

MAIR: A Massive Benchmark for Evaluating Instructed Retrieval

Weiwei Sun¹ Zhengliang Shi² Jiulong Wu³ Lingyong Yan⁴
Xinyu Ma⁴ Yiding Liu⁴ Min Cao³ Dawei Yin^{4*} Zhaochun Ren^{5*}

¹Carnegie Mellon University ²Shandong University

³Soochow University ⁴Baidu Inc. ⁵Leiden University

{sunweiwei, zhengliang.shi, lingyongy, xinyuma2016}@gmail.com

yindawei@acm.org, z.ren@liacs.leidenuniv.nl

Abstract

Recent information retrieval (IR) models are pre-trained and instruction-tuned on massive datasets and tasks, enabling them to perform well on a wide range of tasks and potentially generalize to unseen tasks with instructions. However, existing IR benchmarks focus on a limited scope of tasks, making them insufficient for evaluating the latest IR models. In this paper, we propose MAIR (Massive Instructed Retrieval Benchmark), a heterogeneous IR benchmark that includes 126 distinct IR tasks across 6 domains, collected from existing datasets. We benchmark state-of-the-art instruction-tuned text embedding models and re-ranking models. Our experiments reveal that instruction-tuned models generally achieve superior performance compared to non-instruction-tuned models on MAIR. Additionally, our results suggest that current instruction-tuned text embedding models and re-ranking models still lack effectiveness in specific long-tail tasks. MAIR is publicly available at <https://github.com/sunweiwei/Mair>.

1 Introduction

Large Language Models (LLMs) have demonstrated impressive capabilities in performing a wide range of natural language processing (NLP) tasks by being first pre-trained on large-scale corpora and then instruction-tuned on numerous downstream tasks (Chung et al., 2024; Wang et al., 2023b, 2022b; Wei et al., 2022). This advancement has garnered significant attention in other fields, including Information Retrieval (IR) (Gao and Callan, 2021; Neelakantan et al., 2022; Wang et al., 2022a; Izacard et al., 2021; Wang et al., 2023a; Asai et al., 2022; Su et al., 2022).

IR techniques aim to retrieve a set of relevant candidates from a large corpus based on semantic relevance or other query-specific criteria (Yates et al., 2021; Fan et al., 2022). These techniques

are critical components of many AI applications, from web search (Zhu et al., 2023) to various retrieval-augmented tasks (Gao et al., 2023; Shi et al., 2024b). However, most traditional IR models, typically trained on a single task, often exhibit poor generalization to other IR tasks or domains (Thakur et al., 2021; Zhao et al., 2024). Inspired by the success of LLMs, recent research explores training models for the general-purpose IR (Sachan et al., 2022; Sun et al., 2023; Oh et al., 2024; Weller et al., 2024). These models, instruction-tuned on multiple retrieval tasks, show significant improvements in aligning with the user intent across different IR tasks.

To evaluate the generalization capabilities of newly emerged IR models, several benchmarks such as BEIR (Thakur et al., 2021), KILT (Petroni et al., 2020), and MTEB (Muennighoff et al., 2022) have been established recently, compiling a variety of IR tasks. However, as shown in Table 1, these benchmarks either (i) contain a relatively small number of tasks or (ii) feature tasks that are too similar, thus providing limited coverage of the broad IR landscape. In contrast, instruction-tuned LLMs have been evaluated on hundreds or thousands of diverse NLP tasks (Wang et al., 2022b; Chung et al., 2024).

To fill this gap, this paper introduces MAIR (Massive Instructed Retrieval Benchmark), a large-scale IR benchmark consisting of 126 diverse retrieval tasks with 805 distinct instructions to evaluate model generalization on unseen tasks. In MAIR, we collect tasks from existing IR datasets, such as those found in (i) SIGIR resource track papers, (ii) tasks in existing benchmarks, (iii) publicly accessible shared tasks in TREC (trec.nist.gov), and (iv) recent LLM benchmarks. Tasks in MAIR are elaborately selected to cover various information retrieval requirements in practice, including diverse types of queries and document types, as well as specific relevance criteria. As a result, a

*Corresponding authors

	MAIR (this work)	BEIR Thakur et al. (2021)	KILT Petroni et al. (2020)	FollowIR Weller et al. (2024)	InstructIR Oh et al. (2024)
Number of tasks	126	18	11	3	1
Number of domains	6	4	1	1	1
Number of instructions	805	14	5	104	9,906
Number of test queries	10,038	47,233	50,736	104	9,906
Number of collections	426	18	1	3	1
Total number of docs	4,274,916	38,506,129	5,903,530	98,312	16,072

Table 1: Data statistics of MAIR and other relevant IR benchmarks.

total of 126 unique tasks are selected to constitute the benchmark MAIR.

Subsequently, we clean and merge the data of these tasks, and then sample them, resulting in a dataset that includes 10,038 queries and 4,274,916 documents. We perform the sampling mainly due to the data distribution imbalance. We validate that the sampled benchmark produces results highly correlated with full-scale testing, achieving a balance between evaluation accuracy and cost.

Finally, we manually annotate retrieval instructions based on the queries and corresponding documents for these IR tasks. Following previous work (Weller et al., 2024), these instructions specify the types of queries, documents, and relevance criteria. We ultimately annotate 805 distinct instructions, representing 805 different query-document-relevance combinations. Some tasks include query-level instructions, meaning each query has a different relevance criterion. These annotations make MAIR particularly challenging and suitable for evaluating instruction-tuned retrievers in terms of their ability to follow instructions in completing unseen tasks.

Based on MAIR, we benchmark various different types of retrieval models, including (i) sparse retriever, (ii) single-task text embedding models (Ni et al., 2021; Izacard et al., 2021), (iii) non-instruction-tuned multi-task text embedding models (Li et al., 2023b; Xiao et al., 2023; Wang et al., 2022a), (iv) instruction-tuned embedding models (Wang et al., 2023a; Lee et al., 2024), and (v) re-ranking models (Sun et al., 2023). We evaluate on both *no instruction* and *with instruction*, and found instruction-tuned embedding models show clear improvement when instructions are added. GritLM-7B achieves the best overall score, with an average nDCG@10 of 55.20. Furthermore, both e5-mistral-7b-instruct and GritLM-7B show notable improvement when instructions are added.

2 Data Construction

This section illustrates the process of data collection and benchmark construction. As aforementioned, the benchmark for instruction-tuned IR models should be capable of evaluating IR baselines on a variety of tasks across different domains and using as realistic instruction as possible. To this end, we construct our benchmark based on the following criteria: **(a) Task and domain variety:** To assess the generalization of different IR baselines on various tasks from different domains, we collect data from a large range of domains and tasks. Specifically, the data is collected from 126 distinct IR tasks across 6 domains. Each task is manually filtered to avoid duplication. **(b) Instruction diversity:** For each task, we annotate and review a large number of detailed instructions. These instructions can assist in thoroughly evaluating models’ instruction-following capabilities when searching queries.

2.1 Data Collection

To build a comprehensive IR benchmark for instruction-following evaluation, we collect data from the following four well-known sources:

- **Existing IR Resources:** The specialized resource tracks in some information retrieval conferences (e.g. SIGIR). Specifically, we collect released papers from 2021-2024 SIGIR conferences and finally collect 10 tasks across 2 domains.
- **Other IR Benchmarks:** We leverage tasks from existing IR benchmarks such as BEIR (Thakur et al., 2021) and KILT (Petroni et al., 2020), as well as domain-specific benchmarks for evaluating Voyage embedding¹. These benchmarks have gained attention in the IR community and provide diverse tasks

¹<https://blog.voyageai.com/>

tailored to specific applications. Finally, We collect 55 tasks across 6 domains from the IR benchmarks.

- **TREC Tracks:** The *Text REtrieval Conference* (TREC) is a long-standing and well-established IR conference series organized by the National Institute of Standards and Technology (NIST). TREC provides unique and realistic use cases along with rigorously annotated data. Finally, We collect 34 tasks and across 5 domains from TREC tracks.
- **LLM Evaluation Datasets:** In addition to the above datasets, we also integrate several public LLM evaluation datasets released in the LLM era. Including those datasets can help extending our benchmark to diverse potential instruction-following scenarios. Finally, We collect 27 tasks across 4 domains from the LLM Evaluation Datasets.

Finally, we collect 126 tasks. There is a simple process for us to collect data: We first review the documentation of various conference tracks, merging tasks with identical corpus and settings to avoid duplication. Then, based on task requirements, we download publicly available datasets from official channels. For tasks with incomplete corpus, we use web crawling and other techniques to supplement the data as much as possible. Since the data sources are diverse and the original formats vary substantially, we perform necessary data cleaning operations such as deduplication, keyword extraction, and text normalization on the data. Thus, we unify these task data into a consistent format.

2.2 Data Sampling

After filtering out 126 tasks, we find the data distributions are quite imbalanced among different tasks. Besides, since the scales of some datasets are extremely large with amounts of redundant content, model evaluation of the whole dataset is inefficient and unnecessary. To alleviate the influence of task data imbalance and improve evaluation efficiency, we lightweight our benchmark while preserving its evaluative capability via effective data sampling.

Specifically, for each task, we reduce its sample size following the following two steps:

- **Query Sampling:** First, for all search queries in the task, we perform the K-means algorithm over query embeddings to cluster the queries.

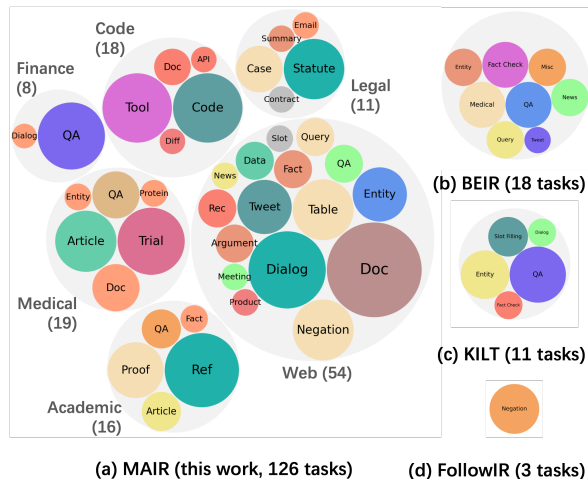


Figure 1: Compared to other datasets, MAIR covers a more diverse types of task. Bubble size represents the number of tasks of each type.

And we randomly sample one single query in each cluster to ensure query diversity while avoiding redundancy. And we set the cluster number as 100 for each task. In this way, we can finally sample 100 queries for each task.

- **Document Sampling:** Second, for each sampled query above, we use all originally annotated documents for evaluation (for example, each document is labeled related or not). Since there are some queries will few annotated documents, we randomly sample some unlabeled documents as negative documents for these queries.

It is noteworthy that some tasks may be originally evaluated on too small size of corpus making it difficult to sample enough documents in the above step 2. To address this issue, we manually combine the documents together if a small task is similar to a larger task. As a result, both tasks can share the same large-scale corpus.

The resulting benchmark comprises 10,038 queries and 4,274,916 documents (about 2 billion tokens based on OpenAI cl32k tokenizer), providing a computationally efficient yet representative testbed for IR models.

2.3 Instruction Annotation

After collecting the tasks above, we manually write the retrieval instructions for each task. All the instructions are written and reviewed by the experts in information retrieval. Following [Asai et al. \(2022\)](#), the basic format of the instruction describes

the query, the passages, and the relevance criterion. For example, instruction of ClinicalTrials is “Given a patient descriptions, retrieve clinical trials that suitable for that patient. The patient description (query) is a questionnaire that is filled by the patient or their clinician. Trials are relevant if the patient met the inclusion criteria and did not meet any exclusion criteria; Trials are partially relevant if the patient met the inclusion criteria but was excluded by one or more exclusion criteria.”.

For tasks where different search queries require unique instructions, we provide query-level annotations. The following are some examples of instructions besides the basic format: For example, task Genomics-AdHoc_2007 (Hersh et al., 2007) aims to retrieve passages that contain a specific type of biomedical entity eg. antibodies, proteins, and strains. We annotate the entity type and its definition in instructions for each query.

After the above steps, we obtain the final datasets, which consist of 10,038 queries from 126 IR tasks, with 805 instructions annotated for each task. There are in total 426 document collections and 4,274,916 documents across various domains, such as news articles, scientific papers, web pages, and code repositories. These datasets contain well-designed instructions and various document collections. They can serve as a comprehensive benchmark to evaluate the ability of retrieval systems in understanding natural language instructions and retrieving relevant information from different sources. See

3 Dataset Analysis

The Table 3 and 4 lists the full list of the tasks in MAIR and their input / output, and task type and domain. Tasks in MAIR mainly come from six domains: (i) Web refers to retrieving information from the general web. 47 IR tasks are included in the web domain. (ii) Academic domain focus on retrieving academic literature or retrieval for academic applications. (iii) The code domain consists of 18 tasks, ranging from code retrieval to tool retrieval, and to code agent. (iv) Legal domain focuses on legal-related tasks, such as case retrieval and statute retrieval. (v) Finance domain consists of 8 tasks that concern IR application in Finance task. (vi) Medical domain mainly consists of the shared tasks in TREC-CDS and GENOMICS, which wide range of retrieval and matching tasks in biomedical.

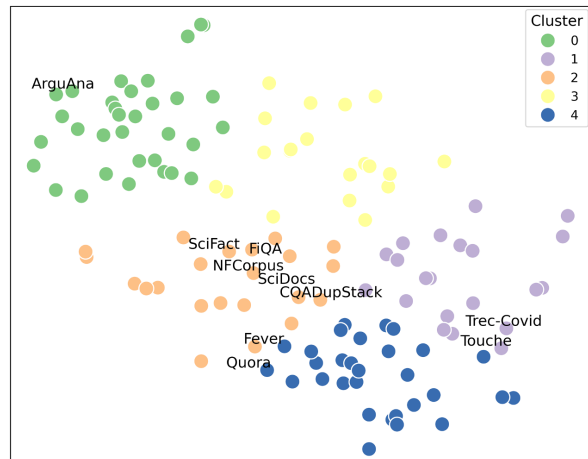


Figure 2: Visualization of the correlation among 126 tasks in MAIR, with annotations for tasks from BEIR. MAIR includes more diverse tasks. Task similarity is determined based on the performance correlation of all baseline models. We employ KMeans for clustering and t-SNE for visualization.

Figure 1 highlights the domain and task categories of MAIR, and compares them with previous datasets. MAIR covers more diverse domains and types of retrieval tasks. Next, we analyze the correlation between different tasks to further study the task diversity of MAIR, and validate the effectiveness of the data sampling approaches.

Task Correlation To measure the task diversity of MAIR, we calculate the similarity of two tasks using the Pearson correlation coefficient of the different models’ performance on two tasks. Based on the results of all baseline models on the MAIR, we calculated the correlation between all tasks, and got a correlation matrix in $\mathbf{M} = \mathbb{R}^{126 \times 126}$, $\mathbf{M}[i, j]$ denote correlation between task i and j . We plot this matrix in Figure 7. To better visualize the matrix, we use t-SNE to visualize the task correlation matrix in 2D.

Figure 2 presents the t-SNE visualization of the correlation matrix, where each task is represented as a point, and the distance between two tasks reflects their correlation. We have annotated the tasks from BEIR and observed that most of these tasks cluster together in the orange group, indicating that the performance of different models on these tasks is highly correlated. In contrast, tasks in MAIR show greater diversity, covering a larger area and thus providing a more comprehensive evaluation of the models. The specific pairwise correlations are displayed in the heatmap in Figure 7, revealing that many tasks exhibit negative correlations. This

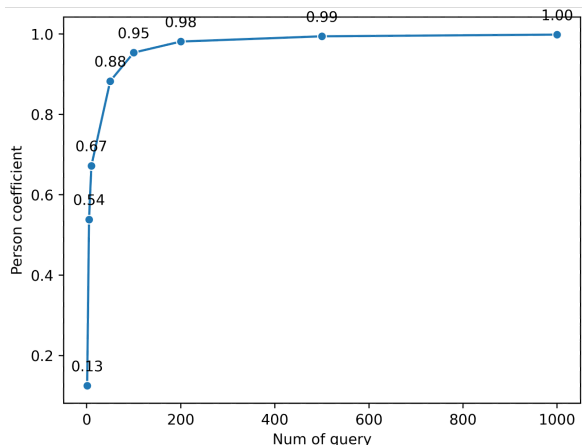


Figure 3: The performance correlation of baseline models with different sampled numbers of queries. Sampling 100 queries achieves a good trade-off between correlation and cost.

suggests that these tasks have significantly different definitions, leading models that perform well on one task to perform poorly on another.

Sampling Effectiveness In our data construction process, we sample the data to obtain a lightweight test set. However, the sampled data may lead to bias in the evaluation results. To measure the effectiveness of our data sampling process, we build a test set with different maximum test sizes cut off, ranging from [1, 5, 10, 50, 100, 200, 500, 10000]. Then, we run the baseline models on these test sets, can compute the task correlation between the full test set (i.e., unsampled test set), and the sampled test set. A high correlation means that evaluation results on the sampled test set are very similar to the evaluation on the full set. Figure 3 illustrates the correlation of different test sets cut off. We can see that retaining more instances leads to a higher correlation, i.e., the evaluation is more robust. Notably, a cut-off of 100 instances achieves over 0.95 Pearson correlation coefficient between the evaluation results, indicating that our sampling method could achieve good reliability while using minimal cost.

4 Experimental Setup

4.1 Models

Sparse Retrieval These measure relevance by computing term overlap. We benchmark BM25, implemented in BM25S (Lù, 2024).

Single-Task Text Embedding These models are trained on a single IR dataset. We benchmark gtr-t5-base, gtr-t5-large, and

contriever-msmarco, all of which are trained on MS MARCO (Ni et al., 2021; Izcard et al., 2021).

Non-Instruction-Tuned Multi-Task Text Embedding These models are typically trained on various annotated IR training datasets combined with massive weekly supervised data. We benchmark (i) the GTE series (Li et al., 2023b): gte-base-en-v1.5, gte-large-en-v1.5; (ii) the BGE series (Xiao et al., 2023): bge-base-en-v1.5, bge-large-en-v1.5; (iii) the E5 series (Wang et al., 2022a): e5-small-v2, e5-base-v2, e5-large-v2; and (iv) all-MiniLM-L6-v2 from Sentence Transformer. We also evaluated text-embedding-3-small by the OpenAI API².

Instruction-Tuned Text Embedding These models are fine-tuned on instruction datasets, where the model input is a query paired with an instruction describing the retrieval task. We benchmark (i) e5-mistral-instruct, which optimizes instruction-following ability using LLM-generated data (Wang et al., 2023a); (ii) NV-Embed-v1, which utilizes bidirectional attention with an additional latent attention layer to enhance text embedding model (Lee et al., 2024); (iii) GritLM-7B (Muennighoff et al., 2024), a unified text embedding and generation model; and (iv) gte-Qwen2-1.5B-instruct (Li et al., 2023b), general text embedding based on Qwen2-1.5B.

Cross-Encoder Re-Rankers These models measure the relevance of paired queries and passages using bidirectional or unidirectional Transformers. We benchmark (i) monoT5-Base, a T5 encoder model trained on MS MARCO, (ii) mxbai-rerank-large-v1 and jina-reranker-v2-base, multi-task reranker developed by Mxbai and Jina.ai, respectively, and (iii) bge-reranker-v2-m3 and bge-reranker-v2-gemma, trained on massive ranking data with XLM-R and Gemma-2B as their respective backbones.

LLM-based Re-Rankers These models prompt general-purpose LLMs to perform re-ranking in a zero-shot setting. We benchmark RankGPT with the gpt-3.5-turbo (Sun et al., 2023).

²<https://platform.openai.com/docs/guides/embeddings>

Model	Avg		Web		Academic		Legal		Medical		Finance		Code	
	✗	✓	✗	✓	✗	✓	✗	✓	✗	✓	✗	✓	✗	✓
Instruction?														
BM25	40.75	-	40.21	-	38.50	-	47.60	-	46.01	-	51.22	-	32.72	-
contriever-msmarco	39.88	-	45.19	-	31.82	-	31.74	-	36.60	-	50.59	-	35.06	-
gtr-t5-base	37.79	-	40.55	-	31.36	-	33.77	-	32.95	-	47.78	-	37.18	-
gtr-t5-large	41.43	-	43.97	-	35.62	-	37.35	-	33.93	-	51.83	-	42.29	-
all-MiniLM-L6-v2	38.66	26.90	40.05	25.56	42.90	35.85	29.20	18.46	35.24	22.32	44.54	28.15	35.90	29.14
e5-small-v2	41.36	34.47	43.24	34.12	38.83	32.83	32.20	28.55	37.89	30.82	55.32	50.49	39.70	35.26
e5-base-v2	43.45	38.80	43.78	38.95	42.00	38.37	35.76	33.06	40.22	30.68	58.33	54.18	43.48	40.13
e5-large-v2	44.75	39.45	44.60	38.93	45.44	40.01	36.96	32.86	41.90	35.13	58.92	54.27	44.28	40.00
gte-base-en-v1.5	44.06	34.25	44.72	30.97	40.15	35.34	37.89	34.73	45.63	37.18	61.64	49.23	40.50	33.11
gte-large-en-v1.5	46.58	33.87	46.38	30.97	47.00	40.23	41.04	34.25	47.75	32.91	63.47	41.57	41.58	32.66
bge-base-en-v1.5	42.58	30.63	43.33	27.70	38.85	31.66	33.56	26.89	43.78	31.58	55.54	38.95	42.46	35.05
bge-base-en-v1.5	44.34	28.01	44.91	24.78	41.20	32.05	37.35	24.09	45.31	25.71	58.69	31.38	42.58	34.77
text-embedding-3-small	47.89	45.28	48.10	46.65	44.43	41.11	45.71	41.06	42.74	37.81	61.36	57.99	48.87	46.74
gte-Qwen2-1.5B-instruct [♠]	49.14	51.81	46.81	51.30	48.87	49.98	50.91	51.34	45.97	45.20	62.85	64.61	50.40	53.48
NV-Embed-v1 [♠]	50.29	51.19	49.76	51.07	50.70	51.14	44.07	44.57	40.71	38.96	68.27	69.59	52.62	54.52
e5-mistral-7b-instruct [♠]	50.85	54.43	48.52	54.57	48.78	50.03	52.94	52.50	48.35	51.10	66.52	68.76	52.29	54.83
GritLM-7B [♠]	48.46	55.26	43.01	54.50	51.46	53.45	50.70	53.14	45.38	48.22	64.28	70.83	53.62	57.52
monot5-base-msmarco	43.51	38.27	47.03	42.24	35.44	27.00	41.45	38.41	36.19	31.31	59.34	56.72	40.27	34.25
mxbai-rerank-large-v1	42.84	22.91	44.33	24.94	32.98	15.53	42.50	19.14	43.25	24.64	56.54	26.23	41.78	23.78
jina-reranker-v2-base	50.88	48.59	52.31	50.74	45.22	42.96	49.32	46.10	46.78	42.87	61.35	59.28	51.03	48.24
bge-reranker-v2-m3	46.59	40.45	48.63	43.64	40.01	33.38	49.18	42.38	41.12	30.85	64.17	60.36	41.47	34.67
bge-reranker-v2-gemma	52.20	38.16	51.26	42.40	52.09	30.89	47.74	30.47	45.69	28.96	68.34	46.41	54.00	39.86
gpt-3.5-turbo [♠]	47.84	49.02	47.80	49.70	41.59	41.02	47.47	48.38	42.21	41.97	64.45	64.52	49.93	52.27

Table 2: **Main Results (nDCG@10)** Rows marked with ✓ indicate results with instruction input, while rows marked with ✗ indicate results without instruction input. Models labeled with ♠ are instruction fine-tuned. The last group represents the re-ranking models, all of which re-rank the top 100 results from text-embedding-3-small.

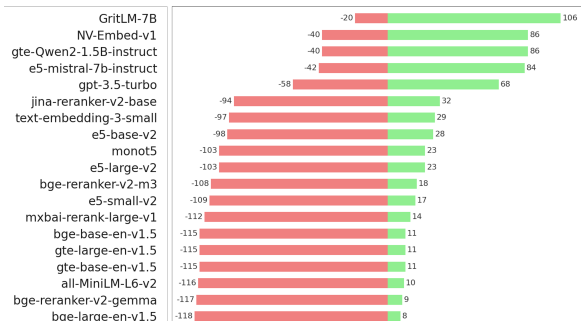


Figure 4: With the addition of instruction, the number of tasks that obtain performance improvement (green part) and reduction (red part). We can see that instruction-tuned models show more improvements while non-instruction-tuned models reduce on most tasks.

4.2 Evaluation

Following previous work, we use **nDCG@10** as the evaluation metric. The overall score is defined as the average score across all queries. We also report the average nDCG@10 for each of the following domains: Web, Academic, Code, Medical, Legal, and Finance. The specific tasks included in each field can be found in Table 5 - 9.

For all models, we consider two settings: (i) *no instruction*, which retrieves passages without task

instruction or with a simple web search instruction when required, and (ii) *+ instruction*, which retrieves passages with instruction input paired with the query. The performance changes after using instructions indicate the model’s ability to understand them. Note that for non-instruction-tuned models, we also test their performance under the instruction setting for reference, even though they may not be optimized to understand instructions.

All re-ranking models use text-embedding-3-small as the first-stage retriever and re-rank the top-100 passages. Passages are truncated to the maximum input length of each model.

5 Evaluation Results

5.1 Main Results

Table 2 reports the evaluation results of the tested models. Table 5-9 reports the detailed results of each individual task. We observe that instruction-tuned embedding models with instruction input achieve the best overall performance. GritLM-7B achieves an average score of 48.40, with a 6.80 nDCG improvement when instructions are added.

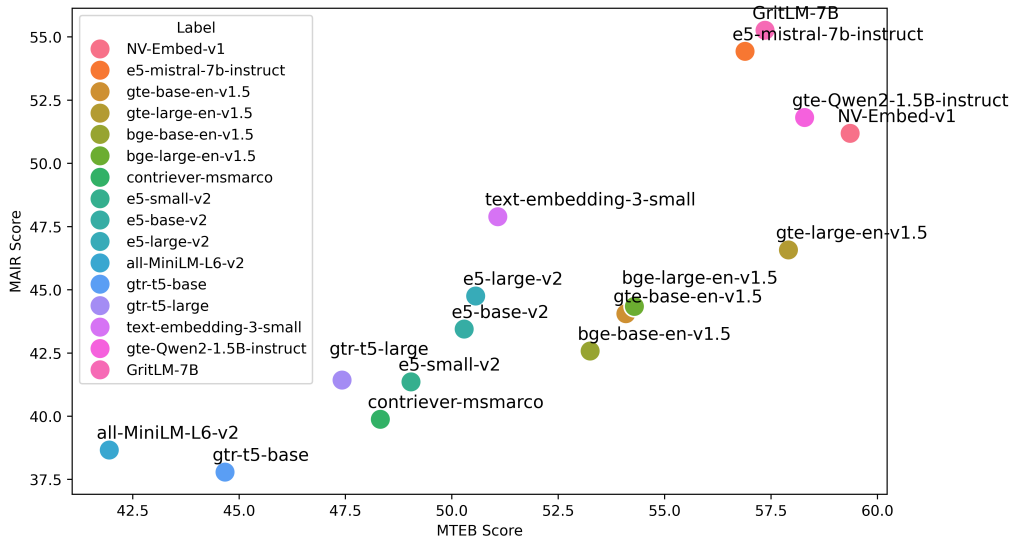


Figure 5: Score between MTEB (Retrieval) and MAIR.

NV-Embed also shows a clear improvement with the addition of instructions, though not as significant as GritLM-7B. This difference is likely due to the LLM-generated instruction-tuning data enhancing the models’ ability to understand instructions (Wang et al., 2023a).

Non-instruction-tuned models experience a performance drop when instructions are added. Some, like text-embedding-3-small, show a slight decrease, while others, such as bge-base-en-v1.5, exhibit a more significant decline.

For the re-ranker, bge-reranker-v2-gemma achieves the best results, outperforming all embedding models when no instructions are provided, but it shows a notable decline in performance when instructions are added. RankGPT based on gpt-3.5-turbo achieves results close to bge-reranker-v2-m3 without instruction input but demonstrates a 1.21 nDCG improvement when instructions are included. This suggests that prompted LLMs can intuitively transfer general instruction-following capabilities to ranking tasks.

5.2 Gain of Instruction

Figure 4 shows that with the addition of instructions, the model’s performance improves on a number of the 126 tasks while decreasing on others. We observe that the instruction-tuned models show improvement on more than half of the tasks, with GritLM-7B performing the best, achieving improvements on 106 out of 126 tasks. NV-Embed-v1, gte-Qwen2-1.5B-instruct, and e5-mistral-7b-instruct achieve similar results

and outperform gpt-3.5-turbo. In contrast, non-instruction-tuned models show a performance decrease on more tasks when instructions are added.

5.3 Compare with MTEB

Figure 5 shows the relationship between models’ performance on MTEB (Retrieval) and MAIR. We can see that the two benchmarks share a similar trend. Among these models, text-embedding-3-small achieves better results on MAIR than on MTEB. For example, on MTEB, gte-large-en-v1.5 outperforms text-embedding-3-small by about 7 points and outperforms e5-mistral-7b-instruct by 1 point. However, it performs worse than them on MAIR. This is probably because text-embedding-3-small has better generalization, as it performs better on massive unseen tasks in MAIR, while gte-large-en-v1.5 is more optimized towards tasks in MTEB. We also observe that single-task models such as contriever-msmarco perform poorly on MAIR, which indicates that MAIR requires more generalization ability.

5.4 Analysis on IFEval

To evaluate the model on challenging instruction-following tasks, we designed the IFEval task (Zhou et al., 2023) within MAIR. IFEval consists of 8 different instruction-following subtasks, such as format (selecting responses in a specific format), keywords (including specific words), and length (adhering to length restrictions). The retrieval task



Figure 6: Results (nDCG@10) on IFEval sub-tasks.

in IFEval is to select the answer that correctly follows the instructions from among the 100 candidates. Specifically, building on the original IFEval data (Zhou et al., 2023), we used gpt-3.5-turbo to generate 100 candidate answers for each question, ensuring that only 10 fully follow the given instructions. IFEval is challenging for retrieval models because (i) the instructions are out-of-training-distribution; (ii) all candidate answers are semantically relevant to the question, thus, the model must focus on the instructions to identify the correct answer.

Figure 6 demonstrates the performance of GritLM-7B, bge-reranker-v2-gemma, and RankGPT with gpt-3.5-turbo, gpt-4o-mini, and gpt-4o on the 8 subtasks of IFEval. The results show that existing instruction-tuned retrieval models still perform poorly on challenge instruction-followed ranking task. For example, the GritLM-7B achieves an nDCG@10 score of less than 60 on 7 out of the 8 tasks. In contrast, advanced LLM gpt-4o achieves an nDCG@10 score of over 80 on 6 out of the 8 tasks. It indicates again the shortcomings of instruction-tuned retrieval models in handling complex information requirements, and utilizing advanced language models as supervisors might be an effective strategy. For full results on IFEval, refer to Tables 13 and 10.

6 Related Work and Background

6.1 Instruction tuning for Retrieval

Instruction tuning has been an effective technique in the development of LLMs (Chung et al., 2024; Wang et al., 2023b, 2022b). In this process, models are trained on diverse instruction-response pairs, empowering them with the ability to adaptively perform unseen tasks based on instructions (Chen et al., 2024; Wei et al., 2022). This has inspired emergent research on training instruction-tuned retrieval models. Asai et al. (2022); Su et al. (2022) are the earliest works in this direction, where they propose training text embedding models with instructions alongside the query, enabling the models to perform well on unseen tasks using instructions. Recent research has started to finetune larger text embedding models with instructions, such as Mistral-7B and Mixtral-7x8B (Muennighoff et al., 2024), while Wang et al. (2023a) further proposes using LLMs to generate instruction-tuning data for training retrievers. These instruction-tuned embeddings have climbed to the top of existing IR leaderboards like MTEB (Muennighoff et al., 2022). Meanwhile, some paper explore prompt-based method to instruct general-purpose LLMs for re-ranking tasks (Sun et al., 2023).

6.2 IR Benchmark

Benchmarking has been a crucial part of IR system development. Traditional IR benchmarks typically evaluate models on narrow tasks, like question-answering, or on a certain domain (Campos et al., 2016; Karpukhin et al., 2020). Recently, with the advance of pre-trained language models and language model-based retrievers, such as monoBERT (Nogueira and Cho, 2019) and GTR (Ni et al., 2021), increasing attention has been paid to constructing multi-task IR benchmarks and evaluating retrievers on out-of-domain tasks. A typical effort is BEIR (Thakur et al., 2021), which consists of 18 tasks across 4 domains. MTEB (Muennighoff et al., 2022) further extends BEIR by adding other embedding-related non-retrieval tasks like clustering and classification. KILT (Petroni et al., 2020) collects 11 knowledge-intensive language tasks. However, given the recent progress of instruction-tuned retrievers, existing benchmarks cannot comprehensively evaluate models' abilities since the number of tasks and the scope of tasks are limited. Some recent papers, such as FollowIR (Weller et al., 2024)

and InstructIR (Oh et al., 2024), propose benchmarks that specifically focus on evaluating models’ instruction-following abilities. FollowIR rewrites the original narratives in three TREC shared tasks to include some exclusionary instructions, making some relevant passages become irrelevant. InstructIR utilizes LLM to generate synthetic user backgrounds and further rewrites the relevant passages using LLM to fit the background. However, these tasks are synthetic and only include limited tasks, making it difficult to robustly evaluate the performance of models on real-world unseen instructions. In comparison, MAIR is a large-scale multi-task IR benchmark consisting of 126 distinct tasks. Most tasks are long-tail and without training data, and detailed instructions for each task are manually annotated.

7 Conclusion

In this paper, we introduce a novel massive instructed IR benchmark called MAIR for evaluating instruction-tuned retrieval models. Compared with existing related benchmarks, MAIR has a more comprehensive coverage of various IR tasks, with different types of queries, documents and relevance criteria. Specifically, MAIR comprises 126 distinct IR tasks, with 805 newly annotated instructions. All tasks and instructions are manually selected and annotated to ensure the benchmark can assess the instruction-following capabilities of different IR models in practice.

Based on MAIR, we benchmark various retrieval models, including both instruction-tuned and non-instruction-tuned models. Our results demonstrate that MAIR poses greater challenges than existing benchmarks, and provides a more comprehensive evaluation.

Limitation

The limitations of this work include the lack of study of retrieval in a multilingual setting. Our benchmark focuses only on the English language and considers only text retrieval. Therefore, we plan to explore multilingual IR settings in the future. Another limitation is the lack of study on prompt sensitivity. It is well known that LLMs are sensitive to prompt words. We plan to annotate more instructions in the future to study how LLM performance is impacted by prompt words.

Ethics Statement

We acknowledge the importance of the ACM Code of Ethics and fully agree with it. We ensure that this work is compatible with the provided code in terms of publicly accessible datasets and models. Risks and harms associated with large language models include the generation of harmful, offensive, or biased content. The new benchmark is composed of various previous datasets and is therefore licensed under their respective data licenses.

References

- Vaibhav Adlakha, Shehzaad Dhuliawala, Kaheer Suleman, Harm de Vries, and Siva Reddy. 2021. Topi-ocqa: Open-domain conversational question answering with topic switching. *Transactions of the Association for Computational Linguistics*, 10:468–483.
- Jafar Afzali, Aleksander Mark Drzewiecki, and Krisztian Balog. 2021. Pointrec: A test collection for narrative-driven point of interest recommendation. *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*.
- Anirudh Ajith, Mengzhou Xia, Alexis Chevalier, Tanya Goyal, Danqi Chen, and Tianyu Gao. 2024. Lit-search: A retrieval benchmark for scientific literature search. *ArXiv*, abs/2407.18940.
- Akari Asai, Timo Schick, Patrick Lewis, Xilun Chen, Gautier Izacard, Sebastian Riedel, Hannaneh Hajishirzi, and Wen tau Yih. 2022. Task-aware retrieval with instructions. In *Annual Meeting of the Association for Computational Linguistics*.
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie J. Cai, Michael Terry, Quoc V. Le, and Charles Sutton. 2021. Program synthesis with large language models. *ArXiv*, abs/2108.07732.
- Paheli Bhattacharya, Kripabandhu Ghosh, Saptarshi Ghosh, Arindam Pal, Parth Mehta, Arnab Bhattacharya, and Prasenjit Majumder. 2019. Overview of the fire 2019 aila track: Artificial intelligence for legal assistance. In *Fire*.
- Alexander Bondarenko, Maik Fröbe, Meriem Beloucif, Lukas Gienapp, Yamen Ajjour, Alexander Panchenko, Christian Biemann, Benno Stein, Henning Wachsmuth, Martin Potthast, and Matthias Hagen. 2020. Overview of touché 2020: Argument retrieval. In *Conference and Labs of the Evaluation Forum*.
- Vera Boteva, Demian Gholipour Ghalandari, Artem Sokolov, and Stefan Riezler. 2016. A full-text learning to rank dataset for medical information retrieval. In *European Conference on Information Retrieval*.
- Daniel Fernando Campos, Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, Li Deng, and Bhaskar Mitra. 2016. Ms marco: A human generated machine reading comprehension dataset. *ArXiv*, abs/1611.09268.
- Arun Tejasvi Chaganty, Megan Leszczynski, Shu Zhang, Ravi Ganti, Krisztian Balog, and Filip Radlinski. 2023. Beyond single items: Exploring user preferences in item sets with the conversational playlist curation dataset. *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*.
- Ilias Chalkidis, Manos Fergadiotis, Nikolaos Manginas, Eva Katakalous, and Prodromos Malakasiotis. 2021. Regulatory compliance through doc2doc information retrieval: A case study in eu/uk legislation where text similarity has limitations. In *Conference of the European Chapter of the Association for Computational Linguistics*.
- Lichang Chen, Shiyang Li, Jun Yan, Hai Wang, Kalpa Gunaratna, Vikas Yadav, Zheng Tang, Vijay Srivasan, Tianyi Zhou, Heng Huang, and Hongxia Jin. 2024. Alpapasus: Training a better alpaca with fewer data. In *ICLR*.
- Xiuying Chen, Hind Alamro, Li Mingzhe, Shen Gao, Rui Yan, Xin Gao, and Xiangliang Zhang. 2022a. Target-aware abstractive related work generation with contrastive learning. *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*.
- Zhiyu Chen, Wenhui Chen, Charese Smiley, Sameena Shah, Iana Borova, Dylan Langdon, Reema N Moussa, Matthew I. Beane, Ting-Hao 'Kenneth' Huang, Bryan R. Routledge, and William Yang Wang. 2021a. Finqa: A dataset of numerical reasoning over financial data. *ArXiv*, abs/2109.00122.
- Zhiyu Chen, SHIYANG LI, Charese Smiley, Zhiqiang Ma, Sameena Shah, and William Yang Wang. 2022b. Convfinqa: Exploring the chain of numerical reasoning in conversational finance question answering. In *Conference on Empirical Methods in Natural Language Processing*.
- Zhiyu Chen, Shuo Zhang, and Brian D. Davison. 2021b. Wtr: A test collection for web table retrieval. *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2024. Scaling instruction-finetuned language models. *JMLR*, pages 1–53.
- Arman Cohan, Sergey Feldman, Iz Beltagy, Doug Downey, and Daniel S. Weld. 2020. Specter: Document-level representation learning using citation-informed transformers. *ArXiv*, abs/2004.07180.
- Emily Dinan, Stephen Roller, Kurt Shuster, Angela Fan, Michael Auli, and Jason Weston. 2018. Wizard of wikipedia: Knowledge-powered conversational agents. *ArXiv*, abs/1811.01241.

- Hady ElSahar, Pavlos Vougiouklis, Arslan Remaci, Christophe Gravier, Jonathon S. Hare, Frédérique Laforest, and Elena Paslaru Bontas Simperl. 2018. T-rex: A large scale alignment of natural language with knowledge base triples. In *International Conference on Language Resources and Evaluation*.
- Angela Fan, Yacine Jernite, Ethan Perez, David Grangier, Jason Weston, and Michael Auli. 2019. Eli5: Long form question answering. *ArXiv*, abs/1907.09190.
- Yixing Fan, Xiaohui Xie, Yinqiong Cai, Jia Chen, Xinyu Ma, Xiangsheng Li, Ruqing Zhang, Jiafeng Guo, et al. 2022. Pre-training methods in information retrieval. *Foundations and Trends® in Information Retrieval*, 16(3):178–317.
- Luyu Gao and Jamie Callan. 2021. Unsupervised corpus aware language model pre-training for dense passage retrieval. *ArXiv*, abs/2108.05540.
- Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. 2023. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997*.
- Lotem Golany, Filippo Galgani, Maya Mamo, Nimrod Parasol, Omer Vandsburger, Nadav Bar, and Ido Dagan. 2024. Efficient data generation for source-grounded information-seeking dialogs: A use case for meeting transcripts. *ArXiv*, abs/2405.01121.
- Zhaochen Guo and Denilson Barbosa. 2018. Robust named entity disambiguation with random walks. *Semantic Web*, 9:459–479.
- Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin Burns, Samir Puranik, Horace He, Dawn Xiaodong Song, and Jacob Steinhardt. 2021a. Measuring coding challenge competence with apps. *ArXiv*, abs/2105.09938.
- Dan Hendrycks, Collin Burns, Anya Chen, and Spencer Ball. 2021b. Cuad: An expert-annotated nlp dataset for legal contract review. *ArXiv*, abs/2103.06268.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Xiaodong Song, and Jacob Steinhardt. 2021c. Measuring mathematical problem solving with the math dataset. *ArXiv*, abs/2103.03874.
- William R. Hersh, Aaron M. Cohen, Jianji Yang, Ravi Teja Bhupatiraju, Phoebe M. Roberts, and Marti A. Hearst. 2007. Trec 2007 genomics track overview. In *Text Retrieval Conference*.
- Johannes Hoffart, Mohamed Amir Yosef, Ilaria Bordino, Hagen Fürstenau, Manfred Pinkal, Marc Spaniol, Bilyana Taneva, Stefan Thater, and Gerhard Weikum. 2011. Robust disambiguation of named entities in text. In *Conference on Empirical Methods in Natural Language Processing*.
- Doris Hoogeveen, Karin M. Verspoor, and Timothy Baldwin. 2015. Cqadupstack: A benchmark data set for community question-answering research. *Proceedings of the 20th Australasian Document Computing Symposium*.
- Christoph Hoppe, David Pelkmann, Nico Migenda, Daniel Hötte, and Wolfram Schenck. 2021. Towards intelligent legal advisors for document retrieval and question-answering in german legal documents. *2021 IEEE Fourth International Conference on Artificial Intelligence and Knowledge Engineering (AIKE)*, pages 29–32.
- Hamel Husain, Hongqiu Wu, Tiferet Gazit, Miltiadis Allamanis, and Marc Brockschmidt. 2019. Code-searchnet challenge: Evaluating the state of semantic code search. *ArXiv*, abs/1909.09436.
- Pranab Islam, Anand Kannappan, Douwe Kiela, Rebecca Qian, Nino Scherrer, and Bertie Vidgen. 2023. Financebench: A new benchmark for financial question answering. *ArXiv*, abs/2311.11944.
- Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2021. Unsupervised dense information retrieval with contrastive learning. *Trans. Mach. Learn. Res.*, 2022.
- Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. 2023. Swe-bench: Can language models resolve real-world github issues? *ArXiv*, abs/2310.06770.
- Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Yu Wu, Sergey Edunov, Danqi Chen, and Wen tau Yih. 2020. Dense passage retrieval for open-domain question answering. *ArXiv*, abs/2004.04906.
- Makoto P. Kato, Hiroaki Ohshima, Ying-Hsang Liu, and Hsin-Liang Chen. 2021. A test collection for ad-hoc dataset retrieval. *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*.
- Anastassia Kornilova and Vladimir Eidelman. 2019. Billsum: A corpus for automatic summarization of us legislation. *ArXiv*, abs/1910.00523.
- Chankyu Lee, Rajarshi Roy, Mengyao Xu, Jonathan Raiman, Mohammad Shoeybi, Bryan Catanzaro, and Wei Ping. 2024. Nv-embed: Improved techniques for training llms as generalist embedding models. *arXiv preprint arXiv:2405.17428*.
- Omer Levy, Minjoon Seo, Eunsol Choi, and Luke Zettlemoyer. 2017. Zero-shot relation extraction via reading comprehension. *ArXiv*, abs/1706.04115.
- Haitao Li, Yunqiu Shao, Yueyue Wu, Qingyao Ai, Yixiao Ma, and Yiqun Liu. 2023a. Lecardv2: A large-scale chinese legal case retrieval dataset. *ArXiv*, abs/2310.17609.

- Zehan Li, Xin Zhang, Yanzhao Zhang, Dingkun Long, Pengjun Xie, and Meishan Zhang. 2023b. Towards general text embeddings with multi-stage contrastive learning. *ArXiv*, abs/2308.03281.
- Tengteng Lin, Qiaosheng Chen, Gong Cheng, Ahmet Soyly, Basil Ell, Ruoqi Zhao, Qing Shi, Xiaxia Wang, Yu Gu, and Evgeny Kharlamov. 2022. Acordar: A test collection for ad hoc content-based (rdf) dataset retrieval. *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*.
- Tianyang Liu, Canwen Xu, and Julian McAuley. 2023. Repobench: Benchmarking repository-level code auto-completion systems. *ArXiv*, abs/2306.03091.
- Antoine Louis, Gerasimos Spanakis, and G. van Dijk. 2021. A statutory article retrieval dataset in french. In *Annual Meeting of the Association for Computational Linguistics*.
- Xing Han Lù. 2024. Bm25s: Orders of magnitude faster lexical search via eager sparse scoring. *ArXiv*, abs/2407.03618.
- 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. *Companion Proceedings of the The Web Conference 2018*.
- Niklas Muennighoff, Qian Liu, Qi Liu, Armel Zebaze, Qinkai Zheng, Binyuan Hui, Terry Yue Zhuo, Swayam Singh, Xiangru Tang, Leandro von Werra, and S. Longpre. 2023. Octopack: Instruction tuning code large language models. *ArXiv*, abs/2308.07124.
- Niklas Muennighoff, Hongjin Su, Liang Wang, Nan Yang, Furu Wei, Tao Yu, Amanpreet Singh, and Douwe Kiela. 2024. Generative representational instruction tuning. *ArXiv*, abs/2402.09906.
- Niklas Muennighoff, Nouamane Tazi, Loic Magne, and Nils Reimers. 2022. Mteb: Massive text embedding benchmark. In *Conference of the European Chapter of the Association for Computational Linguistics*.
- Arvind Neelakantan, Tao Xu, Raul Puri, Alec Radford, Jesse Michael Han, Jerry Tworek, Qiming Yuan, Nikolas A. Tezak, Jong Wook Kim, Chris Hallacy, Johannes Heidecke, Pranav Shyam, Boris Power, Tyna Eloundou Nekoul, Girish Sastry, Gretchen Krueger, David P. Schnurr, Felipe Petroski Such, Kenny Sai-Kin Hsu, Madeleine Thompson, Tabarak Khan, Toki Sherbakov, Joanne Jang, Peter Welinder, and Lilian Weng. 2022. Text and code embeddings by contrastive pre-training. *ArXiv*, abs/2201.10005.
- Jianmo Ni, Chen Qu, Jing Lu, Zhuyun Dai, Gustavo Hernández Abrego, Ji Ma, Vincent Zhao, Yi Luan, Keith B. Hall, Ming-Wei Chang, and Yinfei Yang. 2021. Large dual encoders are generalizable retrievers. *ArXiv*, abs/2112.07899.
- Rodrigo Nogueira and Kyunghyun Cho. 2019. Passage re-ranking with bert. *ArXiv*, abs/1901.04085.
- Hanseok Oh, Hyunji Lee, Seonghyeon Ye, Haebin Shin, Hansol Jang, Changwook Jun, and Minjoon Seo. 2024. Instructir: A benchmark for instruction following of information retrieval models. *arXiv preprint arXiv:2402.14334*.
- Shishir G Patil, Tianjun Zhang, Xin Wang, and Joseph E Gonzalez. 2023. Gorilla: Large language model connected with massive apis. *arXiv preprint arXiv:2305.15334*.
- Fabio Petroni, Aleksandra Piktus, Angela Fan, Patrick Lewis, Majid Yazdani, Nicola De Cao, James Thorne, Yacine Jernite, Vassilis Plachouras, Tim Rocktaschel, and Sebastian Riedel. 2020. Kilt: a benchmark for knowledge intensive language tasks. In *North American Chapter of the Association for Computational Linguistics*.
- Bhawna Piryani, Jamshid Mozafari, and Adam Jatowt. 2024. Chroniclingamericaqa: A large-scale question answering dataset based on historical american newspaper pages. *ArXiv*, abs/2403.17859.
- Yujia Qin, Shi Liang, Yining Ye, Kunlun Zhu, Lan Yan, Ya-Ting Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Runchu Tian, Ruobing Xie, Jie Zhou, Marc H. Gerstein, Dahai Li, Zhiyuan Liu, and Maosong Sun. 2023. Toollm: Facilitating large language models to master 16000+ real-world apis. *ArXiv*, abs/2307.16789.
- Kirk Roberts, Tasmee Alam, Steven Bedrick, Dina Demner-Fushman, Kyle Lo, Ian Soboroff, Ellen M. Voorhees, Lucy Lu Wang, and William R. Hersh. 2020. Trec-covid: rationale and structure of an information retrieval shared task for covid-19. *Journal of the American Medical Informatics Association : JAMIA*, 27:1431 – 1436.
- Devendra Singh Sachan, Mike Lewis, Mandar Joshi, Armen Aghajanyan, Wen tau Yih, Joëlle Pineau, and Luke Zettlemoyer. 2022. Improving passage retrieval with zero-shot question generation. In *Conference on Empirical Methods in Natural Language Processing*.
- Chris Samarinas and Hamed Zamani. 2024. Procis: A benchmark for proactive retrieval in conversations. *ArXiv*, abs/2405.06460.
- Zhengliang Shi, Shen Gao, Xiuyi Chen, Yue Feng, Lingyong Yan, Haibo Shi, Dawei Yin, Zhumin Chen, Suzan Verberne, and Zhaochun Ren. 2024a. Chain of tools: Large language model is an automatic multi-tool learner. *arXiv preprint arXiv:2405.16533*.
- Zhengliang Shi, Shuo Zhang, Weiwei Sun, Shen Gao, Pengjie Ren, Zhumin Chen, and Zhaochun Ren. 2024b. Generate-then-ground in retrieval-augmented generation for multi-hop question answering. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*.

- Ivan Srba, Branislav Pecher, Matús Tomlein, Róbert Móra, Elena Stefancova, Jakub Simko, and Mária Bieliková. 2022. Monant medical misinformation dataset: Mapping articles to fact-checked claims. *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*.
- Hongjin Su, Weijia Shi, Jungo Kasai, Yizhong Wang, Yushi Hu, Mari Ostendorf, Wen tau Yih, Noah A. Smith, Luke Zettlemoyer, and Tao Yu. 2022. One embedder, any task: Instruction-finetuned text embeddings. *ArXiv*, abs/2212.09741.
- Weiwei Sun, Lingyong Yan, Xinyu Ma, Shuaiqiang Wang, Pengjie Ren, Zhumin Chen, Dawei Yin, and Zhaochun Ren. 2023. Is ChatGPT good at search? investigating large language models as re-ranking agents. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*.
- Nandan Thakur, Nils Reimers, Andreas Ruckl'e, Abhishek Srivastava, and Iryna Gurevych. 2021. Beir: A heterogenous benchmark for zero-shot evaluation of information retrieval models. *ArXiv*, abs/2104.08663.
- James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. Fever: a large-scale dataset for fact extraction and verification. *ArXiv*, abs/1803.05355.
- Venktesh V, Abhijit Anand, Avishek Anand, and Vinay Setty. 2024. Quantemp: A real-world open-domain benchmark for fact-checking numerical claims. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*.
- Henning Wachsmuth, Shahbaz Syed, and Benno Stein. 2018. Retrieval of the best counterargument without prior topic knowledge. In *Annual Meeting of the Association for Computational Linguistics*.
- David Wadden, Kyle Lo, Lucy Lu Wang, Shanchuan Lin, Madeleine van Zuylen, Arman Cohan, and Hannaneh Hajishirzi. 2020. Fact or fiction: Verifying scientific claims. In *Conference on Empirical Methods in Natural Language Processing*.
- Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. 2022a. Text embeddings by weakly-supervised contrastive pre-training. *ArXiv*, abs/2212.03533.
- Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. 2023a. Improving text embeddings with large language models. *ArXiv*, abs/2401.00368.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023b. Self-instruct: Aligning language models with self-generated instructions. In *ACL*, pages 13484–13508.
- Yizhong Wang, Swaroop Mishra, Pegah Alipoor-molabashi, Yeganeh Kordi, Amirreza Mirzaei, Anjana Arunkumar, Arjun Ashok, Arut Selvan Dhanasekaran, Atharva Naik, David Stap, Eshaan Pathak, Giannis Karamanolakis, Haizhi Gary Lai, Ishan Purohit, Ishani Mondal, Jacob Anderson, Kirby Kuznia, Krma Doshi, Maitreya Patel, Kuntal Kumar Pal, M. Moradshahi, Mihir Parmar, Mirali Purohit, Neeraj Varshney, Phani Rohitha Kaza, Pulkit Verma, Ravsehaj Singh Puri, Rushang Karia, Shailaja Keyur Sampat, Savan Doshi, Siddhartha Mishra, Sujan Reddy, Sumanta Patro, Tanay Dixit, Xudong Shen, Chitta Baral, Yejin Choi, Noah A. Smith, Hannaneh Hajishirzi, and Daniel Khashabi. 2022b. Super-naturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks. In *Conference on Empirical Methods in Natural Language Processing*.
- Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V Le. 2022. Finetuned language models are zero-shot learners. In *ICLR*.
- Sean Welleck, Jiacheng Liu, Ronan Le Bras, Hannaneh Hajishirzi, Yejin Choi, and Kyunghyun Cho. 2021. Naturalproofs: Mathematical theorem proving in natural language. *ArXiv*, abs/2104.01112.
- Orion Weller, Benjamin Chang, Sean MacAvaney, Kyle Lo, Arman Cohan, Benjamin Van Durme, Dawn Lawrie, and Luca Soldaini. 2024. Followir: Evaluating and teaching information retrieval models to follow instructions. *arXiv preprint arXiv:2403.15246*.
- Orion Weller, Dawn J Lawrie, and Benjamin Van Durme. 2023. Nevir: Negation in neural information retrieval. In *Conference of the European Chapter of the Association for Computational Linguistics*.
- Marco Wrzalik and Dirk Krechel. 2021. Gerdalir: A german dataset for legal information retrieval. In *NLLP*.
- Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas Muennighoff. 2023. C-pack: Packaged resources to advance general chinese embedding. *ArXiv*, abs/2309.07597.
- Jason Yang, Ariane Mora, Shengchao Liu, Bruce J. Wittmann, Anima Anandkumar, Frances H. Arnold, and Yisong Yue. 2024. Care: a benchmark suite for the classification and retrieval of enzymes. In *NeurIPS*.
- Andrew Yates, Rodrigo Nogueira, and Jimmy Lin. 2021. Pretrained transformers for text ranking: Bert and beyond. In *Proceedings of the 14th ACM International Conference on web search and data mining*, pages 1154–1156.
- Tao Yu, Rui Zhang, Kai-Chou Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Z Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir R. Radev. 2018. Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-sql task. *ArXiv*, abs/1809.08887.

- Tao Yu, Rui Zhang, Michihiro Yasunaga, Yi Chern Tan, Xi Victoria Lin, Suyi Li, He Yang Er, Irene Z Li, Bo Pang, Tao Chen, Emily Ji, Shreya Dixit, David Proctor, Sungrok Shim, Jonathan Kraft, Vincent Zhang, Caiming Xiong, Richard Socher, and Dragomir R. Radev. 2019. Sparc: Cross-domain semantic parsing in context. *ArXiv*, abs/1906.02285.
- Wenhao Zhang, Mengqi Zhang, Shiguang Wu, Jiahuan Pei, Zhaochun Ren, Maarten de Rijke, Zhumin Chen, and Pengjie Ren. 2024. Excluir: Exclusionary neural information retrieval. *ArXiv*, abs/2404.17288.
- Wayne Xin Zhao, Jing Liu, Ruiyang Ren, and Ji-Rong Wen. 2024. Dense text retrieval based on pretrained language models: A survey. *ACM Transactions on Information Systems*, 42(4):1–60.
- Qinkai Zheng, Xiao Xia, Xu Zou, Yuxiao Dong, Shanshan Wang, Yufei Xue, Zihan Wang, Lei Shen, Andi Wang, Yang Li, Teng Su, Zhilin Yang, and Jie Tang. 2023. Codegeex: A pre-trained model for code generation with multilingual evaluations on humaneval-x. *ArXiv*, abs/2303.17568.
- Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. 2023. Instruction-following evaluation for large language models. *ArXiv*, abs/2311.07911.
- Shuyan Zhou, Uri Alon, Frank F. Xu, Zhiruo Wang, Zhengbao Jiang, and Graham Neubig. 2022. Docprompting: Generating code by retrieving the docs. In *International Conference on Learning Representations*.
- Fengbin Zhu, Wenqiang Lei, Fuli Feng, Chao Wang, Haozhou Zhang, and Tat seng Chua. 2022. Towards complex document understanding by discrete reasoning. *Proceedings of the 30th ACM International Conference on Multimedia*.
- Yutao Zhu, Huaying Yuan, Shuting Wang, Jiongnan Liu, Wenhan Liu, Chenlong Deng, Zhicheng Dou, and Ji-Rong Wen. 2023. Large language models for information retrieval: A survey. *arXiv preprint arXiv:2308.07107*.

Table 3: List of tasks in MAIR.

Task	Reference	Query	Doc	Task Type	Domain
Competition-Math	Hendrycks et al. (2021c)	Math question	Answer	Question Answering	Academic
ProofWiki_Proof	Welleck et al. (2021)	Math Statement	Proof	Proof Retrieval	Academic
Stacks_Proof	Welleck et al. (2021)	Math Statement	Proof	Proof Retrieval	Academic
Stein_Proof	Welleck et al. (2021)	Math Statement	Proof	Proof Retrieval	Academic
Trench_Proof	Welleck et al. (2021)	Math Statement	Proof	Proof Retrieval	Academic
ProofWiki_Reference	Welleck et al. (2021)	Math Statement	Reference	Reference Retrieval	Academic
Stacks_Reference	Welleck et al. (2021)	Math Statement	Reference	Reference Retrieval	Academic
Stein_Reference	Welleck et al. (2021)	Math Statement	Reference	Reference Retrieval	Academic
Trench_Reference	Welleck et al. (2021)	Math Statement	Reference	Reference Retrieval	Academic
TAD	Chen et al. (2022a)	Paper Abstract	Citation	Reference Retrieval	Academic
TAS2	Chen et al. (2022a)	Paper Abstract	Citation	Reference Retrieval	Academic
StackMathQA	Chen et al. (2022a)	Math Question	Answer	Question Answering	Academic
SciDocs	Cohan et al. (2020)	Paper Title	Citation	Reference Retrieval	Academic
SciFact	Wadden et al. (2020)	Claim	Document	Fact Checking	Academic
FairRanking_2020	Fair-TREC	Key Words	Academic Articles	Web Search	Academic
LitSearch	Ajith et al. (2024)	Question	Articles	Literature Search	Academic
APPS	Hendrycks et al. (2021a)	Code Problem	Solution	Code Retrieval	Code
CodeEditSearch	Muennighoff et al. (2023)	Commit	Code Diff	Code Retrieval	Code
CodeSearchNet	Husain et al. (2019)	Function Header	Function	Code Retrieval	Code
HumanEval-X	Zheng et al. (2023)	Code Problem	Solution	Code Retrieval	Code
LeetCode	LeetCode	Code Problem	Solution	Code Retrieval	Code
MBPP	Austin et al. (2021)	Code Problem	Solution	Code Retrieval	Code
Conala	Zhou et al. (2022)	NL Command	Command Doc	Code Retrieval	Code
TLDR	Zhou et al. (2022)	CNL Command	Command Doc	Code Retrieval	Code
RepoBench	Liu et al. (2023)	Code Context	Next Function	GitHub Agent	Code
SWE-Bench	Jimenez et al. (2023)	GitHub Issue	Related File	GitHub Agent	Code
FoodAPI	Shi et al. (2024a)	Question	Food API	Tool Retrieval	Code
TensorAPI	Patil et al. (2023)	Question	Tensor API	Tool Retrieval	Code
HF-API	Patil et al. (2023)	Question	HuggingFace API	Tool Retrieval	Code
PyTorchAPI	Patil et al. (2023)	Question	PyTorch API	Tool Retrieval	Code
SpotifyAPI	Shi et al. (2024a)	Question	Spotify API	Tool Retrieval	Code
TMDB	Shi et al. (2024a)	Question	TMDB API	Tool Retrieval	Code
WeatherAPI	Shi et al. (2024a)	Question	Weather API	Tool Retrieval	Code
ToolBench	Qin et al. (2023)	Question	Tool	Tool Retrieval	Code
Apple	FinQabench	Question	Paragraph	Question Answering	Finance
ConvFinQA	Chen et al. (2022b)	Dialog	Table & Paragraph	Dialog Retrieval	Finance
FinQA	Chen et al. (2021a)	Question	Paragraph	Question Answering	Finance
FinanceBench	Islam et al. (2023)	Question	Pages	Question Answering	Finance
HC3Finance	HC3Finance	Question	Answer	Question Answering	Finance
TAT-DQA	Zhu et al. (2022)	Question	Table & Paragraph	Question Answering	Finance
Trade-the-event	Trade-the-event	Title	Article	Summary Retrieval	Finance
FiQA	Maia et al. (2018)	Question	Answer	Question Answering	Finance
AILA2019-Case	Bhattacharya et al. (2019)	Situation	Prior Case	Case Retrieval	Legal
AILA2019-Statutes	Bhattacharya et al. (2019)	Situation	Statute	Statute Retrieval	Legal
BSARD	Louis et al. (2021)	Question	Statute	Statute Retrieval	Legal
BillSum	Kornilova and Eidelman (2019)	Summary	Article	Summary Retrieval	Legal
CUAD	Hendrycks et al. (2021b)	Instruction	Highlight	Contract Review	Legal
GerDaLIR	Wrzalik and Krechel (2021)	Legal Case	Prior Cases	Case Retrieval	Legal
LeCaRDv2	Li et al. (2023a)	Legal Case	Prior Cases	Case Retrieval	Legal
LegalQuAD	Hoppe et al. (2021)	Question	Statute	Statute Retrieval	Legal
REGIR-EU2UK	Chalkidis et al. (2021)	EU Directive	UK Legislation	Statute Retrieval	Legal
REGIR-UK2EU	Chalkidis et al. (2021)	UK Legislation	EU Directive	Statute Retrieval	Legal
TREC-Legal_2011	TREC Legal	Request	Communications	Case Retrieval	Legal
NFCorpus	Boteva et al. (2016)	Question	Medical Document	Medical QA	Medical
Trec-Covid	Roberts et al. (2020)	Medical Question	Answer	Medical QA	Medical
Monant	Srba et al. (2022)	Medical Claim	Document	Question Answering	Medical
CliniDS_2014	TREC CDS 2014	Medical Case	Articles	Medical IR	Medical
CliniDS_2015	TREC CDS 2015	Medical Case	Articles	Medical IR	Medical
CliniDS_2016	TREC CDS 2016	Health Record	Articles	Medical IR	Medical

Table 4: List of tasks in MAIR.

Task	Reference	Query	Doc	Task Type	Domain
ClinicalTrials_2021	TREC CDS 2021	Health Record	Clinical Trials	Medical IR	Medical
ClinicalTrials_2022	TREC CDS 2022	Health Record	Clinical Trials	Medical IR	Medical
ClinicalTrials_2023	TREC CDS 2023	Patient Description	Clinical Trials	Medical IR	Medical
Genomics-AdHoc_2004	Genomics	Question	Document	Medical IR	Medical
Genomics-AdHoc_2005	Genomics	Question	Document	Medical IR	Medical
Genomics-AdHoc_2006	Genomics	Question	Document	Medical IR	Medical
Genomics-AdHoc_2007	Genomics	Question	Entity & Passages	Medical IR	Medical
PM_(17,18,19)	Precision Medicine	Patient Data	Clinical Trials	Medical IR	Medical
PM-Article_(19,20)	Precision Medicine	Patient Data	Articles	Medical IR	Medical
CARE	Yang et al. (2024)	Reaction	Proteins Documents	Medical IR	Medical
ELI5	Fan et al. (2019)	Question	Passages	Question Answering	Web
Fever	Thorne et al. (2018)	Claim	Passages	Fact Checking	Web
AY2	Hoffart et al. (2011)	Entity Mention	Entity Page	Entity Linking	Web
WnCw	Guo and Barbosa (2018)	Entity Mention	Entity Page	Entity Linking	Web
WnWi	Guo and Barbosa (2018)	Entity Mention	Entity Page	Entity Linking	Web
TREx	ElSahar et al. (2018)	Entity & Relation	Entity Page	Slot Filling	Web
zsRE	Levy et al. (2017)	Entity & Relation	Entity Page	Slot Filling	Web
WoW	Dinan et al. (2018)	Dialog	Passage	Dialog Retrieval	Web
ArguAna	Wachsmuth et al. (2018)	Claim	Document	Argument Retrieval	Web
CQADupStack	Hoogeveen et al. (2015)	Question	Duplicate Question	Query Retrieval	Web
Quora	Quora	Question	Duplicate Question	Query Retrieval	Web
TopiOCQA	Adlakha et al. (2021)	Dialog	Passage	Dialog Retrieval	Web
Touche	Bondarenko et al. (2020)	Question	Passages	Argument Retrieval	Web
ACORDAR	Lin et al. (2022)	Request	Dataset	Dataset Retrieval	Web
CPCD	Chaganty et al. (2023)	Dialog	Music	Recommendation	Web
ChroniclingAmericaQA	Piryani et al. (2024)	Question	News	News Retrieval	Web
NTCIR	Kato et al. (2021)	Key Words	Dataset	Dataset Retrieval	Web
PointRec	Afzali et al. (2021)	Question	POI	Recommendation	Web
ProCIS-Dialog	Samarinas and Zamani (2024)	Dialog	Passage	Dialog Retrieval	Web
ProCIS-Turn	Samarinas and Zamani (2024)	Utterance	Passage	Dialog Retrieval	Web
QuanTemp	V et al. (2024)	Numerical Claim	Document	Fact Checking	Web
WebTableSearch	Chen et al. (2021b)	Question	Table	Table Retrieval	Web
CAst_(19,20,21,22)	TREC CAst	Dialog	Passage	Dialog Retrieval	Web
DD_2015	Domain2015	Topic	Document	Web Search	Web
DD_2016	Domain2016	Topic	Document	Web Search	Web
DD_2017	Domain2017	Topic	Document	Web Search	Web
FairRanking_2021	Fair-TREC	WikiProject	Wikipedia Page	Web Search	Web
FairRanking_2022	Fair-TREC	WikiProject	Wikipedia Page	Web Search	Web
NeuCLIR-Tech	NeuCLIR	Question	Document	Web Search	Web
NeuCLIR_2022	NeuCLIR	Question	Document	Web Search	Web
NeuCLIR_2023	NeuCLIR	Question	Document	Web Search	Web
ProductSearch_2023	TREC Product Search	Key Words	Product	Product Search	Web
ToT_2023	TREC ToT	Description	Wikipedia Page	Tip-of-the-Tongue	Web
ToT_2024	TREC ToT	Description	Wikipedia Page	Tip-of-the-Tongue	Web
Microblog	TREC Microblog	Topic	Tweet	Tweet Retrieval	Web
MISeD	Golany et al. (2024)	Dialog	Meeting Transcript	Meeting Retrieval	Web
SParC	Yu et al. (2019)	Question	Table	Table Retrieval	Web
SParC-SQL	Yu et al. (2019)	SQL Statement	Table	Table Retrieval	Web
Spider	Yu et al. (2018)	Question	Table	Table Retrieval	Web
Spider-SQL	Yu et al. (2018)	SQL Statement	Table	Table Retrieval	Web
ExcluIR	Zhang et al. (2024)	Question	Passage	Negative Retrieval	Web
FollowIR_Core17	Weller et al. (2024)	Question	Document	Negative Retrieval	Web
FollowIR_News21	Weller et al. (2024)	Question	News	Negative Retrieval	Web
FollowIR_Robust04	Weller et al. (2024)	Question	News	Negative Retrieval	Web
InstructIR	Oh et al. (2024)	Question	Passage	Web Search	Web
NevIR	Weller et al. (2023)	Question	Passage	Negative Retrieval	Web
IFEval	Zhou et al. (2023)	Question	Answer	Question Answering	Web

Table 9: Results (nDCG@10) on MAIR tasks. The last group represents the re-ranking models, which all re-rank the top 100 results of text-embedding-3-small.

Model	Core17		News21		Robust04		InstructIR		NevIR		IFEval	
Instruction?	X	✓	X	✓	X	✓	X	✓	X	✓	X	✓
BM25	11.51	-	20.79	-	18.57	-	36.31	-	66.76	-	28.28	-
contriever-msmarco	18.18	-	27.33	-	26.64	-	50.45	-	62.97	-	25.51	-
gtr-t5-base	22.73	-	22.49	-	26.90	-	49.04	-	63.14	-	22.59	-
gtr-t5-large	22.80	-	26.62	-	29.07	-	50.87	-	59.95	-	23.46	-
all-MiniLM-L6-v2	20.46	17.50	26.86	25.76	25.96	23.92	48.28	74.06	57.43	41.58	26.06	37.39
e5-small-v2	22.84	20.45	28.35	24.02	23.85	19.90	49.82	80.25	62.58	60.99	21.75	36.96
e5-base-v2	21.62	21.03	22.92	22.17	24.32	21.09	48.31	84.64	63.96	55.68	24.47	40.33
e5-large-v2	22.97	29.12	25.27	26.02	29.03	23.89	49.85	83.76	66.60	62.25	22.49	39.28
gte-base-en-v1.5	21.75	21.44	25.63	22.60	27.57	25.52	47.75	86.08	58.48	41.88	25.96	42.18
gte-large-en-v1.5	27.61	26.58	25.16	26.22	30.68	29.14	48.93	87.13	58.95	45.56	25.30	41.23
bge-base-en-v1.5	19.00	20.63	21.23	20.60	23.22	19.50	52.36	87.25	63.74	48.69	24.42	36.48
bge-large-en-v1.5	24.05	18.26	21.41	22.37	27.65	19.16	53.16	85.91	64.42	36.28	22.43	37.65
text-embedding-3-small	27.57	30.04	26.22	25.60	32.19	28.15	53.01	86.80	69.06	61.44	21.50	31.91
e5-mistral-7b-instruct	21.85	33.41	29.35	31.74	30.44	36.05	42.38	85.30	64.59	62.30	24.98	32.15
NV-Embed-v1	26.49	32.81	33.11	30.94	33.04	35.13	49.66	66.14	69.55	69.54	22.45	27.68
GritLM-7B	14.47	34.94	21.87	31.03	29.93	38.27	29.89	79.68	59.49	71.53	22.03	38.73
gte-Qwen2-1.5B-instruct	21.52	32.65	24.36	31.77	30.06	33.79	50.32	82.87	65.28	65.38	21.85	42.76
monot5-base-msmarco	28.99	18.98	26.21	14.77	32.51	24.70	48.47	55.21	76.94	75.70	24.71	31.17
bge-reranker-v2-m3	28.17	26.80	27.31	27.06	33.91	31.80	50.55	78.45	81.92	79.06	20.86	21.23
bge-reranker-v2-gemma	28.28	29.02	31.41	27.62	39.11	30.19	52.99	90.40	80.84	65.00	19.42	53.73
jina-reranker-v2-base	26.68	31.53	29.13	30.20	32.99	34.68	48.87	78.45	79.29	75.07	11.17	20.80
mbai-rerank-large-v1	32.60	25.82	30.64	14.58	33.05	19.53	49.45	86.81	83.53	48.47	17.13	27.74
gpt-3.5-turbo	29.09	33.43	29.62	30.27	29.06	29.73	52.22	78.64	76.96	72.50	15.15	29.04

Table 11: Results of each sub task in SWE-Bench, CQADupStack, and IFEval.

Task Sub-Task	CQADupStack										IFEval									
	Unix		WebMasters		Wordpress		format		keywords		punctuation		change_case		length_constraints		combination		content	
Instruction?	X	✓	X	✓	X	✓	X	✓	X	✓	X	✓	X	✓	X	✓	X	✓	X	✓
BM25	32.80	-	31.60	-	16.31	-	26.81	-	26.67	-	42.69	-	42.85	-	20.91	-	33.11	-	34.10	-
contriever-msmarco	46.18	-	19.83	-	25.00	-	19.87	-	22.40	-	29.02	-	39.50	-	22.20	-	24.50	-	32.82	-
gtr-t5-base	41.33	-	24.31	-	10.00	-	18.85	-	21.27	-	32.62	-	22.13	-	19.07	-	33.82	-	31.01	-
gtr-t5-large	47.52	-	18.72	-	26.57	-	19.33	-	17.93	-	34.26	-	25.36	-	21.29	-	32.31	-	31.73	-
all-MiniLM-L6-v2	55.63	33.63	32.50	30.00	9.20	6.31	22.42	40.41	19.33	57.56	30.13	33.27	39.31	27.34	21.75	24.56	26.03	33.95	32.69	42.10
e5-small-v2	28.74	32.18	40.59	32.64	25.00	17.20	18.22	31.97	10.43	63.87	36.18	36.34	36.28	30.03	18.37	34.38	23.91	26.51	26.14	38.95
e5-base-v2	59.31	33.33	44.78	32.64	24.31	27.46	21.75	36.57	13.01	64.93	36.19	33.80	39.29	30.31	19.09	32.13	19.82	26.08	28.14	44.45
e5-large-v2	53.87	28.18	38.32	38.33	18.87	16.88	17.23	33.94	19.73	65.39	37.78	33.50	36.38	37.36	17.19	32.84	23.32	25.22	33.19	45.90
gte-base-en-v1.5	51.31	44.08	35.70	46.13	22.20	16.31	23.83	51.98	19.08	49.93	32.32	20.50	28.70	27.03	20.37	24.15	25.52	26.84	31.71	49.29
gte-large-en-v1.5	57.62	35.00	40.41	19.31	26.62	5.00	22.58	40.72	17.85	50.29	32.78	40.37	28.12	34.55	21.62	25.67	23.12	25.45	33.35	53.54
bge-base-en-v1.5	44.38	22.88	45.18	40.73	31.07	19.64	21.11	44.40	22.68	54.94	33.71	22.29	28.48	23.68	20.45	23.74	23.39	26.75	30.20	46.03
bge-large-en-v1.5	48.48	35.30	35.59	29.95	25.81	12.62	17.88	40.04	18.76	57.36	31.77	20.72	31.58	20.38	18.89	27.19	19.56	25.70	30.54	51.96
text-embedding-3-small	67.62	61.55	36.88	33.87	23.51	22.35	23.47	34.09	23.11	48.87	32.44	33.47	20.28	17.98	16.94	23.08	13.41	20.78	28.73	38.65
e5-mistral-7b-instruct	60.77	53.93	38.20	48.69	21.31	30.18	23.04	32.74	11.51	58.66	31.98	33.90	31.34	21.96	14.60	22.76	38.96	24.51	30.24	32.65
NV-Embed-v1	68.93	72.62	46.23	53.39	37.64	30.09	21.65	25.74	22.51	42.90	28.55	37.48	27.25	27.58	17.89	19.91	17.77	24.17	31.32	36.96
GritLM-7B	62.88	68.80	45.06	59.81	12.20	16.74	20.75	34.36	5.51	65.11	42.05	33.82	29.92	40.99	12.87	29.17	32.33	49.73	23.13	45.11
gte-Qwen2-1.5B-instruct	48.79	56.49	40.44	36.95	8.56	12.87	18.13	44.52	14.95	58.09	32.27	32.59	28.90	34.94	16.40	23.86	24.68	42.86	25.95	56.47
monot5-base-msmarco	49.93	41.33	35.03	28.18	29.46	25.62	22.18	24.75	20.54	44.52	20.97	26.82	32.71	44.53	25.06	26.08	29.90	22.54	33.89	32.23
bge-reranker-v2-m3	50.21	43.91	20.00	23.01	24.48	22.74	19.62	13.89	19.95	38.65	32.48	27.92	21.82	13.10	21.96	28.75	17.89	12.51	31.01	32.39
bge-reranker-v2-gemma	54.48	41.43	35.80	29.64	30.77	18.05	22.09	46.06	16.68	74.87	23.97	33.63	16.45	22.64	18.02	39.14	17.29	79.48	26.14	89.00
jina-reranker-v2-base	46.09	51.32	29.66	32.37	32.80	24.32	11.55	24.03	4.01	31.94	12.07	20.72	15.91	8.69	17.05	26.22	17.32	18.45	7.53	21.08
mxbai-reranker-large-v1	47.49	29.64	36.05	28.02	23.77	3.87	21.38	34.19	9.76	35.70	15.37	16.52	6.42	20.66	20.11	32.83	25.37	12.26	23.88	27.02
gpt-3.5-turbo	61.19	39.64	26.49	25.71	17.20	14.46	17.59	30.77	18.99	46.48	15.21	31.70	10.23	21.13	18.85	30.91	5.85	16.79	16.45	32.19

Table 12: Results of each sub task in SWE-Bench, CQADupStack, and IFEval.

Task Sub-Task	IFEval startend	
Instruction?	X	✓
BM25	14.55	-
contriever-msmarco	26.57	-
gtr-t5-base	20.55	-
gtr-t5-large	22.68	-
all-MiniLM-L6-v2	26.62	40.72
e5-small-v2	18.86	43.20
e5-base-v2	27.78	59.70
e5-large-v2	14.72	48.01
gte-base-en-v1.5	34.70	62.87
gte-large-en-v1.5	31.31	61.94
bge-base-en-v1.5	27.17	33.74
bge-large-en-v1.5	23.24	48.42
text-embedding-3-small	14.01	37.18
e5-mistral-7b-instruct	35.03	35.23
NV-Embed-v1	17.17	19.75
GritLM-7B	28.59	31.65
gte-Qwen2-1.5B-instruct	28.39	53.10
monot5-base-msmarco	16.90	38.37
bge-reranker-v2-m3	7.83	22.97
bge-reranker-v2-gemma	11.29	72.85
jina-reranker-v2-base	2.29	7.08
mxbai-reranker-large-v1	7.27	16.82
gpt-3.5-turbo	7.88	16.35

Table 13: Results (nDCG@10) on IFEval sub task. AVG means average results.

Model	AVG		detectable_format		keywords		punctuation		change_case		length_constraints		combination		detectable_content		startend		
	X	✓	X	✓	X	✓	X	✓	X	✓	X	✓	X	✓	X	✓	X	✓	
BM25	28.28	-	26.81	-	26.67	-	42.69	-	42.85	-	20.91	-	33.11	-	34.10	-	14.55	-	
text-embedding-3-small	21.50	31.91	23.47	34.09	23.11	48.87	32.44	33.47	20.28	17.98	16.94	23.08	13.41	20.78	28.73	38.65	14.01	37.18	
e5-mistral-7b-instruct	24.98	32.15	23.04	32.74	11.51	58.66	31.98	33.90	31.34	21.96	14.60	22.76	38.96	24.51	30.24	32.65	35.03	35.23	
NV-Embed-v1	22.45	27.68	21.65	25.74	22.51	42.90	28.55	37.48	27.25	27.58	17.89	19.91	17.77	24.17	31.32	36.96	17.17	19.75	
GritLM-7B	22.03	38.73	20.75	34.36	5.51	65.11	42.05	33.82	29.92	40.99	12.87	29.17	32.33	49.73	23.13	45.11	28.59	31.65	
gte-Qwen2-1.5B-instruct	21.85	42.76	18.13	44.52	14.95	58.09	32.27	32.59	28.90	34.94	16.40	23.86	24.68	42.86	25.95	56.47	28.39	53.10	
bge-reranker-v2-gemma	19.42	53.73	22.09	46.06	16.68	74.87	23.97	33.63	16.45	22.64	18.02	39.14	17.29	79.48	26.14	89.00	11.29	72.85	
RankGPT																			
gpt-3.5-turbo	15.15	29.04	17.59	30.77	18.99	46.48	15.21	31.70	10.23	21.13	18.85	30.91	5.85	16.79	16.45	32.19	7.88	16.35	
gpt-4o-mini	-	52.68	-	61.60	-	65.96	-	38.24	-	52.70	-	37.70	-	35.94	-	54.82	-	53.24	-
gpt-4o	-	86.09	-	95.43	-	82.29	-	51.99	-	84.83	-	70.70	-	87.04	-	92.40	-	98.66	-