

MixGR: Enhancing Retriever Generalization for Scientific Domain through Complementary Granularity

Anonymous ACL submission

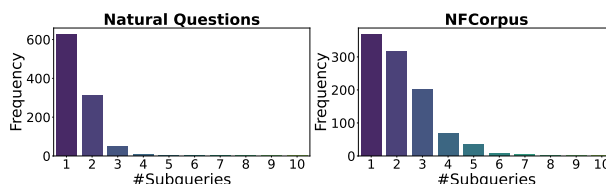
Abstract

Recent studies show the growing significance of document retrieval in the generation of LLMs within the scientific domain by bridging their knowledge gap. However, dense retrievers often struggle with domain-specific retrieval and complex query-document relationships, particularly when query segments correspond to various parts of a document. To alleviate such prevalent challenges, this paper introduces MixGR, which improves dense retrievers' awareness of query-document matching across various levels of granularity in queries and documents using a zero-shot approach. MixGR fuses various metrics based on these granularities to a united score that reflects a comprehensive query-document similarity. Our experiments demonstrate that MixGR outperforms previous document retrieval by 22.6% and 10.4% on nDCG@5 with unsupervised and supervised retrievers, respectively, averaged on queries containing multiple subqueries from four scientific retrieval datasets. Moreover, the efficacy of two downstream scientific question-answering tasks highlights the advantage of MixGR to boost the application of LLMs in the scientific domain.

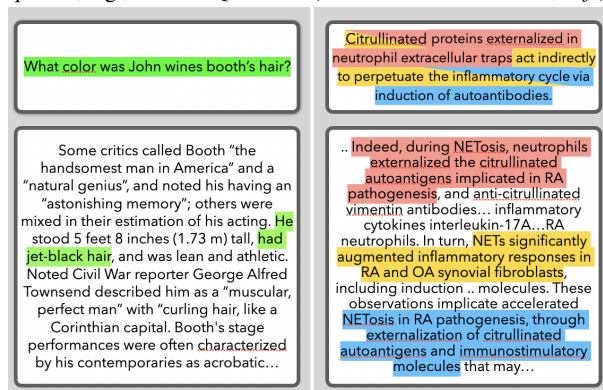
1 Introduction

Recent advances in Large Language Models (LLMs) have significantly impacted various scientific domains (Zhang et al., 2022; Touvron et al., 2023; Birhane et al., 2023; Grossmann et al., 2023). However, LLMs are notorious for their tendency to produce hallucinations, producing unreliable outputs (Ji et al., 2023). To address this, Retrieval-Augmented Generation (RAG; Lewis et al. 2020) has been developed to address this issue by incorporating external knowledge during the generation.

Though notable for accessing external and relevant knowledge, dense retrievers face specific challenges in the scientific domain: (1) *Domain-specific nature*: dense retrievers are typically



(a) Subquery distribution of general and scientific queries: scientific queries, e.g., NFCorpus (Boteva et al. 2016, Right), demonstrate a more diverse range of subqueries per query than general queries, e.g., Natural Questions (Kwiatkowski et al. 2019, Left).



(b) Comparison between general and scientific query-doc retrieval: compared with the general query-doc retrieval exemplified by NQ (Kwiatkowski et al. 2019, Left), the scientific query-doc retrieval exemplified by SciFact (Wadden et al. 2020, Right) demonstrates that one query can be decomposed to multiple subqueries, which can be mapped to different parts of documents.

Figure 1: Scientific document retrieval is shown to be more complicated than general domains.

trained on the general corpus such as Natural Questions (NQ; Kwiatkowski et al. 2019). However, scientific domains differ notably, e.g., the terminology and the pattern of queries as shown in Figure 1a. (2) *Complexity* of scientific documents: they are long, structured (Erera et al., 2019) and contain complex relationships between arguments (Stab et al., 2014). Figure 1a demonstrates that scientific queries tend to contain more subqueries than those in general domains. This indicates that subqueries within a single query may align with different parts of a document (doc), resulting in complex interactions between queries and documents (Figure 1b).

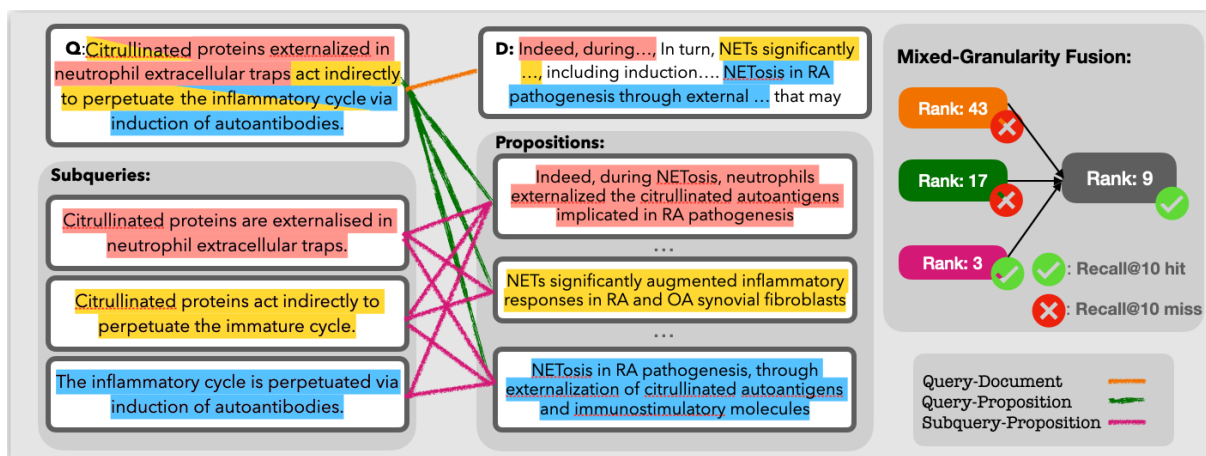


Figure 2: The illustration of $MiXGR$: Both queries and documents (e.g., the query-doc pair from SciFact in Figure 1b) are decomposed into subqueries and propositions, respectively, each containing distinct semantic components. Starting from the original queries and documents along with their decomposed elements, metrics from various granularity combinations are fused into a single integrated score.

Such complexity poses significant challenges for dense retrievers (Lupart et al., 2023). Addressing these challenges requires specific training on the scientific corpus. However, this is often hindered by the necessity of extensive annotations (Wadden et al., 2020) and extra computation (Wang et al., 2021a).

In this study, we introduce a novel **zero-shot** approach that effectively adapts dense retrievers to scientific domains. This method specifically addresses the complexities arising from the composition of scientific queries and their consequent intricate relationships with documents. Inspired by Chen et al. (2023), showing that finer units improve retrievers’ generalization to rare entities, we incorporate more granular retrieval units, specifically propositions (prop), to address domain-specific challenges as shown in Figure 2. Given the complexity between scientific queries and documents (Figure 1b), we also consider finer units within queries—subqueries—to measure query-doc similarity at a finer granularity. This metric captures the similarity between subqueries and propositions, moving beyond simple point similarity between query-doc vectors. Given a query, the distribution of corresponding information within a document is unknown. Additionally, our empirical analysis reveals that similarities at various granularities provide complementary insights. Therefore, for each query-doc pair, we fuse the metrics from these granularities to a unified score, termed **Mixed-Granularity Retrieval** as $MiXGR$, as depicted in Figure 2.

We conducted document retrieval experiments

on four scientific datasets using six dense retrievers, comprising two unsupervised and four supervised models. Our results demonstrate that $MiXGR$ markedly surpasses previous query-doc retrieval methods. Notably, we recorded an average improvement of 22.6% for unsupervised retrievers and 10.4% for supervised retrievers in terms of $nDCG@5$ for queries involving multiple subqueries. Furthermore, documents retrieved via $MiXGR$ substantially enhance the performance of downstream scientific QA tasks, underscoring their potential utility for RAG within scientific domains.

Our contributions are three-fold:

- We identify the challenges within scientific document retrieval, i.e., domain shift and query-doc complexity. We initiate retrieval with mixed granularity within queries and documents to address these issues;
- We propose $MiXGR$, which further incorporates finer granularities within queries and documents, computes query-doc similarity over various granularity combinations, and fuses them as a united score. Our experiments across four datasets and six dense retrievers empirically reveal that $MiXGR$ significantly enhances existing retrievers on the scientific document retrieval and downstream QA tasks;
- Further analysis demonstrates the complementarity of metrics based on different granularities and the generalization of $MiXGR$ in retrieving units finer than documents.

2 Preliminary and Related works

Generalization of Dense Retrievers Dense retrievers generally employ a dual-encoder framework (Yih et al., 2011; Reimers and Gurevych, 2019) to separately encode queries and documents into compact vectors and measure relevance using a non-parametric similarity function (Musmann and Ermon, 2016). However, the simplicity of the similarity function (e.g., cosine similarity) can restrict expressiveness, leading to suboptimal generalization in new domains such as scientific fields that differ from original training data (Thakur et al., 2021). To improve dense retrievers’ adaptability across tasks, researchers have used data augmentation (Wang et al., 2022; Lin et al., 2023; Dai et al., 2023), continual learning (Chang et al., 2020; Sachan et al., 2021; Oguz et al., 2022), and task-aware training (Xin et al., 2022; Cheng et al., 2023). However, these methods still require training on domain-specific data, incurring additional computational costs. This work focuses on *zero-shot* generalization of dense retrievers to scientific fields by incorporating multi-granularity similarities within queries and documents.

Granularity in Retrieval For dense retrieval, the selection of the retrieval unit needs to balance the trade-off between completeness and compactness. Coarser units, like documents or fixed-length passages, theoretically encompass more context but may introduce extraneous information, adversely affecting retrievers and downstream tasks (Shi et al., 2023; Wang et al., 2023). Conversely, finer units like sentences are not always self-contained and may lose context, thereby hindering retrieval (Akkalyoncu Yilmaz et al., 2019; Yang et al., 2020). Additionally, some studies extend beyond complete sentences; for example, Lee et al. (2021a) use phrases as learning units to develop corresponding representations. Meanwhile, ColBERT (Khattab and Zaharia, 2020) addresses token-level query-doc interaction but is hampered by low efficiency.

Chen et al. (2023) propose using *propositions* as retrieval units, defined as atomic expressions of meaning (Min et al., 2023). These units are contextualized and self-contained, including necessary context through decontextualization, e.g., coreference resolution (Zhang et al., 2021). Proposition retrieval improves retrieval of documents with long-tail information, potentially benefiting domain-specific tasks. This motivates the use of

propositions as retrieval units for scientific document retrieval. Furthermore, we extend fine granularity to queries and enhance the query-doc similarity measurement, moving from a point-wise assessment between two vectors to integrating multiple query-doc granularity combinations.

Fusion within Retrieval Each type of retriever, sparse or dense, has its own strength and can be complementary with each other. Based on this insight, previous studies have explored the fusion of searches conducted by different retrievers as a zero-shot solution for domain adaptation (Thakur et al., 2021). A common method involves the convex combination, which linearly combines similarity scores (Karpukhin et al., 2020; Wang et al., 2021b; Ma et al., 2021). However, this approach is sensitive to the weighting of different metrics and score normalization, which complicates configuration across different setups (Chen et al., 2022).

In this work, we enhance retrieval by integrating searches across various query and document granularity levels for a given retriever. To avoid the limitations of convex combination, we use Rank Reciprocal Fusion (RRF; Cormack et al. 2009), a robust, non-parametric method (Chen et al., 2022), to aggregate these searches.

3 MixGR: Mix-Granularity Retrieval

3.1 Finer Units in Queries and Documents

We first decompose queries and documents into atomic units, i.e., subqueries and propositions, respectively. A proposition (or subquery) should meet the following three principal criteria (Min et al., 2023):

- Each proposition conveys a distinct semantic unit, collectively expressing the complete meaning.
- Propositions should be atomic and indivisible.
- According to Choi et al. (2021), propositions should be contextualized and self-contained, including all necessary text information such as resolved coreferences for clear interpretation.

Here, we employ an off-the-shelf model, *propositioner*,¹ for decomposing queries and documents (Chen et al., 2023). This model is developed by distilling the decomposition capacities of GPT-4 (Achiam et al., 2023) to a Flan-T5-Large model

¹<https://huggingface.co/chentong00/propositionizer-wiki-flan-t5-large>

	Query	Document
Accuracy (%)	96.3	94.7
IAA (%)	92.0	89.0

Table 1: Human-evaluated accuracy of query/document decomposition by *propositioner* (Chen et al., 2023).

(Chung et al., 2024) using Wikipedia as the corpus. We sample decomposition results from 100 queries and 100 documents from the datasets in §4.1 and manually label the correctness of decomposition as shown in Table 1. This model is shown to effectively decompose queries and documents into atomic units within scientific domains. Please see Appendix B for further details.

3.2 Multi-Granularity Similarity Calculation

Given these various granularities including queries, subqueries, documents and propositions, we extend the query-doc similarity metrics to include measurements across different combinations of granularities as depicted in Figure 2.

Notations The sets of queries and documents are denoted as \mathcal{Q} and \mathcal{D} , respectively. Given a retriever s , the similarity between a query $q \in \mathcal{Q}$ and a document $d \in \mathcal{D}$ is denoted as $s(q, d)$. A document d can be decomposed to N propositions, i.e., $d = [d_1, \dots, d_N]$. And a query q can be decomposed to M subqueries, i.e., $q = [q_1, \dots, q_M]$.

Query-doc s_{q-d} : The direct and original similarity between q and d is $s_{q-d}(q, d) \equiv s(q, d)$.

Query-prop s_{q-p} : Recent works (Chen et al., 2023) determine query-doc similarity by calculating the maximum similarity between the **query** and individual **propositions** within the document (Lee et al., 2021b; Chen et al., 2023). The computation of this metric, denoted as s_{q-p} , is as follows:

$$s_{q-p}(q, d) = \max_{i=1, \dots, N} \{s(q, d_i)\}. \quad (1)$$

Subquery-prop s_{s-p} : Considering that different parts of a query may be captured by various propositions within a document shown in Figure 1b, we further assess query-doc similarity by analyzing the relationships between **subqueries** and individual **propositions**. The similarity between a query and a document can be defined as the average similarity across subqueries, calculated by identifying the maximum similarity between one subquery and

each proposition, in analogy to MaxSim in ColBERT (Khattab and Zaharia, 2020). This metric, represented by s_{s-p} , is calculated as:

$$s_{s-p}(q, d) = \frac{1}{M} \sum_{i=1}^M \max_{j=1, \dots, N} \{s(q_i, d_j)\}. \quad (2)$$

3.3 Rank Reciprocal Fusion

We then use RRF to fuse these metrics across different query and document granularities:

$$s_f(q, d) = \frac{1}{1 + r_{q-d}(q, d)} + \frac{1}{1 + r_{q-p}(q, d)} + \frac{1}{1 + r_{s-p}(q, d)}, \quad (3)$$

where r_{q-d} , r_{q-p} , $r_{s-p} \in \mathbb{R}_{\geq 0}$ signify the rank of the retrieve results by s_{q-d} , s_{q-p} , and s_{s-p} , respectively. Technically, we retrieve the top- k results R_{q-d}^k , R_{q-p}^k , and R_{s-p}^k by s_{q-d} , s_{q-p} , and s_{s-p} , respectively, where k is set 200 empirically. When a query-doc pair (q', d') in one retrieval result does not exist in the other sets (e.g., $(q', d') \in R_{q-d}^k$ but $(q', d') \notin R_{q-p}^k$), we will calculate the missing similarity (e.g., $s_{q-p}(q', d')$) before aggregation.

4 Experimental Setting

4.1 Scientific Retrieval Datasets

We evaluate our approach on four different scientific retrieval tasks, including NFCorpus (Boteva et al., 2016), SciDocs (Cohan et al., 2020), SciFact (Wadden et al., 2020), and SciQ (Welbl et al., 2017), as shown in Table 4 in Appendix A. We employ the *propositioner* released by Chen et al. (2023) mentioned in §3.1 to break down both queries and documents to atomic units. As we focus with priority on query-doc complexity in scientific domains, we report the experiments and analysis on the subset of the queries which contain multiple subqueries.

4.2 Dense Retrievers

We evaluate the performance of six off-the-shelf dense retrievers, both supervised and unsupervised. Supervised retrievers are trained using human-labeled query-doc pairs in general domains,² while unsupervised models do not require labeled data. These retrievers encode the queries and index the corpus at both document and proposition levels:

²The supervised retrievers used in our experiment have not been trained on these four datasets.

Retriever	Setup	NFCorpus		SciDocs		SciFact		SciQ		Avg.	
		ND@5	ND@20	ND@5	ND@20	ND@5	ND@20	ND@5	ND@20	ND@5	ND@20
Unsupervised Dense Retrievers											
SimCSE	s_{q-d}	16.2	13.3	7.6	9.7	27.1	31.2	62.3	67.3	28.3	30.4
	s_{q-p}	20.0	16.4	<u>8.2</u>	<u>11.1</u>	<u>32.8</u>	<u>37.2</u>	75.6	78.5	34.1	35.8
	s_{s-p}	22.8	18.3	7.3	10.5	32.7	36.9	80.9	<u>83.2</u>	35.9	<u>37.2</u>
	MixGR	<u>22.3</u>	<u>18.1</u>	9.1	12.2	34.8	39.8	84.0	85.5	37.5 (+32.5%)	38.9 (+28.0%)
Contriever	s_{q-d}	42.2	34.9	13.5	18.5	<u>64.5</u>	68.5	67.2	70.0	46.9	48.0
	s_{q-p}	<u>43.0</u>	<u>35.5</u>	<u>14.5</u>	<u>19.4</u>	64.0	<u>68.9</u>	79.7	81.0	50.3	51.2
	s_{s-p}	41.4	34.9	13.5	18.3	63.2	67.5	83.6	<u>84.6</u>	50.4	<u>51.3</u>
	MixGR	44.0	37.1	15.5	20.7	66.4	71.0	85.2	86.7	52.8 (+12.6%)	53.9 (+12.3%)
Supervised Dense Retrievers											
DPR	s_{q-d}	25.1	20.7	7.3	10.4	31.8	37.7	60.6	64.1	31.2	33.2
	s_{q-p}	25.2	20.6	<u>7.8</u>	<u>10.6</u>	36.1	40.5	63.6	67.9	33.2	34.9
	s_{s-p}	<u>26.5</u>	<u>21.4</u>	6.4	10.0	<u>37.1</u>	<u>41.3</u>	<u>67.7</u>	<u>70.7</u>	<u>34.4</u>	<u>35.9</u>
	MixGR	27.7	22.9	8.2	11.5	39.4	43.6	73.6	76.1	37.2 (+19.2%)	38.5 (+16.0%)
ANCE	s_{q-d}	29.9	24.4	<u>9.3</u>	<u>13.1</u>	41.5	45.3	<u>66.4</u>	<u>69.1</u>	36.8	38.0
	s_{q-p}	29.4	24.0	9.2	12.9	43.3	46.4	62.3	66.4	36.0	37.4
	s_{s-p}	<u>30.3</u>	<u>24.5</u>	7.5	11.9	<u>43.5</u>	47.3	66.1	<u>69.1</u>	<u>36.9</u>	<u>38.2</u>
	MixGR	31.9	25.9	9.6	14.1	46.8	49.9	74.4	76.8	40.7 (+10.6%)	41.7 (+9.7%)
TAS-B	s_{q-d}	42.3	34.1	13.8	19.3	60.1	<u>65.6</u>	84.8	86.3	50.2	<u>51.3</u>
	s_{q-p}	<u>42.5</u>	<u>34.4</u>	<u>14.3</u>	<u>18.1</u>	60.7	64.4	<u>85.6</u>	86.3	<u>50.8</u>	50.8
	s_{s-p}	40.9	33.1	12.6	17.2	<u>61.7</u>	65.0	85.3	<u>86.6</u>	50.1	50.5
	MixGR	43.6	35.2	14.0	19.6	62.7	66.9	90.5	91.0	52.7 (+5.0%)	53.2 (+3.7%)
GTR	s_{q-d}	42.1	34.1	13.6	<u>18.9</u>	58.3	62.2	83.3	84.4	49.3	49.9
	s_{q-p}	<u>42.3</u>	<u>34.4</u>	13.2	18.0	60.6	<u>63.3</u>	85.8	86.5	<u>50.5</u>	<u>50.6</u>
	s_{s-p}	41.5	33.6	11.6	16.2	<u>58.4</u>	62.0	<u>88.5</u>	<u>89.0</u>	50.0	50.2
	MixGR	43.3	35.6	13.6	19.2	60.9	64.5	92.9	93.0	52.7 (+6.9%)	53.1 (+6.4%)

Table 2: Document Retrieval Performance (nDCG@ $k = 5, 20$ in percentage, abbreviated as ND@ k): We evaluated four distinct scientific retrieval datasets using two unsupervised and four supervised retrievers. The retrieval results were compared among various metrics: s_{q-d} (previous query-doc similarity), s_{q-p} (Chen et al., 2023), s_{s-p} , and MixGR, as detailed in §3.2. **Bold** presents the best performance across the metrics, while underline denotes the second-best performance. MixGR outperforms all three other metrics, where the percentage in parentheses indicates the relative improvement compared with s_{q-d} .

- SimCSE (Gao et al., 2021) employs a BERT-base (Devlin et al., 2019) encoder trained on randomly selected unlabeled Wikipedia sentences.
- Contriever (Izacard et al., 2022) is an unsupervised retriever evolved from a BERT-base encoder, contrastively trained on segments from unlabelled web and Wikipedia documents.
- DPR (Karpukhin et al., 2020) is built with a dual-encoder BERT-base architecture, finetuned on a suite of open-domain datasets with labels, such as SQuAD (Rajpurkar et al., 2016).
- ANCE (Xiong et al., 2021) mirrors the configuration of DPR but incorporates a training scheme of Approximate Nearest Neighbor Negative Contrastive Estimation (ANCE).
- TAS-B (Hofstätter et al., 2021) is a dual-encoder BERT-base model distilled from ColBERT on MS MARCO (Nguyen et al., 2016).
- GTR (Ni et al., 2022) is a T5-base encoder, focusing on generalization, pre-trained on unlabeled

QA pairs, and fine-tuned on labeled data including MS MARCO.

More details on retrievers and experimental setups are presented in Appendices C and D.

4.3 Document Retrieval Evaluation

We assess the performance of MixGR in the task of document retrieval. Due to input length limitations for retrievers (Karpukhin et al., 2020), we divide each document into fixed-length chunks of up to 128 words. In practice, for MixGR and baselines, we identify the retrieved chunks, map them back to their original documents, and return the top- k documents. We use Normalised Cumulative Discount Gain (nDCG@ k) as the evaluation metrics for document retrieval. Unlike Recall@ k , which only indicates the presence of golden documents in the retrieved list, nDCG@ k also accounts for both the ranking of retrievals and the relevance judgment of golden documents (Thakur et al., 2021). The baselines will be the metrics containing the ho-

332 homogeneous granularity introduced in the previous
333 section, i.e., s_{q-d} , s_{q-p} and s_{s-p} .

334 4.4 Downstream QA Evaluation

335 As previously mentioned, scientific documents are
336 vital for LLMs due to the rapid advancements in
337 science and the limited availability of such con-
338 tent in training datasets. To better understand how
339 $MixGR$ enhances downstream QA tasks, we im-
340 plement the *retrieval-then-read* approach on two
341 datasets SciQ and SciFact. We retrieve and rank the
342 top- k documents based on scores, s_{q-d} and $MixGR$,
343 then concatenate them to form the context. During
344 our evaluations, we limit the number of document
345 chunks retrieved to 1 and 3—thus, only the top
346 k documents are injected into the reader model.
347 We assess the performance by measuring the Ex-
348 act Match (EM) rate—the proportion of responses
349 where the predicted answer perfectly aligns with
350 the reference answer (Kamalloo et al., 2023), de-
351 noted as $EM@k$. Specifically, we utilize Llama-3-
352 8B-Instruct³ (Touvron et al., 2023) as the reader
353 model. We take the original query-doc retrieval
354 setup, i.e., retrieval based on s_{q-d} , as the baseline.
355 Please refer to Appendix E for more details.

356 5 Results

357 This section analyzes the impact of mixed-
358 granularity retrieval on document retrieval and
359 downstream applications. We highlight the effec-
360 tiveness of our proposed fine-grained and mixed-
361 granularity approaches in enhancing performance
362 across various metrics.

363 5.1 Document Retrieval

364 Table 2 reports the results of document retrieval.
365 We observe that retrieval by $MixGR$ outperforms
366 all single-granularity retrieval with both unsuper-
367 vised and supervised dense retrievers in most cases.

368 With unsupervised retrievers, $MixGR$ signifi-
369 cantly outperforms the query-doc similarity, s_{q-d} ,
370 across all four datasets. There is an average
371 $nDCG@5$ improvement of +9.2 and +5.9 (32.5%
372 and 12.6% relatively) for SimCSE and Contriever,
373 respectively.

374 With supervised retrievers, improvements associ-
375 ated with $MixGR$ are also observed, although they
376 are not as significant as with unsupervised retriev-
377 ers. This indicates that $MixGR$ effectively narrows

³<https://huggingface.co/meta-llama/Meta-Llama-3-8B-Instruct>

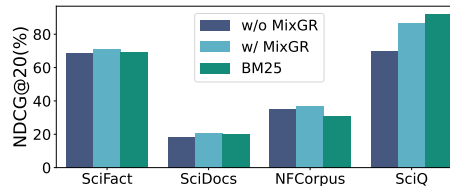


Figure 3: Comparison between BM25 and Contriever (w/ and w/o $MixGR$) on $nDCG@20$: Contriever w/ $MixGR$ outperforms BM25 in three out of four datasets.

378 the distributional gap between dense retrievers and
379 scientific domains.

380 Unsupervised retrievers benefit more from 381 $MixGR$ than supervised ones.

382 Remarkably, with $MixGR$, the unsupervised retriever Contriever out-
383 performs supervised models, as evidenced by its superior average results across four datasets. This
384 result is particularly significant given that Contriever typically underperforms compared to TAS-
385 B and GTR when evaluated using traditional query-
386 document similarity measures. Additionally, the
387 study (Thakur et al., 2021) reveals that sparse re-
388 trievers like BM25 often excel over dense retrievers
389 in domain-specific retrieval tasks. As shown in Fig-
390 ure 3, Contriever outperforms BM25 in three out
391 of four datasets when applied with $MixGR$. Simi-
392 larly, SimCSE also outperforms DPR under the
393 $MixGR$ scheme. These findings emphasize the sub-
394 stantial enhancements that $MixGR$ contributes to
395 unsupervised retrievers within scientific domains.
396
397

398 Finer granularity helps retrieval more.

399 Among three metrics within $MixGR$, the subquery-
400 proposition measurement s_{s-p} shows a distinct
401 advantage over the other two, as highlighted by
402 the underlined results in Table 2. The original
403 query-doc metric, s_{q-d} , outperforms the subquery-
404 proposition measurement only when using the re-
405 triever TAS-B. These findings corroborate and ex-
406 pand upon Chen et al. (2023), suggesting that finer
407 query-doc similarity measurement significantly im-
408 proves document retrieval performance.

409 5.2 Downstream QA Tasks

410 Table 3 reports the results of scientific question an-
411 swering when the documents retrieved by $MixGR$
412 are fed into LLMs, i.e. the readers. It is observed
413 that EM scores achieved with $MixGR$ generally
414 surpass those of the baseline across two datasets,
415 six dense retrievers, and multiple numbers of in-
416 put documents. This underscores the effectiveness

Setup		SciFact		SciQ	
		EM@1	EM@3	EM@1	EM@3
Unsupervised Dense Retrievers					
SimCSE	s_{q-d}	50.0	61.6	54.7	58.2
	MixGR	48.3	62.8	61.3	66.4
Contriever	s_{q-d}	63.4	75.6	53.9	63.3
	MixGR	64.0	70.9	61.7	66.0
Supervised Dense Retrievers					
DPR	s_{q-d}	51.2	59.9	52.0	57.4
	MixGR	51.7	65.7	57.4	62.5
ANCE	s_{q-d}	51.7	65.1	52.7	59.4
	MixGR	57.6	69.2	54.7	62.9
TAS-B	s_{q-d}	62.8	74.4	60.5	66.4
	MixGR	62.2	70.3	64.5	67.6
GTR	s_{q-d}	61.0	72.1	59.8	64.8
	MixGR	62.8	73.8	64.1	66.0

Table 3: Scientific Question Answering on SciFact and SciQ using Llama-3-8B-Instruct (Touvron et al., 2023): the top-1 and 3 document chunks retrieved by retrievers, following the metrics s_{q-d} and MixGR, were fed into the reader. **Bold** indicates the better performance.

of MixGR in enhancing the performance of downstream QA tasks.

6 Analysis

In this section, we explore the complementary advantages of various similarity metrics across multiple granularities within MixGR through an ablation study. Although the finer-granularity metric, s_{s-p} , generally enhances performance as previously discussed, it can occasionally result in degradation when compared to original query-document similarity s_{q-d} . We identify specific conditions under which the finer-granularity metric offers greater benefits. Previous works (Chen et al., 2023) primarily explored multiple granularities in *documents*. We conduct a control experiment to highlight the significance of incorporating multiple granularities in *queries* in the MixGR framework, which also validate the generalization of MixGR on the retrieval units finer than documents.

6.1 Ablation Study

In our ablation study, we conducted a systematic evaluation of the impact of various granularity measures— s_{q-d} (query-doc similarity), s_{q-p} (query-prop similarity), and s_{s-p} (subquery-prop similarity)—on the performance of six retrievers. By individually omitting each of these measures from the calculation of MixGR as defined in Equation 3, we assessed the significance of each granularity level. Specifically, the extent of performance degradation upon removal of a measure indicates

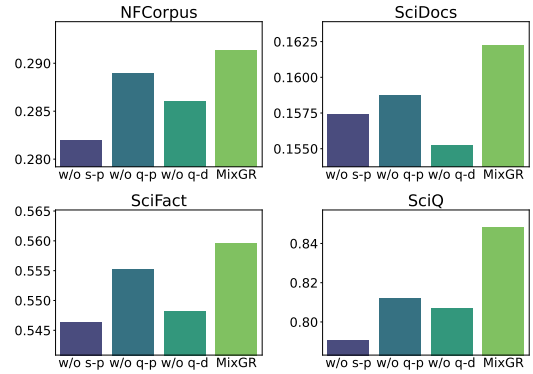


Figure 4: Ablation study of MixGR on the nDCG@20 metrics averaged on six retrievers: MixGR achieves optimal performance when combining these three metrics, indicating their complementary nature.

its importance; greater degradation suggests higher importance of that particular granularity metric.

As illustrated in Figure 4, the nDCG@20 performance declined across all three setups and datasets, demonstrating that the metrics are complementary to each other. The degree of performance degradation varied across different configurations, highlighting the importance of each granularity measure. Notably, the most significant declines in performance consistently occurred in configurations excluding s_{q-d} and s_{s-p} . This observation suggests that s_{q-p} , while beneficial, is the *least* critical measure for retrieval tasks in scientific domains. Please refer to Table 6 in Appendix F.1 for detailed results.

6.2 When is finer granularity beneficial?

Therefore, to more effectively compare the impacts of s_{q-d} and s_{s-p} , we categorized the *correctly* retrieved pairs (complex query, ⁴ doc) by MixGR in SciFact, using SimCSE, into two distinct groups:

- $r_{q-d} \succ r_{s-p}$: The query-doc rank of s_{q-d} is higher than the subquery-prop rank of s_{s-p} ;
- $r_{q-d} \prec r_{s-p}$: The query-doc rank of s_{q-d} is lower than the subquery-prop rank of s_{s-p} .

Upon analyzing the number of propositions in documents, a significant pattern emerges: based on the distributions present in Figure 5, the number of propositions in $r_{q-d} \prec r_{s-p}$ is generally higher than in $r_{q-d} \succ r_{s-p}$. This underscores the importance of incorporating finer units within documents, especially for those containing more propositions, and suggests potential degradation in dense retrievers

⁴We refer *complex query* as the query containing no fewer than three subqueries.

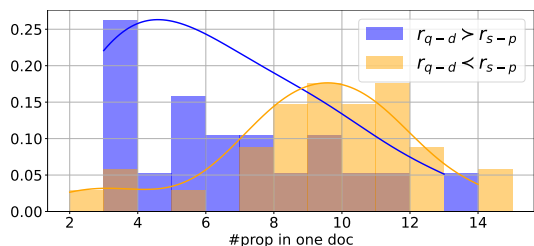


Figure 5: Distribution of proposition number within documents in two sets. There are more propositions within document when $r_{q-d} < r_{s-p}$ than $r_{q-d} > r_{s-p}$.

when handling such documents. For other retrievers’ results, please refer to Appendix F.3.

6.3 MixGR on Proposition Retrieval

Previous sections present the effectiveness of MixGR on scientific document retrieval. While previous works (Chen et al., 2023) focus on finer document granularity, we specifically assess MixGR on the proposition as the retrieval units. This controlled study highlights the benefits of MixGR, which incorporates different granularities within queries and documents, in general text retrieval beyond document-level granularity.

For a given query q and a proposition p , the conventional similarity is denoted by $s_{q-p}^p \equiv s(q, p)$. When the query is further broken down into multiple sub-queries, we introduce a finer granularity measure, s_{s-p}^p , which is defined as the maximum similarity between these sub-queries and the proposition. s_{s-p}^p is mathematically defined as follows:

$$s_{s-p}^p(q, p) = \max_{i=1, \dots, M} \{s(q_i, p)\}. \quad (4)$$

Therefore, the merged score by RRF, $s_f^p(q, p)$, is calculated as:

$$s_f^p(q, p) = \frac{1}{1 + r_{q-p}^p(q, p)} + \frac{1}{1 + r_{s-p}^p(q, p)}, \quad (5)$$

where r_{q-p}^p and r_{s-p}^p signify the rank of the retrieve results by s_{q-p}^p and s_{s-p}^p , respectively.

Following $s_{q-p}^p(q, p)$ and $s_f^p(q, p)$, we input the first 50 and 200 words in propositions retrieved with SimCSE on SciFact and SciQ into the reader Llama-3-8B-Instruct. This process adheres to the same setups outlined in §4.4. As shown in Figure 6, the performance advance observed with mixed-granularity retrieval on propositions, compared to the original query-prop similarity, demonstrates the effectiveness of using mixed-granularity in retrieval. This substantiates the generalizability of MixGR beyond document-level granularity. Please refer to Appendix F.2 for details.

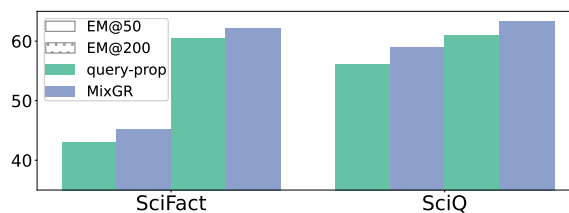


Figure 6: Proposition retrieval with MixGR: We evaluate Exact Match of Llama-3-8B-Instruct on SciFact and SciQ with the first 50 and 200 words of propositions, i.e., EM@50 and EM@200, retrieved by SimCSE as the context. Please refer to Table 7 for other retrievers in Appendix F.2.

6.4 Prospect: Adaptive MixGR

Here, we outline potential future research directions. In §6.1, we observed the complementary nature of retrieval results achieved using different granularities. Additionally, as discussed in §6.2, we noted a distinct pattern where retrieval guided by a specific granularity outperforms others. These findings indicate that metrics based on different granularities each have relatively distinct strengths in specific contexts, presenting opportunities for further exploration. Unlike the non-parametric method of fusion by RRF, which overlooks the relative importance of components, an adaptive approach could enhance fusion and, consequently, improve retrieval performance with dense retrievers—a prospect we aim to explore in future research.

7 Conclusion

In this work, we identify key challenges for dense retrievers in scientific document retrieval, namely domain shift and query-document complexity. In response, we propose a zero-shot approach, MixGR, that utilizes atomic components in queries and documents to calculate their similarity with greater nuance. We then use Rank Reciprocal Fusion (RRF) to integrate these metrics, modeling query-doc similarity at different granularities into a unified score that enhances document retrieval.

Our experiments demonstrate that MixGR significantly enhances the existing dense retriever on document retrieval within the scientific domain. Moreover, MixGR has proven beneficial for downstream applications such as scientific QA. The analysis reveals a synergistic relationship among the components of MixGR, and suggests evolving our non-parametric fusion framework into a more general method as a future research direction.

551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599

Limitation

Our work explores retrieval guided by an integral metric that incorporates various levels of granularity. We identify several limitations in our approach: (1) *Coverage of Retrievers*: Our study categorizes dense retrievers into supervised and unsupervised models, yet all utilize a dual-encoder structure. Future studies could include a more diverse array of retriever architectures. (2) *Coverage of Domains*: While our main focus is on the scientific domain, and we extend to three additional domains in Appendix G, there are still many domains we have not explored. (3) *Languages*: Our research is limited to an English corpus. The applicability of MixGR in multilingual contexts also deserves further validation and exploration.

Ethical Statements

We foresee no ethical concerns and potential risks in our work. All of the retrieval models and datasets are open-sourced, as shown in Table 10 in Appendix H. The LMs we applied are also publicly available. Given our context, the outputs of LLMs should be insensitive.

References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Zeynep Akkalyoncu Yilmaz, Wei Yang, Haotian Zhang, and Jimmy Lin. 2019. [Cross-domain modeling of sentence-level evidence for document retrieval](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3490–3496, Hong Kong, China. Association for Computational Linguistics.

Abeba Birhane, Atoosa Kasirzadeh, David Leslie, and Sandra Wachter. 2023. Science in the age of large language models. *Nature Reviews Physics*, 5(5):277–280.

Vera Boteva, Demian Gholipour, Artem Sokolov, and Stefan Riezler. 2016. A full-text learning to rank dataset for medical information retrieval. In *Advances in Information Retrieval: 38th European Conference on IR Research, ECIR 2016, Padua, Italy, March 20–23, 2016. Proceedings 38*, pages 716–722. Springer.

Wei-Cheng Chang, X Yu Felix, Yin-Wen Chang, Yiming Yang, and Sanjiv Kumar. 2020. [Pre-training tasks for embedding-based large-scale retrieval](#). In *International Conference on Learning Representations*. 600–602.

Tao Chen, Mingyang Zhang, Jing Lu, Michael Bendersky, and Marc Najork. 2022. Out-of-domain semantics to the rescue! zero-shot hybrid retrieval models. In *European Conference on Information Retrieval*, pages 95–110. Springer. 604–608.

Tong Chen, Hongwei Wang, Sihao Chen, Wenhao Yu, Kaixin Ma, Xinran Zhao, Dong Yu, and Hongming Zhang. 2023. Dense x retrieval: What retrieval granularity should we use? *arXiv preprint arXiv:2312.06648*. 609–613.

Hao Cheng, Hao Fang, Xiaodong Liu, and Jianfeng Gao. 2023. [Task-aware specialization for efficient and robust dense retrieval for open-domain question answering](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 1864–1875, Toronto, Canada. Association for Computational Linguistics. 614–620.

Eunsol Choi, Jennimaria Palomaki, Matthew Lamm, Tom Kwiatkowski, Dipanjan Das, and Michael Collins. 2021. Decontextualization: Making sentences stand-alone. *Transactions of the Association for Computational Linguistics*, 9:447–461. 621–623.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2024. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70):1–53. 626–630.

Arman Cohan, Sergey Feldman, Iz Beltagy, Doug Downey, and Daniel Weld. 2020. [SPECTER: Document-level representation learning using citation-informed transformers](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2270–2282, Online. Association for Computational Linguistics. 631–637.

Gordon V Cormack, Charles LA Clarke, and Stefan Buettcher. 2009. Reciprocal rank fusion outperforms condorcet and individual rank learning methods. In *Proceedings of the 32nd international ACM SIGIR conference on Research and development in information retrieval*, pages 758–759. 638–642.

Zhuyun Dai, Vincent Y Zhao, Ji Ma, Yi Luan, Jianmo Ni, Jing Lu, Anton Bakalov, Kelvin Guu, Keith Hall, and Ming-Wei Chang. 2023. [Promptagator: Few-shot dense retrieval from 8 examples](#). In *The Eleventh International Conference on Learning Representations*. 644–649.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *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)*, pages 650–652. 651–655.

657	4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.	
658		
659	Shai Erera, Michal Shmueli-Scheuer, Guy Feigenblat, Ora Peled Nakash, Odellia Boni, Haggai Roitman, Doron Cohen, Bar Weiner, Yosi Mass, Or Rivlin, Guy Lev, Achiya Jerbi, Jonathan Herzig, Yufang Hou, Charles Jochim, Martin Gleize, Francesca Bonin, Francesca Bonin, and David Konopnicki. 2019. A summarization system for scientific documents . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP): System Demonstrations</i> , pages 211–216, Hong Kong, China. Association for Computational Linguistics.	
660		
661		
662		
663		
664		
665		
666		
667		
668		
669		
670		
671		
672	Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. SimCSE: Simple contrastive learning of sentence embeddings . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 6894–6910, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	
673		
674		
675		
676		
677		
678		
679	Igor Grossmann, Matthew Feinberg, Dawn C Parker, Nicholas A Christakis, Philip E Tetlock, and William A Cunningham. 2023. Ai and the transformation of social science research. <i>Science</i> , 380(6650):1108–1109.	
680		
681		
682		
683		
684	Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. 2020. Array programming with NumPy . <i>Nature</i> , 585(7825):357–362.	
685		
686		
687		
688		
689		
690		
691		
692		
693		
694		
695	Sebastian Hofstätter, Sheng-Chieh Lin, Jheng-Hong Yang, Jimmy Lin, and Allan Hanbury. 2021. Efficiently teaching an effective dense retriever with balanced topic aware sampling. In <i>Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , pages 113–122.	
696		
697		
698		
699		
700		
701		
702	John D Hunter. 2007. Matplotlib: A 2d graphics environment. <i>Computing in science & engineering</i> , 9(03):90–95.	
703		
704		
705	Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2022. Unsupervised dense information retrieval with contrastive learning . <i>Transactions on Machine Learning Research</i> .	
706		
707		
708		
709		
710	Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. <i>ACM Comput. Surv.</i> , 55(12):248:1–248:38.	
711		
712		
713		
714		
	Ehsan Kamalloo, Nouha Dziri, Charles Clarke, and Davood Rafiei. 2023. Evaluating open-domain question answering in the era of large language models . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 5591–5606, Toronto, Canada. Association for Computational Linguistics.	715
		716
		717
		718
		719
		720
		721
	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.	722
		723
		724
		725
		726
		727
		728
	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.	729
		730
		731
		732
		733
		734
	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. 2019. Natural questions: a benchmark for question answering research. <i>Transactions of the Association for Computational Linguistics</i> , 7:453–466.	735
		736
		737
		738
		739
		740
		741
	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In <i>Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles</i> .	742
		743
		744
		745
		746
		747
		748
	Jinhyuk Lee, Mujeen Sung, Jaewoo Kang, and Danqi Chen. 2021a. Learning dense representations of phrases at scale . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 6634–6647, Online. Association for Computational Linguistics.	749
		750
		751
		752
		753
		754
		755
		756
	Jinhyuk Lee, Alexander Wettig, and Danqi Chen. 2021b. Phrase retrieval learns passage retrieval, too . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 3661–3672, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	757
		758
		759
		760
		761
		762
	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. <i>Advances in Neural Information Processing Systems</i> , 33:9459–9474.	763
		764
		765
		766
		767
		768
	Sheng-Chieh Lin, Akari Asai, Minghan Li, Barlas Oguz, Jimmy Lin, Yashar Mehdad, Wen-tau Yih, and Xilun Chen. 2023. How to train your dragon: Diverse augmentation towards generalizable dense retrieval . In	769
		770
		771
		772

773	<i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 6385–6400, Singapore. Association for Computational Linguistics.	830
774		831
775		832
776	Simon Lupart, Thibault Formal, and Stéphane Clinchant.	833
777	2023. Ms-shift: An analysis of ms marco distribution shifts on neural retrieval. In <i>European Conference on Information Retrieval</i> , pages 636–652. Springer.	834
778		835
779		836
780	Xueguang Ma, Kai Sun, Ronak Pradeep, and Jimmy Lin.	837
781	2021. A replication study of dense passage retriever. <i>arXiv preprint arXiv:2104.05740</i> .	838
782		839
783	Macedo Maia, Siegfried Handschuh, André Freitas, Brian Davis, Ross McDermott, Manel Zarrouk, and Alexandra Balahur. 2018. Www’18 open challenge: financial opinion mining and question answering. In <i>Companion proceedings of the the web conference 2018</i> , pages 1941–1942.	840
784		841
785		842
786		843
787		844
788		845
789	Sewon Min, Kalpesh Krishna, Xinxu Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. FActScore: Fine-grained atomic evaluation of factual precision in long form text generation . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 12076–12100, Singapore. Association for Computational Linguistics.	846
790		847
791		848
792		849
793		850
794		851
795		852
796		853
797	Stephen Mussmann and Stefano Ermon. 2016. Learning and inference via maximum inner product search. In <i>International Conference on Machine Learning</i> , pages 2587–2596. PMLR.	854
798		855
799		856
800		857
801	Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. 2016. Ms marco: A human-generated machine reading comprehension dataset.	858
802		859
803		860
804		861
805	Jianmo Ni, Chen Qu, Jing Lu, Zhuyun Dai, Gustavo Hernandez Abrego, Ji Ma, Vincent Zhao, Yi Luan, Keith Hall, Ming-Wei Chang, and Yinfei Yang. 2022. Large dual encoders are generalizable retrievers . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 9844–9855, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	862
806		863
807		864
808		865
809		866
810		867
811		868
812		869
813	Barlas Oguz, Kushal Lakhotia, Anshit Gupta, Patrick Lewis, Vladimir Karpukhin, Aleksandra Piktus, Xilun Chen, Sebastian Riedel, Scott Yih, Sonal Gupta, and Yashar Mehdad. 2022. Domain-matched pre-training tasks for dense retrieval . In <i>Findings of the Association for Computational Linguistics: NAACL 2022</i> , pages 1524–1534, Seattle, United States. Association for Computational Linguistics.	870
814		871
815		872
816		873
817		874
818		875
819		876
820		877
821	Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. Pytorch: An imperative style, high-performance deep learning library . In <i>Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada</i> , pages 8024–8035.	878
822		879
823		880
824		881
825		882
826		883
827		884
828		885
829		886

887	Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller,	Johannes Welbl, Nelson F. Liu, and Matt Gardner. 2017.	944
888	Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-	Crowdsourcing multiple choice science questions .	945
889	thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan	In <i>Proceedings of the 3rd Workshop on Noisy User-</i>	946
890	Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,	<i>generated Text</i> , pages 94–106, Copenhagen, Den-	947
891	Isabel Kloumann, Artem Korenev, Punit Singh Koura,	mark. Association for Computational Linguistics.	948
892	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-		
893	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien	949
894	tinet, Todor Mihaylov, Pushkar Mishra, Igor Moly-	Chaumond, Clement Delangue, Anthony Moi, Pier-	950
895	bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-	ric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz,	951
896	stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,	et al. 2019. Huggingface’s transformers: State-of-	952
897	Ruan Silva, Eric Michael Smith, Ranjan Subrama-	the-art natural language processing . <i>ArXiv preprint</i> ,	953
898	nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-	abs/1910.03771.	954
899	lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,		
900	Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,	Ji Xin, Chenyan Xiong, Ashwin Srinivasan, Ankita	955
901	Melanie Kambadur, Sharan Narang, Aurélien Ro-	Sharma, Damien Jose, and Paul Bennett. 2022. Zero-	956
902	driguez, Robert Stojnic, Sergey Edunov, and Thomas	shot dense retrieval with momentum adversarial do-	957
903	Scialom. 2023. Llama 2: Open foundation and fine-	main invariant representations . In <i>Findings of the As-</i>	958
904	tuned chat models. <i>CoRR</i> , abs/2307.09288.	<i>sociation for Computational Linguistics: ACL 2022</i> ,	959
		pages 4008–4020, Dublin, Ireland. Association for	960
		Computational Linguistics.	961
905	Henning Wachsmuth, Shahbaz Syed, and Benno Stein.		
906	2018. Retrieval of the best counterargument without	Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang,	962
907	prior topic knowledge . In <i>Proceedings of the 56th</i>	Jialin Liu, Paul N. Bennett, Junaid Ahmed, and	963
908	<i>Annual Meeting of the Association for Computational</i>	Arnold Overwijk. 2021. Approximate nearest neigh-	964
909	<i>Linguistics (Volume 1: Long Papers)</i> , pages 241–251,	bor negative contrastive learning for dense text re-	965
910	Melbourne, Australia. Association for Computational	trieval . In <i>International Conference on Learning</i>	966
911	Linguistics.	<i>Representations</i> .	967
912	David Wadden, Shanchuan Lin, Kyle Lo, Lucy Lu	Yinfei Yang, Daniel Cer, Amin Ahmad, Mandy Guo,	968
913	Wang, Madeleine van Zuylen, Arman Cohan, and	Jax Law, Noah Constant, Gustavo Hernandez Abrego,	969
914	Hannaneh Hajishirzi. 2020. Fact or fiction: Verifying	Steve Yuan, Chris Tar, Yun-Hsuan Sung, et al. 2020.	970
915	scientific claims . In <i>Proceedings of the 2020 Con-</i>	Multilingual universal sentence encoder for semantic	971
916	<i>ference on Empirical Methods in Natural Language</i>	retrieval. In <i>Proceedings of the 58th Annual Meet-</i>	972
917	<i>Processing (EMNLP)</i> , pages 7534–7550, Online. As-	<i>ing of the Association for Computational Linguistics:</i>	973
918	sociation for Computational Linguistics.	<i>System Demonstrations</i> , pages 87–94.	974
919	Kexin Wang, Nils Reimers, and Iryna Gurevych. 2021a.	Wen-tau Yih, Kristina Toutanova, John C. Platt, and	975
920	TSDAE: Using transformer-based sequential denois-	Christopher Meek. 2011. Learning discriminative	976
921	ing auto-encoder for unsupervised sentence embed-	projections for text similarity measures . In <i>Proceed-</i>	977
922	ding learning . In <i>Findings of the Association for</i>	<i>ings of the Fifteenth Conference on Computational</i>	978
923	<i>Computational Linguistics: EMNLP 2021</i> , pages	<i>Natural Language Learning</i> , pages 247–256, Port-	979
924	671–688, Punta Cana, Dominican Republic. Associa-	land, Oregon, USA. Association for Computational	980
925	tion for Computational Linguistics.	Linguistics.	981
926	Kexin Wang, Nandan Thakur, Nils Reimers, and Iryna	Hongming Zhang, Xinran Zhao, and Yangqiu Song.	982
927	Gurevych. 2022. GPL: Generative pseudo labeling	2021. A brief survey and comparative study of recent	983
928	for unsupervised domain adaptation of dense retrieval .	development of pronoun coreference resolution in	984
929	In <i>Proceedings of the 2022 Conference of the North</i>	English . In <i>Proceedings of the Fourth Workshop on</i>	985
930	<i>American Chapter of the Association for Computa-</i>	<i>Computational Models of Reference, Anaphora and</i>	986
931	<i>tional Linguistics: Human Language Technologies</i> ,	<i>Coreference</i> , pages 1–11, Punta Cana, Dominican	987
932	pages 2345–2360, Seattle, United States. Association	Republic. Association for Computational Linguistics.	988
933	for Computational Linguistics.		
934	Shuai Wang, Shengyao Zhuang, and Guido Zuccon.	Susan Zhang, Stephen Roller, Naman Goyal, Mikel	989
935	2021b. Bert-based dense retrievers require interpo-	Artetxe, Moya Chen, Shuohui Chen, Christopher	990
936	lation with bm25 for effective passage retrieval. In	Dewan, Mona T. Diab, Xian Li, Xi Victoria Lin,	991
937	<i>Proceedings of the 2021 ACM SIGIR international</i>	Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shus-	992
938	<i>conference on theory of information retrieval</i> , pages	ter, Daniel Simig, Punit Singh Koura, Anjali Srid-	993
939	317–324.	har, Tianlu Wang, and Luke Zettlemoyer. 2022.	994
		OPT: open pre-trained transformer language mod-	995
		els. <i>CoRR</i> , abs/2205.01068.	996
940	Zhiruo Wang, Jun Araki, Zhengbao Jiang, Md Rizwan		
941	Parvez, and Graham Neubig. 2023. Learning to filter		
942	context for retrieval-augmented generation. <i>arXiv</i>		
943	<i>preprint arXiv:2311.08377</i> .		

Appendix

997

A Datasets

998

Different from the setup of the original dataset, we split one document into several chunks with a maximum of 128 words. This is because some dense retrievers such as DPR (Karpukhin et al., 2020) have the requirement of maximum input. Too long inputs will be overflow, leading to the loss of information. The chunk selected can be used to locate the document in the original dataset during the evaluation. Specifically, for SciQ, we reformulate the dataset from a QA task to a retrieval task. Originally, this task aims to answer scientific questions given the context. We collect the contexts in training, validation and test sets as the corpus.

999

1000

1001

1002

1003

1004

1005

Also, we will explain our motivation of focusing the queries containing subqueries:

1006

- Chen et al. (2023) have studied the advantage of using propositions, i.e., the atomic units within documents, as the retrieval units given a complete query. And MixGR will not affect the retrieval results of single-subquery queries.
- In this work, we highlight the advantages of mixed-granularity retrieval that incorporates finer units in both queries and documents. Queries containing multiple subqueries are particularly well-suited to our research problem, as they will have different combinations with the documents.

1007

1008

1009

1010

1011

1012

Statistic	NFCorpus (Boteva et al., 2016)	SciDocs (Cohan et al., 2020)	SciFact (Wadden et al., 2020)	SciQ (Welbl et al., 2017)
#Query	1 016	1 000	1 109	884
#Multi-semantics queries	647	206	283	256
#Subqueries	3 337	522	614	874
#Documents	3 633	25 657	5 183	12 241
#Propositions	67 110	351 802	87 190	91 635

Table 4: Statistics for the NFCorpus, SciDocs, SciFact, and SciQ datasets.

B Query and Document Decomposition

1013

Here, we will complement the necessary information regarding the query and document decomposition.

1014

B.1 Subquery and Proposition Examples

1015

Here, we present examples of subqueries and propositions decomposed from the documents. The example is the decomposition of the example in Figure 1.

1016

1017

Query: Citrullinated proteins externalized in neutrophil extracellular traps act indirectly to perpetuate the inflammatory cycle via induction of autoantibodies.

- Subquery-0: Citrullinated proteins are externalized in neutrophil extracellular traps.
- Subquery-1: Citrullinated proteins act indirectly to perpetuate the inflammatory cycle.
- Subquery-2: The inflammatory cycle is perpetuated via induction of autoantibodies.

Document: RA sera and immunoglobulin fractions from RA patients with high levels of ACPA and/or rheumatoid factor significantly enhanced NETosis, and the NETs induced by these autoantibodies displayed distinct protein content. Indeed, during NETosis, neutrophils externalized the citrullinated autoantigens implicated in RA pathogenesis, and anti-citrullinated vimentin antibodies potentially induced NET formation. Moreover, the inflammatory cytokines interleukin-17A (IL-17A) and tumor necrosis factor- α (TNF- α) induced NETosis in RA neutrophils. In turn, NETs significantly augmented inflammatory responses in RA and OA synovial fibroblasts, including induction of IL-6, IL-8, chemokines, and adhesion molecules. These observations implicate accelerated NETosis in RA pathogenesis, through externalization of citrullinated autoantigens and immunostimulatory molecules that may promote aberrant adaptive and innate immune responses in the joint and in the periphery, and perpetuate pathogenic mechanisms in this disease.

- Proposition-0: RA sera and immunoglobulin fractions from RA patients with high levels of ACPA and/or rheumatoid factor significantly enhanced NETosis.
- Proposition-1: NETs induced by these autoantibodies displayed distinct protein content.
- Proposition-2: During NETosis, neutrophils externalized the citrullinated autoantigens implicated in RA pathogenesis.
- Proposition-3: Anti-citrullinated vimentin antibodies potentially induced NET formation.
- Proposition-4: Interleukin-17A (IL-17A) and tumor necrosis factor- (TNF-) induced NETosis in RA neutrophils.
- Proposition-5: NETs significantly augmented inflammatory responses in RA and OA synovial fibroblasts.
- Proposition-6: NETs inducing IL-6, IL-8, chemokines, and adhesion molecules occurred in RA and OA synovial fibroblasts.
- Proposition-7: These observations implicate accelerated NETosis in RA pathogenesis.
- Proposition-8: NETosis externalizes citrullinated autoantigens and immunostimulatory molecules.
- Proposition-9: NETosis may promote aberrant adaptive and innate immune responses in the joint and in the periphery.
- Proposition-10: NETosis may perpetuate pathogenic mechanisms in RA.

B.2 Remarks on *Propositioner*

During our manual check on the decomposition results of *propositioner* (Chen et al., 2023), we find the following potential flaws.

- (1) Wrong logic during decomposition:

Query: Identification of Design Elements for a Maturity Model for Interorganizational Integration: A Comparative Analysis

→ *Subqueries:* ['Identification of Design Elements for a Maturity Model for Interorganizational Integration.', 'A Comparative Analysis is used for identifying design elements.']

(2) Hallucination:

Query: Bigger ocean waves and waves that carry more sediment cause a greater extent of what?

→ *Subqueries:* ['Bigger ocean waves cause a greater extent of erosion.', 'Waves that carry more sediment cause a greater extent of erosion.']

(3) Information loss:

Query: The reduction was 1.6 ± 1.6 in controls. ...

→ *Subqueries:* ['The reduction in migraine headache was 1.6 1.6 in controls.', ...]

We find that the proposition will convert the questions to declarative sentences during decomposition. This may stem from the fact that its training corpus is Wikipedia, where a small portion of sentences are questions. Still, we find that *propositioner* can still decompose question-style queries, as shown in the following example:

Query: What is the purpose of bright colors on a flower's petals?

→ *Subqueries:* ["The purpose of bright colors on a flower's petals is unknown."]

B.3 Human Evaluation on Query and Document Decomposition

As mentioned in §3.1, we evaluate the decomposition outputs by *propositioner* (Chen et al., 2023), 100 samples for both query and document decomposition. Concretely, we ask three students at the post-graduate levels to evaluate the results, who are paid above the local minimum hourly wage. The instruction is shown below:

Propositions in documents (or subqueries in queries) are defined as follows:

- Each proposition conveys a distinct semantic unit, collectively expressing the complete meaning.
- Propositions should be atomic and indivisible.
- According to Choi et al. (2021), propositions should be contextualized and self-contained, including all necessary text information such as coreferences for clear interpretation.

Given the document (query) and the corresponding propositions (subqueries) generated by the model, please check whether the document (query) has been correctly decomposed. Please write *1* as correct, and *0* as incorrect.

C Retrievers Models

Table 5 presents the dense retrievers applied in the experimental section, i.e., §4.

D Offline Indexing

The *pyserini* and *faiss* libraries were employed to convert retrieval units into embeddings. We leveraged GPUs for encoding these text units in batches with a batch size of 64 and a floating precision

Model	HuggingFace Checkpoint
SimCSE (Gao et al., 2021)	princeton-nlp/unsup-simcse-bert-base-uncased
Contriever (Izacard et al., 2022)	facebook/contriever
DPR (Karpukhin et al., 2020)	facebook/dpr-ctx_encoder-multiset-base facebook/dpr-question_encoder-multiset-base
ANCE (Xiong et al., 2021)	castorini/ance-dpr-context-multi castorini/ance-dpr-question-multi
TAS-B (Hofstätter et al., 2021)	sentence-transformers/msmarco-distilbert-base-tas-b
GTR (Ni et al., 2022)	sentence-transformers/gtr-t5-base

Table 5: Model checkpoints released on HuggingFace. For DPR and ANCE, two different models encode the context and query.

f16. Following the preprocessing of these embeddings, all experiments conducted involved the utilization of an exact search method for inner products using `faiss.IndexFlatIP`,

E Downstream Tasks

The templates of LLama for downstream QA tasks, i.e., SciFact and SciQ, are listed as follows. For SciQ, we convert it from multiple choice question answering to open question answering.

Given the knowledge source: *context* \n Question: *query* \n Reply with one phrase. \n Answer:

As SciFact is a fact-checking task, we here check whether LLMs can predict the relationship between the context and the claim. The template of SciFact is shown as follows:

Context: {*context*} \n Claim: {*query*} \n For the claim, the context is supportive, contradictory, or not related? \n Options: (A) Supportive (B) Contradictory (C) Not related \n Answer:"

F Detailed Results

F.1 Ablation Study

As discussed in §6.1, we remove the component, i.e., query-doc similarity, query-prop similarity, or subquery-prop similarity, and assess the corresponding performance compared with `MixGR`. In Table 6, it is observed that `MixGR` outperforms all its components.

F.2 `MixGR` for Propositional Retrieval

Here, we evaluate `MixGR` on the retrieval units beyond documents, e.g., propositions, which Table 7 present. We observe that `MixGR` can outperform the previous document retrieval based on the similarity between query and proposition, on proposition retrieval, as discussed in §6.3.

F.3 Advantageous pattern for finer granularity measurement

In Table 8, we can notice the average number of propositions in $r_{q-d} \prec r_{s-p}$ is more than $r_{q-d} \succ r_{s-p}$. This shows that the finer granularity can better deal with the documents with more propositions than the original query-document similarity.

G `MixGR` for Other Domains

Our work provides a comprehensive analysis of the impact of `MixGR` on scientific text retrieval, considering both the variety of datasets and the use of dense retrievers. The applicability of `MixGR` to other domains remains an open question. We explore this by conducting document retrieval experiments on three distinct datasets: ConditionalQA (Sun et al., 2022), FiQA (Maia et al., 2018), and Arguana (Wachsmuth et al., 2018), which belong to the domains of law, finance, and argumentation, respectively.

Retriever	Setup	NFCorpus		SciDocs		SciFact		SciQ		Avg.	
		ND@5	ND@20	ND@5	ND@20	ND@5	ND@20	ND@5	ND@20	ND@5	ND@20
Unsupervised Dense Retrievers											
SimCSE	w/o s_{s-p}	19.6	16.0	8.7	11.5	32.3	37.0	76.1	78.0	34.2	35.6
	w/o s_{q-p}	21.4	17.4	8.5	11.6	33.1	37.4	77.9	79.6	35.2	36.5
	w/o s_{q-d}	22.8	18.6	8.5	11.9	33.9	39.0	80.7	82.2	36.5	37.9
	MixGR	22.3	18.1	9.1	12.2	34.8	39.8	84.0	85.5	37.5	38.9
Contriever	w/o s_{s-p}	43.6	36.2	14.8	20.0	65.6	69.9	78.0	80.1	50.5	51.5
	w/o s_{q-p}	43.0	36.6	14.6	20.1	66.3	70.8	81.6	83.3	51.4	52.7
	w/o s_{q-d}	43.2	36.3	14.7	20.0	65.0	69.5	83.3	84.8	51.6	52.6
	MixGR	44.0	37.1	15.5	20.7	66.4	71.0	85.2	86.7	52.8	53.9
Supervised Dense Retrievers											
DPR	w/o s_{s-p}	26.5	21.9	8.2	11.2	35.0	40.8	66.6	69.9	34.1	35.9
	w/o s_{q-p}	27.5	22.8	7.5	11.2	38.3	42.4	71.0	73.1	36.1	37.4
	w/o s_{q-d}	26.6	22.2	8.0	11.2	38.0	42.1	69.5	72.2	35.5	36.9
	MixGR	27.7	22.9	8.2	11.5	39.4	43.6	73.6	76.1	37.2	38.5
ANCE	w/o s_{s-p}	30.7	25.2	10.0	13.7	45.8	48.9	69.0	72.0	38.9	40.0
	w/o s_{q-p}	32.0	26.2	9.0	13.4	46.8	50.4	71.3	73.9	39.8	41.0
	w/o s_{q-d}	30.8	25.1	8.8	13.4	44.9	48.6	67.8	70.1	38.1	39.3
	MixGR	31.9	25.9	9.6	14.1	46.8	49.9	74.4	76.8	40.7	41.7
TAS-B	w/o s_{s-p}	42.9	34.7	13.8	19.2	61.4	66.7	86.7	87.0	51.2	51.9
	w/o s_{q-p}	42.9	34.9	13.8	19.6	63.2	67.3	88.3	88.8	52.1	52.7
	w/o s_{q-d}	42.7	34.5	13.6	18.8	62.1	65.3	85.2	85.9	50.9	51.1
	MixGR	43.6	35.2	14.0	19.6	62.7	66.9	90.5	91.0	52.7	53.2
GTR	w/o s_{s-p}	43.2	35.2	13.4	18.9	60.9	64.5	87.2	87.5	51.2	51.5
	w/o s_{q-p}	43.0	35.5	13.8	19.5	60.6	64.7	88.4	88.5	51.4	52.0
	w/o s_{q-d}	42.4	34.9	12.6	18.0	61.5	64.4	89.0	89.3	51.4	51.6
	MixGR	43.3	35.6	13.6	19.2	60.9	64.5	92.9	93.0	52.7	53.1

Table 6: Ablation study (nDCG@ $k = 5, 20$ in percentage, abbreviated as ND@ k): We evaluated four distinct scientific retrieval datasets using two unsupervised and four supervised retrievers. The retrieval results were compared using various metrics: MixGR w/o s_{s-q} , MixGR w/o s_{q-p} , MixGR w/o s_{s-p} , and MixGR, as detailed in §3.

The results are detailed in Table 9. We observe that MixGR’s benefits are considerably more limited, or even negative, outside the scientific context. This disparity may be attributed to the varying degrees of alignment between the domain-specific characteristics of each field and the training corpus of the dense retrievers. Or, *propositioner* can not perform well in these domains. Such findings further underscore the potentially distinct domain-specific nature of scientific document retrieval.

H Licences of Scientific Artifacts

Setup	SciFact		SciQ		
	EM@50	EM@200	EM@50	EM@200	
Unsupervised Dense Retrievers					
SimCSE	s_{q-d}	43.0	60.5	56.2	60.9
	Mi×GR	45.3	62.2	59.0	63.3
Contriever	s_{q-d}	49.4	67.4	56.2	62.9
	Mi×GR	47.7	71.5	57.4	62.5
Supervised Dense Retrievers					
DPR	s_{q-d}	49.4	56.4	55.5	60.2
	Mi×GR	52.3	59.9	59.0	60.9
ANCE	s_{q-d}	47.1	61.6	53.9	60.5
	Mi×GR	45.9	66.9	55.5	59.8
TAS-B	s_{q-d}	50.0	69.8	56.2	60.9
	Mi×GR	52.3	68.0	58.2	62.9
GTR	s_{q-d}	41.9	66.3	60.2	63.7
	Mi×GR	45.9	63.4	60.9	65.2

Table 7: Scientific Question Answering (Exact Match) was conducted using LLama-3 (Touvron et al., 2023) with propositions retrieved by six retrievers. Here, EM@50 and EM@200 have been reported, where the first 50 and 200 words are fed into the reader models. **Bold** indicates superior performance, and it is observed that retrieval using Mi×GR on proposition units generally outperforms the baseline.

Model	Avg. #prop in $r_{q-d} \prec r_{s-p}$	Avg. #prop in $r_{q-d} \succ r_{s-p}$
SimCSE	9.06	6.32
Contriever	8.25	7.24
ANCE	8.12	8.15
DPR	8.54	7.88
GTR	8.45	6.79
TAS-B	8.00	7.52

Table 8: Average number of propositions in two sets of document for different retrievers, i.e., $r_{q-d} \prec r_{s-p}$ and $r_{q-d} \succ r_{s-p}$. We can notice the average number of propositions in $r_{q-d} \prec r_{s-p}$ is more than $r_{q-d} \succ r_{s-p}$. This shows that the finer granularity can better deal with the documents with more propositions.

Retriever	Setup	Arguana		ConditionalQA		FiQA		Avg.	
		ND@5	ND@20	ND@5	ND@20	ND@5	ND@20	ND@5	ND@20
Unsupervised Dense Retrievers									
SimCSE	s_{q-d}	16.4	25.9	52.3	58.0	8.4	10.9	25.7	31.6
	s_{q-p}	12.5	20.9	53.7	59.5	7.6	9.7	24.6	30.0
	s_{s-p}	6.3	12.3	42.8	50.8	9.3	11.6	19.5	24.9
	MixGR	12.7	22.4	57.7	63.3	10.6	13.8	27.0	33.2
Contriever	s_{q-d}	25.9	36.0	82.5	83.9	25.0	29.9	44.5	49.9
	s_{q-p}	24.8	35.9	81.8	83.5	18.8	23.1	41.8	47.5
	s_{s-p}	24.1	34.5	63.3	67.2	18.6	22.9	35.3	41.5
	MixGR	28.7	39.2	83.5	84.5	24.7	29.8	45.6	51.2
Supervised Dense Retrievers									
DPR	s_{q-d}	9.0	16.6	58.5	63.6	12.0	14.6	26.5	31.6
	s_{q-p}	8.4	16.9	60.1	64.7	8.4	10.9	25.6	30.8
	s_{s-p}	6.1	12.2	34.8	41.8	9.2	11.8	16.7	21.9
	MixGR	8.2	16.3	59.9	65.4	11.2	14.9	26.4	32.2
ANCE	s_{q-d}	12.0	20.5	64.2	68.0	14.6	18.2	30.3	35.6
	s_{q-p}	11.7	21.3	64.0	68.2	8.5	10.9	28.1	33.5
	s_{s-p}	10.1	18.6	41.4	48.1	8.4	11.3	20.0	26.0
	MixGR	12.4	21.8	66.2	69.8	12.8	16.2	30.5	36.0
TAS-B	s_{q-d}	27.9	37.8	75.3	77.9	26.7	31.5	43.3	49.0
	s_{q-p}	18.8	30.5	76.4	78.7	15.3	19.7	36.8	43.0
	s_{s-p}	12.9	20.8	60.8	65.2	13.9	17.8	29.2	34.6
	MixGR	22.6	33.6	77.7	79.2	22.8	27.9	41.1	46.9
GTR	s_{q-d}	31.4	40.7	79.8	82.3	34.4	39.6	48.5	54.2
	s_{q-p}	25.6	36.9	80.1	82.0	22.8	27.4	42.8	48.8
	s_{s-p}	20.4	30.0	62.9	67.7	19.6	24.2	34.3	40.6
	MixGR	29.4	39.4	82.4	84.1	30.8	36.1	47.5	53.2

Table 9: Comparison between MixGR and its components on ConditionalQA, Arguana, and FiQA. We can find that the similarity based on the finer granularity s_{s-p} and MixGR won't bring as many benefits as their performance in the scientific domains, even the degradation.

Artifacts/Packages	Citation	Link	License
<i>Artifacts(datasets/benchmarks).</i>			
SciFact	(Wadden et al., 2020)	https://huggingface.co/datasets/BeIR/scifact	cc-by-sa-4.0
SciDocs	(Cohan et al., 2020)	https://huggingface.co/datasets/BeIR/scidocs	cc-by-sa-4.0
SciQ	(Welbl et al., 2017)	https://huggingface.co/datasets/bigbio/sciq	cc-by-nc-3.9
NFCorpus	(Boteva et al., 2016)	https://huggingface.co/datasets/BeIR/nfcorpus	cc-by-sa-4.0
<i>Packages</i>			
PyTorch	(Paszke et al., 2019)	https://pytorch.org/	BSD-3 License
transformers	(Wolf et al., 2019)	https://huggingface.co/transformers/v2.11.0/index.html	Apache License 2.0
numpy	(Harris et al., 2020)	https://numpy.org/	BSD License
matplotlib	(Hunter, 2007)	https://matplotlib.org/	BSD compatible License
vllm	(Kwon et al., 2023)	https://github.com/vllm-project/vllm	Apache License 2.0
<i>Models</i>			
LLaMA-3	(Touvron et al., 2023)	https://huggingface.co/meta-llama/Meta-Llama-3-8B-Instruct	LICENSE
SimCSE	(Gao et al., 2021)	https://huggingface.co/princeton-nlp/unsup-simcse-bert-base-uncased	MIT license
Contriever	(Izacard et al., 2022)	https://huggingface.co/facebook/contriever	License
DPR	(Karpukhin et al., 2020)	https://huggingface.co/facebook/dpr-ctx_encoder-multiset-base	cc-by-nc-4.0
ANCE	(Xiong et al., 2021)	https://huggingface.co/castorini/ance-dpr-context-multi	MIT license
TAS-B	(Hofstätter et al., 2021)	https://huggingface.co/sentence-transformers/msmarco-distilbert-base-tas-b	Apache License 2.0
GTR	(Ni et al., 2022)	https://huggingface.co/sentence-transformers/gtr-t5-base	Apache License 2.0

Table 10: Details of datasets, major packages, and existing models we use. The datasets we reconstructed or revised and the code/software we provide are under the MIT License.