

# Joint Inference of Retrieval and Generation for Passage Re-ranking

Anonymous ACL submission

## Abstract

Passage retrieval is a crucial component of modern open-domain question answering (QA) systems, providing information for downstream QA components to generate accurate and transparent answers. In this study we focus on passage re-ranking, proposing a simple yet effective method, *Joint Passage Re-ranking* (JPR), that optimizes the mutual information between query and passage distributions, integrating both cross-encoders and generative models in the re-ranking process. Experimental results demonstrate that JPR outperforms conventional re-rankers and language model scorers in both open-domain QA retrieval settings and diverse retrieval benchmarks under zero-shot settings.<sup>1</sup>

## 1 Introduction

Passage retrieval is a crucial component in open-domain question answering (QA) (Chen and Yih, 2020), a task that requires answering questions from a wide range of domains and could be applied in systems that fulfill user’s information needs (Voorhees et al., 1999). Retrieval offers downstream QA systems grounding information, which not only improves accuracy in a lot of cases but also provides transparency to how systems generate answers, similar to how articles provide references and citations, such that model hallucinations can be checked with ease. Furthermore, the set of documents to be retrieved from, or knowledge base, can be quickly updated with new documents and knowledge such that models can adapt to temporal changes, and do not need to be continuously re-trained nor require online training paradigms for continual learning.

Early retrieval methods are typically based on term-matching, such as BM25 (Robertson et al., 2009) or TF-IDF (Salton et al., 1975). Such methods, called sparse retrievers, perform keyword

matching efficiently with an inverted index to find relevant contexts. Sparse retrievers often achieve reasonable performance while being computationally efficient and does not require training, but are shown to have limited abilities beyond lexical matching.

Recently, dense retrievers that encode text with continuous embeddings have been heavily studied and utilized in contemporary QA systems, often outperforming their sparse counterparts on high resource evaluation settings (Karpukhin et al., 2020). There are a few drawbacks however, such as higher computational demands during both training and inference, inability to handle large contexts (Luan et al., 2021), and difficulty in generalizing to new domains especially those with limited data (Reddy et al., 2021). Hybrid methods have been explored to get the best of both worlds, generally utilizing an efficient sparse method to retrieve a larger number of possibly relevant contexts, and then perform passage re-ranking with a more computationally-intensive dense model for refined scoring (Nogueira and Cho, 2019).

In this work, we focus on passage re-ranking and explore the use of generative models alongside conventional re-rankers. Previous work have explored pre-trained language models (LM) as the re-ranking scorer (Sachan et al., 2022), however we find that it underperforms conventional re-rankers for both supervised and zero-shot settings. Starting from maximizing mutual information for inference, which had been explored in several other tasks such as dialogue generation (Li et al., 2016), machine translation (Li and Jurafsky, 2016), and QA (Tang et al., 2017; Luo et al., 2022), we show how a small generative model can be effectively used with conventional cross-encoding re-rankers for improved performance. Experiments on a supervised setting for open-domain QA retrieval and a zero-shot setting across a suite of diverse retrieval benchmarks validate our approach. Our contributions can be

<sup>1</sup>Source code is available. See Appx. A.

summarized as follows:

- We propose *Joint Passage Re-ranking* (JPR), a method utilizing both a cross-encoder and a generative model in the retrieval re-ranking process, optimizing the mutual information between query and document distributions.
- We demonstrate that JPR outperforms conventional re-rankers and generative scorers in open-domain QA retrieval evaluation and diverse zero-shot retrieval datasets.

## 2 Joint Passage Re-ranking (JPR)

Consider the two distributions  $p(\mathbf{x})$  and  $p(\mathbf{z})$  over all queries  $\mathbf{x} \in \mathcal{X}$  and all passages  $\mathbf{z} \in \mathcal{Z}$ . The conditional distributions  $p(\mathbf{z}|\mathbf{x})$  and  $p(\mathbf{x}|\mathbf{z})$  can be used to infer one domain based on the other. The joint distribution  $p(\mathbf{x}, \mathbf{z})$  characterizes the combined structure of both domains, where  $p(\mathbf{x}, \mathbf{z}) = p(\mathbf{x})p(\mathbf{z}|\mathbf{x}) = p(\mathbf{z})p(\mathbf{x}|\mathbf{z})$ .

Here  $p_\phi(\mathbf{z}|\mathbf{x})$  defines a passage retrieval model, which we parametrize by  $\phi$ , generally trained with maximum likelihood estimation (MLE):  $\mathcal{L}_{\text{retrieval}}(\phi) \triangleq -\mathbb{E}_{\mathbf{x}, \mathbf{z} \sim p(\mathbf{x}, \mathbf{z})} [\log p_\phi(\mathbf{z}|\mathbf{x})]$ . During inference, finding the most probable relevant passage can be written as:

$$\hat{z} = \arg \max_z \log p_\phi(\mathbf{z}|\mathbf{x}). \quad (1)$$

Since we focus on passage re-ranking, we treat  $p_\phi(\mathbf{z}|\mathbf{x})$  in Eq. 1 as re-ranking scores.

### 2.1 Inference with Maximum Mutual Information

In our work, we approach inference by finding passage that maximizes the (pairwise) mutual information (MMI) between both domains instead:

$$\hat{z} = \arg \max_z \left( \log p(\mathbf{z}|\mathbf{x}) - \log p(\mathbf{z}) \right). \quad (2)$$

We see that maximizing MI adds a penalizing term compared to MLE in Eq. 1, which avoids favoring passages that unconditionally have a higher probability, and biases the model towards those that are specific to the given query. A hyperparameter  $\lambda$  is used to control the regularization term. Using Bayes' theorem, we can rewrite Eq. 2 as:

$$\begin{aligned} \hat{z} &= \arg \max_z \left( \log p(\mathbf{z}|\mathbf{x}) - \lambda \log p(\mathbf{z}) \right) \\ &= \arg \max_z \left( (1 - \lambda) \log p(\mathbf{z}|\mathbf{x}) + \lambda \log p(\mathbf{x}|\mathbf{z}) \right). \end{aligned} \quad (3)$$

The MMI objective is equivalent to the convex combination of the terms  $\log p(\mathbf{z}|\mathbf{x})$  and  $\log p(\mathbf{x}|\mathbf{z})$ . Notice that the latter term can be viewed as a condi-

tional generation model that gives the probability of generating a query given a passage. We denote the generative model by  $p_\theta(\mathbf{x}|\mathbf{z})$  with parameters  $\theta$ . This term was previously explored as the sole inference objective in Sachan et al. (2022), in which an LM was used as a question generator for rescoring. Instead of using either the retrieval model or the generative model only, as explored in prior work, Eq. 3 provides a simple way to use both models jointly for inference, which we refer to as *Joint Passage Re-ranking* (JPR).

### 2.2 Joint Fine-tuning

A straightforward way to obtain the two models that can be used for the aforementioned MMI-based inference is to train both models using MLE separately. The retrieval model can be trained with  $\mathcal{L}_{\text{retrieval}}(\phi)$ , while the generative model can be trained with a simple LM loss  $\mathcal{L}_{\text{generation}}(\theta)$ .

However, the terms in Eq. 3 are derived when the distributions are matched, that is, when  $p(\mathbf{x})p_\phi(\mathbf{z}|\mathbf{x}) = p(\mathbf{z})p_\theta(\mathbf{x}|\mathbf{z})$ . When the two models are optimized independently, we cannot ensure that this holds. We therefore attempt to enforce this constraint during fine-tuning. Similar to previous work on dual supervised learning, we approach this by adding a regularization term, defined as the symmetric KL divergence between the two distributions:  $\mathcal{L}_{\text{match}}(\phi, \theta) \triangleq D_{\text{sym-KL}}(p_\phi(\mathbf{x}, \mathbf{z}) || p_\theta(\mathbf{x}, \mathbf{z}))$ . The joint training objective is obtained by combining all three losses:  $\mathcal{L}(\phi, \theta) \triangleq \mathcal{L}_{\text{retrieval}} + \mathcal{L}_{\text{generation}} + \alpha \mathcal{L}_{\text{match}}$ , where  $\alpha$  is a regularization hyperparameter.

## 3 Experiments

### 3.1 Open-Domain QA Retrieval

#### 3.1.1 Data

First, we evaluate on two standard open-domain QA retrieval benchmark datasets: Natural Questions (NQ; Kwiatkowski et al., 2019) and TriviaQA (Joshi et al., 2017). Wikipedia passages used in DPR (Karpukhin et al., 2020) were used in these experiments, which consists of 21M 100-word passages from the English Wikipedia dump of Dec. 20, 2018 (Lee et al., 2019). Additional dataset information can be found in Appx. B.

#### 3.1.2 Setup and Baselines

We adopt the setting from prior work using standard dataset splits, retrieving the top 100 passages for re-ranking. We use Pyserini (Lin et al., 2021) for

Re-ranking Method	Cross-Encoder? $\log p_\phi(z x)$	Generative? $\log p_\theta(x z)$	Natural Questions			TriviaQA		
			Top-1	Top-5	Top-10	Top-1	Top-5	Top-10
BM25	$\times$	$\times$	22.1	43.8	54.5	46.3	66.3	71.7
BERT-FT	$\checkmark$	$\times$	49.4	66.4	71.4	66.7	77.6	80.2
T5-FT	$\times$	$\checkmark$	34.3	59.6	66.7	56.8	74.1	78.0
UPR (T0-3B)	$\times$	$\checkmark$	36.8	61.6	68.2	57.7	75.4	78.5
JPR	$\checkmark$	$\checkmark$	51.0	<b>68.0</b>	<u>72.3</u>	68.3	78.3	<b>80.5</b>
JPR-FT	$\checkmark$	$\checkmark$	<u>51.4</u>	67.5	71.9	<b>69.2</b>	<b>78.5</b>	<b>80.5</b>
UPR (LLaMA-33B)	$\times$	$\checkmark$	35.0	61.5	69.0	57.2	76.7	79.5
JPR (LLaMA-33B)	$\checkmark$	$\checkmark$	48.2	66.9	71.5	<u>70.1</u>	<u>79.3</u>	<u>80.8</u>

Table 1: Top- $K$  retrieval accuracy (%) on the Natural Questions and TriviaQA test sets. All non-BM25 methods re-rank the top-100 passages retrieved by BM25. Best overall are in **bold** while best non-LLM are underlined.

BM25 as the initial retriever, with default Lucene parameters of  $k = 0.9$  and  $b = 0.4$ . We report top- $K$  retrieval accuracy, the standard metric.

We compare JPR against several baselines: 1) cross-encoding re-ranker (BERT-FT), a fine-tuned BERT-based (Devlin et al., 2019) re-ranker, running inference with Eq. 1; 2) generative re-ranker (T5-FT), a fine-tuned T5 conditional generation model (Raffel et al., 2020) with the second term of Eq. 3 as inference objective; and 3) UPR (Sachan et al., 2022), a generator-only re-ranker using the larger pre-trained T0-3B model (Sanh et al., 2022).

For our approach, we report one setting with joint inference (JPR), and another with joint fine-tuning followed by the MMI-based inference (JPR-FT). Joint inference uses the separately fine-tuned retrieval re-ranker and generative re-ranker described above directly. For joint fine-tuning, we bootstrap with the two models, and further fine-tune with our proposed objective to match the retriever and generator distributions.  $\lambda$  and  $\alpha$  are chosen by performance on the development set. Additional details can be found in Appx. C.

Furthermore, we aim to explore the effects of scaling generative re-rankers up. We experiment with a large language model (LLM), the 33B-parameter LLaMA (Touvron et al., 2023), as our generative re-ranker for both UPR and JPR.

### 3.1.3 Results and Discussion

Open-domain QA retrieval results are shown in Table 1. Using the conventional cross-encoder BERT-FT on initial BM25 results yields decent improvements. UPR, not fine-tuned but being much larger, significantly underperforms BERT-FT. The fine-tuned generator T5-FT, 15 $\times$  smaller than the T0-3B model in UPR, nearly matches the performance of UPR. When using JPR, which corresponds to scoring with Eq. 3 using the re-ranker BERT-FT

and the generator T5-FT, surpasses all baselines. The generator, although used by itself underperforms BERT-FT, boosts performance especially for the top retrieved passages. Matching distributions (JPR-FT) by fine-tuning for a small amount of steps further improves performance, albeit more modestly. For LLM generative re-ranking, despite being multitudes larger, LLaMA-33B surprisingly underperforms against T5-FT and T0-3B on NQ for both UPR and JPR, however on TriviaQA JPR with LLaMA-33B achieves best overall results. Appx. D shows further results for different model pairs.

## 3.2 Zero-Shot Retrieval

### 3.2.1 Data

We further evaluate in a transfer learning setting on BEIR (Thakur et al., 2021), a commonly used benchmark consisting of a suite of information retrieval datasets that span multiple tasks and domains. Datasets in the benchmark contain queries and passages of a variety of styles and lengths, and no training data is provided, making it considerably difficult for models to perform well across all datasets. See Appx. E for more details.

### 3.2.2 Setup and Baselines

We follow BEIR’s zero-shot evaluation on all tasks, using MS MARCO (Nguyen et al., 2017) as training data. Pyserini is used for BM25 to retrieve 100 passages, with default parameters and indexing title and passage as separate fields<sup>23</sup>. The Normalized Cumulative Discount Gain (nDCG@ $K$ ) (Wang et al., 2013) is used for evaluation, with  $K = 10$ , computed by the official TREC evaluation tool (Van Gysel and de Rijke, 2018).

We compare against the three baselines used

<sup>2</sup>Pyserini reproductions for BEIR can be found at <https://castorini.github.io/pyserini/2cr/beir.html>.

<sup>3</sup>We follow BEIR and retrieve 100, which is more practical.

Dataset	BM25	Re-ranking Method					
		BERT- FT	T5-FT	UPR	JPR	UPR (LLM)	JPR (LLM)
TREC-DL 2019	50.8	<u>74.9</u>	<u>65.6</u>	-	<u>75.0</u>	-	-
TREC-COVID	65.6	75.7	75.7	76.5	78.2	76.5	77.2
NFCorpus	32.6	35.0	33.2	34.8	35.3	33.5	35.7
NQ	32.9	53.3	43.8	44.5	52.1	45.3	54.0
HotpotQA	60.3	70.7	68.5	70.9	72.4	72.3	72.1
FiQA-2018	23.6	34.7	35.7	42.0	38.5	40.3	36.6
ArguAna	<i>41.4</i>	<i>41.8</i>	50.2	<i>50.9</i>	49.3	28.5	43.3
Touché-2020	36.7	27.1	25.0	21.0	26.8	18.5	25.7
CQADupStack	29.9	37.1	37.7	40.2	39.7	42.9	39.0
Quora	78.9	82.5	81.2	83.6	84.8	84.4	84.1
DBPedia	31.3	40.9	34.6	35.5	40.5	35.1	41.6
SCIDOCS	15.8	16.6	16.9	17.6	18.3	18.1	17.1
FEVER	75.3	81.8	75.7	61.3	82.5	62.5	79.7
Climate-FEVER	21.3	25.3	18.4	14.6	25.2	11.2	24.9
SciFact	66.5	68.8	69.3	70.4	72.7	65.7	70.3
Average	43.7	49.4	47.6	47.4	<b>51.2</b>	45.3	50.1

Table 2: Zero-shot results on BEIR, scores denote **nDCG@10**. All methods re-rank the top-100 passages retrieved by BM25, except for TREC-DL 2019 to compare to prior work. Best overall are in **bold**. Underlined indicate in-domain performance, and *italicized* are based on Pyserini reproductions, differing from those reported in prior work.

previously with slight differences: 1) conventional re-ranker (BERT-FT), using a BERT-based re-ranker pre-trained on MS MARCO with the same configuration (Reimers and Gurevych, 2019); 2) generative re-ranker (T5-FT), using the same t5-base-lm-adapt but fine-tuned on MS MARCO; and 3) UPR, but re-ranked over 100 instead of 1000. For our proposed approach, we only evaluate the joint inference method (JPR), as the MS MARCO pre-trained re-ranker from SBERT<sup>4</sup> is already at a saddle point, and using it to bootstrap leads to degraded performance. Detailed training hyperparameters can be found in Appx. F.

### 3.2.3 Results and Discussion

Zero-shot results on BEIR are presented in Table 2. JPR attains roughly 2% absolute gain on average simply by utilizing both re-ranker and generator for inference, which is more prominent when compared against in-domain performances in Sec. 3.1 and on TREC-DL 2019. JPR surpasses BERT-FT on 10 out of the 14 tasks and is roughly equal on the other 4, and eclipses T5-FT on 13 of 14. Notably, for two tasks, FEVER and Climate-FEVER, generative re-rankers struggle and exhibit degraded performance, whereas JPR avoids this issue and outperforms BERT-FT. When using the comparatively huge LLaMA, we see that UPR worsens on

average, mostly due to major underperformance on tasks such as ArguAna, Touché-2020, FEVER, and Climate-FEVER. On most other tasks it outperforms UPR, suggesting that larger models’ effects may scale both ways, positively on familiar tasks, such as CQADupStack which LLaMA had exposure during LM training, and negatively on a few out-of-domain ones. JPR (LLM) can mitigate the worst cases, however it mostly does not outperform JPR that uses the considerably smaller generator.

## 4 Related Work

Passage re-ranking seeks to combine the advantages of sparse retrieval methods, such as efficiency, precise matching, and low-resource generalizability (Sciavolino et al., 2021; Reddy et al., 2021), with the superior performance of dense methods in the presence of extensive annotated data (Karpukhin et al., 2020; Guu et al., 2020). Early work by Nogueira and Cho (2019) examined BERT-based supervised re-rankers, while later research proposed reader prediction-based re-ranking (Mao et al., 2021) and attempted to use LMs as re-rankers (Sachan et al., 2022), although with limitations. Concurrent to our study, Sun et al. (2023) explored using the exceptionally larger ChatGPT models for re-ranking<sup>5</sup>.

MMI-based objectives, originally introduced in speech recognition to measure input-output dependence (Bahl et al., 1986; Woodland and Povey, 2002), have been applied to different tasks such as dialogue (Li et al., 2016), machine translation (Li and Jurafsky, 2016), and QA (Luo et al., 2022). MMI-based joint inference and learning have been explored in question answering and generation (Tang et al., 2017), language understanding and generation (Su et al., 2020), and various vision and language tasks (Xia et al., 2017).

## 5 Conclusion

In this study, we introduce a simple and effective approach to enhance re-ranking for passage retrieval. By jointly utilizing a conventional cross-encoding re-ranker and a conditional query generator for inference, we optimize the mutual information between the query and passage distributions, achieving improvements in open-domain QA retrieval, and more significantly in zero-shot information retrieval tasks.

<sup>5</sup>Sun et al. (2023) reported results only on a subset of BEIR and uses BM25 “flat” (cf. “multifield”).

<sup>4</sup>[https://www.sbert.net/docs/pretrained\\_cross-encoders.html](https://www.sbert.net/docs/pretrained_cross-encoders.html)

## 319 Limitations

320 First, improvements under the supervised setting  
321 for open-domain QA retrieval are diminished as  
322  $K$  increases, and roughly equals out with using  
323 conventional re-rankers at  $K = 20$ ; however, there  
324 are still many use cases especially for large models  
325 with limited context that can benefit from the im-  
326 provements of our approach. Additionally, in this  
327 work we tackle passage re-ranking for retrieval, fo-  
328 cusing on the second stage re-ranking scores using  
329 dense cross encoders and generative models. We  
330 have not explored approaching the retrieval process  
331 without passage re-ranking, that is, directly apply-  
332 ing the MMI objective to train a dense retrieval  
333 model, which could potentially lead to larger im-  
334 provements but comes with much higher computa-  
335 tional costs. We leave this for future work.

## 336 Ethics Statement

337 In this work, we used publicly available models and  
338 datasets for training and evaluation, and did not col-  
339 lect data or any personal information. The trained  
340 models could however potentially be misused and  
341 pose ethical risks typical of large language models  
342 when deployed in real-world applications, if not  
343 thoroughly audited.

## 344 References

- 345 L. Bahl, P. Brown, P. de Souza, and R. Mercer. 1986.  
346 [Maximum mutual information estimation of hidden  
347 markov model parameters for speech recognition.](#)  
348 In *ICASSP '86. IEEE International Conference on  
349 Acoustics, Speech, and Signal Processing*, volume 11,  
350 pages 49–52.
- 351 Alexander Bondarenko, Maik Fröbe, Meriem Be-  
352 louchif, Lukas Gienapp, Yamen Ajjour, Alexander  
353 Panchenko, Chris Biemann, Benno Stein, Henning  
354 Wachsmuth, Martin Potthast, and Matthias Hagen.  
355 2020. Overview of touché 2020: Argument retrieval.  
356 In *Experimental IR Meets Multilinguality, Multi-  
357 modality, and Interaction*, pages 384–395, Cham.  
358 Springer International Publishing.
- 359 Vera Boteva, Demian Gholipour, Artem Sokolov, and  
360 Stefan Riezler. 2016. A full-text learning to rank  
361 dataset for medical information retrieval. In *Ad-  
362 vances in Information Retrieval*, pages 716–722,  
363 Cham. Springer International Publishing.
- 364 Danqi Chen and Wen-tau Yih. 2020. [Open-domain  
365 question answering.](#) In *Proceedings of the 58th An-  
366 nual Meeting of the Association for Computational  
367 Linguistics: Tutorial Abstracts*, pages 34–37, Online.  
368 Association for Computational Linguistics.

- Arman Cohan, Sergey Feldman, Iz Beltagy, Doug  
Downey, and Daniel Weld. 2020. [SPECTER:  
Document-level representation learning using  
citation-informed transformers.](#) In *Proceedings  
of the 58th Annual Meeting of the Association  
for Computational Linguistics*, pages 2270–2282,  
Online. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and  
Kristina Toutanova. 2019. [BERT: Pre-training of  
deep bidirectional transformers for language under-  
standing.](#) In *Proceedings of the 2019 Conference of  
the North American Chapter of the Association for  
Computational Linguistics: Human Language Tech-  
nologies, Volume 1 (Long and Short Papers)*, pages  
4171–4186, Minneapolis, Minnesota. Association for  
Computational Linguistics.
- Thomas Diggelmann, Jordan Boyd-Graber, Jannis Bu-  
lian, Massimiliano Ciaramita, and Markus Leip-  
pold. 2020. Climate-fever: A dataset for verifica-  
tion of real-world climate claims. *arXiv preprint  
arXiv:2012.00614*.
- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasu-  
pat, and Ming-Wei Chang. 2020. Realm: Retrieval-  
augmented language model pre-training. In *Proceeed-  
ings of the 37th International Conference on Machine  
Learning, ICML'20. JMLR.org*.
- Faegheh Hasibi, Fedor Nikolaev, Chenyan Xiong, Krisz-  
tian Balog, Svein Erik Bratsberg, Alexander Kotov,  
and Jamie Callan. 2017. [Dbpedia-entity v2: A test  
collection for entity search.](#) In *Proceedings of the  
40th International ACM SIGIR Conference on Re-  
search and Development in Information Retrieval,  
SIGIR '17*, page 1265–1268, New York, NY, USA.  
Association for Computing Machinery.
- Doris Hoogeveen, Karin M. Verspoor, and Timothy  
Baldwin. 2015. [Cquadupstack: A benchmark data  
set for community question-answering research.](#) In  
*Proceedings of the 20th Australasian Document Com-  
puting Symposium (ADCS), ADCS '15*, pages 3:1–  
3:8, New York, NY, USA. ACM.
- Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke  
Zettlemoyer. 2017. [TriviaQA: A large scale distantly  
supervised challenge dataset for reading comprehen-  
sion.](#) In *Proceedings of the 55th Annual Meeting of  
the Association for Computational Linguistics (Vol-  
ume 1: Long Papers)*, pages 1601–1611, Vancouver,  
Canada. Association for Computational Linguistics.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick  
Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and  
Wen-tau Yih. 2020. [Dense passage retrieval for open-  
domain question answering.](#) In *Proceedings of the  
2020 Conference on Empirical Methods in Natural  
Language Processing (EMNLP)*, pages 6769–6781,  
Online. Association for Computational Linguistics.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Red-  
field, Michael Collins, Ankur Parikh, Chris Alberti,  
Danielle Epstein, Illia Polosukhin, Jacob Devlin, Ken-  
ton Lee, Kristina Toutanova, Llion Jones, Matthew

427	Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob	<i>Companion Proceedings of the The Web Conference</i>	484
428	Uszkoreit, Quoc Le, and Slav Petrov. 2019. <a href="#">Natural questions: A benchmark for question answering research</a> . <i>Transactions of the Association for Computational Linguistics</i> , 7:452–466.	2018, WWW '18, page 1941–1942, Republic and Canton of Geneva, CHE. International World Wide Web Conferences Steering Committee.	485
429			486
430			487
431			
432	Kenton Lee, Ming-Wei Chang, and Kristina Toutanova.	Yuning Mao, Pengcheng He, Xiaodong Liu, Yelong	488
433	2019. <a href="#">Latent retrieval for weakly supervised open domain question answering</a> . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 6086–6096, Florence, Italy. Association for Computational Linguistics.	Shen, Jianfeng Gao, Jiawei Han, and Weizhu Chen.	489
434		2021. <a href="#">Reader-guided passage reranking for open-domain question answering</a> . In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 344–350, Online. Association for Computational Linguistics.	490
435			491
436			492
437			493
438	Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2016. <a href="#">A diversity-promoting objective function for neural based models</a> . In <i>Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 110–119, San Diego, California. Association for Computational Linguistics.	Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. <a href="#">Distributed representations of words and phrases and their compositionality</a> . In <i>Advances in Neural Information Processing Systems</i> , volume 26. Curran Associates, Inc.	495
439			496
440			497
441			498
442			499
443			
444		Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng.	500
445		2017. <a href="#">MS MARCO: A human-generated MACHine reading COMprehension dataset</a> .	501
446	Jiwei Li and Dan Jurafsky. 2016. Mutual information and diverse decoding improve neural machine translation. <i>arXiv preprint arXiv:1601.00372</i> .		502
447			503
448		Rodrigo Nogueira and Kyunghyun Cho. 2019. <a href="#">Passage re-ranking with bert</a> . <i>arXiv preprint arXiv:1901.04085</i> .	504
449	Jimmy Lin, Xueguang Ma, Sheng-Chieh Lin, Jheng-Hong Yang, Ronak Pradeep, and Rodrigo Nogueira.		505
450	2021. <a href="#">Pyserini: A Python toolkit for reproducible information retrieval research with sparse and dense representations</a> . In <i>Proceedings of the 44th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2021)</i> , pages 2356–2362.	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, Peter J Liu, et al. 2020. <a href="#">Exploring the limits of transfer learning with a unified text-to-text transformer</a> . <i>J. Mach. Learn. Res.</i> , 21(140):1–67.	507
451			508
452			509
453			510
454			511
455			
456		Revanth Gangi Reddy, Vikas Yadav, Md Arafat Sultan, Martin Franz, Vittorio Castelli, Heng Ji, and Avirup Sil. 2021. <a href="#">Towards robust neural retrieval models with synthetic pre-training</a> . <i>arXiv preprint arXiv:2104.07800</i> .	512
457	Ilya Loshchilov and Frank Hutter. 2018. <a href="#">Decoupled weight decay regularization</a> . In <i>International Conference on Learning Representations</i> .		513
458			514
459			515
460	Yi Luan, Jacob Eisenstein, Kristina Toutanova, and Michael Collins. 2021. <a href="#">Sparse, dense, and attentional representations for text retrieval</a> . <i>Transactions of the Association for Computational Linguistics</i> , 9:329–345.	Nils Reimers and Iryna Gurevych. 2019. <a href="#">Sentence-BERT: Sentence embeddings using Siamese BERT-networks</a> . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.	516
461			517
462			518
463			519
464			520
465	Hongyin Luo, Shang-Wen Li, Mingye Gao, Seunghak Yu, and James Glass. 2022. <a href="#">Cooperative self-training of machine reading comprehension</a> . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 244–257, Seattle, United States. Association for Computational Linguistics.		521
466			522
467			523
468			524
469		Stephen Robertson, Hugo Zaragoza, et al. 2009. The probabilistic relevance framework: Bm25 and beyond. <i>Foundations and Trends® in Information Retrieval</i> , 3(4):333–389.	525
470			526
471			527
472			528
473	Zhuang Ma and Michael Collins. 2018. <a href="#">Noise contrastive estimation and negative sampling for conditional models: Consistency and statistical efficiency</a> . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 3698–3707, Brussels, Belgium. Association for Computational Linguistics.	Devendra Sachan, Mike Lewis, Mandar Joshi, Armen Aghajanyan, Wen-tau Yih, Joelle Pineau, and Luke Zettlemoyer. 2022. <a href="#">Improving passage retrieval with zero-shot question generation</a> . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 3781–3797, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	529
474			530
475			531
476			532
477			533
478			534
479			535
480	Macedo Maia, Siegfried Handschuh, André Freitas, Brian Davis, Ross McDermott, Manel Zarrouk, and Alexandra Balahur. 2018. <a href="#">Www'18 open challenge: Financial opinion mining and question answering</a> . In	G. Salton, A. Wong, and C. S. Yang. 1975. <a href="#">A vector space model for automatic indexing</a> . <i>Commun. ACM</i> , 18(11):613–620.	537
481			538
482			539
483			

540	Victor Sanh, Albert Webson, Colin Raffel, Stephen	Christophe Van Gysel and Maarten de Rijke. 2018.	598
541	Bach, Lintang Sutawika, Zaid Alyafeai, Antoine	<a href="#">Py trec_eval: An extremely fast python interface to</a>	599
542	Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey,	<a href="#">trec_eval</a> . In <i>The 41st International ACM SIGIR Con-</i>	600
543	M Saiful Bari, Canwen Xu, Urmish Thakker,	<i>ference on Research &amp; Development in Information</i>	601
544	Shanya Sharma Sharma, Eliza Szczechla, Taewoon	<i>Retrieval, SIGIR '18</i> , page 873–876, New York, NY,	602
545	Kim, Gunjan Chhablani, Nihal Nayak, Debajyoti	USA. Association for Computing Machinery.	603
546	Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	604
547	Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong,	Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz	605
548	Harshit Pandey, Rachel Bawden, Thomas Wang, Tr-	Kaiser, and Illia Polosukhin. 2017. <a href="#">Attention is all</a>	606
549	ishala Neeraj, Jos Rozen, Abheesht Sharma, An-	<a href="#">you need</a> . In <i>Advances in Neural Information Pro-</i>	607
550	drea Santilli, Thibault Fevry, Jason Alan Fries, Ryan	<i>cessing Systems</i> , volume 30. Curran Associates, Inc.	608
551	Teehan, Teven Le Scao, Stella Biderman, Leo Gao,	Ellen Voorhees, Tasmee Alam, Steven Bedrick, Dina	609
552	Thomas Wolf, and Alexander M Rush. 2022. <a href="#">Multi-</a>	Demner-Fushman, William R. Hersh, Kyle Lo, Kirk	610
553	<a href="#">task prompted training enables zero-shot task gener-</a>	Roberts, Ian Soboroff, and Lucy Lu Wang. 2021.	611
554	<a href="#">alization</a> . In <i>International Conference on Learning</i>	<a href="#">Trec-covid: Constructing a pandemic information</a>	612
555	<i>Representations</i> .	<a href="#">retrieval test collection</a> . <i>SIGIR Forum</i> , 54(1).	613
556	Christopher Sciavolino, Zexuan Zhong, Jinhyuk Lee,	Ellen M Voorhees et al. 1999. The trec-8 question	614
557	and Danqi Chen. 2021. <a href="#">Simple entity-centric ques-</a>	answering track report. In <i>Trec</i> , volume 99, pages	615
558	<a href="#">tions challenge dense retrievers</a> . In <i>Proceedings of</i>	77–82. Citeseer.	616
559	<i>the 2021 Conference on Empirical Methods in Natu-</i>	Henning Wachsmuth, Shahbaz Syed, and Benno Stein.	617
560	<i>ral Language Processing</i> , pages 6138–6148, Online	2018. <a href="#">Retrieval of the best counterargument without</a>	618
561	and Punta Cana, Dominican Republic. Association	<a href="#">prior topic knowledge</a> . In <i>Proceedings of the 56th</i>	619
562	for Computational Linguistics.	<i>Annual Meeting of the Association for Computational</i>	620
563	Shang-Yu Su, Yung-Sung Chuang, and Yun-Nung Chen.	<i>Linguistics (Volume 1: Long Papers)</i> , pages 241–251,	621
564	2020. <a href="#">Dual inference for improving language under-</a>	Melbourne, Australia. Association for Computational	622
565	<a href="#">standing and generation</a> . In <i>Findings of the Associ-</i>	<i>Linguistics</i> .	623
566	<i>ation for Computational Linguistics: EMNLP 2020</i> ,	David Wadden, Shanchuan Lin, Kyle Lo, Lucy Lu	624
567	pages 4930–4936, Online. Association for Computa-	Wang, Madeleine van Zuylen, Arman Cohan, and	625
568	tional Linguistics.	Hannaneh Hajishirzi. 2020. <a href="#">Fact or fiction: Verifying</a>	626
569	Weiwei Sun, Lingyong Yan, Xinyu Ma, Pengjie	<a href="#">scientific claims</a> . In <i>Proceedings of the 2020 Con-</i>	627
570	Ren, Dawei Yin, and Zhaochun Ren. 2023. Is	<i>ference on Empirical Methods in Natural Language</i>	628
571	chatgpt good at search? investigating large lan-	<i>Processing (EMNLP)</i> , pages 7534–7550, Online. As-	629
572	guage models as re-ranking agent. <i>arXiv preprint</i>	sociation for Computational Linguistics.	630
573	<i>arXiv:2304.09542</i> .	Yining Wang, Liwei Wang, Yuanzhi Li, Di He, and Tie-	631
574	Duyu Tang, Nan Duan, Tao Qin, Zhao Yan, and Ming	Yan Liu. 2013. A theoretical analysis of ndcg type	632
575	Zhou. 2017. Question answering and question gener-	ranking measures. In <i>Conference on learning theory</i> ,	633
576	ation as dual tasks. <i>arXiv preprint arXiv:1706.02027</i> .	pages 25–54. PMLR.	634
577	Nandan Thakur, Nils Reimers, Andreas Rücklé, Ab-	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien	635
578	hishek Srivastava, and Iryna Gurevych. 2021. <a href="#">BEIR:</a>	Chaumond, Clement Delangue, Anthony Moi, Pier-	636
579	<a href="#">A heterogeneous benchmark for zero-shot evaluation</a>	ric Cistac, Tim Rault, Remi Louf, Morgan Funtow-	637
580	<a href="#">of information retrieval models</a> . In <i>Thirty-fifth Con-</i>	icz, Joe Davison, Sam Shleifer, Patrick von Platen,	638
581	<i>ference on Neural Information Processing Systems</i>	Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu,	639
582	<i>Datasets and Benchmarks Track (Round 2)</i> .	Teven Le Scao, Sylvain Gugger, Mariama Drame,	640
583	James Thorne, Andreas Vlachos, Christos	Quentin Lhoest, and Alexander Rush. 2020. <a href="#">Trans-</a>	641
584	Christodoulopoulos, and Arpit Mittal. 2018.	<a href="#">formers: State-of-the-art natural language processing</a> .	642
585	<a href="#">FEVER: a large-scale dataset for fact extraction</a>	In <i>Proceedings of the 2020 Conference on Empirical</i>	643
586	<a href="#">and VERification</a> . In <i>Proceedings of the 2018</i>	<i>Methods in Natural Language Processing: System</i>	644
587	<i>Conference of the North American Chapter of</i>	<i>Demonstrations</i> , pages 38–45, Online. Association	645
588	<i>the Association for Computational Linguistics:</i>	for Computational Linguistics.	646
589	<i>Human Language Technologies, Volume 1 (Long</i>	P.C. Woodland and D. Povey. 2002. <a href="#">Large scale dis-</a>	647
590	<i>Papers)</i> , pages 809–819, New Orleans, Louisiana.	<a href="#">criminative training of hidden markov models for</a>	648
591	Association for Computational Linguistics.	<a href="#">speech recognition</a> . <i>Computer Speech &amp; Language</i> ,	649
592	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	16(1):25–47.	650
593	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	Yingce Xia, Tao Qin, Wei Chen, Jiang Bian, Nenghai	651
594	Baptiste Rozière, Naman Goyal, Eric Hambro,	Yu, and Tie-Yan Liu. 2017. Dual supervised learning.	652
595	Faisal Azhar, et al. 2023. Llama: Open and effi-	In <i>Proceedings of the 34th International Conference</i>	653
596	cient foundation language models. <i>arXiv preprint</i>	<i>on Machine Learning-Volume 70</i> , pages 3789–3798.	654
597	<i>arXiv:2302.13971</i> .		

Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. [HotpotQA: A dataset for diverse, explainable multi-hop question answering](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2369–2380, Brussels, Belgium. Association for Computational Linguistics.

## A Source Code

We provide anonymized source code at <https://www.dropbox.com/s/n6zzvs01yutliwp/jpr.zip>. More details can be found in `README.md`.

## B Open-Domain QA Retrieval Datasets

We show the number of train/dev/test examples in NQ and TriviaQA in Table 3. Please refer to [Kwiatkowski et al. \(2019\)](#) and [Joshi et al. \(2017\)](#) for more details. Note that NQ is licensed under Apache License 2.0, which we follow, and TriviaQA does not provide dataset licenses.

Dataset	Train	Dev	Test
Natural Questions	58,880	8,757	3,610
TriviaQA	60,413	8,837	11,313

Table 3: Dataset splits for NQ and TriviaQA.

## C Open-Domain QA Retrieval Training and Inference Details

### C.1 Training

Generally, conventional cross-encoders are trained to minimize the negative likelihood  $\mathcal{L}_{\text{retrieval}}(\phi) \triangleq -\mathbb{E}_{x,z \sim p(x,z)} [\log p_\phi(z|x)]$ , where  $p_\phi(z|x)$  is usually calculated from the retrieval score of question-passage pairs, with the partition function approximated by a noise contrastive approach trained either with a classification or a ranking objective ([Ma and Collins, 2018](#)). We choose to fine-tune our cross-encoder, BERT-FT, using a 6-layer transformer model ([Vaswani et al., 2017](#)), which takes the concatenated input of a query and a passage, with the binary classification objective for noise contrastive learning ([Mikolov et al., 2013](#)). The 6-layer SBERT model MiniLM-L-6-v2 we use was previously pre-trained on MS MARCO, which we fine-tune for 2 epochs using the top 32 passages from BM25 on the NQ/TriviaQA training set. We train with a batch size of 128, learning rate of  $5e-5$ , linear warmup and decay with ratio of 0.1.

For training of T5-FT, we fine-tune with  $\mathcal{L}_{\text{generation}}(\theta)$  using the `t5-base-lm-adapt` model, a 12-layer encoder-decoder configuration with

Hyper-parameter	NQ		TriviaQA	
	BERT-FT	T5-FT	BERT-FT	T5-FT
learning rate	1e-5	2e-5	1e-5	1.5e-5
batch size	96	64	64	64
$\alpha$	0.0005	0.0005	0.005	0.005

Table 4: Training hyperparameters for NQ and TriviaQA selected by performance on the dev set.

220M parameters initialized from T5-base v1.1 and trained for an additional 100k steps with an LM objective. It takes a ground truth passage as input with its corresponding query as the decoder target. Ground truth query-passage pairs from the training set was used to fine-tune the model for 2 epochs. We use a batch size of 64, learning rate of  $5e-5$ , and linear warmup and decay ratio of 0.1. Hyperparameters were chosen by performance on the dev set.

UPR uses the pre-trained T0-3B directly without any fine-tuning.

JPR uses BERT-FT and T5-FT, described earlier, directly during inference (see Sec. C.2 below). JPR-FT requires further fine-tuning, which we train for another epoch. Training hyperparameters were searched with the dev set, with one run for each hyperparameter setting, shown in Table 4. We report results for the model with the best-performing run on the dev set.

All models were trained with HuggingFace’s Transformers library ([Wolf et al., 2020](#)), using the AdamW optimizer ([Loshchilov and Hutter, 2018](#)) with default parameters. The maximum sequence lengths for queries and passages were set to 128 and 512, respectively, for generative models. For the cross-encoding BERT-FT, we set the maximum concatenated length to be 512. Training was done with four Nvidia A6000 GPUs, with around 2.5 GPU hours per epoch, equating to around 250 GPU hours in total.

### C.2 Inference

For the conventional cross-encoding re-ranker (BERT-FT), we re-rank with Eq. 1 by directly ranking the retrieval scores. When using BERT-FT in JPR, we approximate  $\log p_\phi(z|x)$  by taking SoftMax over the scores for the 100 retrieved passages. For generative re-rankers T5-FT and UPR, we follow [Sachan et al. \(2022\)](#) and estimate  $\log p_\theta(x|z)$  with length-normalized conditional likelihood followed by taking SoftMax over the passages. For JPR, the preceding two terms are weight-averaged according to Eq. 3.



Cross-encoder	Generative Model	#params	Top-1	Top-5	Top-10
TinyBERT	✗	4.4M	37.8	60.3	67.0
MiniLM-L-4	✗	19.2M	47.5	65.9	70.9
MiniLM-L-6 (BERT-FT)	✗	22.7M	49.4	66.4	71.4
BERT-base	✗	109.5M	49.2	66.0	70.8
BERT-large	✗	335.1M	49.8	67.5	71.7
✗	T5-tiny	15.6M	25.7	51.4	62.0
✗	T5-small	77.0M	30.7	57.1	65.2
✗	T5-base (T5-FT)	247.6M	34.4	59.7	66.9
MiniLM-L-6	T5-tiny	38.3M	49.6	67.0	71.6
MiniLM-L-6	T5-small	99.7M	50.4	67.3	71.7
MiniLM-L-6	T5-base	270.3M	50.4	67.3	71.8

Table 5: Top- $K$  retrieval accuracy (%) on NQ for different model combinations with the proposed JPR.

## D Results on Open-Domain QA Retrieval with Different Cross-encoding and Generative Model Pairs

We further show the efficacy of JPR on NQ by conducting additional evaluations on NQ with various model combinations. We experiment with BERT models of different sizes for the cross-encoders, and for generative models we chose T5 models of multiple models sizes. All cross-encoding models were previously pre-trained on MS MARCO, which we fine-tune on NQ, and the T5 models were fine-tuned on NQ, all following training procedures reported in Sec. C. For inference, we use  $\lambda = 0.5$  and follow the inference steps outlined in Sec. C.2. The results are shown in Table 5.

From the results, notice that when T5-small is paired with MiniLM-L-6 for JPR, it aligns with the performance of T5-base paired with MiniLM-L-6. This observation underscores that the additional parameters of T5-base may be superfluous in our application. When comparing JPR (MiniLM-L-6 & T5-small) with the standalone BERT-base, which is in the same parameter ballpark, and the larger BERT-large, it’s evident that the gains from JPR are not solely attributable to model size.

## E BEIR Benchmark

The BEIR benchmark contains 18 datasets from a variety of text retrieval tasks and domains, 14 of which are publicly available. In this work we evaluate baselines and our approach on the publicly available datasets in BEIR: TREC-COVID (Voorhees et al., 2021), NFCorpus (Boteva et al., 2016), NQ (Kwiatkowski et al., 2019), HotpotQA (Yang et al., 2018), FiQA-2018 (Maia

et al., 2018), ArguAna (Wachsmuth et al., 2018), Touché-2020 (Bondarenko et al., 2020), CQADupStack (Hoogeveen et al., 2015), Quora<sup>6</sup>, DB-Pedia (Hasibi et al., 2017), SCIDOCS (Cohan et al., 2020), FEVER (Thorne et al., 2018), Climate-FEVER (Diggelmann et al., 2020), and SciFact (Wadden et al., 2020). For details on dataset statistics, links, and licenses please refer to BEIR (Thakur et al., 2021). Note that datasets in BEIR that are under copyright were not used in this study, and 4 out of the 14 publicly available datasets do not report dataset licenses. We follow the intended uses for each dataset license.

## F Zero-shot Retrieval Training and Inference Details

For BEIR, since the SBERT model was already pre-trained on MS MARCO, we directly use it for BERT-FT. On the other hand, T5-FT stills requires fine-tuning, which we train for 3 epochs on query-passage pairs in the training set, with batch size of 16 and learning rate of  $5e-5$  with no warmup. The inference process is the same as open-domain QA retrieval, described earlier in Sec. C.2.

<sup>6</sup><https://quoradata.quora.com/First-Quora-Dataset-Release-Question-Pairs>