

PromptReps: Prompting Large Language Models to Generate Dense and Sparse Representations for Zero-Shot Document Retrieval

Anonymous ACL submission

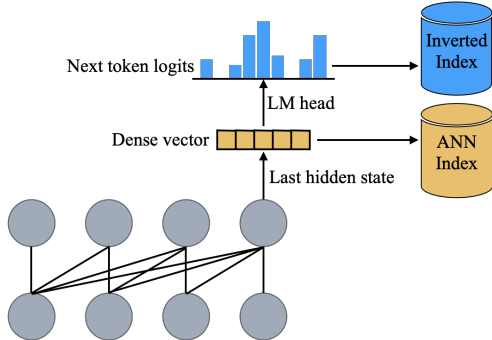
Abstract

Utilizing large language models (LLMs) for zero-shot document ranking is done in one of two ways: 1) prompt-based re-ranking methods, which require no further training but are only feasible for re-ranking a handful of candidate documents due to computational costs; and 2) unsupervised contrastive trained dense retrieval methods, which can retrieve relevant documents from the entire corpus but require a large amount of paired text data for contrastive training. In this paper, we propose PromptReps, which combines the advantages of both categories: no need for training and the ability to retrieve from the whole corpus. Our method only requires prompts to guide an LLM to generate query and document representations for effective document retrieval. Specifically, we prompt the LLMs to represent a given text using a single word, and then use the last token’s hidden states and the corresponding logits associated with the prediction of the next token to construct a hybrid document retrieval system. The retrieval system harnesses both dense text embedding and sparse bag-of-words representations given by the LLM. We further explore variations of this core idea that consider the generation of multiple words, and representations that rely on multiple embeddings and sparse distributions. Our experimental evaluation on the MSMARCO, TREC deep learning and BEIR zero-shot document retrieval datasets illustrates that this simple prompt-based LLM retrieval method can achieve a similar or higher retrieval effectiveness than state-of-the-art LLM embedding methods that are trained with large amounts of unsupervised data, especially when using a larger LLM.¹

1 Introduction

Large Language Models (LLMs) such as GPT4 and LLaMA, which are pretrained on massive cor-

¹Code for fully reproducing the results is available at <https://anonymous.4open.science/r/PromptReps-anonymous-58DE>.



<System> You are an AI assistant that can understand human language.
<User> Passage: “[text]”. Use one word to represent the passage in a retrieval task. Make sure your word is in lowercase.
<Assistant> The word is: “

Figure 1: Overview of PromptReps. LLMs are prompted to simultaneously generate dense and sparse representations, which are then used to build search indices.

pura and finetuned to follow user instructions, have strong zero-shot natural language understanding capabilities (OpenAI, 2023; Touvron et al., 2023). Via prompting, LLMs excel in various text generation tasks such as question answering, writing assistance, and conversational agent (Hendrycks et al., 2021; Liu et al., 2023). Inspired by the success of LLMs on natural language understanding tasks, research has explored the potential of using LLMs to perform unsupervised document ranking.

One line of work focuses on directly prompting LLMs to infer document relevance to a given query (Sachan et al., 2022; Zhuang et al., 2023a; Ma et al., 2023b; Sun et al., 2023; Pradeep et al., 2023; Zhuang et al., 2023b; Qin et al., 2024). For instance, RankGPT (Sun et al., 2023) casts document re-ranking as a permutation generation task, prompting LLMs to generate re-ordered document identifiers according to the document’s relevance to the query. These methods leverage LLMs for document ranking in a complete zero-shot setting where no further training is required. However,

these methods can only serve as a second-stage re-ranker on a handful of candidate documents. This is because each prompt requires one full LLM inference: for example, in the case of a corpus with 1M documents, a pointwise approach would require the construction of 1M prompts and thus the execution of 1M (costly) LLM inferences – making it unfeasible for an online search engine.

Another line of research leverages LLMs as a text embedding model for dense document retrieval (Lee et al., 2024; Wang et al., 2024a,b; BehnamGhader et al., 2024). For example, E5-mistral (Wang et al., 2024b) employs LLMs to create synthetic datasets of query-document pairs. These paired text data are then used to perform unsupervised contrastive training for a Mistral LLM-based dense retriever. Since the queries and documents are encoded with LLMs separately; i.e., using a bi-encoder architecture, these methods could serve as a first-stage document retriever. However, all existing LLM-based retrievers require an unsupervised contrastive training step to transform a generative LLM into a text-embedding model. Even with parameter-efficient training techniques such as LoRA (Hu et al., 2022), this extra training is still very expensive. For example, the contrastive training of E5-mistral using a large batch size (2048) and LoRA took \approx 18 hours on 32 V100 GPUs (Wang et al., 2024b).

In this work, we propose a new zero-shot LLM-based document retrieval method called PromptReps. We demonstrate that LLMs can be directly prompted to produce query and document embeddings, which can serve as effective text representations for neural retrieval systems. Specifically, we prompt an LLM by asking it to use a single word to represent a query or a document. Then, we extract the last layer’s hidden state of the last token in the prompt as the dense representation of the input text. Simultaneously, we utilize the logits associated with predicting the subsequent token to form a sparse representation. As illustrated in Figure 1, through a single forward pass, we generate text representations for a document that can be indexed for dense, sparse, or hybrid search architectures. We also explore alternative representations in addition to the core idea in this paper, where we generate multiple words, and use multiple embeddings to represent an item (Figures 3 and 4).

Our empirical evaluation on multiple datasets show that PromptReps can achieve a similar or higher zero-shot retrieval effectiveness than previ-

ous trained LLM-based embedding methods, especially when a large LLM is utilized. Of key importance is that our method is the first LLM-based method that can effectively perform full corpus retrieval while at the same time not requiring contrastive training, demonstrating that prompt engineering for generative LLMs is capable of generating robust representations for retrieval.

2 Related Work

2.1 Supervised Neural Retrievers

Neural retrievers based on the bi-encoder architecture bring significant improvements over traditional best-match retrievers such as BM25. Dense retrievers such as DPR (Karpukhin et al., 2020), ANCE (Xiong et al., 2021), ColBERT (Khattab and Zaharia, 2020), are based on encoder-only language models and encode text into low-dimensional dense vectors, conducting search with (approximate) nearest neighbor search. On the other hand, sparse neural retrievers such as DeepImpact (Mallia et al., 2021), uniCOIL (Lin and Ma, 2021), TILDE (Zhuang and Zuccon, 2021c,b), and SPLADE (Formal et al., 2021), also based on encoder-only language models, encode text into high-dimensional sparse vectors as bag-of-words representations, conducting search in an inverted index. Recent work has also explored fine-tuning generative LLMs as dense retrievers such as ReplLaMA (Ma et al., 2023a) and LLaRA (Liao et al., 2024). A hybrid neural retrieval system refers to a system that combines the rankings provided by both dense and sparse retrievers, often resulting in an enhanced final ranking (Lin and Ma, 2021; Wang et al., 2021a).

All these retrievers are trained with supervised relevance judgment data (e.g., MS MARCO (Bajaj et al., 2018)) using contrastive learning. Our work instead focuses on building a hybrid neural retrieval system with zero-shot dense and sparse document representations without supervised contrastive learning and based on generative LLMs. This capability has two implications: (1) no contrastive training is required, which is expensive when applied to LLMs with several billions parameters, and (2) no human-labelled training data is required, which may be laborious and expensive to obtain. With regards to the first point, Wang et al. (2024b) reported that the training of E5-mistral (7B parameters) took about 18 hours on 32 V100

GPUs, for an approximate cost of USD \$2,300², emissions of ≈ 5.6 kgCO₂e and consumption of ≈ 37.7 L of water for the associated cooling activities³. Scaling this training to more and larger LLMs, and more data, will consequently further increase costs. Our proposed method does not incur these additional contrastive pre-training costs. With regards to the second point, dense retrievers have shown to have poor generalisability when applied to data out-of-domain or out-of-task compared to the data used for contrastive training (Thakur et al., 2021; Zhuang and Zuccon, 2021a, 2022; Ren et al., 2023; Lin et al., 2023; Lupart et al., 2022). In presence of shift in data between training and deployment, retrieval losses can be significant: dense retrieval effectiveness can plummet far below that of best-match models like BM25 (Khramtsova et al., 2023, 2024). The acquisition of in-domain/in-task training data can be costly, laborious and often impractical/impossible especially in domain-specific applications and when dealing with sensitive, private data.

2.2 Unsupervised Neural Retrievers

There have also been attempts at training effective neural retrievers without relying on human relevance judgments. Methods such as Contriever (Izacard et al., 2022) and E5 (Wang et al., 2024a), train a dense retriever with large-scale pseudo query-document pairs to build unsupervised (synthetic) training data. LLMs have also been adapted as unsupervised text embedding models for first-stage document retrieval. For instance, HyDE (Gao et al., 2023a) enhances query representations for an unsupervised retriever by replacing the original query with LLM-generated hypothetical documents. More recent work has focused on directly converting generative LLMs into a text-embedding model with unsupervised contrastive pre-training. Methods like E5-Mistral-Inst (Wang et al., 2024b) and Gecko (Lee et al., 2024) use large-scale weakly supervised paired text data or LLM-generated query-document pair data to perform contrastive training on top of LLMs. LLM2Vec (BehnamGhader et al., 2024), on the other hand, conducts further masked next token prediction pre-training with bidirectional attention, and SimCSE (Gao et al., 2021) trains on raw text data to transform LLMs into text encoders. Al-

²Based on 4 On-Demand p3dn.24xlarge instances, June 2024.

³Emissions and water consumption estimates obtained using the frameworks of Scells et al. (2022); Zuccon et al. (2023).

though no labeled data is used, these methods require synthetic or unsupervised paired text data to perform contrastive pre-training (thus still experiencing training costs in terms of computations; and further computational costs may be associated with the generation of synthetic training data). Our method instead relies solely on prompt engineering to transform LLM into a robust text encoder for document retrieval without any extra training.

2.3 Prompting LLMs for document ranking

Inspired by the prompt-following capacity of LLMs, recent studies have explored prompting LLMs for document re-ranking. For instance, UPR (Sachan et al., 2022) ranks documents pointwise by prompting the LLM to generate a relevant query for a given document and rank documents based on the likelihood of generating the query. RankGPT (Sun et al., 2023) and LRL (Ma et al., 2023b) propose to re-rank a list of documents at once and generate permutations for the reordered list. Pairwise (Qin et al., 2024) and Setwise (Zhuang et al., 2023b) prompting methods have also been explored to improve effectiveness and efficiency in the LLM re-ranking pipeline. These methods are only feasible for re-ranking a handful of candidate documents, thus limited to second-stage document re-ranking. In contrast, our approach utilizes prompts to construct the first-stage retrievers.

2.4 Prompting LLM for sentence embeddings

The methods most similar to ours prompt LLMs to generate sentence embeddings for semantic textual similarity (STS) tasks (Jiang et al., 2023b; Lei et al., 2024; Zhang et al., 2024). These previous methods also used an Explicit One-word Limitation (EOL) prompt, which also instructs LLMs to represent a sentence with one word. However, these methods only evaluate such prompts on STS datasets, and their effectiveness on information retrieval datasets with large document corpora is unknown. Additionally, these methods only represent text with dense embeddings from the hidden states; our method instead generates dense and sparse representations simultaneously to build a hybrid retrieval system. Our empirical results show that dense embeddings alone perform poorly for document retrieval tasks with some LLMs, but sparse representations are much more robust, and the best retrieval effectiveness is achieved with the hybrid retrieval system with scaled model size.

3 PromptReps

Previous work that leverages LLMs for document ranking are limited to document re-ranking tasks with prompts or rely on contrastive learning to transform a generative LLM into an embedding model for document retrieval. Unlike these previous works, here we aim to directly prompt LLMs to generate both dense embedding representations and sparse bag-of-words representations for document retrieval without any form of extra training effort. To achieve this, we devise the prompt as illustrated in Figure 1 as the input text for LLMs, where `<System>` `<User>` and `<Assistant>` are LLM pre-defined conversational prefix tokens and `[text]` is the placeholder for passage text.

When using this prompt for text generation, the language model needs to find a single word in its token vocabulary that can best represent the given passage to generate. However, since there could be multiple words to represent the passage, there might be multiple tokens in the vocabulary that have a high probability of being sampled by the language model. Such a distribution over the vocabulary, which is often refers to as "logits", could provide a good representation of the given passage. In addition, since the logits are computed by the last layer hidden state⁴ of the last input token (' '), which is a dense vector embedding, it could also serve as a dense representation of the passage.

Based on the above intuition, we develop a sparse + dense hybrid document retrieval system by utilizing both the next token logits and the last layer hidden states outputted by the LLM with our designed prompt.

Specifically, during the document indexing phase, we pass all the documents (one at the time) with our prompt into the LLM to get output hidden states and logits. To build a sparse retrieval pipeline with logits, we first need to sparsify the logits representation to be able to perform efficient sparse retrieval. This is because logits originally had values for all tokens in the vocabulary, essentially forming dense vectors with dimensions equal to the vocabulary size. To sparsify the logit representations for sparse retrieval, we perform the following steps:

1. Lowercase the input document text to align with the phrase "Make sure your word is in lowercase." in the prompt since this phrase skewed

the sampling distribution towards lowercase tokens (a "sparser" distribution). We then utilize the NLTK toolkit (Bird and Loper, 2004) to extract all words in the document, filtering out standard English stopwords and punctuation.

2. Next, we use the LLM's tokenizer to tokenize each extracted word and obtain their token IDs⁵. We retain only the values corresponding to the obtained token IDs in the logits and set the rest of the dimensions to zero, thereby considering only tokens present in the documents, thus enabling exact term matching in retrieval.
3. Next, we follow the SPLADE recipe (Formal et al., 2021), using the ReLU function to remove dimensions with negative values and applying log-saturation to the logits to prevent certain tokens from dominating. To further enhance the sparsity of logits, we only keep tokens within the top 128 values if the logits had more than 128 non-zero values after the previous steps.
4. Finally, the logits are quantized by multiplying the original values by 100 and taking the integer operation on that, and the obtained values represent the weights of corresponding tokens.

With these adjustments, the logits representations of documents are heavily sparsified, allowing for efficient sparse retrieval with an inverted index.

For dense retrieval, we directly use the hidden states as the embeddings of the documents. For indexing these embeddings, we simply normalize all the embeddings and add them into an Approximate Nearest search (ANN) vector index.

At query time, we process the queries exactly the same as the documents, with the only exception being that the term "passage" in the prompt is replaced with "query". The dense representation of the query is utilized for semantic search via the ANN index, while the sparse representation of the query is employed for exact term matching via the inverted index. Following previous work (Wang et al., 2021b), we compute the final document scores by applying min-max normalization to both dense and sparse document scores. These normalized scores are then linearly interpolated with equal weights to produce the final document scores.

⁴Often through dot product between the last hidden state with all token embeddings.

⁵Note that many words may be split into sub-tokens, resulting in multiple token IDs, all of which are considered in the logits

Table 1: nDCG@10 scores of BEIR 13 publicly available datasets. The best scores of methods without interpolating with BM25 are highlighted in bold.

LLM	-	Unsup Contrastive training		PromptReps (ours)						
		BERT-330M	Llama3-8B-I	Llama3-8B-I			Llama3-70B-I			
Dataset	BM25	E5-PT _{large}	LLM2Vec	Dense	Sparse	Hybrid	Dense	Sparse	Hybrid	+BM25
arguana	39.70	44.4	51.73	29.70	22.85	32.98	31.65	24.66	35.27	39.53
climatefever	16.51	15.7	23.58	19.92	9.98	21.38	19.95	12.14	22.18	23.34
dbpedia	31.80	37.1	26.78	31.53	28.84	37.71	31.12	28.30	37.59	41.63
fever	65.13	68.6	53.42	56.28	52.35	71.11	42.06	51.75	63.97	74.06
fiqa	23.61	43.2	28.56	27.11	20.33	32.40	30.80	22.16	34.66	35.35
hotpotqa	63.30	52.2	52.37	19.64	44.75	47.05	24.32	42.12	48.51	65.29
nfcopus	32.18	33.7	26.28	29.56	28.18	32.98	33.84	29.74	36.08	37.64
nq	30.55	41.7	37.65	34.43	29.55	43.14	38.25	30.37	46.97	48.30
quora	78.86	86.1	84.64	72.55	68.27	80.45	76.14	68.77	82.56	85.83
scidocs	14.90	21.8	10.39	18.51	11.57	17.59	20.59	13.25	19.10	18.82
scifact	67.89	72.3	66.36	52.68	58.48	65.71	63.12	61.53	70.34	73.58
trec-covid	59.47	61.8	63.34	59.52	54.59	69.25	67.64	63.00	76.85	80.29
touche	44.22	19.8	12.82	14.85	18.47	21.78	15.56	18.65	22.35	34.15
avg	43.70	44.61	41.38	35.87	34.48	44.13	38.08	35.88	45.88	50.60

4 Experimental setup

Dataset and evaluation: We evaluate the document ranking effectiveness of both baseline methods and our proposed PromptReps using MS-MARCO (Nguyen et al., 2016) passage retrieval, TREC deep learning (Craswell et al., 2020) and BEIR (Thakur et al., 2021). These datasets encompass various IR tasks, providing a heterogeneous evaluation environment. For MSMARCO we report MRR@10 and for TREC deep learning and BIER we report nDCG@10 scores, the commonly employed evaluation measures for these datasets.

Baselines: We compare PromptReps with strong unsupervised first-stage retrievers including BM25, a classic term frequency-based sparse retrieval method, and E5-PT_{large} (Wang et al., 2024a), a state-of-the-art BERT-based dense embedding method trained on 1.3B carefully crafted unsupervised text pairs. LLM2Vec (BehnamGhader et al., 2024), a Llama3-8B-Instruct LLM-based dense embedding method trained with bi-directional attention, masked next token prediction, and SimCSE (Gao et al., 2021) on the Wikipedia corpus.

Implementation of PromptReps: PromptReps is implemented using four base LLMs: Mistral-7b-Instruct-v0.2⁶ (Jiang et al., 2023a), Phi-3-mini-4k-instruct⁷ (Abdin et al., 2024), Llama3-8B-Instruct⁸, and Llama3-70B-Instruct⁹ (AI@Meta, 2024). Dense and sparse document and query encodings are implemented using the Huggingface

⁶<https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.2>

⁷<https://huggingface.co/microsoft/Phi-3-mini-4k-instruct>

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

⁹<https://huggingface.co/meta-llama/Meta-Llama-3-70B-Instruct>

Transformers library (Wolf et al., 2020) and the Tevatron toolkit (Gao et al., 2023b). The Faiss library (Douze et al., 2024) is used to build the ANN index with cosine similarity as the embedding distance metric, and Pyserini (Lin et al., 2021) is utilized to construct the inverted index for sparse retrieval. For the dense and sparse ranking hybrid, the Ranx library (Bassani and Romelli, 2022) is employed. In our experiments, we report dense only, sparse only, and the full hybrid results.

5 Results

We start by showing our overall results on the BEIR dataset, which we treated as test set; we then analyse choices in instantiation of PromptReps, including different variations in the prompt using the MS MARCO and TREC deep learning datasets, which we used as development datasets to inform the choices we made to run PromptReps on BEIR.

5.1 Zero-shot retrieval effectiveness on BEIR

We present our results on BEIR in Table 1. The first observation highlights that BM25 is a very strong zero-shot retrieval method, capable of outperforming LLM2Vec, based on the Llama3-8B-Instruct LLM, across numerous datasets, achieving a higher average nDCG@10 score. This outcome implies that even with a large-size LLM, bi-directional attention enabled, additional pre-training, and SimCSE-based unsupervised contrastive training, there remains a gap in transforming a decoder-only LLM into an effective retrieval method.

On the other hand, E5-PT_{large}, based on the BERT-large model, is the first method that can

Table 2: Investigated prompts. The systems prompt and any text string before the prompts in this table are the same as Figure 1, thus omitted. <A> denotes the model-specific assistant special token.

ID	Prompts
1	Use one word to represent the passage in a retrieval task.<A>The word is: "
2	Use one word to represent the passage.<A>The word is: "
3	Use one most important word to represent the passage in a retrieval task. Make sure your word is in lowercase.<A>The word is: "
4	Use one word to represent the passage in a retrieval task.<A>
5	Use one most important word to represent the passage in a retrieval task.<A>The word is: "
6	Use one word to represent the passage in a retrieval task. Make sure your word is in lowercase.<A>The word is: "

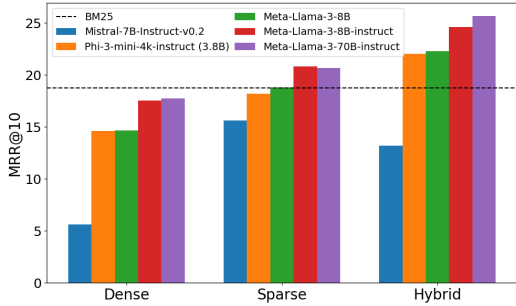


Figure 2: MRR@10 scores on MS MARCO of PromptReps with different LLMs.

outperform BM25 without any supervised training data. However, it has been trained on a massive, carefully mined text pair dataset with a large batch size, which may require more data-collecting efforts and computational resources than LLM2Vec.

Our proposed PromptReps with Llama3-8B-Instruct LLM has lower nDCG@10 scores when only using dense or sparse retrieval. However, the hybrid system (combining dense and sparse) contributes notable retrieval effectiveness improvements, surpassing BM25 and approaching the state-of-the-art E5-PT_{large}. Notably, this is achieved without any form of extra training but solely relying on prompts.

The scaling law of LLM (Kaplan et al., 2020) also applies here. When changing Llama3-8B-Instruction to Llama3-70B-Instruction, the dense and sparse retrieval effectiveness of our PromptReps further improves, with the hybrid approach surpassing E5-PT_{large}.

We further note that when interpolating dense, sparse and BM25, the average nDCG@10 achieved a remarkable score of 50.60. These results demonstrate that it is possible to build a strong retrieval system with LLMs and BM25 without the need for any unsupervised or supervised training.

5.2 Sensitivity to different prompts

In the previous experiments, we always use the prompt illustrated in Figure 1. In this section, we

Table 3: Retrieval effectiveness of different prompts on TREC deep learning and MS MARCO. The ID correspond to the prompt IDs list in Table 2.

ID	Methods	DL2019	DL2020	MSMARCO
-	BM25	49.73	48.76	18.75
-	LLM2Vec	-	-	13.61
PromptReps Llama3-8B-Instruct (ours)				
1	Dense	49.26	40.28	16.26
2	Dense	43.32	31.60	12.52
3	Dense	49.20	43.90	17.49
4	Dense	0.00	0.00	0.00
5	Dense	47.19	40.17	16.02
6	Dense	50.62	43.81	17.54
1	Sparse	41.77	44.81	20.12
2	Sparse	39.90	43.10	19.13
3	Sparse	43.50	44.87	20.42
4	Sparse	21.77	20.49	7.22
5	Sparse	42.18	44.17	19.78
6	Sparse	42.25	45.60	20.85
1	Hybrid	53.67	54.35	23.68
2	Hybrid	50.65	49.25	21.76
3	Hybrid	55.64	53.83	23.86
4	Hybrid	13.47	11.81	5.06
5	Hybrid	54.16	52.06	23.25
6	Hybrid	55.58	56.66	24.62

study how different prompts impact the retrieval effectiveness. Particularly, we design six different prompts¹⁰, listed in Table 2, and conduct experiments on TREC deep learning 2019 and 2020 datasets, and MS MARCO passage retrieval dev sub-dataset. We use Llama3-8B-Instruction as the base LLM for PromptReps. The results are listed in Table 3. We also report results of Recall@1000 and other base LLMs in Appendix A.

The results demonstrate that PromptReps can achieve a similar level of retrieval effectiveness as BM25 and surpass LLM2Vec with most of the prompts. The only prompt that does not work well is prompt #4, which does not include the phrase “The word is: “” to force the LLM to generate the representative word as the next token. This is expected because, without this phrase, the first generated token would be a general token such as

¹⁰The prompt in Figure 1 is the prompt number 6 in Table 2.

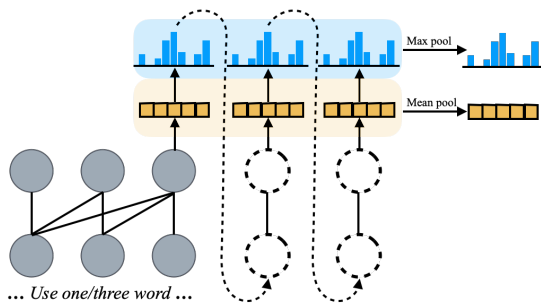


Figure 3: First-word single-representations or Multi-token single-representation.

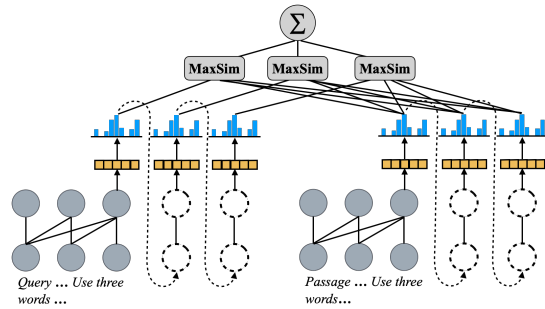


Figure 4: Multi-representations with ColBERT scoring.

466 “The” which is not representative of the input text.

467 Interestingly, our results also show that LLMs
 468 have instruction-following ability in this representa-
 469 tion generation task. For instance, comparing
 470 prompts #1 and #2, the only difference is the phrase
 471 “in a retrieval task”, and the prompt with this
 472 phrase yields higher retrieval effectiveness across
 473 all datasets. Additionally, comparing prompts #1
 474 and #6, the difference is the phrase “Make sure your
 475 word is in lowercase”, which matches our sparse
 476 exact matching method where we first lowercase
 477 the input text. This phrase can further improve the
 478 retrieval effectiveness. Finally, using the adjective
 479 phrase “most important” in the prompt does not
 480 significantly impact the results.

481 5.3 Impact of different LLMs

482 In this section, we explore how different base
 483 LLMs impact PromptReps. For this study, we in-
 484 vestigate five state-of-the-art open-sourced decoder-
 485 only LLMs, covering different model sizes and
 486 models with or without instruction tuning. We use
 487 prompt #6 for all LLMs¹¹ and report MRR@10
 488 scores on the MS MARCO datasets. The results
 489 are illustrated in Figure 2; more detailed results
 490 including on TREC deep learning datasets are re-
 491 ported in Appendix A.

492 The results show that the hybrid retrieval effec-
 493 tiveness of our PromptReps consistently outper-
 494 forms BM25, regardless of which LLM is used,
 495 with the only exception of Mistral-7B-Instruct.
 496 When using the Mistral-7B-Instruct LLM, the
 497 dense-only retriever performs poorly. Surpris-
 498 ingly, implementing PromptReps with Phi-3-mini-
 499 4k-instruct achieved much higher retrieval effec-
 500 tiveness than that of Mistral-7B-Instruct, despite
 501 having far less parameters (3.8B).

¹¹Only the model specific conversational special tokens are changed.

502 Meta-Llama-3 models are generally very effec-
 503 tive for our method. For 8B models, the instruc-
 504 tion-tuned model performs significantly better than the
 505 pretrained-only model, indicating that the instruc-
 506 tion fine-tuning is helpful to further improve our
 507 method. The 70B instruction-tuned model achieved
 508 the best hybrid retrieval results, but the dense-only
 509 and sparse-only retrieval effectiveness are similar
 510 to the 8B instruction-tuned model. These results
 511 agree with the BEIR results presented in Table 1.

512 6 Alternative representations and scoring

513 In the previous sections, we only considered us-
 514 ing the representations (dense and sparse) yielded
 515 from the last token in the prompt for document
 516 retrieval. These representations, in the context of
 517 generative LLMs, are responsible for predicting
 518 the first generated token. We define this setting as
 519 **First-token single-representation**. We have demon-
 520 strated that this simple way of generating represen-
 521 tations is effective for document retrieval; however,
 522 these representations might be sub-optimal. For
 523 example, LLMs use sub-word tokenization algo-
 524 rithms such as SentencePiece (Kudo and Richard-
 525 son, 2018). This tokenization might split a word
 526 into sub-words, meaning that the first generated
 527 token might just be a sub-word. Using the rep-
 528 resentation of the whole word might be a better
 529 representation than the first token representation.
 530 Additionally, previous works in multi-vector dense
 531 retrieval such as ColBERT (Khattab and Zaharia,
 532 2020) demonstrated that using multiple represen-
 533 tations could be beneficial for document retrieval.
 534 How can we use PromptReps to also generate
 535 single-word representations or multiple represen-
 536 tations that can potentially enhance the retrieval
 537 effectiveness? In this section, we explore these
 538 alternative representations.

539 **First-word single-representation** and **Multi-**
 540 **token single-representation**. Instead of just using

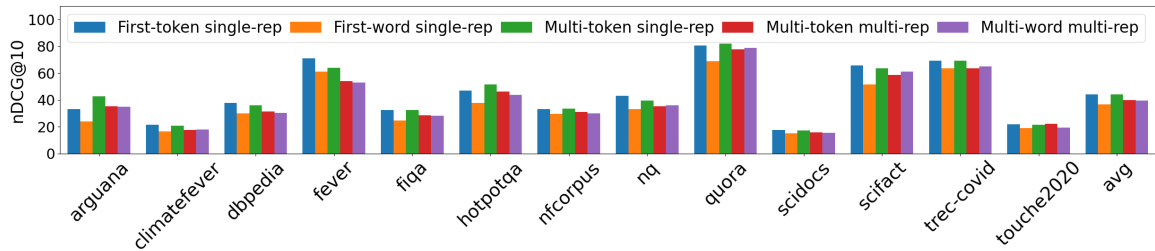


Figure 5: Hybrid retrieval results of different representation methods on BEIR.

541 the representations for the first generated token, 542 these two methods let the LLM finish the genera- 543 tion¹² of the whole word or multiple words, con- 544 trolled by the given prompt (“Use one word” or 545 “Use three words”), as illustrated in Figure 3. The 546 end of generation is detected by the token ‘”’. We 547 then pool all the representations of the generated 548 tokens to form a single dense and sparse representa- 549 tion to index the input text. For the dense representa- 550 tion we use mean pooling and for the sparse representa- 551 tion we use max pooling. Once representa- 552 tions are obtained, the scoring is the same as 553 *First-token single-representation*.

554 *Multi-token multi-representation* and *Multi-* 555 *word multi-representation*. Instead of using a single 556 representation for retrieval, these two methods 557 prompt the LLM to generate multiple words 558 and then index each generated representation sepa- 559 rately. The difference between the two is that 560 *Multi-token multi-representation* keeps all the to- 561 ken representations in the index, while *Multi-word* 562 *multi-representation* first groups tokens into words 563 by using space, and then creates a single repre- 564 sentation for each word by using max pooling 565 for sparse representations and mean pooling for 566 dense representations. During retrieval, we follow 567 the ColBERT scoring method where the relevance 568 score of a document is computed by the sum of the 569 maximum similarity of each query representation 570 against each document representation (Figure 4).

571 Hybrid retrieval results are shown in Figure 5, 572 and full dense and sparse retrieval results in Ap- 573 pendix B. Results show that all the explored 574 methods are able to perform document retrieval. 575 The *First-token single-representation* and *Multi-* 576 *token single-representation* generally perform the 577 best. However, we note that *Multi-token single-* 578 *representation* requires more token generation 579 steps and thus has higher query latency. The *First-* 580 *word single-representation* performs the worst, 581 suggesting that sub-word representations hurt the

582 retrieval performance for single-word generation 583 prompts. On the other hand, multi-representation 584 methods with ColBERT scoring methods do not 585 seem beneficial. Thus, we conclude that the sim- 586 plest *First-token single-representation* is sufficient 587 to represent the input text for document retrieval.

588 7 Conclusion

589 We introduced PromptReps, a simple yet effec- 590 tive method that prompts LLMs to generate dense 591 and sparse representations for zero-shot document 592 retrieval without any further unsupervised or su- 593 pervised training. Our work reveals that modern 594 LLMs are effective text encoders by themselves, 595 and prompt engineering is sufficient to stimulate 596 their text encoding ability.

597 For future works, techniques like few-shot in- 598 context learning (Brown et al., 2020), chain-of- 599 thought prompting (Wei et al., 2022), and auto- 600 prompt optimization methods (Yang et al., 2024; 601 Fernando et al., 2023), which have proven to be ef- 602 fective in text-generation tasks, could potentially be 603 leveraged here to enhance embedding generation.

604 Moreover, it has been shown that the instruction- 605 following ability of LLMs could be transferred 606 to embedding models with synthetic instruction 607 fine-tuning data (Wang et al., 2024b). In our 608 work, we always keep the instruction prompt con- 609 sistent across different IR tasks, which could be 610 sub-optimal. It is interesting to investigate how to 611 customize instructions for PromptReps to generate 612 embeddings specific to different domains, tasks, or 613 even to multi-lingual and cross-lingual IR settings.

614 Finally, our prompting method could be seen as a 615 simple approach to obtaining a better initialization 616 of LLM-based embedding models, which is much 617 more cost-effective than methods requiring further 618 pre-training (BehnamGhader et al., 2024; Li et al., 619 2023). All the previous contrastive pre-training 620 with paired text data and synthetically generated 621 data could be applied on top of our method and 622 could potentially yield improved LLM-based em- 623 bedding models.

¹²We simply use greedy generation.

8 Limitations

PromptReps has higher query latency than other LLM-based dense retrievers if no further optimization is implemented. This limitation comes from two aspects.

First, although the computation of document representations happens offline thus will not affect query latency, the query representations are created online. PromptReps adds extra prompt texts on top of the query text thus has a longer input length – and LLM inference time is proportional to prompt length. However, we believe this limitation can be mitigated by leveraging recent works on prompt compression to compress the fixed prompt tokens into few or even a single latent token (Ge et al., 2024; Cheng et al., 2024).

Secondly, the highest effectiveness for our PromptReps is achieved in the hybrid retrieval setting. Compared to previous works which use dense representations only, the hybrid setting requires both dense and sparse retrieval, thus the extra sparse retrieval introduces extra query latency (and requires additional disk/memory space for the inverted index). However, PromptReps actually only requires a limited query latency overhead if dense and sparse retrieval are implemented in parallel. In our method, obtaining both dense and sparse representations only requires a single LLM forward inference; the only extra computation is the dot product of the dense vector with the token embeddings, which is very fast on GPU. For document search, since we heavily sparsified the sparse representation, in our experiments, our sparse retriever is much faster than BM25, and the bottleneck is the dense retriever. Since the dense and sparse search could be run in parallel and the hybrid operation is a simple linear interpolation of both rankings (very fast on CPU), the query latency of the hybrid process only depends on the dense retrieval latency, and it is thus very close to previous methods.

9 Ethical considerations

In our experiments, we use PromptReps coupled with LLMs with a large number of parameters (up to 70B in our experiments) to encode the BEIR and MS MARCO datasets, which contain millions of documents. Although no LLM training was conducted, we are aware that our experiments might still have consumed significant energy, thus contributing to CO2 emissions (Scells et al., 2022) and water consumption (Zuccon et al., 2023).

In addition, since we leverage LLMs in a black-box manner and LLMs’ generation might contain biases (Gallegos et al., 2024), the representations generated by LLMs may be biased towards certain contents or topics. Future work could consider how to mitigate biases in PromptReps via prompt engineering.

References

- Marah Abdin et al. 2024. [Phi-3 technical report: A highly capable language model locally on your phone](#). *Preprint*, arXiv:2404.14219.
- AI@Meta. 2024. [Llama 3 model card](#).
- Payal Bajaj, Daniel Campos, Nick Craswell, Li Deng, Jianfeng Gao, Xiaodong Liu, Rangan Majumder, Andrew McNamara, Bhaskar Mitra, Tri Nguyen, Mir Rosenberg, Xia Song, Alina Stoica, Saurabh Tiwary, and Tong Wang. 2018. [Ms marco: A human generated machine reading comprehension dataset](#). *Preprint*, arXiv:1611.09268.
- Elias Bassani and Luca Romelli. 2022. [ranx.fuse: A python library for metasearch](#). In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management, CIKM ’22*, pages 4808–4812, New York, NY, USA. Association for Computing Machinery.
- Parishad BehnamGhader, Vaibhav Adlakha, Marius Mosbach, Dzmitry Bahdanau, Nicolas Chapados, and Siva Reddy. 2024. [Llm2vec: Large language models are secretly powerful text encoders](#). *Preprint*, arXiv:2404.05961.
- Steven Bird and Edward Loper. 2004. [NLTK: The natural language toolkit](#). In *Proceedings of the ACL Interactive Poster and Demonstration Sessions*, pages 214–217, Barcelona, Spain. Association for Computational Linguistics.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Xin Cheng, Xun Wang, Xingxing Zhang, Tao Ge, Si-Qing Chen, Furu Wei, Huishuai Zhang, and Dongyan Zhao. 2024. [xrag: Extreme context compression for retrieval-augmented generation with one token](#). *Preprint*, arXiv:2405.13792.

728	Nick Craswell, Bhaskar Mitra, Emine Yilmaz, Daniel Campos, and Ellen M. Voorhees. 2020. Overview of the trec 2019 deep learning track . <i>Preprint</i> , arXiv:2003.07820.	783
729		784
730		785
731		786
732	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.	787
733		788
734		789
735		790
736		791
737		792
738		793
739		794
740		795
741	Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. 2024. The faiss library . <i>Preprint</i> , arXiv:2401.08281.	796
742		797
743		798
744		799
745	Chrisantha Fernando, Dylan Banarse, Henryk Michalewski, Simon Osindero, and Tim Rocktäschel. 2023. Promptbreeder: Self-referential self-improvement via prompt evolution . <i>Preprint</i> , arXiv:2309.16797.	800
746		801
747		802
748		803
749		804
750	Thibault Formal, Benjamin Piwowarski, and Stéphane Clinchant. 2021. Splade: Sparse lexical and expansion model for first stage ranking . In <i>Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '21</i> , pages 2288–2292, New York, NY, USA. Association for Computing Machinery.	805
751		806
752		807
753		808
754		809
755		810
756		811
757	Isabel O Gallegos, Ryan A Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K Ahmed. 2024. Bias and fairness in large language models: A survey. <i>Computational Linguistics</i> , pages 1–79.	812
758		813
759		814
760		815
761		816
762	Luyu Gao, Xueguang Ma, Jimmy Lin, and Jamie Callan. 2023a. Precise zero-shot dense retrieval without relevance labels . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1762–1777, Toronto, Canada. Association for Computational Linguistics.	817
763		818
764		819
765		820
766		821
767		822
768		823
769	Luyu Gao, Xueguang Ma, Jimmy Lin, and Jamie Callan. 2023b. Tevatron: An efficient and flexible toolkit for neural retrieval . In <i>Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '23</i> , pages 3120–3124, New York, NY, USA. Association for Computing Machinery.	824
770		825
771		826
772		827
773		828
774		829
775		830
776	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.	831
777		832
778		833
779		834
780		835
781		836
782		837
		838
		839
		840
	Tao Ge, Hu Jing, Lei Wang, Xun Wang, Si-Qing Chen, and Furu Wei. 2024. In-context autoencoder for context compression in a large language model . In <i>The Twelfth International Conference on Learning Representations</i> .	841
		842
		843
		844
		845
		846
		847
		848
		849
		850
		851
		852
		853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900
		901
		902
		903
		904
		905
		906
		907
		908
		909
		910
		911
		912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

841	Ekaterina Khramtsova, Shengyao Zhuang, Mahsa Bak- tashmotlagh, and Guido Zuccon. 2024. Leveraging llms for unsupervised dense retriever ranking. <i>arXiv preprint arXiv:2402.04853</i> .	Simon Lupart, Thibault Formal, and Stéphane Clinchant. 2022. <i>Ms-shift: An analysis of ms marco distribution shifts on neural retrieval</i> . In <i>European Conference on Information Retrieval</i> .	898 899 900 901
845	Taku Kudo and John Richardson. 2018. <i>SentencePiece: A simple and language independent subword tok- enizer and detokenizer for neural text processing</i> . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 66–71, Brussels, Belgium. Association for Computational Linguistics.	Xueguang Ma, Liang Wang, Nan Yang, Furu Wei, and Jimmy Lin. 2023a. <i>Fine-tuning llama for multi-stage text retrieval</i> . <i>Preprint</i> , arXiv:2310.08319.	902 903 904
848	Jinhyuk Lee, Zhuyun Dai, Xiaoqi Ren, Blair Chen, Daniel Cer, Jeremy R. Cole, Kai Hui, Michael Boratto, Rajvi Kapadia, Wen Ding, Yi Luan, Sai Meher Karthik Duddu, Gustavo Hernandez Abrego, Weiqiang Shi, Nithi Gupta, Aditya Kusupati, Pra- teek Jain, Siddhartha Reddy Jonnalagadda, Ming- Wei Chang, and Iftexhar Naim. 2024. <i>Gecko: Ver- satile text embeddings distilled from large language models</i> . <i>Preprint</i> , arXiv:2403.20327.	Xueguang Ma, Xinyu Zhang, Ronak Pradeep, and Jimmy Lin. 2023b. <i>Zero-shot listwise document reranking with a large language model</i> . <i>Preprint</i> , arXiv:2305.02156.	905 906 907 908
852	Yibin Lei, Di Wu, Tianyi Zhou, Tao Shen, Yu Cao, Chongyang Tao, and Andrew Yates. 2024. <i>Meta-task prompting elicits embedding from large language models</i> . <i>Preprint</i> , arXiv:2402.18458.	Antonio Mallia, Omar Khattab, Torsten Suel, and Nicola Tonello. 2021. <i>Learning passage impacts for in- verted indexes</i> . In <i>Proceedings of the 44th Inter- national ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '21</i> , pages 1723–1727, New York, NY, USA. Association for Computing Machinery.	909 910 911 912 913 914 915
861	Chaofan Li, Zheng Liu, Shitao Xiao, and Yingxia Shao. 2023. <i>Making large language models a better founda- tion for dense retrieval</i> . <i>Preprint</i> , arXiv:2312.15503.	Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. 2016. Ms marco: A human-generated machine read- ing comprehension dataset.	916 917 918 919
866	Jiayi Liao, Sihang Li, Zhengyi Yang, Jiancan Wu, Yancheng Yuan, Xiang Wang, and Xiangnan He. 2024. <i>Llara: Large language-recommendation as- sistant</i> . <i>Preprint</i> , arXiv:2312.02445.	OpenAI. 2023. GPT-4 technical report. <i>arXiv:2303.08774</i> .	920 921
872	Jimmy Lin and Xueguang Ma. 2021. <i>A few brief notes on deepimpact, coil, and a conceptual frame- work for information retrieval techniques</i> . <i>Preprint</i> , arXiv:2106.14807.	Ronak Pradeep, Sahel Sharifymoghaddam, and Jimmy Lin. 2023. <i>Rankvicuna: Zero-shot listwise document reranking with open-source large language models</i> . <i>Preprint</i> , arXiv:2309.15088.	922 923 924 925
876	Jimmy Lin, Xueguang Ma, Sheng-Chieh Lin, Jheng- Hong Yang, Ronak Pradeep, and Rodrigo Nogueira. 2021. <i>Pyserini: A python toolkit for reproducible information retrieval research with sparse and dense representations</i> . In <i>Proceedings of the 44th Inter- national ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '21</i> , pages 2356–2362, New York, NY, USA. Association for Computing Machinery.	Zhen Qin, Rolf Jagerman, Kai Hui, Honglei Zhuang, Junru Wu, Le Yan, Jiaming Shen, Tianqi Liu, Jialu Liu, Donald Metzler, Xuanhui Wang, and Michael Bendersky. 2024. <i>Large language models are effec- tive text rankers with pairwise ranking prompting</i> . <i>Preprint</i> , arXiv:2306.17563.	926 927 928 929 930 931
885	Sheng-Chieh Lin, Akari Asai, Minghan Li, Barlas Oguz, Jimmy Lin, Yashar Mehdad, Wen-tau Yih, and Xilun Chen. 2023. <i>How to train your dragon: Diverse augmentation towards generalizable dense retrieval</i> . <i>arXiv preprint arXiv:2302.07452</i> .	Ruiyang Ren, Yingqi Qu, Jing Liu, Xin Zhao, Qifei Wu, Yuchen Ding, Hua Wu, Haifeng Wang, and Ji- Rong Wen. 2023. <i>A thorough examination on zero- shot dense retrieval</i> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 15783–15796, Singapore. Association for Computa- tional Linguistics.	932 933 934 935 936 937 938
890	Yiheng Liu, Tianle Han, Siyuan Ma, Jiayue Zhang, Yuanyuan Yang, Jiaming Tian, Hao He, Antong Li, Mengshen He, Zhengliang Liu, Zihao Wu, Lin Zhao, Dajiang Zhu, Xiang Li, Ning Qiang, Dingang Shen, Tianming Liu, and Bao Ge. 2023. <i>Summary of chatgpt-related research and perspective towards the future of large language models</i> . <i>Meta-Radiology</i> , 1(2):100017.	Devendra Sachan, Mike Lewis, Mandar Joshi, Armen Aghajanyan, Wen-tau Yih, Joelle Pineau, and Luke Zettlemoyer. 2022. <i>Improving passage retrieval with zero-shot question generation</i> . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natu- ral Language Processing</i> , pages 3781–3797, Abu Dhabi, United Arab Emirates. Association for Com- putational Linguistics.	939 940 941 942 943 944 945 946
897		Harrison Scells, Shengyao Zhuang, and Guido Zuccon. 2022. <i>Reduce, reuse, recycle: Green information retrieval research</i> . In <i>Proceedings of the 45th Inter- national ACM SIGIR Conference on Research and Development in Information Retrieval</i> , pages 2825– 2837.	947 948 949 950 951 952

953	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 <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 14918–14937, Singapore. Association for Computational Linguistics.	
961	Nandan Thakur, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, and Iryna Gurevych. 2021. BEIR: A heterogeneous benchmark for zero-shot evaluation of information retrieval models . In <i>Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)</i> .	
967	Hugo Touvron et al. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv:2307.09288</i> .	
969	Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. 2024a. Text embeddings by weakly-supervised contrastive pre-training . <i>Preprint</i> , arXiv:2212.03533.	
974	Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. 2024b. Improving text embeddings with large language models . <i>Preprint</i> , arXiv:2401.00368.	
978	Shuai Wang, Shengyao Zhuang, and Guido Zuccon. 2021a. Bert-based dense retrievers require interpolation with bm25 for effective passage retrieval. In <i>Proceedings of the 2021 ACM SIGIR international conference on theory of information retrieval</i> , pages 317–324.	
984	Shuai Wang, Shengyao Zhuang, and Guido Zuccon. 2021b. Bert-based dense retrievers require interpolation with bm25 for effective passage retrieval . In <i>Proceedings of the 2021 ACM SIGIR International Conference on Theory of Information Retrieval, ICTIR '21</i> , pages 317–324, New York, NY, USA. Association for Computing Machinery.	
991	Jason Wei, Xuezhong Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models . In <i>Advances in Neural Information Processing Systems</i> , volume 35, pages 24824–24837. Curran Associates, Inc.	
998	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 38–45, Online. Association for Computational Linguistics.	
	Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang, Jialin Liu, Paul N. Bennett, Junaid Ahmed, and Arnold Overwijk. 2021. Approximate nearest neighbor negative contrastive learning for dense text retrieval . In <i>International Conference on Learning Representations</i> .	1010 1011 1012 1013 1014 1015
	Chengrun Yang, Xuezhong Wang, Yifeng Lu, Hanxiao Liu, Quoc V Le, Denny Zhou, and Xinyun Chen. 2024. Large language models as optimizers . In <i>The Twelfth International Conference on Learning Representations</i> .	1016 1017 1018 1019 1020
	Bowen Zhang, Kehua Chang, and Chunping Li. 2024. Simple techniques for enhancing sentence embeddings in generative language models . <i>Preprint</i> , arXiv:2404.03921.	1021 1022 1023 1024
	Shengyao Zhuang, Bing Liu, Bevan Koopman, and Guido Zuccon. 2023a. Open-source large language models are strong zero-shot query likelihood models for document ranking . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 8807–8817, Singapore. Association for Computational Linguistics.	1025 1026 1027 1028 1029 1030 1031
	Shengyao Zhuang, Honglei Zhuang, Bevan Koopman, and Guido Zuccon. 2023b. A setwise approach for effective and highly efficient zero-shot ranking with large language models . <i>Preprint</i> , arXiv:2310.09497.	1032 1033 1034 1035
	Shengyao Zhuang and Guido Zuccon. 2021a. Dealing with typos for BERT-based passage retrieval and ranking. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 2836–2842, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	1036 1037 1038 1039 1040 1041 1042
	Shengyao Zhuang and Guido Zuccon. 2021b. Fast passage re-ranking with contextualized exact term matching and efficient passage expansion . <i>Preprint</i> , arXiv:2108.08513.	1043 1044 1045 1046
	Shengyao Zhuang and Guido Zuccon. 2021c. Tilde: Term independent likelihood model for passage re-ranking . In <i>Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '21</i> , pages 1483–1492, New York, NY, USA. Association for Computing Machinery.	1047 1048 1049 1050 1051 1052 1053
	Shengyao Zhuang and Guido Zuccon. 2022. Characterbert and self-teaching for improving the robustness of dense retrievers on queries with typos. SIGIR '22, pages 1444–1454, New York, NY, USA. Association for Computing Machinery.	1054 1055 1056 1057 1058
	Guido Zuccon, Harris Scells, and Shengyao Zhuang. 2023. Beyond co2 emissions: The overlooked impact of water consumption of information retrieval models. In <i>Proceedings of the 2023 ACM SIGIR International Conference on Theory of Information Retrieval</i> , pages 283–289.	1059 1060 1061 1062 1063 1064

1065 **A Full results on TREC deep learning**
1066 **and MS MARCO**

1067 In Table 4 we present the full results we abstained
1068 on TREC deep learning datasets and MS MARCO
1069 passage retrieval dataset. The prompt ID is refer to
1070 Table 2.

1071 **B Full results of different representation**
1072 **methods**

1073 In Table 5 we present the full results of different
1074 representation and scoring methods discussed in
1075 Section 6.

Table 4: TREC deep learning and MS MARCO performance of different prompts and LLMs.

Prompt ID	Methods	DL2019		DL2020		MS MARCO Dev	
		nDCG@10	Recall@1000	nDCG@10	Recall@1000	MRR@10	Recall@1000
-	BM25	49.73	74.50	48.76	80.31	18.75	85.73
-	LLM2Vec	-	-	-	-	13.61	94.70
Phi-3-mini-4k-instruct (3.8B)							
1	Dense	46.78	70.10	42.84	67.60	15.45	82.68
2	Dense	34.64	55.15	30.62	50.47	10.85	66.04
3	Dense	49.62	75.79	43.21	71.87	15.78	86.24
4	Dense	39.12	60.77	28.33	57.20	9.26	72.31
5	Dense	43.94	72.51	39.00	70.57	13.50	83.08
6	Dense	40.77	62.05	37.20	58.39	14.64	79.29
1	Sparse	41.51	69.56	40.95	69.70	16.89	84.72
2	Sparse	40.67	60.59	39.36	61.58	16.38	75.43
3	Sparse	42.28	74.33	40.72	72.14	18.16	87.09
4	Sparse	38.05	65.15	34.33	64.13	14.75	78.59
5	Sparse	40.68	70.84	39.26	69.02	16.04	84.41
6	Sparse	41.98	71.20	41.99	69.66	18.19	86.55
1	Hybrid	53.04	79.99	52.76	77.64	21.61	92.21
2	Hybrid	50.51	69.75	43.92	66.47	19.22	81.35
3	Hybrid	55.53	81.68	51.35	79.49	21.76	93.53
4	Hybrid	48.53	76.29	40.37	73.64	18.23	87.18
5	Hybrid	52.08	80.16	50.52	79.30	20.30	92.37
6	Hybrid	51.10	75.84	49.24	73.98	22.06	91.41
Meta-Llama-3-8B-Instruct							
1	Dense	49.26	73.03	40.28	68.77	16.26	81.96
2	Dense	43.32	64.77	31.60	61.35	12.52	73.89
3	Dense	49.20	71.69	43.90	69.96	17.49	84.50
4	Dense	0.00	0.00	0.00	0.00	0.00	0.04
5	Dense	47.19	72.00	40.17	66.71	16.02	82.56
6	Dense	50.62	73.01	43.81	68.39	17.54	82.91
1	Sparse	41.77	67.28	44.81	71.36	20.12	85.71
2	Sparse	39.90	66.00	43.10	69.08	19.13	83.74
3	Sparse	43.50	66.74	44.87	72.93	20.42	85.14
4	Sparse	21.77	41.94	20.49	50.51	7.22	56.35
5	Sparse	42.18	67.18	44.17	71.94	19.78	85.37
6	Sparse	42.25	66.58	45.60	72.82	20.85	85.57
1	Hybrid	53.67	83.52	54.35	78.42	23.68	92.84
2	Hybrid	50.65	80.31	49.25	76.64	21.76	90.12
3	Hybrid	55.64	81.90	53.83	79.15	23.86	92.99
4	Hybrid	13.47	37.81	11.81	45.22	5.06	50.50
5	Hybrid	54.16	82.06	52.06	78.70	23.25	92.77
6	Hybrid	55.58	83.44	56.66	79.14	24.62	93.11
Meta-Llama-3-8B							
6	Dense	43.90	67.38	35.50	63.34	14.67	79.61
6	Sparse	38.41	64.83	43.34	67.57	18.82	82.63
6	Hybrid	51.13	77.07	46.34	75.42	22.31	90.87
Mistral-7B-Instruct-v0.2							
6	Dense	13.96	27.26	16.77	26.69	5.61	40.27
6	Sparse	39.84	58.05	37.29	63.53	15.62	77.55
6	Hybrid	32.58	57.00	32.95	63.12	13.18	77.98
Meta-Llama-3-70B-Instruct							
6	Dense	51.95	77.30	45.01	73.66	17.76	85.65
6	Sparse	44.07	68.60	44.14	70.99	20.70	86.42
6	Hybrid	58.39	86.22	59.17	81.57	25.66	93.75
6	+ BM25	63.18	88.56	62.55	86.28	27.63	95.83

Table 5: Full results of different representation and scoring methods on BEIR.

Dataset	First token single rep			First-word single rep			Multi token single rep			Multi-token multi-rep			Multi-word multi-rep		
	Dense	Sparse	Hybrid	Dense	Sparse	Hybrid	Dense	Sparse	Hybrid	Dense	Sparse	Hybrid	Dense	Sparse	Hybrid
arguana	29.70	22.85	33.32	20.54	24.59	23.80	41.78	24.46	42.61	36.69	23.03	35.19	36.47	24.13	34.96
climatefever	19.92	9.98	21.38	13.88	11.28	16.67	22.19	9.29	20.90	19.40	6.72	17.56	18.75	8.10	18.09
dbpedia	31.53	28.84	37.71	22.71	28.70	30.08	31.83	26.33	36.03	27.46	18.18	31.35	24.50	21.66	30.32
fever	56.28	52.35	71.11	40.97	57.10	61.20	50.49	51.36	64.13	44.53	30.31	54.04	38.81	37.95	52.97
fiqa	27.11	20.33	32.40	17.61	19.60	24.74	28.94	20.73	32.44	25.26	19.41	28.38	26.28	19.50	28.30
hotpotqa	19.64	44.75	47.05	10.35	46.25	37.79	29.94	46.50	51.50	23.80	39.39	46.38	21.93	40.25	43.68
nfcopus	29.56	28.18	32.98	21.98	29.15	29.49	28.97	28.65	33.65	25.68	25.39	31.20	22.95	25.37	30.05
nq	34.43	29.55	43.14	22.83	29.25	33.18	35.09	25.88	39.36	31.36	23.28	35.38	30.55	22.78	35.95
quora	72.55	68.27	80.45	51.68	66.68	69.09	78.90	66.83	81.93	72.71	63.21	77.91	73.57	63.38	78.77
scidocs	18.51	11.57	17.59	12.73	12.05	14.97	18.12	11.83	17.39	16.13	11.08	15.68	15.85	11.40	15.44
scifact	52.68	58.48	65.71	26.66	58.75	51.59	52.55	59.32	63.75	45.53	54.23	58.53	47.05	53.21	61.18
trec-covid	59.52	54.59	69.25	51.00	55.04	63.53	63.28	51.73	69.16	60.97	49.88	63.54	61.17	46.24	65.10
touche2020	14.85	18.47	21.65	12.23	21.58	19.13	15.59	19.24	21.44	15.86	17.81	22.10	15.41	18.09	19.26
avg	35.87	34.48	44.13	25.01	35.39	36.56	38.28	34.01	44.18	34.26	29.38	39.79	33.33	30.16	39.54