited to the full-paper-to-full-paper setting, but generalizes effectively to the A2A scenario, yielding meaningful performance gains.

Comparison against Metadata-Augmented Baselines

We also consider a setting where additional metadata is provided to the baselines by concatenating author information to both queries and can-

Table 20: Performance of COR on the SciDocs benchmark under the F2F setup.

Retrieval Method	Recall@5	Recall@10	Average
A2A Retrieval			
JEmb-v2 (A2A)	63.61 \pm 0.00	69.86 \pm 0.00	66.74
BGE-M3 (A2A)	55.56 \pm 0.00	65.97 \pm 0.00	60.77
Inf-Ret-v1 (A2A)	65.28 \pm 0.00	76.11 \pm 0.00	70.70
F2F Retrieval			
JEmb-v2 (F2F)	65.97 \pm 0.00	71.94 \pm 0.00	68.96
BGE-M3 (F2F)	53.75 \pm 0.00	62.92 \pm 0.00	58.34
Inf-Ret-v1 (F2F)	65.42 \pm 0.00	73.89 \pm 0.00	69.66
Ours (COR)			
GPT-4.1 + JEmb-v2	66.25 \pm 1.58	71.53 \pm 1.10	68.89
GPT-4.1 + BGE-M3	58.33 \pm 0.64	66.81 \pm 0.96	62.57
GPT-4.1 + Inf-Ret-v1	66.25 \pm 1.46	77.50 \pm 0.96	71.88

didate documents. As shown in Table 15, CoR remains effective under this comparison and continues to outperform the baselines even when they are augmented with author metadata.

Comparison with a Single Aspect-Aware Comprehensive Query To examine the effectiveness of our multiple aspect-aware query design, we compare CoR against a baseline that uses a single comprehensive query covering all three aspects simultaneously. As shown in Table 18, CoR with multiple aspect-specific queries outperforms the single-query baseline, indicating the benefit of separating queries by aspect in paper-to-paper retrieval.

Effectiveness on the Existing Benchmark (A2A)

While our main experiments are conducted on SCIFULLBENCH, which provides a full-paper-to-full-paper retrieval setting with more recent publications, we also examine whether our findings hold under previously established benchmarks. To this end, we additionally evaluate CoR on the SciDocs-Cocite split (Cohan et al., 2020), which is originally defined under an A2A retrieval setting, using 993 queries whose abstract content is available²². As shown in Table 19, CoR remains effective under the existing benchmark with A2A setting, yielding statistically meaningful improvements over baseline methods.

Effectiveness on the Existing Benchmark (F2F)

In addition to evaluating CoR on the existing A2A benchmark setting of SciDocs, we further extend this benchmark to a F2F configuration. For this purpose, we construct a SciDocs-Full variant from the Cocite, CoRead, and CoView splits of SciDocs (Cohan et al., 2020), where both queries and candidate documents contain full-paper content, following the same preprocessing pipeline as

SCIFULLBENCH introduced in Appendix A. As shown in Table 20, CoR yields consistent improvements over baselines on this extended version of the benchmark, indicating that our method remains also effective across different retrieval configurations within the established benchmarks.

Case Study We provide the qualitative examples of the retrieved results from SCIFULLBENCH and PATENTFULLBENCH in Table 23 and Table 24, respectively, as well as the examples of the generated queries from them, in Appendix G.

²²<https://huggingface.co/datasets/allenai/scirepeval>

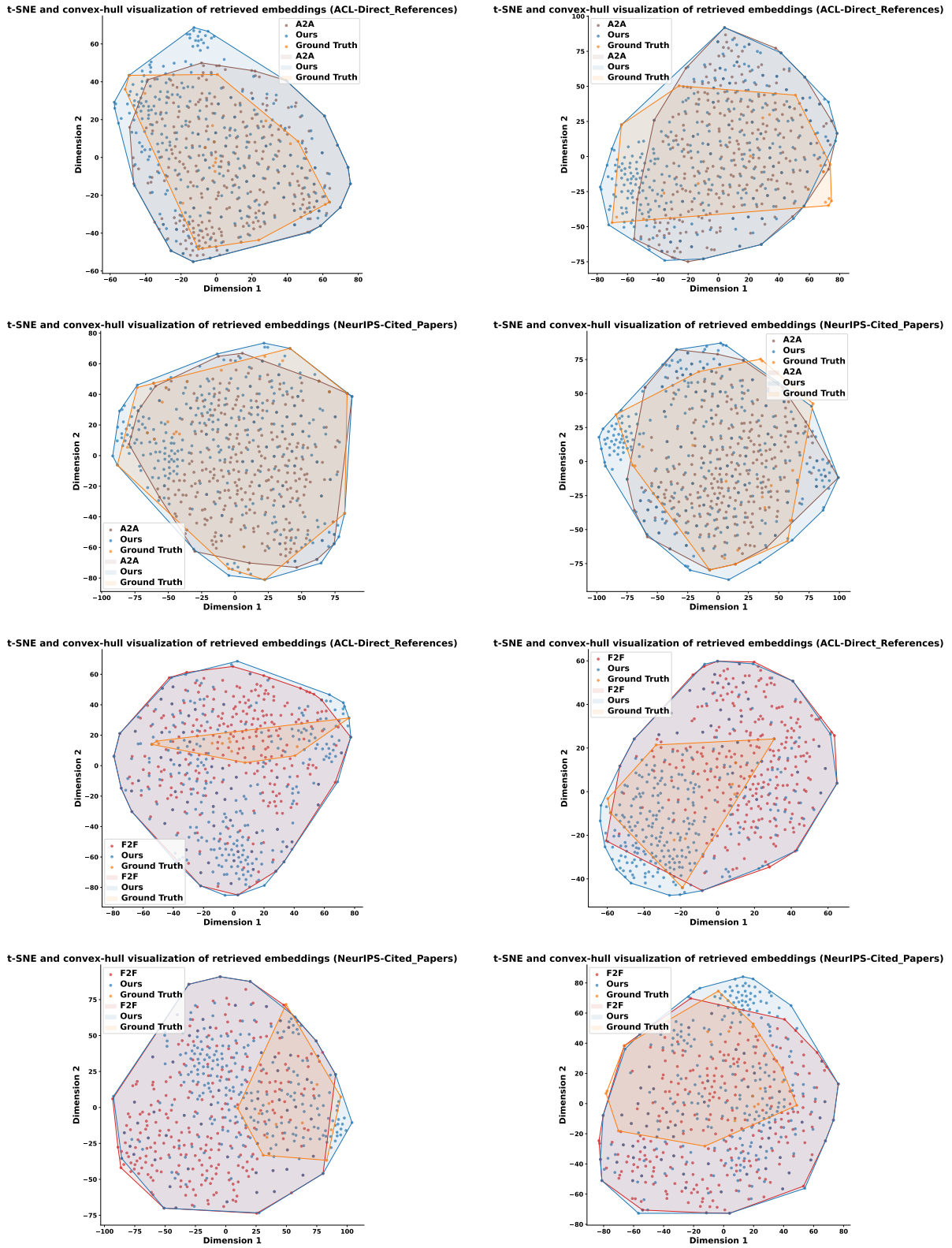


Figure 10: t-SNE and convex hull boundary visualization of Top-300 retrieved document embeddings across eight queries from respective benchmark splits. We report the results with Jina-Embeddings-v2-Base-EN as the embedding model.

Table 21: Main results with additional metrics.

IR Method	ICLR				NeurIPS			
	References		Citations		References		Citations	
	nDCG@300	Recall@300	nDCG@300	Recall@300	nDCG@300	Recall@300	nDCG@300	Recall@300
Lexical-Based Retrievers								
BM-25 (A2A)	24.78	34.69	23.35	32.59	31.79	41.73	31.68	42.22
BM-25 (F2F)	31.09	44.54	34.57	46.28	21.59	29.49	38.38	48.65
Domain-Specific Retriever								
SciNCL-A2A	34.04	50.51	29.57	44.13	38.44	53.08	36.57	51.67
SPECTER2-Base-A2A	32.71	48.99	30.97	45.23	37.77	51.82	37.99	52.95
SPECTER2-Adapter-MTL CTRL-A2A	33.91	49.98	29.49	43.29	38.54	52.16	36.75	51.46
SciMult-MHAAExpert-A2A	28.26	42.56	24.61	36.88	33.64	47.52	32.03	45.25
Jina-Embeddings-v2-BASE-EN								
A2A (Abstract-to-Abstract)	34.43	48.38	30.47	44.49	39.12	51.45	37.22	51.28
F2F (Full-to-Full)	35.27	49.58	33.71	47.85	38.82	51.62	37.89	51.41
CoR w/ Llama-3.2-3B-Instruct (w/ DPO)	37.19	54.60	36.23	53.08	40.56	56.40	40.93	59.33
CoR w/ QWEN-2.5-3B-Instruct (w/ DPO)	36.82	54.37	36.55	53.59	41.00	57.19	42.17	60.35
BGE-M3								
A2A (Abstract-to-Abstract)	29.51	41.95	27.15	39.61	34.65	46.22	33.48	46.08
F2F (Full-to-Full)	31.78	44.10	29.06	41.13	35.63	47.57	35.38	46.89
CoR w/ Llama-3.2-3B-Instruct (w/ DPO)	31.21	47.24	31.96	46.53	36.40	51.53	37.12	52.01
CoR w/ QWEN-2.5-3B-Instruct (w/ DPO)	32.10	48.16	32.94	47.97	37.58	52.69	38.50	53.57
Inf-Retriever-v1-1.5B								
A2A (Abstract-to-Abstract)	43.03	59.66	35.40	51.56	48.64	63.77	41.81	57.70
F2F (Full-to-Full)	39.01	53.85	34.06	47.89	23.17	30.69	37.09	49.69
CoR w/ Llama-3.2-3B-Instruct (w/ DPO)	44.65	63.19	40.66	58.41	49.96	67.43	46.66	64.66
CoR w/ QWEN-2.5-3B-Instruct (w/ DPO)	44.65	62.93	41.36	59.24	49.77	66.62	47.06	64.88
ACL								
IR Method	References		Citations		References		Citations	
	nDCG@300	Recall@300	nDCG@300	Recall@300	nDCG@300	Recall@300	nDCG@300	Recall@300
Lexical-Based Retrievers								
BM-25 (A2A)	21.47	31.80	21.90	30.74	20.57	29.78	19.82	28.00
BM-25 (F2F)	22.77	34.64	27.47	37.95	24.00	35.55	28.69	39.21
Domain-Specific Retriever								
SciNCL-A2A	26.72	41.63	26.29	39.79	25.52	39.78	24.86	38.02
SPECTER2-Base-A2A	25.82	40.48	27.61	41.19	24.32	37.55	26.09	39.22
SPECTER2-Adapter-MTL CTRL-A2A	25.76	40.13	26.31	39.61	24.87	37.59	24.65	37.47
SciMult-MHAAExpert-A2A	24.19	37.49	23.76	35.31	22.02	34.75	21.38	32.26
Jina-Embeddings-v2-BASE-EN								
A2A (Abstract-to-Abstract)	25.50	38.80	27.66	40.52	24.41	37.00	25.23	37.32
F2F (Full-to-Full)	27.27	41.66	28.80	40.95	27.83	41.73	27.87	40.34
CoR w/ Llama-3.2-3B-Instruct (w/ DPO)	28.52	44.92	32.76	48.05	28.15	43.96	30.84	46.00
CoR w/ QWEN-2.5-3B-Instruct (w/ DPO)	28.62	45.67	33.40	49.00	27.80	43.99	31.68	47.60
BGE-M3								
A2A (Abstract-to-Abstract)	23.65	35.81	25.95	38.40	22.25	33.29	22.87	34.15
F2F (Full-to-Full)	24.39	36.67	25.78	37.01	23.88	35.08	23.92	34.73
CoR w/ Llama-3.2-3B-Instruct (w/ DPO)	25.12	40.19	30.18	44.46	24.12	38.24	27.55	40.87
CoR w/ QWEN-2.5-3B-Instruct (w/ DPO)	25.77	40.71	30.44	45.29	24.66	38.97	28.22	42.14
Inf-Retriever-v1-1.5B								
A2A (Abstract-to-Abstract)	34.17	51.51	32.20	47.33	32.68	48.68	29.68	43.72
F2F (Full-to-Full)	33.35	49.59	30.61	43.72	32.37	46.40	30.02	43.57
CoR w/ Llama-3.2-3B-Instruct (w/ DPO)	35.09	55.14	36.98	54.13	33.86	52.11	35.37	51.77
CoR w/ QWEN-2.5-3B-Instruct (w/ DPO)	35.16	54.68	37.74	54.78	33.65	51.91	35.99	52.40

Table 22: Main results with additional metrics.

IR Method	ICLR				NeurIPS			
	References		Citations		References		Citations	
	Recall@100	Recall@200	nDCG@200	Recall@100	Recall@100	Recall@200	nDCG@200	Recall@100
Lexical-Based Retrievers								
BM-25 (A2A)	24.51	30.77	21.83	22.57	30.86	37.60	30.08	31.17
BM-25 (F2F)	31.66	39.59	32.72	33.01	20.19	25.07	36.84	37.57
Domain-Specific Retriever								
SciNCL-A2A	34.92	44.55	27.33	29.49	38.10	47.37	34.40	36.90
SPECTER2-Base-A2A	33.47	43.43	28.66	30.32	36.65	46.15	36.05	37.74
SPECTER2-Adapter-MTL CTRL-A2A	34.51	43.72	27.43	29.13	37.30	46.60	34.70	37.24
SciMult-MHAAExpert-A2A	28.44	37.13	22.57	23.51	33.01	41.98	30.14	31.66
Jina-Embeddings-v2-BASE-EN								
A2A (Abstract-to-Abstract)	34.71	43.35	28.38	30.09	37.91	46.30	35.16	36.85
F2F (Full-to-Full)	35.17	44.29	31.60	33.25	37.25	46.15	35.88	37.43
CoR w/ Llama-3.2-3B-Instruct (w/ DPO)	38.36	48.71	33.94	36.50	41.16	50.57	38.57	42.55
CoR w/ QWEN-2.5-3B-Instruct (w/ DPO)	38.52	48.33	34.16	36.52	41.36	51.07	39.91	43.81
BGE-M3								
A2A (Abstract-to-Abstract)	29.69	37.23	25.44	26.66	33.28	41.65	31.69	33.40
F2F (Full-to-Full)	32.46	39.41	27.12	28.50	34.81	42.98	33.57	34.31
CoR w/ Llama-3.2-3B-Instruct (w/ DPO)	32.88	42.47	29.99	32.13	36.47	46.17	35.29	37.58
CoR w/ QWEN-2.5-3B-Instruct (w/ DPO)	33.76	43.28	30.95	33.07	37.34	47.40	36.50	39.19
Inf-Retriever-v1-1.5B								
A2A (Abstract-to-Abstract)	44.49	53.91	33.03	35.03	48.20	57.89	39.88	42.97
F2F (Full-to-Full)	39.21	48.37	32.15	34.65	22.95	28.06	35.35	37.46
CoR w/ Llama-3.2-3B-Instruct (w/ DPO)	47.15	57.14	38.26	41.44	50.87	61.38	44.45	48.48
CoR w/ QWEN-2.5-3B-Instruct (w/ DPO)	46.81	57.22	38.83	41.99	50.78	60.62	44.91	49.11
ACL								
IR Method	References		Citations		References		Citations	
	Recall@100	Recall@200	mAP@30	Recall@200	Recall@100	Recall@200	mAP@30	Recall@200
Lexical-Based Retrievers								
BM-25 (A2A)	23.24	28.45	4.82	27.61	21.50	26.44	4.36	24.64
BM-25 (F2F)	23.68	30.25	7.24	33.81	24.92	31.12	7.44	34.90
Domain-Specific Retriever								
SciNCL-A2A	29.18	36.62	5.67	34.69	27.23	34.55	4.87	33.05
SPECTER2-Base-A2A	28.28	35.81	6.21	36.38	25.47	33.12	5.33	34.14
SPECTER2-Adapter-MTL CTRL-A2A	28.42	35.73	5.81	34.77	26.28	33.14	4.88	32.36
SciMult-MHAAExpert-A2A	25.96	33.07	5.11	30.78	23.37	30.32	4.13	27.35
Jina-Embeddings-v2-BASE-EN								
A2A (Abstract-to-Abstract)	27.75	34.49	6.20	35.33	26.08	32.64	5.24	32.43
F2F (Full-to-Full)	28.78	36.46	7.03	35.88	29.51	36.98	6.32	35.51
CoR w/ Llama-3.2-3B-Instruct (w/ DPO)	31.96	39.81	7.67	42.93	30.76	39.21	6.57	40.31
CoR w/ QWEN-2.5-3B-Instruct (w/ DPO)	32.02	40.41	7.92	43.75	30.83	38.96	6.66	41.86
BGE-M3								
A2A (Abstract-to-Abstract)	25.88	32.01	5.56	33.76	23.50	29.36	4.57	29.63
F2F (Full-to-Full)	25.15	32.21	5.98	32.64	24.81	30.84	5.24	30.08
CoR w/ Llama-3.2-3B-Instruct (w/ DPO)	28.43	35.75	6.84	39.48	26.19	33.31	5.78	35.63
CoR w/ QWEN-2.5-3B-Instruct (w/ DPO)	28.65	36.33	6.75	40.22	26.77	34.28	6.01	37.01
Inf-Retriever-v1-1.5B								
A2A (Abstract-to-Abstract)	37.55	46.24	7.65	41.77	34.73	43.40	6.35	39.00
F2F (Full-to-Full)	36.35	44.37	7.81	38.86	34.65	41.83	6.98	38.92
CoR w/ Llama-3.2-3B-Instruct (w/ DPO)	39.44	49.3	9.13	48.23	37.26	46.81	8.05	46.29
CoR w/ QWEN-2.5-3B-Instruct (w/ DPO)	39.73	49.31	9.60	48.82	36.80	46.37	8.38	46.51

2124 **F Algorithm Details**

2125 In this section, we provide further details on the
2126 functional roles of our COR framework with formal
2127 algorithms. Specifically, Algorithm 2 defines
2128 an aspect-aware retrieval process from Section 3.2,
2129 where the query optimizers decompose the input
2130 document into aspect-aware queries for exploration
2131 and then the papers are retrieved from multi-vector
2132 corpora. Algorithm 3 describes the aspect-aware
2133 next query selection process in Section 3.3. Algo-
2134 rithm 4 shows the overall exploration process of
2135 CoR, using Algorithm 2 and Algorithm 3. Finally,
2136 Algorithm 5 explains the Post-Order Aggregation
2137 algorithm of CoR explained in Section 3.3.

Algorithm 2 Aspect-Aware Multi-Vector Retrieval

```
1: Require: Input paper  $D$ ; Input paper abstract  $D_{\text{abstract}}$ ; Top- $k$  per query  $K$ ;  
2: Name of the Parent Retrieved Results PARENT; Corpora  $\mathcal{C} \leftarrow \{\mathcal{C}_{\text{abstract}}, \mathcal{C}_{\text{chunked}}\}$ ;  
3:  
4: Initialize: Functional aspect-aware LLM query optimizer agents  $\mathcal{F} = \{f_M, f_E, f_R\}$   
5:  
6: function ONEHOPRETRIEVAL( $D$ , PARENT,  $\mathcal{C}$ ,  $K$ )  
7:    $R \leftarrow []$  ▷ Memory to save aspect-aware retrieved results per document  
8:   for all  $f_i \in \mathcal{F}$  do  
9:     NAME  $\leftarrow$  PARENT  $\parallel f_i$  ▷ Update name for aspect-aware retrieval for document D  
10:     $q \leftarrow f_i(D)$   
11:     $T \leftarrow$  Retrieve( $q$ ,  $\mathcal{C}_{\text{chunked}}$ ,  $K$ ) ▷ Top-K Retrieved Results from chunked corpus  
12:     $\mathcal{R}_q \leftarrow [h(x) \mid x \in T]$  ▷ Retrieved Results mapped back to paper  
13:     $R.\text{APPEND}((\mathcal{R}_q, \text{PARENT}, \text{NAME}, f_i))$   
14:  NAME  $\leftarrow$  PARENT  $\parallel$  abstract  
15:   $q \leftarrow D_{\text{abstract}}$   
16:   $T \leftarrow$  Retrieve( $q$ ,  $\mathcal{C}_{\text{abstract}}$ ,  $K$ ) ▷ Top-K Retrieved Results from abstract corpus  
17:   $\mathcal{R}_q \leftarrow [h(x) \mid x \in T]$  ▷ Retrieved Results mapped to paper  
18:   $R.\text{APPEND}((\mathcal{R}_q, \text{PARENT}, \text{NAME}, \text{abstract}))$   
19:  return  $R$ 
```

Algorithm 3 Aspect-Aware Next Query Selection

```
1: Require: Root paper  $D$ ; Max depth  $R$ ; selection starting index  $\gamma$ ; Top- $k$  per query  $K$   
2: Aspect Aware Cache CACHE; memory to save query paper for next round retrieval  $Q_{\text{next}}$   
3:  
4: function INITIALASPECT(NAME) ▷ Obtain Initial Branching Aspect of the root document  
5:   Parse NAME as (ROOT  $\parallel a_1 \parallel a_2 \parallel \dots$ )  
6:   return  $a_1$   
7:  
8: function SELECTNEXTQUERY( $P$ , CACHE,  $\gamma$ )  
9:   for all  $r \in P$  do  
10:    TOPK, PARENT, NAME, PREVIOUS ASPECT  $\leftarrow r$   
11:    if PREVIOUS ASPECT = abstract then  
12:      CONTINUE ▷ Skip selection if previous retrieved results were from abstract query  
13:      
14:    INITIAL ASPECT  $\leftarrow$  INITIALASPECT(NAME)  
15:    for  $j = \gamma$  to  $K$  do  
16:      DOC  $\leftarrow$  TOPK[ $j$ ]  
17:      if DOC  $\notin$  CACHE[INITIAL ASPECT] then  
18:        CACHE[INITIAL ASPECT].append(DOC)  
19:         $Q_{\text{next}}.\text{APPEND}((\text{DOC}, \text{NAME}))$   
20:        break  
21:      
22:      
23:    ▷ Select most similar document not in aspect-aware cache starting from the  $\gamma$  index.  
24:
```

Algorithm 4 Chain-of-Retrieval

```
1: Require: Root paper  $D$ ; Max depth  $R$ ; selection starting index  $\gamma$ ; Top-k per query  $K$ ; Corpora  
    $\mathcal{C} \leftarrow \{\mathcal{C}_{\text{abstract}}, \mathcal{C}_{\text{chunked}}\}$   
2:  
3: Initialize:  
4: query queue  $Q \leftarrow [(D, \text{ROOT})]$   
5: memory to save retrieved results  $\mathcal{M} \leftarrow []$   
6:  $\text{CACHE}[a] \leftarrow \emptyset \quad \forall a \in \{\text{M}, \text{E}, \text{R}\}$   $\triangleright$  Initialize Aspect-Aware Cache for Method, Experiment, and  
   Research Question Aspect  
7:  
8: for  $r = 0$  to  $R - 1$  do  
9:    $\mathcal{M}[r] \leftarrow [], Q_{\text{next}} \leftarrow []$   
10:  for all  $(d, \text{PARENT})$  in  $Q$  do  
11:     $P \leftarrow \text{ONEHOPRETRIEVAL}(d, \text{PARENT}, \mathcal{C}, K)$   
12:     $\mathcal{M}[r].\text{EXTEND}(P)$   
13:     $NQ \leftarrow \text{SELECTNEXTQUERY}(NQ, P, \text{CACHE}, \gamma)$   
14:     $Q \leftarrow NQ$ 
```

Algorithm 5 Recursive Post-Order Merging

```
1: Require: Root paper  $D$ ; Max depth  $R$ ; All Retrieved Results after depth  $R$  of retrieval  $\mathcal{M}$ ;  
2:  
3: Ensure: Final top- $k$  candidates for  $D$   
4:  
5: function  $\text{GROUP}(L)$   $\triangleright$  Group list of per round retrieved results with same parent nodes  
6:   return  $\{p \mapsto [x \in L : \text{Parent}(x) = p]\}$   
7:  
8: function  $\text{MERGESIBLINGS}(L)$   $\triangleright$  Merge retrieved results with same parent nodes  
9:    $S \leftarrow \{\}$   
10:  for  $(\text{PARENT}, G) \in \text{GROUP}(L)$  do  
11:     $S[\text{PARENT}] \leftarrow \text{RRF}(G)$   
12:  return  $S$   
13:  
14: function  $\text{MERGEEDGES}(P, S)$   $\triangleright$  Merge retrieved results from parent nodes and merged child nodes  
15:    $O \leftarrow []$   
16:  for  $(\text{TopK}, \text{PARENT}, \text{NAME}) \in P$  do  
17:    if  $\text{PARENT} \in \text{Keys}(S)$  then  
18:       $O.\text{append}(\text{RRF}(\text{TopK}, S[\text{PARENT}])))$   
19:  return  $O$   
20:  
21:  $r \leftarrow R - 1$   
22:  
23: while  $r > 0$  do  $\triangleright$  Post-Order Merge, starting from the bottom level nodes  
24:    $S \leftarrow \text{MERGESIBLINGS}(\mathcal{M}[r])$   
25:    $\mathcal{M}[r - 1] \leftarrow \text{MERGEEDGES}(\mathcal{M}[r - 1], S)$   
26:    $r \leftarrow r - 1$   
27: return  $\text{RRF}(\mathcal{M}[0])$ 
```

G Case Study

Table 23: Example of the retrieved documents for the input document with CoR from a EMNLP-Citations split in SciFULLBENCH. Due to the length of input documents, we provide only the title and abstract for the input, as well as the titles for the retrieved documents. Retrieved documents that belong to ground truth are highlighted in blue.

2139

2140

Input Document Meta Data	<p>[Title] Stance Detection on Social Media with Background Knowledge</p> <p>[Abstract] Identifying users' stances regarding specific targets/topics is a significant route to learning public opinion from social media platforms. Most existing studies of stance detection strive to learn stance information about specific targets from the context, in order to determine the user's stance on the target. However, in real-world scenarios, we usually have a certain understanding of a target when we express our stance on it. In this paper, we investigate stance detection from a novel perspective, where the background knowledge of the targets is taken into account for better stance detection. To be specific, we categorize background knowledge into two categories: episodic knowledge and discourse knowledge, and propose a novel Knowledge-Augmented Stance Detection (KASD) framework. For episodic knowledge, we devise a heuristic retrieval algorithm based on the topic to retrieve the Wikipedia documents relevant to the sample. Further, we construct a prompt for ChatGPT to filter the Wikipedia documents to derive episodic knowledge. For discourse knowledge, we construct a prompt for ChatGPT to paraphrase the hashtags, references, etc., in the sample, thereby injecting discourse knowledge into the sample. Experimental results on four benchmark datasets demonstrate that our KASD achieves state-of-the-art performance in in-target and zero-shot stance detection.</p>
Ground-Truth Document Meta Data	<p>[Title] A More Advanced Group Polarization Measurement Approach Based on LLM-Based Agents and Graphs</p> <p>[Title] A Survey of Stance Detection on Social Media: New Directions and Perspectives</p> <p>[Title] Chain of Stance: Stance Detection with Large Language Models</p> <p>[Title] A Challenge Dataset and Effective Models for Conversational Stance Detection</p> <p>[Title] Mitigating Biases of Large Language Models in Stance Detection with Counterfactual Augmented Calibration</p> <p>[Title] Multi-modal Stance Detection: New Datasets and Model</p> <p>[Title] A Logically Consistent Chain-of-Thought Approach for Stance Detection</p> <p>[Title] Stance Detection with Collaborative Role-Infused LLM-Based Agents</p> <p>[Title] Prompting and Fine-Tuning Open-Sourced Large Language Models for Stance Classification</p> <p>[Title] Ladder-of-Thought: Using Knowledge as Steps to Elevate Stance Detection</p> <p>[Title] Investigating Chain-of-thought with ChatGPT for Stance Detection on Social Media</p>
Top@30 Retrieved Document Meta Data	<p>[Title] A Survey of Stance Detection on Social Media: New Directions and Perspectives</p> <p>[Title] Stance Detection with Collaborative Role-Infused LLM-Based Agents</p> <p>[Title] A Survey on Stance Detection for Mis- and Disinformation Identification</p> <p>[Title] Prompting and Fine-Tuning Open-Sourced Large Language Models for Stance Classification</p> <p>[Title] Chain of Stance: Stance Detection with Large Language Models</p> <p>[Title] A Challenge Dataset and Effective Models for Conversational Stance Detection</p> <p>[Title] Enabling Contextual Soft Moderation on Social Media through Contrastive Textual Deviation</p> <p>[Title] Stance Detection on Social Media with Fine-Tuned Large Language Models</p> <p>[Title] Stance Detection in Web and Social Media: A Comparative Study</p> <p>[Title] Multi-modal Stance Detection: New Datasets and Model</p> <p>[Title] TATA: Stance Detection via Topic-Agnostic and Topic-Aware Embeddings</p> <p>[Title] A Benchmark for Cross-Domain Argumentative Stance Classification on Social Media</p> <p>[Title] Mitigating Biases of Large Language Models in Stance Detection with Counterfactual Augmented Calibration</p> <p>[Title] Advancing Annotation of Stance in Social Media Posts: A Comparative Analysis of Large Language Models and Crowd Sourcing</p> <p>[Title] DEEM: Dynamic Experienced Expert Modeling for Stance Detection</p> <p>[Title] FarExStance: Explainable Stance Detection for Farsi</p> <p>[Title] Investigating Chain-of-thought with ChatGPT for Stance Detection on Social Media</p> <p>[Title] Relative Counterfactual Contrastive Learning for Mitigating Pretrained Stance Bias in Stance Detection</p> <p>[Title] A Logically Consistent Chain-of-Thought Approach for Stance Detection</p> <p>[Title] Embeddings-Based Clustering for Target Specific Stances: The Case of a Polarized Turkey</p> <p>[Title] Reinforcement Tuning for Detecting Stances and Debunking Rumors Jointly with Large Language Models</p> <p>[Title] Examining the Influence of Political Bias on Large Language Model Performance in Stance Classification</p> <p>[Title] KCD: Knowledge Walks and Textual Cues Enhanced Political Perspective Detection in News Media</p> <p>[Title] Deciphering Political Entity Sentiment in News with Large Language Models: Zero-Shot and Few-Shot Strategies</p> <p>[Title] Argumentative Stance Prediction: An Exploratory Study on Multimodality and Few-Shot Learning</p> <p>[Title] A More Advanced Group Polarization Measurement Approach Based on LLM-Based Agents and Graphs</p> <p>[Title] Beneath the Tip of the Iceberg: Current Challenges and New Directions in Sentiment Analysis Research</p> <p>[Title] WIBA: What Is Being Argued? A Comprehensive Approach to Argument Mining</p> <p>[Title] KGAP: Knowledge Graph Augmented Political Perspective Detection in News Media</p> <p>[Title] "We Demand Justice!": Towards Social Context Grounding of Political Texts</p>

Table 24: Example of the retrieved documents for the input document with CoR from a Citations split in PATENT-FULLBENCH. Due to the length of input documents, we provide only the title and abstract for the input, as well as the titles for the retrieved documents. Retrieved documents that belong to ground truth are highlighted in blue.

Input Document Meta Data	<p>[Abstract] The present disclosure describes methods and systems directed towards providing scaled engagement and views of an e-sports event. Instead of providing the same distribution of live e-sport event data to all remote viewers of a live e-sports event, features associated with e-sports gaming network could be used to customize the distribution of live e-sport event data to promote immersive viewer experience. The enhanced immersion can also be carried out in a virtual reality or augmented reality setting. The features would be capable of providing additional information, different views, and a variety of different commentators for the e-sports event so that the viewer can be more engaged when viewing the particular e-sports event. With the increased engagement from remote viewers, the distribution of live e-sports event data can also be further modified for monetization by incorporating advertisements as well.</p>
Ground-Truth Document Meta Data	<p>[Title] Statistical driven tournaments [Title] Scaled VR engagement and views in an e-sports event [Title] Player to spectator handoff and other spectator controls [Title] Creation of winner tournaments with fandom influence [Title] User-driven spectator channel for live game play in multi-player games [Title] Methods and systems to increase interest in and viewership of content before, during and after a live event [Title] Discovery and detection of events in interactive content [Title] Integrating commentary content and gameplay content over a multi-user platform [Title] Statistically defined game channels [Title] Online tournament integration [Title] De-interleaving gameplay data</p>
Top@30 Retrieved Document Meta Data	<p>[Title] Methods and systems to increase interest in and viewership of content before, during and after a live event [Title] Scaled VR engagement and views in an e-sports event [Title] Creation of winner tournaments with fandom influence [Title] Player to spectator handoff and other spectator controls [Title] Discovery and detection of events in interactive content [Title] Real-time modifications in augmented reality experiences [Title] User-driven spectator channel for live game play in multi-player games [Title] Integrating commentary content and gameplay content over a multi-user platform [Title] Online tournament integration [Title] Integrating augmented reality experiences with other components [Title] AR-based connected portal shopping [Title] External screen streaming for an eyewear device [Title] Systems and methods for generating and facilitating access to a personalized augmented rendering of a user [Title] Systems, methods and apparatuses of digital assistants in an augmented reality environment and local determination of virtual object placement and apparatuses of single or multi-directional lens as portals between a physical world and a digital world component of the augmented reality environment [Title] Statistically defined game channels [Title] AR/VR enabled contact lens [Title] Shared augmented reality unboxing experience [Title] De-interleaving gameplay data [Title] Statistical driven tournaments [Title] Automated augmented reality experience creation based on sample source and target images [Title] Augmented reality unboxing experience [Title] Real-time video dimensional transformations of video for presentation in mixed reality-based virtual spaces [Title] Systems and methods for pinning content items to locations in an augmented reality display based on user preferences [Title] Avatar customization system [Title] Controlling interactive fashion based on body gestures [Title] Systems, methods, and apparatus for enhanced headsets [Title] User interfaces for wide angle video conference [Title] User interface for multi-user communication session [Title] Real-time upper-body garment exchange [Title] Applying animated 3D avatar in AR experiences</p>

2141

2142

Document Meta-Data

Abstract: We introduce GAIA, a benchmark for General AI Assistants that, if solved, would represent a milestone in AI research. GAIA proposes real-world questions that require a set of fundamental abilities such as reasoning, multi-modality handling, web browsing, and generally tool-use proficiency. GAIA questions are conceptually simple for humans yet challenging for most advanced AIs: we show that human respondents obtain 92% vs. 15% for GPT-4 equipped with plugins. This notable performance disparity contrasts with the recent trend of LLMs outperforming humans on tasks requiring professional skills in e.g. law or chemistry. GAIA's philosophy departs from the current trend in AI benchmarks suggesting to target tasks that are ever more difficult for humans. We posit that the advent of Artificial General Intelligence (AGI) hinges on a system's capability to exhibit similar robustness as the average human does on such questions. Using GAIA's methodology, we devise 466 questions and their answer. We release our questions while retaining answers to 300 of them to power a leader-board hereby accessible.

Method-Focused Query:

GAIA, a benchmark for General AI Assistants, presents a novel approach to evaluating the capabilities of large language models (LLMs) in performing real-world tasks that require a range of fundamental abilities, including reasoning, multi-modality handling, web browsing, and tool-use proficiency. The methodology employed in developing GAIA involves creating a set of 466 carefully crafted questions and their answers, covering various assistant use cases such as daily personal tasks, science, and general knowledge. Each question is designed to admit a short, single correct answer, enabling simple and robust automatic evaluation.

To address the limitations of traditional LLM benchmarks, GAIA focuses on targeting tasks that are conceptually simple for humans yet challenging for most advanced AIs. The questions are intended to be easy to create, yet vary, root in the real world, and pose a unique combination of factual and inferential demands. Furthermore, GAIA's design principles prioritize interpretability, robustness against memorization, and easiness of use, aiming to provide a more nuanced evaluation framework for assessing the capabilities of LLMs.

The development of GAIA involved several stages, starting with the creation of a seed set of questions by the authors themselves, followed by extensive annotation and validation phases involving human curators and evaluators. Two independent annotators assessed each question for ambivalence, and discrepancies were resolved through iterative refinements. This rigorous process ensured that the questions remained unambiguous and suitable for evaluation.

Upon releasing the benchmark, GAIA evaluates LLMs' performance across various capabilities, including web browsing, multi-modality processing, coding, and diverse filetype reading. The evaluation procedure leverages a zero-shot prompt strategy, utilizing a prefix prompt to guide the model toward producing a factual answer. Results indicate significant disparities between human and LLM performances, highlighting the need for more sophisticated evaluation frameworks that account for the complexities inherent in real-world scenarios.

Notably, GAIA's methodology diverges from previous benchmarks by focusing on the overall performance of AI assistants rather than individual tasks or narrow skill sets. This shift underscores the importance of considering the broader cognitive landscape encompassed by general intelligence, particularly in light of the growing emphasis on achieving true Artificial General Intelligence (AGI). Future directions for GAIA entail exploring extensions to accommodate emerging topics, incorporating human-in-the-loop feedback, and investigating novel methodologies for tackling issues such as linguistic and cultural diversity.

Overall, GAIA represents a pioneering endeavor in establishing a standardized benchmark for

evaluating the capabilities of General AI Assistants, offering valuable insights into the strengths and weaknesses of contemporary LLM architectures. By fostering a deeper understanding of the intricacies underlying effective AI assistance, GAIA paves the way for continued innovation in the pursuit of AGI.

Experiment-Focused Query:

Here's a detailed summary of the experiment setup, datasets, evaluation metrics, baselines, and key findings of the paper "GAIA: A Benchmark for General AI Assistants":

Experiment Setup:

The authors propose a novel benchmark called GAIA (General AI Assistant) consisting of 466 carefully crafted questions and their answers. Each question requires a series of steps to arrive at a factual answer, mimicking real-world scenarios where AI assistants interact with humans. The questions cover various domains, including science, literature, and everyday life, and are designed to be conceptually simple yet challenging for current AI systems.

Datasets:

The authors develop a custom dataset of 466 questions, divided into three levels of difficulty: Level 1 (simplest), Level 2 (moderate), and Level 3 (most challenging). The questions are paired with their answers, reasoning traces, and metadata describing the required capabilities and tools.

Evaluation Metrics:

The authors employ a combination of human evaluation and automated scoring to assess the performance of AI systems on the GAIA benchmark. Specifically:

*** Human Score:**

The authors annotate a subset of questions with their expected answers and measure the similarity between the predicted answer and the true answer.

*** Automated Scoring:**

The authors implement a scoring mechanism that evaluates the correctness of the predicted answer based on the presence of certain keywords and phrases in the answer.

Baselines:

The authors compare the performance of several popular large language models (LLMs) on the GAIA benchmark, including GPT-4, AutoGPT, and human evaluators. The baseline performances serve as a reference point for assessing the effectiveness of the proposed GAIA benchmark.

Key Findings:

* Despite achieving impressive results on various NLP tasks, current LLMs struggle significantly on the GAIA benchmark, particularly on Level 3 questions.

* The authors observe that equipping LLMs with external tools and APIs can significantly enhance their performance on the GAIA benchmark, highlighting the importance of integrating external resources into AI systems.

* The GAIA benchmark demonstrates the limitations of traditional evaluation metrics, such as perplexity and BLEU score, in capturing the nuances of human-AI interaction.

The authors conclude that solving the GAIA benchmark represents a significant milestone toward developing truly general-purpose AI assistants that can effectively assist humans in a wide range of contexts.

Overall, the GAIA benchmark offers a comprehensive framework for evaluating the capabilities of AI systems in simulating real-world human-AI interactions, emphasizing the importance of

integrating external resources, improving reasoning and decision-making, and enhancing overall human-AI alignment.

Research-Question Focused Query:

Here's a detailed explanation of the scientific paper, focusing on its motivation, research questions, and contributions:

Motivation:

The paper introduces a novel benchmark called GAIA (General AI Assistant) aimed at assessing the capabilities of General AI Assistants. The author argues that traditional benchmarks, such as GLUE and MMLU, are insufficient for evaluating AI systems' ability to generalize across various tasks and domains. Current benchmarks tend to focus on narrow tasks, whereas real-world applications require AI systems to tackle complex, multifaceted problems. GAIA addresses this limitation by presenting a comprehensive framework for evaluating AI assistants' reasoning, multi-modality handling, web browsing, and tool-use proficiency.

Research Questions:

The primary research question driving this paper is: "Can we develop a benchmark that comprehensively evaluates the capabilities of General AI Assistants?" Specifically, the author seeks to investigate:

- * Can AI systems generalize across various tasks and domains?
- * How effective are current benchmarks in capturing the complexities of real-world applications?

To address these questions, the author designs and implements GAIA, a benchmark comprising 466 carefully crafted questions covering various aspects of AI assistance, including reasoning, multi-modality handling, web browsing, and tool-use proficiency.

Contributions: The paper contributes several significant ideas to the field of Natural Language Processing (NLP):

1. **Comprehensive Benchmark for General AI Assistants:** GAIA offers a standardized framework for evaluating AI assistants' capabilities, filling the gap left by existing benchmarks.
2. **Realistic Task Set:** The benchmark features questions that mimic real-world scenarios, encouraging AI systems to apply their knowledge and skills in context-dependent settings.
3. **Focus on Reasoning and Multi-Modality Handling:** GAIA emphasizes the importance of reasoning and multi-modality handling in AI assistants, highlighting areas where current models fall short.
4. **Evaluation Methodology:** The paper presents a systematic approach to evaluating AI assistants, incorporating human evaluation, model-based evaluation, and automatic scoring techniques.

Overall, the paper demonstrates the feasibility of developing a comprehensive benchmark for General AI Assistants, paving the way for improved evaluation methodologies and more sophisticated AI systems. . . .

Document Meta-Data

Abstract: The present disclosure describes methods and systems directed towards providing scaled engagement and views of an e-sports event. Instead of providing the same distribution of live e-sport event data to all remote viewers of a live e-sports event, features associated with e-sports gaming network could be used to customize the distribution of live e-sport event data to promote immersive viewer experience. The enhanced immersion can also be carried out in a virtual reality or augmented reality setting. The features would be capable of providing additional information, different views, and a variety of different commentators for the e-sports event so that the viewer can be more engaged when viewing the particular e-sports event. With the increased engagement from remote viewers, the distribution of live e-sports event data can also be further modified for monetization by incorporating advertisements as well.

Method-Focused Query: The disclosed method for providing scaled engagement and views of a live electronic sports (e-sport) event involves a systematic approach to enhance viewer immersion through customized data distribution. Initially, the method collects comprehensive information from the live e-sport event, including in-game elements, player statistics, and audience reactions. This data is processed to generate detailed statistics about the players, which are then used to create a dynamic display of the event from an in-game camera perspective. The generated display is further modified based on individual viewer preferences, which may include adjustments for augmented reality (AR) environments, different camera angles, and the integration of advertisements. The modifications are tailored to enhance the viewer's experience by providing relevant statistics, commentary, and interactive features, such as the ability to follow specific players or engage with other viewers. The final modified display is distributed to the viewer's device, ensuring a personalized and immersive viewing experience that can also facilitate monetization through targeted advertisements. The system architecture includes an event server for data collection, a gaming server for processing and display generation, and a public network for distribution, all of which work together to deliver a rich, engaging e-sport viewing experience.

Claim-Focused Query: The claims of this patent outline a method and system for enhancing viewer engagement and experience during live electronic sports (e-sports) events through customized data distribution. The independent claim details a method that includes collecting event data, generating player statistics, and creating a display that presents the event from an in-game camera perspective. This display is then modified based on viewer preferences, which may include augmented reality integration and personalized camera angles, before being distributed to remote viewers. The dependent claims specify refinements such as the inclusion of game-related information from player devices, audience reactions, advertisement incorporation, and user profile-based modifications that tailor the viewing experience to individual preferences. Additional claims address the provision of multiple viewing perspectives, interactive functions for viewers, and audio from other viewers, enhancing the immersive experience. The system claims describe an architecture comprising an event server and a gaming server that collectively manage data collection, processing, and distribution, while also allowing for the integration of advertisements and user interactions. Overall, the claims provide broad protection for methods and systems that facilitate a customizable and engaging viewing experience for remote participants in e-sports events.

Background Focused Query: The present invention pertains to the field of electronic sports (e-sports), specifically focusing on enhancing viewer engagement and experience during live e-sports events through customized data distribution. Existing methods of broadcasting e-sports events typically provide a uniform viewing experience to all remote viewers, lacking personalization and immersive features that could significantly enhance engagement. Current streaming platforms do not adequately leverage the wealth of data generated during live competitions, such as player statistics, audience reactions, and in-game dynamics, which limits the depth of viewer interaction and

understanding of the event. Furthermore, the integration of advanced technologies like virtual reality (VR) and augmented reality (AR) remains underutilized, preventing viewers from experiencing events in a more immersive manner akin to being physically present. There is a pressing need for systems that can dynamically modify the presentation of e-sports data based on individual viewer preferences, including customizable perspectives, additional commentary, and interactive features, while also addressing challenges related to monetization through targeted advertisements. Operating constraints include the need for real-time data processing, compatibility with various user devices, and the ability to handle high viewer throughput without compromising latency or quality. The application context spans competitive gaming leagues, online streaming services, and event organizers, all of whom seek to enhance viewer satisfaction and engagement while navigating the complexities of data management and user interaction in a rapidly evolving digital landscape.

2143

H Human Evaluation Guideline

2144

In this section, we provide the guidelines for human evaluation provided to annotators in Table 6.

2145

Human Evaluation Guideline

1. Evaluation Guidelines

Each task contains Title, Abstract information of Input Query Document, and has 3 different randomly shuffled retrieval variants with 5 retrieved documents with their titles and abstracts respectively. You will annotate each Query–Document pair with **EXACTLY ONE** of the following labels:

Method

Experiment

Research Question/Motivation

Irrelevant

A total of 15 retrieved documents must be labeled per Input paper.

The evaluation unit is one retrieved document for one query.

For each retrieved document, Exactly One label can be assigned.

2. Evaluation Instruction

1. Evaluation Criteria

1.1 Research Question Definition:

Choose this label when the retrieved document addresses the same or a highly similar research problem, task, or scientific question as the query.

Criteria:

Same or very similar task/problem/research objectives Overlapping motivation or problem definition High-level goals align strongly Do not consider overly vague and broad similarities as Motivation-Wise Related.

Examples:

Correct: Both works aim to address the fundamental challenge of retrieving relevant information from long-context documents, focusing on the issues of efficiently utilizing the important content within the query and candidate documents, and mitigate the effect of irrelevant noise information.

Incorrect: Both papers focus on Large Language Models.

1.2 Method Definition:

Choose this label when the retrieved document is relevant mainly because it provides similar methodology, to solve its respective problem. Focus on the similarity between Algorithms, Model Architectures, Training Methods, Optimization Strategies, or Theoretical Frameworks.

Criteria:

- 1) Similar algorithm or model architecture
- 2) Methodological or optimization insights overlap
- 3) Technical design relevant to the query
- 4) Try to Avoid Overly Vague Methodological Similarities

Examples:

Correct Although the query document A focuses on research agents for solution generation and the retrieved document B addresses automated web agents, their core methodological frameworks are highly similar. Both frameworks implement a training-free, self-evolving memory architecture in which past experiences with the highest rewards are stored during an offline stage and later retrieved during inference. This mechanism enables the policy model to generalize to unseen tasks and effectively perform a form of curriculum learning through experience reuse.

Incorrect Both Paper A and Paper B leverage large language models as autonomous agents and employ multi-agent frameworks to address their respective problem settings.

1.3 Experiment Definition:

Choose this label when the retrieved document is primarily relevant due to its experimental framework: datasets, benchmarks, evaluation methodologies, or baseline setups.

Criteria:

- 1) Introduces a similar dataset or benchmark used in the query
- 2) Uses similar evaluation protocol
- 3) Provides relevant baselines or experiment structures
- 4) Try to Avoid Overly Vague Experimental Similarities

Examples:

Correct Both paper A and paper B attempts to solve the long-horizon forgetting problem of automated Agents with different methodologies and research objectives. Both paper A and paper B performs evaluation on OsWorld benchmark and VisualWebArena benchmark, while CoACT and UITars as their strong baselines.

Incorrect Paper A and Paper B both utilize benchmarks to evaluate agents.

1.4 Irrelevant Definition:

Choose this label when the retrieved document does not meaningfully relate to the query in terms of research question, method, or experimental setup.

Criteria:

- 1) Only superficial keyword overlap (e.g., both mention “LLM”)
- 2) Different domain or unrelated task
- 3) No substantial conceptual or technical relevance

Examples:

Vision Language Model layer compression for GUI Agents for efficiency improvement <->
Preference Optimization for Vision Language Models

2. How to Resolve Ambiguous Cases When a retrieved document could belong to multiple categories, use the following priority:

Phase 1: Check the aspect that aligns strongly among three given aspects.

Research Question → if the research problem aligns strongly

Method → if the alignment is mainly methodological

Experiment → if the relevance is primarily datasets/benchmarks/evaluation

Phase 2 Irrelevant → if none of the above apply

3. Quality Checklist Before submitting, verify the following:

Exactly one label selected per document

Strong task alignment → **Research Question**

Methodological overlap → **Method**

Dataset/evaluation relevance → **Experiment**

Unrelated documents → **irrelevant**

There are duplicate candidate documents across variants. Please make sure that the labels are applied consistently across all query, document pairs.

I Prompts

Here, we provide prompts used for experiments.

Method-Focused Agent Prompt for Paper-to-Paper Retrieval

Instruction:

You are a specialized research assistant tasked with generating a structured, detailed explanation of a scientific paper based on its **Methodology**. Your goal is to provide a clear yet comprehensive summary that makes it easy to identify relevant papers by emphasizing the **methodology** of given paper.

IMPORTANT

- Your explanation is going to be used as a query to retrieve similar papers **METHOD** wise.
- Make sure that your explanation can retrieve highly relevant papers easily.
- There are also other agents who are tasked with generating explanation on given paper. Unlike you, they are focused on experiments, and research questions of given paper. You must try to avoid overlap with possible explanations that the other two agents might generate.

Input:

You will be given the full text of a scientific paper. Carefully analyze its content, with a particular focus on the **METHODOLOGY** section, to extract its main approaches.

Key Considerations:

Highlight specific method/approach details, avoiding vague or overly general descriptions. Use precise language to ensure clarity while maintaining depth.

Output Format:

Generate a well-structured, detailed and yet clear paragraph that effectively captures the paper's approach, and key concepts in a concise yet informative manner. Focus on high-level insights rather than excessive detail.

You must not include title and abstract of given paper in your answer, and try to put it into your own words with high level reasoning after reading the paper.

Experiment-Focused Agent Prompt for Paper-to-Paper Retrieval

You are a specialized research assistant tasked with generating a structured, detailed explanation of a scientific paper's experimental setup.

Your goal is to clearly outline the **datasets, evaluation metrics, baselines, and key experimental findings**, making it easy to understand how the paper validates its approach.

IMPORTANT

- Your explanation is going to be used as a query to retrieve similar papers **EXPERIMENT** wise.
- Make sure that your explanation can retrieve highly relevant papers easily.
- There are also other agents who are tasked with generating explanation on given paper. Unlike you, they are focused on methods, and research questions of given paper. You must try to avoid overlap with possible explanations that the other two agents might generate.

Input:

You will be provided with the full text of a scientific paper. Carefully analyze its content, paying particular attention to the **Experiments, Results, and Evaluation** sections to extract the key experimental details.

Key Considerations:

Datasets & Benchmarks: Clearly specify the datasets and benchmarks used for evaluation.

Baselines & Comparisons: Identify what methods or models the paper compares against.

Key Results & Insights: Summarize the main experimental findings without excessive detail.

Output Format:

Generate a clear, well structured and detailed paragraph that highlights the experimental methodology, datasets, evaluation metrics, baselines, and key results. Focus on high-level insights rather than excessive detail.

You must not include title and abstract of given paper in your answer, and try to put it into your own words with high level reasoning after reading the paper.

Research Question-Focused Agent Prompt for Paper-to-Paper Retrieval

Instruction:

You are a specialized research assistant tasked with generating a structured, detailed explanation of a scientific paper based on its **Motivation, Research Questions, and Contributions**. Your goal is to provide a clear yet comprehensive summary that makes it easy to identify relevant papers by emphasizing the **core problem, key contributions, and research objectives**.

IMPORTANT

- Your explanation is going to be used as a query to retrieve similar papers **RESEARCH QUESTION** wise.
- Make sure that your explanation can retrieve highly relevant papers easily.
- There are also other agents who are tasked with generating explanation on given paper. Unlike you, they are focused on experiments, and methods of given paper. You must try to avoid overlap with possible explanations that the other two agents might generate.

Input:

You will be given the full text of a scientific paper. Carefully analyze its content, with a particular focus on the Introduction and Conclusion, to extract its main contributions, research questions, and motivations.

Key Considerations:

Highlight specific motivations, research questions, and contributions, avoiding vague or overly general descriptions.

Use precise language to ensure clarity while maintaining depth.

Output Format: Generate a well-structured, detailed and yet clear paragraph that effectively captures the paper's motivation, problem statement, research questions, and key contributions in a concise yet informative manner. Focus on high-level insights rather than excessive detail

You must not include title and abstract of given paper in your answer, and try to put it into your own words with high level reasoning after reading the paper.

Base Agent Prompt for Paper-to-Paper Retrieval

Instruction: You are a specialized research assistant tasked with generating a structured,detailed explanation of a scientific paper based three different aspects.

1. Method-Specific Queries:

Generate a structured, detailed explanation of a scientific paper based on its Methodology Your goal is to provide a clear yet comprehensive summary that makes it easy to identify relevant papers by emphasizing the methodology of given paper.

IMPORTANT

- Your explanation is going to be used as a query to retrieve similar papers **METHOD** wise.
- Make sure that your explanation can retrieve highly relevant papers easily.

Input:

You will be given the full text of a scientific paper. Carefully analyze its content, with a particular focus on the **METHODOLOGY** section, to extract its main approaches.

Key Considerations:

Highlight specific method/approach details, avoiding vague or overly general descriptions. Use precise language to ensure clarity while maintaining depth.

Output Format:

Generate a well-structured, detailed and yet clear paragraph that effectively captures the paper's approach, and key concepts in a concise yet informative manner. Focus on high-level insights rather than excessive detail You must not include title and abstract of given paper in your answer, and try to put it into your own words with high level reasoning after reading the paper.

2. Experiment-Specific Queries:

Generate a structured, detailed explanation of a scientific paper's experimental setup Your goal is to clearly outline the datasets, evaluation metrics, baselines, and key experimental findings, making it easy to understand how the paper validates its approach.

IMPORTANT

- Your explanation is going to be used as a query to retrieve similar papers **EXPERIMENT** wise.
- Make sure that your explanation can retrieve highly relevant papers easily.

You will be provided with the full text of a scientific paper. Carefully analyze its content, paying particular attention to the Experiments, Results, and Evaluation sections to extract the key experimental details.

Key Considerations:

Datasets & Benchmarks: Clearly specify the datasets and benchmarks used for evaluation.

Baselines & Comparisons: Identify what methods or models the paper compares against.

Key Results & Insights: Summarize the main experimental findings without excessive detail.

Output Format:

Generate a clear, well structured and detailed paragraph that highlights the experimental methodology, datasets, evaluation metrics, baselines, and key results. Focus on high-level insights rather than excessive detail. You must not include title and abstract of given paper in your answer, and try to put it into your own words with high level reasoning after reading the paper.

3. Research Question-Specific Queries:

Generate a structured, detailed explanation of a scientific paper based on its Motivation, Research Questions, and Contributions. Your goal is to provide a clear yet comprehensive summary that

makes it easy to identify relevant papers by emphasizing the core problem, key contributions, and research objectives.

IMPORTANT

- Your explanation is going to be used as a query to retrieve similar papers **RESEARCH QUESTION** wise.
- Make sure that your explanation can retrieve highly relevant papers easily.

Input:

You will be given the full text of a scientific paper. Carefully analyze its content, with a particular focus on the Introduction and Conclusion, to extract its main contributions, research questions, and motivations.

Key Considerations:

Highlight specific motivations, research questions, and contributions, avoiding vague or overly general descriptions. Use precise language to ensure clarity while maintaining depth.

Output Format:

Generate a well-structured, detailed and yet clear paragraph that effectively captures the paper's motivation, problem statement, research questions, and key contributions in a concise yet informative manner. Focus on high-level insights rather than excessive detail. You must not include title and abstract of given paper in your answer, and try to put it into your own words with high level reasoning after reading the paper.

Output:

Return a structured json file for respective Method-Specific Queries, Experiment-Specific Queries, and Research Question-Specific queries, with the respective keys as "method_query", "experiment_query", and "research_question_query". Each key should contain the generated explanation as a string.

Prompt for Single Comprehensive Query for Paper-to-Paper Retrieval

Instruction:

You are a specialized research assistant tasked with generating a structured, detailed explanation of a scientific paper based three different aspects. You should generate a single query covering method, experiment, and research question aspects. Carefully read the below instructions and generate a single comprehensive query that covers below three aspects.

1. Method-Specific Queries:

Generate a structured, detailed explanation of a scientific paper based on its Methodology. Your goal is to provide a clear yet comprehensive summary that makes it easy to identify relevant papers by emphasizing the methodology of given paper.

IMPORTANT

- Your explanation is going to be used as a query to retrieve similar papers **METHOD** wise.
- Make sure that your explanation can retrieve highly relevant papers easily.

Input:

You will be given the full text of a scientific paper. Carefully analyze its content, with a particular focus on the **METHODOLOGY** section, to extract its main approaches.

Key Considerations:

Highlight specific method/approach details, avoiding vague or overly general descriptions. Use precise language to ensure clarity while maintaining depth.

Output Format:

Generate a well-structured, detailed and yet clear paragraph that effectively captures the paper's approach, and key concepts in a concise yet informative manner. Focus on high-level insights rather than excessive detail You must not include title and abstract of given paper in your answer, and try to put it into your own words with high level reasoning after reading the paper.

2. Experiment-Specific Queries:

Generate a structured, detailed explanation of a scientific paper's experimental setup Your goal is to clearly outline the datasets, evaluation metrics, baselines, and key experimental findings, making it easy to understand how the paper validates its approach.

IMPORTANT

- Your explanation is going to be used as a query to retrieve similar papers **EXPERIMENT** wise.
- Make sure that your explanation can retrieve highly relevant papers easily.

You will be provided with the full text of a scientific paper. Carefully analyze its content, paying particular attention to the Experiments, Results, and Evaluation sections to extract the key experimental details.

Key Considerations:

Datasets & Benchmarks: Clearly specify the datasets and benchmarks used for evaluation.

Baselines & Comparisons: Identify what methods or models the paper compares against.

Key Results & Insights: Summarize the main experimental findings without excessive detail.

Output Format:

Generate a clear, well structured and detailed paragraph that highlights the experimental methodology, datasets, evaluation metrics, baselines, and key results. Focus on high-level insights rather than excessive detail. You must not include title and abstract of given paper in your answer, and try to put it into your own words with high level reasoning after reading the paper.

3. Research Question-Specific Queries:

Generate a structured, detailed explanation of a scientific paper based on its Motivation, Research Questions, and Contributions. Your goal is to provide a clear yet comprehensive summary that

makes it easy to identify relevant papers by emphasizing the core problem, key contributions, and research objectives.

IMPORTANT

- Your explanation is going to be used as a query to retrieve similar papers **RESEARCH QUESTION wise**.
- Make sure that your explanation can retrieve highly relevant papers easily.

Input:

You will be given the full text of a scientific paper. Carefully analyze its content, with a particular focus on the Introduction and Conclusion, to extract its main contributions, research questions, and motivations.

Key Considerations:

Highlight specific motivations, research questions, and contributions, avoiding vague or overly general descriptions. Use precise language to ensure clarity while maintaining depth.

Output Format: Generate a well-structured, detailed and yet clear paragraph that effectively captures the paper's motivation, problem statement, research questions, and key contributions in a concise yet informative manner. Focus on high-level insights rather than excessive detail. You must not include title and abstract of given paper in your answer, and try to put it into your own words with high level reasoning after reading the paper.

Output:

Return a comprehensive query that covers the Methodology, Research Questions, and Experimental Details within a single paper.

Method-Focused Agent Prompt for Patent-to-Patent Retrieval

You are a specialized research assistant tasked with generating a structured, detailed explanation of a patent's **METHOD**.

Your goal is to summarize how the invention works in practice: its core technical principles, implementation procedures, system architecture, and functional mechanisms. This explanation will be used as a query to retrieve patents with similar methods or implementation techniques.

IMPORTANT

- Focus **ONLY** on the Detailed Description and Embodiments sections.
- Capture the technical processes, structures, system flows, or algorithms.
- Paraphrase at a high-level while keeping enough technical detail for retrieval.
- Other agents will generate explanations about the invention claims and the background/problems of the given patent. You must avoid overlap with those aspects.

Input

You will be provided with the full text of a patent. Carefully analyze its Detailed Description, Examples, and Figures to extract the methodological details.

Key Considerations

Here are some key aspects that you may focus on.

- **Core technical principle:** what mechanism enables the invention to function?
- **Implementation structure:** key components, modules, or subsystems.
- **Operational flow:** how the method proceeds step by step.
- **Variants or embodiments:** different configurations or modes of execution.
- **Integration context:** how it interacts with external systems or environments.

Output Format

Produce a single, clear, and well-structured paragraph that captures the **METHOD** of the invention. Write in concise, retrieval-friendly technical language that highlights implementation strategies and system functionality.

Claim-Focused Agent Prompt for Patent-to-Patent Retrieval

You are a specialized research assistant tasked with generating a structured, detailed explanation of a patent's **CLAIMS**.

Your goal is to clearly outline the legal protection scope: what is being claimed, how broad the claims are, and what technical features or components are covered.

This explanation will be used as a query to retrieve patents with similar **CLAIM** structures and protection scopes.

IMPORTANT

- Focus **ONLY** on the **CLAIMS** section.
- Emphasize the scope of protection, claimed components, processes, and relationships among them.
- Avoid background information, prior art, or detailed embodiments (those are handled by other agents).
- Do not directly copy claim sentences; paraphrase into concise, high-level, retrieval-friendly language.
- Include both the broad independent claims (core invention scope) and notable dependent claims (specific refinements).
- Other agents will generate explanations about the invention **details/method** and the **background/problems** of the given patent. You must avoid overlap with those aspects.

Input

You will be provided with the full text of a patent. Carefully analyze the **CLAIMS** section.

Output Format

Produce a clear, well-structured paragraph that captures the scope and focus of the **CLAIMS**. Use neutral, technical language suitable for retrieval so that patents with overlapping protection scopes can be surfaced.

Background-Focused Agent Prompt for Patent-to-Patent Retrieval

You are a specialized research assistant tasked with generating a structured, detailed explanation of a patent's **Problem / Background** (i.e., the technical field and the shortcomings of prior art).

Your goal is to clearly state the problem space—technical domain, prior-art limitations, unresolved challenges, operating constraints, and desired (non-solution) performance objectives—so that similar patents can be retrieved by shared problem patterns rather than specific solutions.

IMPORTANT

- Your explanation is going to be used as a query to retrieve patents that are similar in the **PROBLEM SPACE**.
- Include concrete domain terms (components, materials, data types), relevant standards/regulations, operating conditions (e.g., throughput, latency, power, temperature), failure modes, and application contexts that characterize the problem.
- Other agents will generate explanations about the invention details/method and the claims made from the input patent. You must avoid overlap with those aspects.

Input

You will be provided with the full text of a patent.

Key Considerations

Here are some key aspects that you may focus on.

- **Technical Field:** the domain and subdomain (e.g., “wireless edge inference for medical imaging”).
- **Prior Art & Limitations:** concrete bottlenecks, inefficiencies, failure cases, safety/privacy concerns, interoperability issues.
- **Unmet Technical Objectives:** what must be achieved (targets or constraints).
- **Operating Constraints & Edge Cases:** data distributions, environmental/thermal limits, network conditions, power/memory budgets, size/weight/cost constraints, lifecycle/maintenance issues.
- **Application Contexts & Stakeholders:** deployment scenarios, industries, users, and integration boundaries.

Output Format

Generate a single, clear, well-structured paragraph in neutral technical language that captures **ONLY** the **Problem / Background**. Paraphrase in your own words at a high level; do not copy text or include any solution, embodiment, or claim content. The paragraph should be information-dense and retrieval-friendly so that it can surface patents that confront the same technical challenges.