

- general text embeddings with multi-stage contrastive learning. *arXiv preprint arXiv:2308.03281*.
- Yu-An Liu, Ruqing Zhang, Jiafeng Guo, Maarten de Rijke, Yixing Fan, and Xueqi Cheng. 2024. Robust neural information retrieval: An adversarial and out-of-distribution perspective. *arXiv preprint arXiv:2407.06992*.
- Jianmo Ni, Chen Qu, Jing Lu, Zhuyun Dai, Gustavo Hernández Ábrego, Ji Ma, Vincent Y. Zhao, Yi Luan, Keith B. Hall, Ming-Wei Chang, and Yinfei Yang. 2021. [Large dual encoders are generalizable retrievers](#). *Preprint*, arXiv:2112.07899.
- Fabio Petroni, Patrick Lewis, Aleksandra Piktus, Tim Rocktäschel, Yuxiang Wu, Alexander H Miller, and Sebastian Riedel. How context affects language models’ factual predictions. In *Automated Knowledge Base Construction*.
- Juan Ramos et al. 2003. Using tf-idf to determine word relevance in document queries. In *Proceedings of the first instructional conference on machine learning*, volume 242, pages 29–48. Citeseer.
- Stephen Robertson, Hugo Zaragoza, et al. 2009. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends® in Information Retrieval*, 3(4):333–389.
- 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 *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Kexin Wang, Nandan Thakur, Nils Reimers, and Iryna Gurevych. 2021. Gpl: Generative pseudo labeling for unsupervised domain adaptation of dense retrieval. *arXiv preprint arXiv:2112.07577*.
- Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. 2022. Text embeddings by weakly-supervised contrastive pre-training. *arXiv preprint arXiv:2212.03533*.
- Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. 2024. Multilingual e5 text embeddings: A technical report. *arXiv preprint arXiv:2402.05672*.
- Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas Muennighoff. 2023. [C-pack: Packaged resources to advance general chinese embedding](#). *Preprint*, arXiv:2309.07597.
- Renchunzi Xie, Ambroise Odonnat, Vasiliï Feofanov, Ievgen Redko, Jianfeng Zhang, and Bo An. 2024. Characterising gradients for unsupervised accuracy estimation under distribution shift. *arXiv preprint arXiv:2401.08909*.
- 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.
- Yue Yu, Chenyan Xiong, Si Sun, Chao Zhang, and Arnold Overwijk. 2022. Coco-dr: Combating distribution shift in zero-shot dense retrieval with contrastive and distributionally robust learning. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 1462–1479.
- Jingtao Zhan, Jiaxin Mao, Yiqun Liu, Jiafeng Guo, Min Zhang, and Shaoping Ma. 2021. Optimizing dense retrieval model training with hard negatives. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 1503–1512.

ratio $(1 - r(\mathcal{C}))$.

The graph's descending trend indicates that $1 - r(\mathcal{C})$ is proportional to retriever performance, as datasets with higher retrieval performance show greater Non-OOD ratios. GradNormIR clearly demonstrates this relationship, showing high Non-OOD ratios for Quora, Arguana, and Touché, and low ratios for FiQA, Scidocs, NFCorpus, and COVID. While GenQuery also exhibits a descending trend, it shows minimal variation from Quora to NFCorpus, making OOD corpus detection less effective.

To predict OOD corpora, we set γ to 0.5 based on the average performance across all datasets. With this threshold, we identify Scifact, Touch'e, DBPedia, FiQA, Scidocs, NFCorpus, and COVID as OOD corpora.

F Feasibility of GradNormIR

We aim to validate whether GradNormIR can identify the documents that are difficult for the the models to retrieve. To this end, we inspect if there is a consistent relationship between the computed gradient norm and the likelihood of a document successfully retrieved by its associated queries.

Evaluation Metric. To evaluate the effectiveness, we measure the document-to-query (d2q) as the standard metric. In each dataset, annotations are provided in the form of $\{q_i, D_i\}_{i=1}^N$, where q_i is a query and D_i is the set of relevant documents. We reorganize these annotations as $\{d_i, Q_i\}_{i=1}^N$, where Q_i represents the set of relevant queries for each document d_i . For a document to be considered effectively retrievable, it should be retrieved for all its relevant queries.

To quantify this, we define the d2q recall as follows:

$$\text{recall}_{\text{d2q}} = \frac{\sum_{q_i \in Q_i} \mathbb{I}\{d_i \in D^+(q_i)\}}{|Q_i|}, \quad (6)$$

where \mathbb{I} is an indicator function and $D^+(q_i)$ represents the top- k retrieved documents (with $k = 100$).

When the retriever model generalizes well for a document d_i , the d2q recall value will be high. Additionally, if the retriever generalizes effectively on d_i , the gradient norm associated with d_i will be low, as the retriever does not need to make substantial updates based on the contrastive loss for d_i . Therefore, there should be an inverse relationship: higher the d2q recall values correspond to lower the gradient norms.

Results. Figure 7 illustrates the relationship between GradNormIR and d2q recall. We divide the data points into quartiles based on GradNormIR values, sorted in ascending order and labeled as Q1, Q2, Q3, and Q4. The x-axis represents these quartiles, while the y-axis shows the average d2q recall for each group.

The results reveal a strong inverse correlation between GradNormIR and retrieval performance. As GradNormIR values increase from Q1 to Q4, d2q recall decreases. This indicates that higher GradNormIR values (Q4) are associated with documents that are more challenging for the retriever to retrieve consistently. Conversely, lower GradNormIR values (Q1) correspond to higher recall, indicating better retrieval performance. When d2q recall approaches 1, such as Quora and SciFact, this trend becomes less noticeable. This is likely because the datasets have been trained on; nearly all documents are well generalized and easily retrievable.

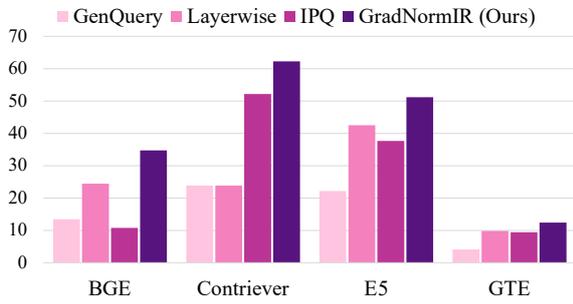


Figure 5: Results of relevance gains via OOD document filtering.

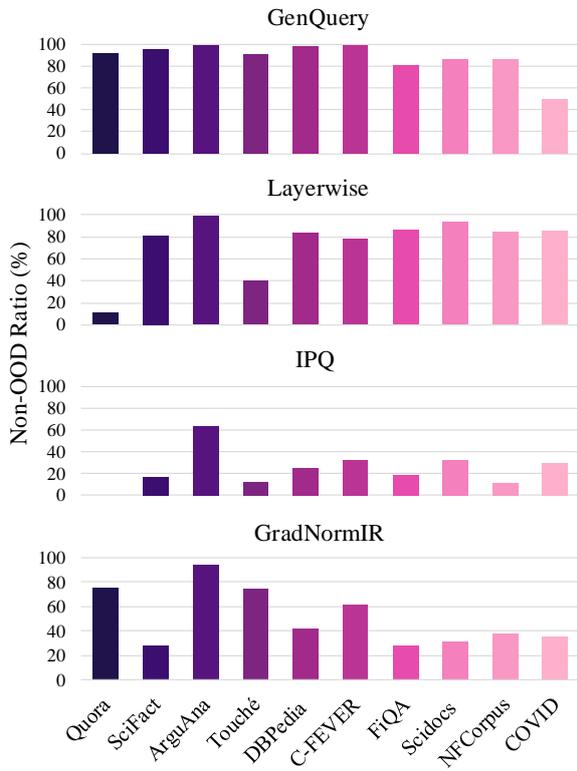


Figure 6: Relation Between OOD Ratio and Performance

