

Factuality and Transparency Are All RAG Needs! Self-Explaining Contrastive Evidence Re-ranking

Francielle Vargas
São Paulo State University
francielle.vargas@unesp.br

Daniel Pedronette
São Paulo State University
pedronette@unesp.br

ABSTRACT

This extended abstract introduces Self-Explaining Contrastive Evidence Re-Ranking (CER), a novel method that restructures retrieval around factual evidence by fine-tuning embeddings with contrastive learning and generating token-level attribution rationales for each retrieved passage. Hard negatives are automatically selected using a subjectivity-based criterion, forcing the model to pull factual rationales closer while pushing subjective or misleading explanations apart. As a result, the method creates an embedding space explicitly aligned with evidential reasoning. We evaluated our method on clinical trial reports, and initial experimental results show that CER improves retrieval accuracy, mitigates the potential for hallucinations in RAG systems, and provides transparent, evidence-based retrieval that enhances reliability, especially in safety-critical domains.

1 INTRODUCTION

Despite recent advances in Retrieval-Augmented Generation (RAG), recent studies have shown that the benefits of RAG can be undermined by noisy retrieval, where passages that truly support a claim are mixed with those that are only topically related, incomplete, or even contradictory, reducing factual correctness and reliability [9, 10]. This limitation is particularly concerning in sensitive domains such as healthcare, where inaccurate or unsupported answers can have severe consequences, for example, compromising patient safety [14] or leading to legal [11] and ethical [2, 3, 5, 8] implications. In addition, most RAG systems are not explicitly parameterized to distinguish evidentially valid passages from texts that are semantically proximal lack factual grounding. Consequently, these systems frequently conflate topical similarity with evidential adequacy, resulting in systematic semantic confounding and the absence of retrieval mechanisms optimized for factual correctness. This limitation, combined with insufficient term-level interpretability, directly compromises the factual robustness of downstream generation—an especially critical concern in biomedical and clinical contexts [4, 12, 14]. For example, in response to the query “Is lemon effective for curing cancer?”, a standard retriever may prioritize passages containing lexical associations with “lemon,” “cancer,” or even misinformation, while failing to surface clinically substantiated evidence. An evidence-based retriever, by contrast, must

elevate passages reporting verified trial outcomes or medical consensus. When generation models are conditioned on retrieval distributions lacking evidential precision, the likelihood of propagating incomplete, biased, or spurious information increases substantially, amplifying hallucinations and degrading factual reliability [6].

Recent work has attempted to address these deficiencies through contrastive learning [7], training retrievers to discriminate between relevant and irrelevant passages by maximizing representational separation between positive and negative samples [13–15]. However, these approaches typically treat all negative instances as homogeneous, disregarding the various semantic functions they may serve with respect to a target claim. For example, providing corroborating evidence, presenting contradictory findings, or being entirely unrelated. Moreover, current contrastive retrievers optimize relevance implicitly through similarity scores, offering limited transparency regarding the retrieval rationale and leaving the underlying evidential reasoning process opaque.

To fill these limitations, we propose Self-Explaining Contrastive Evidence Re-ranking (CER), which integrates triplet-based contrastive learning with explicit evidential attribution. CER fine-tunes Contriever [7] using a cosine and euclidean triplet loss to restructure the embedding space, as factual evidence is systematically drawn closer to the query representation, while subjective, contradictory, or irrelevant content is repelled. Empirical evaluation on large-scale clinical trial corpora demonstrates significant gains in retrieval accuracy.

2 RELATED WORK

Chatzikyriakidis and Natsina [1] generate Modern Greek interwar poetry with GPT-4-turbo and GPT-4o using RAG and contrastive learning. RAG retrieves thematically and stylistically similar poems to guide generation, while the contrastive variant adds opposing examples for style control. Expert evaluations and quantitative metrics (vocabulary density, average words per sentence, readability index) show that RAG improves style and thematic consistency. Sriram et al. [13] create a dense retriever fine-tuned with contrastive learning to improve evidence retrieval for complex fact-checking. Contrastive Fact-Checking Reranker (CFR), extends Contriever [7] using multiple supervision signals: GPT-4 distillation, LERC-based answer equivalence, and human-annotated gold data from the AVeriTeC dataset. The model enhances second-stage retrieval by ranking documents that better support claim verification. Experiments show a 6% improvement in veracity classification accuracy and 9% higher top-document relevance on AVeriTeC, with consistent gains across benchmarks, demonstrating improved reasoning and robustness in real-world fact-checking.

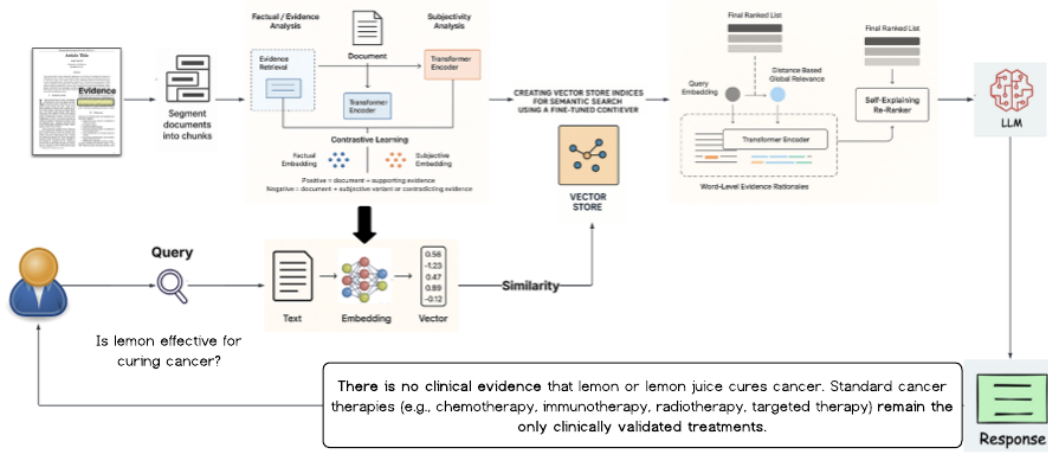


Figure 1: The pipeline for self-Explaining Re-Ranking with Contrastive Evidence Selection for Retrieval-Augmented Generation.

3 THE PROPOSED METHOD

The pipeline is decomposed into two complementary stages: **contrastive retrieval**, which learns an evidence-sensitive embedding space through (i) document chunking and contrastive training and (ii) contextual retrieval; and **self-explaining re-ranking**, which provides fine-grained attribution signals to identify the most query-aligned passages through (iii) self-explaining re-ranking. The selected evidence—augmented with token-level rationales—is then passed to an LLM, which generates a response grounded in verifiable evidence. The system subsequently returns the final answer to the user. In this design, the first stage ensures that retrieval prioritizes factual evidence rather than merely topical content, while the second stage refines ranking via token-level semantic contributions, enabling interpretable and evidence-based selection of supporting passages. Next, we describe the method in detail.

4 EVALUATION AND INITIAL RESULTS

We evaluated only the proposed fine-tuned Contriever model, reporting results using recall@K and precision@K (e.g., $K=5$).

As shown in 3, The UMAP projections show that our fine-tuned model, trained with the proposed subjectivity-based triplet selection, produces a much clearer separation between evidential (positive) and non-evidential (negative) embeddings. While the self-learning baseline displays substantial overlap between classes, the fine-tuned model forms well-separated clusters, indicating that contrastive training effectively distinguishes evidence from non-evidence. We also evaluated the efficiency between applying Euclidean and cosine distance metrics, as shown in Figure 2.

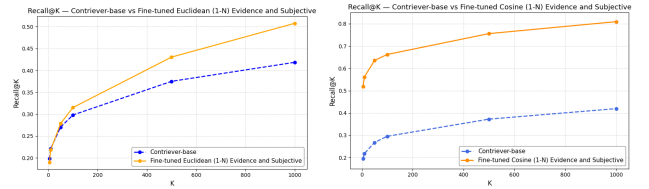


Figure 2: Self-learning and fine-tuned models evaluated using both Euclidean and cosine distance metrics.

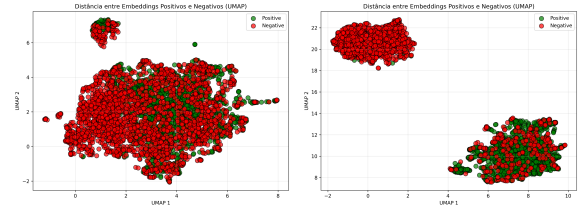


Figure 3: Results for the Contriever self-learning model (left) and the Contriever fine-tuned model (right). For the self-learning model, the average distance among positive pairs is Intra-Pos = 0.5716 and among negative pairs is Intra-Neg = 0.5977; the average distance between positives and negatives is Inter = 0.5901. For the fine-tuned model, the average distance among positive pairs is Intra-Pos = 0.7766 and among negative pairs is Intra-Neg = 0.8141; the average distance between positives and negatives is Inter = 0.8110.

ACKNOWLEDGEMENTS

The authors are grateful to São Paulo Research Foundation - FAPESP (grants #2025/01118-2 and #2024/04890-5) for financial support.

LIMITATIONS

These findings represent only initial results obtained from experiments using the fine-tuned model trained with our proposed criteria for selecting positive and negative examples in the triplet loss for contrastive learning. While the results indicate that the selection strategy is effective in structuring the embedding space and separating evidential from non-evidential content, they should be interpreted as preliminary evidence rather than definitive validation. More robust experiments are still required to thoroughly evaluate the proposed method, both in the contrastive retrieval stage and in the re-ranking stage. Furthermore, no assessment has yet been conducted on the quality of the explanations generated by the self-explaining component, which remains an essential step for understanding the interpretability and reliability of the approach.

REFERENCES

- [1] Stergios Chatzikyriakidis and Anastasia Natsina. 2025. Poetry in RAGs: Modern Greek interwar poetry generation using RAG and contrastive training. In *Proceedings of the 5th International Conference on Natural Language Processing for Digital Humanities*, Mika Härmäläinen, Emily Öhman, Yuri Bizzoni, So Miyagawa, and Khalid Alnajjar (Eds.). Association for Computational Linguistics, Albuquerque, USA, 257–264. <https://doi.org/10.18653/v1/2025.nlp4dh-1.22>
- [2] Yen-Shan Chen, Jing Jin, Peng-Ting Kuo, Chao-Wei Huang, and Yun-Nung Chen. 2025. LLMs are Biased Evaluators But Not Biased for Fact-Centric Retrieval Augmented Generation. In *Findings of the Association for Computational Linguistics: ACL 2025*, Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (Eds.). Association for Computational Linguistics, Vienna, Austria, 26669–26684. <https://doi.org/10.18653/v1/2025.findings-acl.1369>
- [3] Jiali Cheng and Hadi Amiri. 2025. EqualizeIR: Mitigating Linguistic Biases in Retrieval Models. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)*, Luis Chiruzzo, Alan Ritter, and Lu Wang (Eds.). Association for Computational Linguistics, Albuquerque, New Mexico, 889–898. <https://doi.org/10.18653/v1/2025.naacl-short.75>
- [4] Jean-Benoit Delbrouck, Pierre Chambon, Christian Bluethgen, Emily Tsai, Omar Almusa, and Curtis Langlotz. 2022. Improving the Factual Correctness of Radiology Report Generation with Semantic Rewards. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (Eds.). Association for Computational Linguistics, Abu Dhabi, United Arab Emirates, 4348–4360. <https://doi.org/10.18653/v1/2022.findings-emnlp.319>
- [5] Kshitish Ghate, Tessa Charlesworth, Mona T. Diab, and Aylin Caliskan. 2025. Biases Propagate in Encoder-based Vision-Language Models: A Systematic Analysis From Intrinsic Measures to Zero-shot Retrieval Outcomes. In *Findings of the Association for Computational Linguistics: ACL 2025*, Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (Eds.). Association for Computational Linguistics, Vienna, Austria, 18562–18580. <https://doi.org/10.18653/v1/2025.findings-acl.955>
- [6] Wentao Hu, Wengyu Zhang, Yiyang Jiang, Chen Jason Zhang, Xiaoyong Wei, and Li Qing. 2025. Removal of Hallucination on Hallucination: Debate-Augmented RAG. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (Eds.). Association for Computational Linguistics, Vienna, Austria, 15839–15853. <https://doi.org/10.18653/v1/2025.acl-long.770>
- [7] Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2022. Unsupervised Dense Information Retrieval with Contrastive Learning. [arXiv:2112.09118 \[cs.LR\]](https://arxiv.org/abs/2112.09118) <https://arxiv.org/abs/2112.09118>
- [8] Youngwoo Kim, Razieh Rahimi, and James Allan. 2024. Discovering Biases in Information Retrieval Models Using Relevance Thesaurus as Global Explanation. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (Eds.). Association for Computational Linguistics, Miami, Florida, USA, 19530–19547. <https://doi.org/10.18653/v1/2024.emnlp-main.1089>
- [9] Nelson Liu, Tianyi Zhang, and Percy Liang. 2023. Evaluating Verifiability in Generative Search Engines. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, Singapore, 7001–7025. <https://doi.org/10.18653/v1/2023.findings-emnlp.467>
- [10] Fabio Petroni, Patrick Lewis, Aleksandra Piktus, Tim Rocktäschel, Yuxiang Wu, Alexander H. Miller, and Sebastian Riedel. 2020. How Context Affects Language Models’ Factual Predictions. [arXiv:2005.04611 \[cs.CL\]](https://arxiv.org/abs/2005.04611) <https://arxiv.org/abs/2005.04611>
- [11] General Data Protection Regulation. 2016. Regulation (EU) 2016/679 of the European Parliament and of the Council of 27 April 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46. *Official Journal of the European Union (OJ)* 59, 1–88 (2016).
- [12] Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed Chi, Nathanael Schärli, and Denny Zhou. 2023. Large language models can be easily distracted by irrelevant context. In *Proceedings of the 40th International Conference on Machine Learning (Honolulu, Hawaii, USA) (ICML ’23)*. JMLR.org, Article 1291, 18 pages.
- [13] Aniruddh Sriram, Fangyuan Xu, Eunsol Choi, and Greg Durrett. 2024. Contrastive Learning to Improve Retrieval for Real-World Fact Checking. In *Proceedings of the Seventh Fact Extraction and VERification Workshop (FEVER)*, Michael Schlichtkrull, Yulong Chen, Chenxi Whitehouse, Zhenyun Deng, Mubashara Akhtar, Rami Aly, Zhijiang Guo, Christos Christodoulopoulos, Oana Cocarascu, Arpit Mittal, James Thorne, and Andreas Vlachos (Eds.). Association for Computational Linguistics, Miami, Florida, USA, 264–279. <https://doi.org/10.18653/v1/2024.fever-1.28>
- [14] Liwen Sun, Jialun Zhao, Wenjing Han, and Chenyan Xiong. 2025. Fact-Aware Multimodal Retrieval Augmentation for Accurate Medical Radiology Report Generation. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, Luis Chiruzzo, Alan Ritter, and Lu Wang (Eds.). Association for Computational Linguistics, Albuquerque, New Mexico, 643–655. <https://doi.org/10.18653/v1/2025.naacl-long.28>
- [15] Junde Wu, Jiayuan Zhu, Yunli Qi, Jingkun Chen, Min Xu, Filippo Menolascina, Yueming Jin, and Vicente Grau. 2025. Medical Graph RAG: Evidence-based Medical Large Language Model via Graph Retrieval-Augmented Generation. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (Eds.). Association for Computational Linguistics, Vienna, Austria, 28443–28467. <https://doi.org/10.18653/v1/2025.acl-long.1381>