

BioHiCL: Hierarchical Multi-Label Contrastive Learning for Biomedical Retrieval with MeSH Labels

Mengfei Lan, Lecheng Zheng, Halil Kilicoglu*

University of Illinois Urbana-Champaign
{mlan3, lecheng4, halil}@illinois.edu

Abstract

Effective biomedical information retrieval requires modeling domain semantics and hierarchical relationships among biomedical texts. Existing biomedical generative retrievers build on coarse binary relevance signals, limiting their ability to capture semantic overlap. We propose BioHiCL (*Biomedical Retrieval with Hierarchical Multi-Label Contrastive Learning*), which leverages hierarchical MeSH annotations to provide structured supervision for multi-label contrastive learning. Our models, BioHiCL-Base (0.1B)¹ and BioHiCL-Large (0.3B)², achieve promising performance on biomedical retrieval, sentence similarity, and question answering tasks, while remaining computationally efficient for deployment.

1 Introduction

Dense retrievers learn vector representations of text from large-scale unsupervised, supervised, and synthetic corpora, achieving strong performance on general-domain information retrieval (IR) benchmarks (Ni et al., 2022a; Wang et al., 2024; Xiao et al., 2024; Lassance et al., 2024). However, these models often fail to capture biomedical-specific semantics, limiting their effectiveness in tasks involving specialized terminology.

To bridge this gap, various biomedical IR models have been proposed. Encoder-based biomedical language models, such as BiomedBERT (Chakraborty et al., 2020), pretrain transformer encoders on large-scale biomedical corpora to capture domain-specific terminology and semantics, but their retrieval effectiveness is limited by the availability of retrieval task-specific supervision. More recently, large-scale biomedical IR models, including MedCPT and BMRetriever (Jin et al., 2023; Xu et al., 2024), leverage contrastive

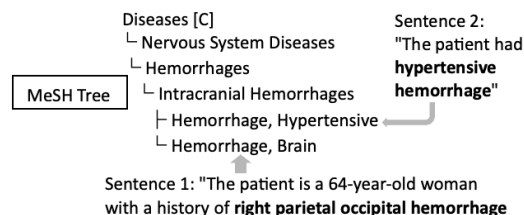


Figure 1: The MeSH labels for a pair of sentences annotated as "neutral" (irrelevant) in MedNLI.

learning on extensive biomedical semantic-related text pairs to achieve improved semantic alignment. Despite their effectiveness, these models depend on computationally expensive training and coarse-grained relevance signals (e.g., MedNLI-derived entailment labels (Xu et al., 2024) or query–article clicks (Jin et al., 2023)), making it challenging to capture the fine-grained and partially overlapping semantics characteristic of biomedical texts.

Semantic relationships in biomedical text are complex. Documents and sentences often share partial semantic overlap that cannot be adequately modeled with binary relevance signals alone, limiting the ability of dense retrievers to learn fine-grained biomedical representations. Medical Subject Headings (MeSH) offer a natural source of multi-aspect supervision. MeSH is a hierarchically organized controlled vocabulary curated by domain experts (U.S. National Library of Medicine, 2024), explicitly encoding both concept overlap and hierarchical structure. For example, Figure 1 shows a sentence pair labeled as *neutral* in MedNLI, while their MeSH annotations share a common parent concept (*Intracranial Hemorrhages*) in the disease hierarchy indicating the relevance between sentences. This hierarchy-aware signal reveals semantic relatedness beyond binary relevance, offering richer supervision for representation learning.

In this work, we propose **BioHiCL** (*Biomedical Retrieval with Hierarchical Multi-Label Contrastive Learning*), a framework that leverages partially overlapping and hierarchical MeSH annotations to supervise dense retrieval. Using expert-

*Corresponding author

¹<https://huggingface.co/LunaLan07/BioHiCL-Base>

²<https://huggingface.co/LunaLan07/BioHiCL-Large>

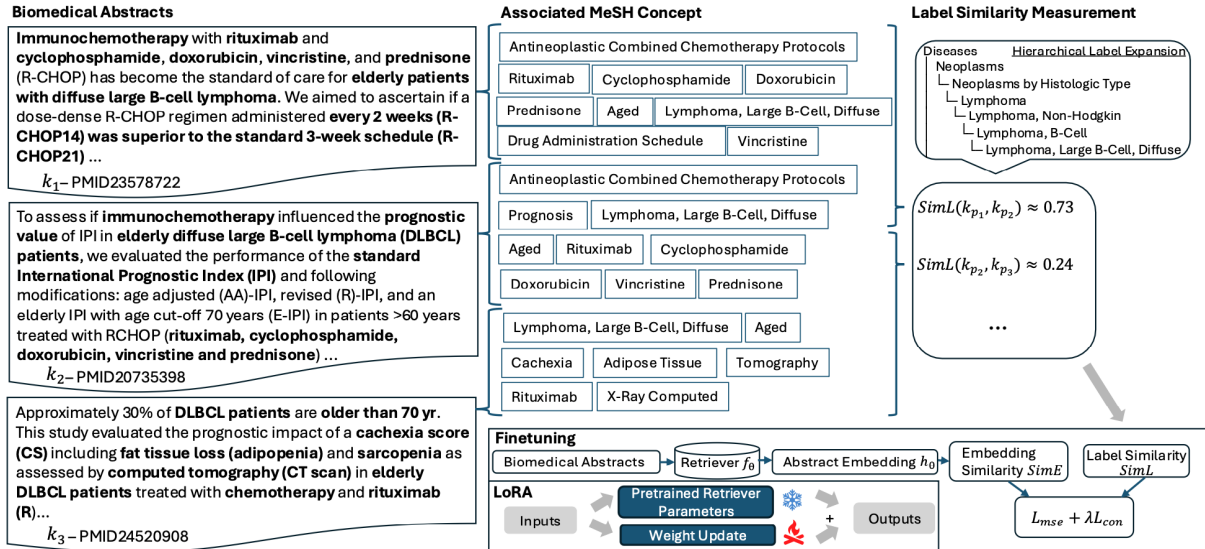


Figure 2: Overview of BioHiCL. During LoRA fine-tuning, embedding similarity (SimE) is aligned with MeSH-based label similarity (SimL), so that abstract pairs with greater MeSH overlap (k_1, k_2) are represented as more similar than pairs with less overlap, such as (k_2, k_3).

curated MeSH labels from the BioASQ Task 1a benchmark (Nentidis et al., 2022), BioHiCL adapts a general-domain dense retriever to the biomedical domain via parameter-efficient LoRA fine-tuning (Hu et al., 2022), transferring broad language understanding while enabling fine-grained biomedical representation learning.

Our contributions are summarized as follows: (a) We introduce a hierarchical multi-label contrastive learning framework that models graded semantic overlap between texts by aligning embedding similarity with depth-aware label similarity, going beyond binary or instance-level contrastive objectives; (b) we develop two models, BioHiCL-Base (0.1B) and BioHiCL-Large (0.3B), trained using structured MeSH supervision; (c) experiments on biomedical retrieval, sentence similarity, and question answering demonstrate improvements over the baselines; and (d) efficiency analysis shows that BioHiCL models run efficiently on a single A100 40GB GPU, supporting practical deployment³.

2 Methodology

In this section, we introduce **BioHiCL**, a biomedical retrieval framework that aligns embedding similarity with depth-weighted MeSH semantic similarity via hierarchical multi-label contrastive learning. Figure 2 provides an overview of the framework.

Problem Statement. Given a collection of biomedical abstracts $\mathcal{K} = \{k_1, k_2, \dots, k_n\}$ with MeSH annotations $\mathcal{M} = \{m_1, m_2, \dots, m_n\}$, our

goal is to learn an encoder $f_\theta : \mathcal{K} \rightarrow \mathbb{R}^d$ that maps each abstract k_i to a d -dimensional embedding $h_i = f_\theta(k_i)$. The embedding space should reflect semantic relationships among abstracts, so that the similarity between embeddings corresponds to MeSH-based label space similarity.

Hierarchical MeSH Label Representation.

Each abstract k_i is annotated with a multi-label set of major MeSH headings m_i , drawn from the hierarchically structured MeSH ontology with 16 branches⁴. Deeper nodes in the branch represent more specific biomedical concepts. To incorporate this hierarchical structure, we expand each label set m_i to include the ancestors of each MeSH label in the multi-label set, resulting in m_i^{hier} . This expanded set is then encoded as a multi-hot vector $y_i \in \{0, 1\}^C$, where C is the total number of MeSH concepts considered. We use $\mathcal{C} = \{c_1, c_2, \dots, c_C\}$ to represent the full set of considered MeSH labels associated in the hot-vectors. Notice that $y_{i,j} = 1$ indicates the presence of the j -th concept c_j in m_i^{hier} . To reflect specificity, each concept $c_j \in \mathcal{C}$ is assigned a depth-based weight:

$$w_j = \log(d(c_j) + 1), \quad (1)$$

where $d(c_j)$ denotes the depth of label c_j in the MeSH hierarchy. The depth-oriented label weighting set $\mathbf{w} = \{w_1, w_2, \dots, w_C\}$ allows the model to emphasize matches on more specific, deeper labels over higher-level, more general concepts.

³Code released at: <https://github.com/MengfeiLan/BioHiCL>

⁴[Link](#) to download the MeSH hierarchy.

Each abstract k_i is encoded using the pretrained general-domain dense retriever BGE (Xiao et al., 2024), yielding a sentence embedding $h_i \in \mathbb{R}^d$ from the model’s pooled output. For any pair of abstracts k_p and k_q , we define **Embedding Similarity** (similarity in embedding space) as $\text{SimE}(k_p, k_q) = \cos(h_p, h_q)$, and **Label Similarity** (similarity in MeSH label space, weighted by \mathbf{w}) as $\text{SimL}(k_p, k_q) = \cos(y_p \odot \mathbf{w}, y_q \odot \mathbf{w})$, where \odot denotes element-wise multiplication. By weighting MeSH labels with \mathbf{w} , the model gives more importance to matches on specific concepts, aligning embedding similarity with meaningful biomedical semantic overlap.

To encourage the embeddings to capture MeSH-informed semantic structure, we train the model so that the embedding similarity $\text{SimE}(\cdot)$ reflects the label-based similarity $\text{SimL}(\cdot)$. We achieve this alignment using a filtered regression objective:

$$\mathcal{L}_{\text{mse}} = \sum_{(i,j) \in \mathcal{P}} \text{MSE}(\text{SimE}(k_i, k_j), \text{SimL}(k_i, k_j)) \quad (2)$$

where $\mathcal{P} = \{(i, j) : i \neq j, \text{SimL}(k_i, k_j) > \beta\}$. The threshold $\text{SimL}(\cdot) > \beta$ ensures that only abstract pairs with meaningful topic overlap in the MeSH label space are selected. This avoids noisy supervision from weakly related or unrelated documents, for which label similarity provides little semantic guidance.

Hierarchical Multi-Label Contrastive Learning. While the regression loss \mathcal{L}_{mse} encourages embedding similarity to reflect MeSH-based label similarity, it alone may lead to trivial solutions where embeddings collapse to a single point (Shen et al., 2022). To address this, we introduce a hierarchy-aware contrastive loss that explicitly separates semantically related and unrelated abstracts in the embedding space (Lan et al., 2024; Zheng et al., 2022, 2024). We define positive and negative pairs based on hierarchical MeSH similarity:

- **Positive pairs** (k_i, k_i^+): abstracts satisfying $\text{SimL}(k_i, k_i^+) > \beta$, where the degree of similarity reflects both shared concepts and hierarchical proximity in the MeSH tree.
- **Negative pairs** (k_i, k_j^-): abstracts sharing no labels, i.e., $\text{SimL}(k_i, k_j^-) = 0$, representing unrelated concepts.

We formulate the hierarchical multi-label contrastive learning as follows:

$$\mathcal{L}_{\text{con}} = -\mathbb{E}_i \log \left[\frac{\text{SimL}(k_i, k_i^+) \exp(\text{SimE}(k_i, k_i^+))}{\sum_{k_j^-} \exp(\text{SimE}(k_i, k_j^-))} \right]. \quad (3)$$

Notably, the positive pairs are weighted by their label similarity $\text{SimL}(k_i, k_i^+)$, which incorporates hierarchical information: matches on deeper, more specific MeSH concepts contribute more to the similarity than matches on higher-level, general concepts. This makes the contrastive learning hierarchical, as the model explicitly considers the structure and specificity of the labels when shaping the embedding space.

Parameter-Efficient Fine-Tuning with LoRA.

The final training objective combines the regression and contrastive losses to jointly align embeddings with hierarchical label similarity while maintaining discriminative power:

$$\mathcal{L} = \mathcal{L}_{\text{mse}} + \lambda \mathcal{L}_{\text{con}}. \quad (4)$$

To adapt the pretrained encoder to biomedical text efficiently, we use Low-Rank Adaptation (LoRA) (Hu et al., 2022). All original parameters are frozen, and trainable low-rank adapters $A^{(l)}, B^{(l)}$ are injected into layers:

$$W_{\text{adapted}}^{(l)} = W^{(l)} + B^{(l)} A^{(l)}, \quad r \ll \min(d, k). \quad (5)$$

3 Experiments and Results

3.1 Datasets and Baseline Models

To adapt the general-domain retrievers, BAAI/bge-base-en-v1.5 and BAAI/bge-large-en-v1.5, to biomedical text, we fine-tune the retrievers on 80K abstracts randomly sampled from the latest BioASQ Task 1a release (v2022)⁵, covering 29,681 unique MeSH terms (12.68 per abstract on average). For model selection, we validate on the TREC 2022 Clinical Trials (TREC-CT) benchmark (Roberts et al., 2022)⁶, which matches patient case descriptions to relevant clinical trials using graded relevance judgments. We select the checkpoint that achieves the lowest relevance loss on the query–document pairs from the benchmark.

We evaluate the tuned BioHiCL models across a diverse set of biomedical benchmarks: **information retrieval**, measured by NDCG@10, including NFCorpus (Boteva et al., 2016), TREC-COVID (Voorhees et al., 2021), SciFact (Wadden et al., 2022), and SCIDOCs (Cohan et al., 2020); **sentence similarity**, evaluated via Spearman correlation between gold labels and embedding similarities, covering BIOSSES (Soğancıoğlu et al.,

⁵<https://participants-area.bioasq.org/datasets/>

⁶<https://www.trec-cds.org/2022.html>

Models		Information Retrieval (nDCG@10)					Sentence Similarity (Spearman Correlation)		Question Answering (Recall@1)
Name	Size	NFCorpus	TREC- COVID	SciFact	SCIDOCS	IR Avg	BIOSSES	SciFact (Sentence)	PubMedQA
General Domain									
Contriever (Lei et al., 2023)	0.1B	0.328	0.596	0.677	0.165	0.442	0.833	0.265	0.479
bge-base-en-v1.5 (Xiao et al., 2024)	0.1B	0.368	0.798	0.735	0.215	0.529	0.860	0.339	0.856
SpladeV3 (Lassance et al., 2024)	0.3B	0.357	0.748	0.710	0.158	0.493	0.827	0.303	0.844
bge-large-en-v1.5 (Xiao et al., 2024)	0.3B	0.369	0.748	0.735	0.209	0.515	0.844	0.346	0.871
e5-large-v2 (Wang et al., 2022)	0.3B	0.371	0.665	0.726	0.201	0.491	0.830	0.294	0.717
gtr-t5-xl (Ni et al., 2022b)	1.2B	0.343	0.580	0.635	0.159	0.429	0.789	0.203	0.685
SGPT-1.3B (Muennighoff, 2022)	1.3B	0.320	0.730	0.682	0.162	0.474	0.830	0.214	0.754
Biomedical Domain									
BiomedBERT (Chakraborty et al., 2020)	0.1B	0.049	0.190	0.262	0.021	0.131	0.791	0.211	0.134
MedCPT-Query (Jin et al., 2023)	0.1B	0.340	0.697	0.724	0.123	0.471	0.837	0.298	0.834
BMRetriever-410M (Xu et al., 2024)	0.4B	0.321	<u>0.831</u>	0.711	0.167	0.501	0.840	0.121	0.879
BMRetriever-1B (Xu et al., 2024)	1B	0.344	0.840	0.760	0.180	0.531	0.858	0.107	0.810
BiCA-small (Sinha et al., 2025)	0.03B	0.347	0.661	0.727	0.214	0.487	0.872	0.318	0.856
BiCA-base (Sinha et al., 2025)	0.1B	0.378	0.684	<u>0.762</u>	<u>0.231</u>	0.514	<u>0.880</u>	0.335	0.868
BioHiCL-base (ours)	0.1B	<u>0.379</u>	0.812	0.757	0.225	0.543	0.896	<u>0.350</u>	<u>0.893</u>
BioHiCL-large (ours)	0.3B	0.385	0.765	0.765	0.228	<u>0.534</u>	0.868	0.359	0.898

Table 1: Overall evaluation results. Bolds and underlines indicate the first- and second-best for each evaluation.

2017) and SciFact sentence-level labeling (Wadden et al., 2022); and **question answering**, using PubMedQA (Jin et al., 2019). We present further details of the evaluation datasets in Appendix A.

We compare our approach against several publicly available dense retrieval models ranging from 0.1B to 1.5B parameters designed for computationally efficient large-scale retrieval, as presented in Table 1. In the experiment, we set $\lambda = 0.1$ and $\beta = 0.3$. Additional hyperparameter analysis over λ and β is presented in Appendix B.3, and further implementation details are included in Appendix B.

3.2 Experiment Results

Our models, BioHiCL-Base and BioHiCL-Large, achieve strong performance across biomedical benchmarks (Table 1). BioHiCL-Base attains the highest IR average of 0.543, with second best scores on NFCorpus (0.379) and SCIDOCS (0.225), while BioHiCL-Large ranks second with an IR average of 0.534 and the best NFCorpus score (0.385) and SciFact score (0.765). In sentence similarity, BioHiCL-Base leads on BIOSSES (0.896), and BioHiCL-Large performs best on SciFact sentences (0.359). On PubMedQA, BioHiCL-Large achieves the best Recall@1 score as 0.898, and BioHiCL-Base reaches the second best as 0.893. Across tasks, BioHiCL remains competitive with both general-domain and biomedical-specific models, highlighting its cross-task robustness.

Computational efficiency. Despite containing only 0.1B parameters, BioHiCL-Base achieves performance comparable to larger models such as BMRetriever-1B (1B parameters). Scaling to BioHiCL-large (0.3B) yields further gains on some

benchmarks, including sentence similarity on SciFact, while preserving strong overall IR performance. Appendix C presents a detailed analysis of runtime and computational cost, demonstrating that the BioHiCL models are efficient to run on a single A100 GPU and suitable for large-scale real-world application with time efficiency.

Challenges of instruction-based retrievers.

Instruction-based retrievers rely on task-specific prompts to generate effective embeddings. This dependence can limit their ability to generalize to unseen tasks, where an appropriate prompt may not be readily available. In our experiments, models such as BMRetriever and gtr-t5-xl show reduced performance on sentence similarity when task-specific prompting is not provided, suggesting sensitivity to prompt design in zero-shot settings.

3.3 Ablation Study

We conduct ablations by removing each component of BioHiCL individually while keeping all other settings fixed, to identify the source of performance gains. Results are shown in Table 2. All four structural components (ancestor label, depth weighting, L_{MSE} , L_{Con}) lead to performance drops when removed, demonstrating their effectiveness. The impact varies across components: removing the hierarchical contrastive loss L_{Con} causes the largest drop (IR Avg 0.528 vs. 0.543), showing its importance in maintaining discriminative structure and preventing embedding collapse. Removing ancestor label expansion also leads to a decrease, indicating that incorporating hierarchical labels provides further signal beyond flat multi-label overlap. We also compare LoRA with full fine-tuning, where

Model	NFC	TREC COVID	SciFact	SCIDOCS	IR Avg
BioHiCL (0.1B)	0.379	0.812	0.757	0.225	0.543
w/o Ancestor Label	0.375	0.804	0.752	0.219	0.538
w/o Depth-based Weight	0.378	0.810	0.754	0.223	0.541
w/o L_{MSE}	0.376	0.798	0.753	0.221	0.537
w/o L_{Con}	0.367	0.782	0.747	0.215	0.528
w/o LoRA	0.380	0.805	0.760	0.224	0.542

Table 2: Ablation study.

Model	NFC	TREC COVID	SciFact	SCIDOCS	IR Avg
e5-large-v2	0.371→ 0.374	0.665→ 0.683	0.726→ 0.746	0.201→ 0.206	0.491→ 0.502
bge-base-en-v1.5	0.368→ 0.379	0.798→ 0.812	0.735→ 0.757	0.215→ 0.225	0.529→ 0.543
BMRetriever-410m	0.321→ 0.105	0.831→ 0.417	0.711→ 0.541	0.167→ 0.053	0.501→ 0.279
BiCA-small	0.347→ 0.348	0.661→ 0.672	0.727→ 0.731	0.214→ 0.214	0.487→ 0.491

Table 3: Additional fine-tuning baselines. “→” denotes performance before and after fine-tuning.

LoRA updates only 0.3% of the model parameters to reduce memory and computational cost while maintaining comparable performance.

3.4 Additional Fine-Tuning Baselines

We fine-tune additional information retrieval models on MeSH labels using the same procedure as BioHiCL, with results summarized in Table 3. Overall, general-domain models benefit from fine-tuning, showing consistent improvements (e.g., e5: IR Avg 0.491 → 0.502; bge: 0.529 → 0.543), indicating that task-specific supervision helps adapt embeddings to the biomedical domain. Biomedical-domain models exhibit limited improvements, with BiCA-small improving only marginally (0.487 → 0.491), suggesting it is already well aligned with domain-specific semantics. BMRetriever-410M has substantial performance degradation (0.501 → 0.279) because of objective mismatch, as replacing its original instruction-based training objective with a hierarchical contrastive loss may override its retrieval-specialized embedding geometry.

4 Related Works

Biomedical IR has been studied over the years, with early methods focusing on rule-based indexing and ranking (Robertson et al., 2009; Jiang and Zhai, 2007; Edinger et al., 2018). To go beyond exact term matching and better capture semantic meaning, neural approaches have been developed to learn dense representations in an unsupervised manner (Lee et al., 2020; Gu et al., 2020; Chen et al., 2019, 2020; Mohan et al., 2018). Supervi-

sion from biomedical IR datasets further improves these models: MedCPT (Jin et al., 2023) uses query–document click data to guide relevance modeling, BMRetriever (Xu et al., 2024) leverages natural language inference and question-answering datasets as supervision, and BiCA (Sinha et al., 2025) incorporates citation-based signals to enhance domain-specific dense retrieval. However, no previous work has tried MeSH-based supervision as textual relatedness signals.

5 Conclusions

We introduce BioHiCL, a biomedical dense retriever leveraging hierarchical MeSH supervision to capture fine-grained semantic overlap. It learns rich biomedical representations while preserving retrieval ability by adapting general-domain retrievers via LoRA fine-tuning. Experiments demonstrate the effectiveness and efficiency of BioHiCL.

Although our evaluation focuses on biomedical data, the underlying principle of BioHiCL, which is to align embedding spaces with hierarchical multi-label structures, is not inherently tied to MeSH. Similar hierarchical supervision signals naturally exist in other domains, such as Wikipedia category graphs (Milne and Witten, 2008) and e-commerce product taxonomies (Garza et al., 2020), where multi-label structures encode graded semantic relationships. Future exploration could extend hierarchical multi-label learning beyond the biomedical domain for dense retrieval.

6 Limitations

BioHiCL relies on the availability and quality of expert-curated hierarchical labels (MeSH), which may not exist or may be incomplete, noisy, or inconsistently applied in other domains, limiting the general applicability of the framework. In addition, the fixed hierarchy and depth-based weighting assume that semantic specificity is well captured by the ontology structure, which may not fully reflect contextual or task-dependent semantic relevance.

Acknowledgements

This work was supported by the National Library of Medicine of the National Institutes of Health under the award number R01LM014079. The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health. The funder had no role in considering the study design or in the collection, analysis, interpretation of data, writing of the report, or decision to submit the article for publication.

References

- Vera Boteva, Demian Gholipour, Artem Sokolov, and Stefan Riezler. 2016. A full-text learning to rank dataset for medical information retrieval. In *European Conference on Information Retrieval*, pages 716–722. Springer.
- Souradip Chakraborty, Ekaba Bisong, Shweta Bhatt, Thomas Wagner, Riley Elliott, and Francesco Mosconi. 2020. Biomedbert: A pre-trained biomedical language model for qa and ir. In *Proceedings of the 28th international conference on computational linguistics*, pages 669–679.
- Qingyu Chen, Kyubum Lee, Shankai Yan, Sun Kim, Chih-Hsuan Wei, and Zhiyong Lu. 2020. Bioconceptvec: creating and evaluating literature-based biomedical concept embeddings on a large scale. *PLoS computational biology*, 16(4):e1007617.
- Qingyu Chen, Yifan Peng, and Zhiyong Lu. 2019. Biosentvec: creating sentence embeddings for biomedical texts. In *2019 IEEE International Conference on Healthcare Informatics (ICHI)*, pages 1–5. IEEE.
- 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.
- Tracy Edinger, Dina Demner-Fushman, Aaron M Cohen, Steven Bedrick, and William Hersh. 2018. Evaluation of clinical text segmentation to facilitate cohort retrieval. In *AMIA Annual Symposium Proceedings*, volume 2017, page 660.
- José Luis Garza, Ralph Peeters, and Christian Bizer. 2020. **Wdc-24 gold standard for product categorization**.
- Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. 2020. **Domain-specific language model pretraining for biomedical natural language processing**.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, and 1 others. 2022. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3.
- Jing Jiang and ChengXiang Zhai. 2007. An empirical study of tokenization strategies for biomedical information retrieval. *Information Retrieval*, 10(4):341–363.
- Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. 2019. Pubmedqa: A dataset for biomedical research question answering. In *Proceedings of the 2019 conference on empirical methods in natural language processing and the 9th international joint conference on natural language processing (EMNLP-IJCNLP)*, pages 2567–2577.
- Qiao Jin, Won Kim, Qingyu Chen, Donald C Comeau, Lana Yeganova, W John Wilbur, and Zhiyong Lu. 2023. Medcpt: Contrastive pre-trained transformers with large-scale pubmed search logs for zero-shot biomedical information retrieval. *Bioinformatics*, 39(11):btad651.
- Mengfei Lan, Lecheng Zheng, Shufan Ming, and Halil Kilicoglu. 2024. Multi-label sequential sentence classification via large language model. In *Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA, November 12–16, 2024*, Findings of EMNLP, pages 16086–16104. Association for Computational Linguistics.
- Carlos Lassance, Hervé Déjean, Thibault Formal, and Stéphane Clinchant. 2024. Splade-v3: New baselines for splade. *arXiv preprint arXiv:2403.06789*.
- Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2020. Biobert: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4):1234–1240.
- Yibin Lei, Liang Ding, Yu Cao, Changtong Zan, Andrew Yates, and Dacheng Tao. 2023. **Unsupervised dense retrieval with relevance-aware contrastive pre-training**. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 10932–10940. Toronto, Canada. Association for Computational Linguistics.
- Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, Benjamin Bossan, and Marian Tietz. 2022. PEFT: State-of-the-art parameter-efficient fine-tuning methods. <https://github.com/huggingface/peft>.
- David Milne and Ian H Witten. 2008. Learning to link with wikipedia. In *Proceedings of the 17th ACM conference on Information and knowledge management*, pages 509–518.
- Sunil Mohan, Nicolas Fiorini, Sun Kim, and Zhiyong Lu. 2018. A fast deep learning model for textual relevance in biomedical information retrieval. In *Proceedings of the 2018 world wide web conference*, pages 77–86.
- Niklas Muennighoff. 2022. Sgpt: Gpt sentence embeddings for semantic search. *arXiv preprint arXiv:2202.08904*.
- Anastasios Nentidis, Georgios Katsimpras, Eirini Vantorou, Anastasia Krithara, Antonio Miranda-Escalada, Luis Gasco, Martin Krallinger, and Georgios Paliouras. 2022. Overview of bioasq 2022: the tenth bioasq challenge on large-scale biomedical semantic indexing and question answering. In *International conference of the cross-language evaluation forum for European languages*, pages 337–361. Springer.

- Jianmo Ni, Chen Qu, Jing Lu, Zhuyun Dai, Gustavo Hernandez Abrego, Ji Ma, Vincent Zhao, Yi Luan, Keith Hall, Ming-Wei Chang, and 1 others. 2022a. Large dual encoders are generalizable retrievers. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9844–9855.
- Jianmo Ni, Chen Qu, Jing Lu, Zhuyun Dai, Gustavo Hernandez Abrego, Ji Ma, Vincent Zhao, Yi Luan, Keith Hall, Ming-Wei Chang, and Yinfei Yang. 2022b. [Large dual encoders are generalizable retrievers](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9844–9855, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Kirk Roberts, Dina Demner-Fushman, Ellen M Voorhees, Steven Bedrick, and William R Hersh. 2022. Overview of the trec 2022 clinical trials track. In *TREC*.
- Stephen Robertson, Hugo Zaragoza, and 1 others. 2009. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends® in Information Retrieval*, 3(4):333–389.
- Kendrick Shen, Robbie M Jones, Ananya Kumar, Sang Michael Xie, Jeff Z HaoChen, Tengyu Ma, and Percy Liang. 2022. Connect, not collapse: Explaining contrastive learning for unsupervised domain adaptation. In *International conference on machine learning*, pages 19847–19878. PMLR.
- Aarush Sinha, Roshan Balaji, Nirav Pravinbhai Bhatt, and 1 others. 2025. Bica: Effective biomedical dense retrieval with citation-aware hard negatives. *AAAI 2026*.
- Gizem Soğancıoğlu, Hakime Öztürk, and Arzucan Özgür. 2017. Biosses: a semantic sentence similarity estimation system for the biomedical domain. *Bioinformatics*, 33(14):i49–i58.
- 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 (Round 2)*.
- U.S. National Library of Medicine. 2024. Medical subject headings (mesh). <https://www.nlm.nih.gov/mesh/meshhome.html>.
- Ellen Voorhees, Tasmee Alam, Steven Bedrick, Dina Demner-Fushman, William R Hersh, Kyle Lo, Kirk Roberts, Ian Soboroff, and Lucy Lu Wang. 2021. Trec-covid: constructing a pandemic information retrieval test collection. In *ACM SIGIR Forum*, volume 54, pages 1–12. ACM New York, NY, USA.
- David Wadden, Kyle Lo, Bailey Kuehl, Arman Cohan, Iz Beltagy, Lucy Lu Wang, and Hannaneh Hajishirzi. 2022. Scifact-open: Towards open-domain scientific claim verification. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 4719–4734.
- 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. [Improving text embeddings with large language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11897–11916, Bangkok, Thailand. Association for Computational Linguistics.
- Shitao Xiao, Zheng Liu, Peitian Zhang, Niklas Muenighoff, Defu Lian, and Jian-Yun Nie. 2024. C-pack: Packed resources for general chinese embeddings. In *Proceedings of the 47th international ACM SIGIR conference on research and development in information retrieval*, pages 641–649.
- Ran Xu, Wenqi Shi, Yue Yu, Yuchen Zhuang, Yanqiao Zhu, May Dongmei Wang, Joyce C Ho, Chao Zhang, and Carl Yang. 2024. Bmretriever: Tuning large language models as better biomedical text retrievers. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 22234–22254.
- Lecheng Zheng, Baoyu Jing, Zihao Li, Hanghang Tong, and Jingrui He. 2024. Heterogeneous contrastive learning for foundation models and beyond. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD 2024, Barcelona, Spain, August 25-29, 2024*, pages 6666–6676. ACM.
- Lecheng Zheng, Jinjun Xiong, Yada Zhu, and Jingrui He. 2022. Contrastive learning with complex heterogeneity. In *KDD '22: The 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Washington, DC, USA, August 14 - 18, 2022*, pages 2594–2604. ACM.

A Evaluation Datasets

A.1 Information Retrieval

We leverage the BEIR framework (Thakur et al., 2021) to conduct a unified evaluation of retrievers across four biomedical IR benchmarks, consistent with the previous biomedical IR tasks (Jin et al., 2023; Xu et al., 2024; Sinha et al., 2025).

A.2 Sentence Similarity

BIOSSES. BIOSSES (Soğancıoğlu et al., 2017) is a benchmark dataset for evaluating biomedical sentence similarity. It consists of sentence pairs sampled from PubMed abstracts, each annotated by biomedical experts with a similarity score ranging

from 0 (completely dissimilar) to 4 (semantically equivalent). Consistent with previous work, we encode each sentence independently using the sentence encoder and compute the cosine similarity between the resulting embeddings for each sentence pair. Model performance is measured using the Spearman correlation between the predicted similarity scores and the gold-standard human annotations, assessing the model’s ability to capture graded semantic similarity between biomedical sentences.

SciFact-Sentence. Beyond the document-level annotations, SciFact (Wadden et al., 2022) provides fine-grained sentence-level retrieval annotations that enable the evaluation of sentence similarity (Wadden et al., 2022). We construct claim–sentence pairs by pairing each claim with all sentences from its annotated evidence abstract, assigning a positive label to sentences appearing in the claim’s evidence attribute and a negative label to all others. Both claims and sentences are encoded with the pretrained models, and cosine similarity is computed for each pair. Performance is quantified using the Spearman correlation between the similarity scores and the binary labels, measuring the model’s ability to assign higher similarity to true evidence sentences.

A.3 Question-Answering

PubMedQA. We evaluate retrieval performance on PubMedQA (Jin et al., 2019) by formulating the task as large-scale question-to-abstract retrieval. We use the `pqa_labeled` split, which contains 1,000 biomedical questions paired with gold-standard evidence abstracts manually annotated by domain experts. These questions serve as retrieval queries.

To construct a challenging retrieval corpus, we augment the evidence set from `pqa_labeled` split with an additional 211K PubMed abstracts drawn from the `pqa_artificial` split. Because the `pqa_artificial` split lacks human-curated question–evidence annotations, we do not use the corresponding questions from the split as retrieval queries. The resulting retrieval corpus contains approximately 212K abstracts.

For each question, we independently encode the question and all collected abstracts using the evaluated retrieval model and compute cosine similarity between the question embedding and abstract embeddings. Abstracts are ranked by similarity, and performance is measured using Recall@1, which

evaluates whether the top-ranked retrieved abstract matches the gold-standard evidence associated with the question.

B Experiment Setup

B.1 Baseline IR Models

We consider both effectiveness and computational practicality when selecting baseline models. Specifically, we include the state-of-the-art retrieval models that demonstrate competitive performance on general-domain and biomedical IR benchmarks, while remaining feasible for real-world deployment. Computational efficiency is a critical factor for retrieval systems operating over large-scale corpora, where excessive latency or hardware requirements can limit applicability.

Some large-scale retrieval models require multiple GPUs and extended runtime even for evaluation on standard benchmarks such as BEIR, making them impractical for large-scale or time-sensitive applications. For example, although `intfloat/e5-small-v2` achieves strong performance on BEIR, its evaluation requires multiple powerful GPUs over several days⁷. Similarly, recent biomedical retrievers with billions of parameters introduce substantial inference latency and computational cost (Xu et al., 2024).

In this work, we therefore restrict our comparisons to state-of-the-art baselines that offer a balance between retrieval quality and efficiency, and that can be evaluated on a single NVIDIA A100 GPU with 40GB memory. The selected baselines span model sizes from 0.1B to 1.5B parameters, reflecting realistic deployment constraints for large-scale retrieval:

General Domain IR Baselines.

- **Contriever** (Lei et al., 2023)⁸: An unsupervised dense retriever that leverages contrastive learning to align queries and documents in the embedding space.
- **SpladeV3** (Lassance et al., 2024)⁹: A sparse lexical-semantic retriever that transforms input text into high-dimensional sparse representations, combining the advantages of classical bag-of-words and dense embeddings.

⁷<https://github.com/microsoft/unilm/blob/master/e5>

⁸<https://huggingface.co/facebook/contriever>

⁹<https://huggingface.co/naver/splade-v3>

- bge-base-en-v1.5 (Xiao et al., 2024)¹⁰: A bi-encoder model trained to generate general-purpose sentence embeddings for semantic search and retrieval tasks.
- bge-large-en-v1.5 (Xiao et al., 2024)¹¹: A larger version of bge-base, offering higher embedding capacity and improved retrieval quality to enable better semantic matching for longer and more complex queries.
- e5-large-v2 (Wang et al., 2022)¹²: A dense text embedding model that produces embeddings optimized for cross-encoder retrieval pipelines.
- gtr-t5-xl (Ni et al., 2022b)¹³: A sequence-to-sequence Transformer model adapted for retrieval tasks by generating embeddings that capture query-document similarity.
- SGPT-1.3B (Muennighoff, 2022)¹⁴: A generative pre-trained Transformer for sentence embeddings.

Biomedical Domain

- BiomedBERT (Chakraborty et al., 2020)¹⁵: A domain-specific BERT model pretrained on biomedical text to capture medical terminology and context. It provides reasonable sentence embeddings but limited document-level retrieval performance.
- MedCPT-Query (Jin et al., 2023)¹⁶: A generative pre-trained biomedical retriever trained on query-document pairs from PubMed.
- BMRetriever-410M (Xu et al., 2024)¹⁷: A 410M-parameter biomedical dense retriever unsupervised pretrained on large biomedical corpora and finetuned on a mixture of human-annotated retrieval datasets.

¹⁰<https://huggingface.co/BAAI/bge-base-en-v1.5>

¹¹<https://huggingface.co/BAAI/bge-large-en-v1.5>

¹²<https://huggingface.co/intfloat/e5-large-v2>

¹³<https://huggingface.co/sentence-transformers/gtr-t5-xl>

¹⁴<https://huggingface.co/Muennighoff/SGPT-1.3B-weightedmean-msmarco-specb-bitfit>

¹⁵<https://huggingface.co/microsoft/BiomedNLP-BiomedBERT-base-uncased-abstract-fulltext>

¹⁶<https://huggingface.co/ncbi/MedCPT-Query-Encoder>

¹⁷<https://huggingface.co/BMRetriever/BMRetriever-410M>

β	λ	NFC	TREC COVID	SciFact	SCIDOCS	IR Avg
0.1	0.05	0.383	0.799	0.758	0.222	0.541
0.1	0.1	0.382	0.796	0.757	0.221	0.539
0.1	0.2	0.383	0.798	0.759	0.221	0.540
0.3	0.05	0.383	0.805	0.758	0.222	0.542
0.3	0.1	0.379	0.812	0.757	0.225	0.543
0.3	0.2	0.382	0.808	0.759	0.222	0.543
0.5	0.05	0.381	0.804	0.758	0.221	0.541
0.5	0.1	0.383	0.806	0.758	0.222	0.542
0.5	0.2	0.381	0.805	0.759	0.221	0.542

Table 4: Hyperparameter analysis.

- BMRetriever-1B (Xu et al., 2024)¹⁸: A larger 1B-parameter variant of BMRetriever, designed to improve retrieval accuracy by increasing model capacity.
- BiCA-small (Sinha et al., 2025)¹⁹: A compact biomedical dense retriever trained using citation-aware hard negative mining, where citation links among PubMed articles are used to construct challenging negative examples that improve discriminative ability.
- BiCA-base (Sinha et al., 2025)²⁰: The larger version of BiCA-small, offering better embedding quality and improved performance.

B.2 Training Parameters

In the experiment, we set $\lambda = 0.1$ and $\beta = 0.3$. We select BAAI/bge-base-en-v1.5 and BAAI/bge-large-en-v1.5 as the backbones due to their state-of-the-art performance on general-domain information retrieval tasks. We use NVIDIA A100 with 40GB GPU to fine-tune and evaluate the model on multiple retrieval datasets. We use PEFT package (Mangrulkar et al., 2022) to tune model using LoRA, and the model is trained with the AdamW optimizer with zero weight decay. We set the LoRA learning rate as $1e-5$, and set the batch size as 32, within which to perform the mean squared error loss and contrastive loss. The created BioHiCL models support an input length of up to 8192 tokens.

B.3 Parameter Analysis

We conduct a grid search over representative values of the hyperparameters ($\beta \in \{0.1, 0.3, 0.5\}$, $\lambda \in \{0.05, 0.1, 0.2\}$), and evaluate model performance on information retrieval benchmarks. As

¹⁸<https://huggingface.co/BMRetriever/BMRetriever-1B>

¹⁹<https://huggingface.co/bisectgroup/BiCA-small>

²⁰<https://huggingface.co/bisectgroup/BiCA-base>

Model	NFC	TREC-COVID	SciFact	SCIDOCS	IR Avg
BioHiCL (0.1B)	0.380 ± 0.003	0.812 ± 0.013	0.757 ± 0.012	0.225 ± 0.003	0.543 ± 0.005
BioHiCL (0.3B)	0.380 ± 0.006	0.759 ± 0.015	0.766 ± 0.014	0.223 ± 0.002	0.532 ± 0.005

Table 5: Variability of Model Training. Mean \pm 95% confidence interval (computed using the t-distribution with degrees of freedom=4) is reported in the table.

shown in Table B.3, the performance remains stable across these settings. For example, the IR average varies slightly from 0.539 to 0.543, sentence similarity metrics on BIOSSES range from 0.889 to 0.897, and QA performance on PubMedQA varies between 0.883 and 0.895. Notably, the originally chosen values ($\beta=0.3$, $\lambda=0.1$) achieve promising performance across most metrics, balancing retrieval and downstream task objectives effectively.

B.4 Variability of Training

We conducted three independent fine-tuning runs of BioHiCL using different random seeds and reported results as a mean \pm 95% confidence interval (computed using the t-distribution with degrees of freedom=4). The results are shown in Table 5. The results for BioHiCL reported in Table 1 are obtained from the publicly released pretrained checkpoint released on HuggingFace.

C Efficiency Analysis

Table 6 reports the latency and GPU memory consumption of the evaluated models. We measure three metrics: average corpus encoding latency (ms per document), average query encoding latency (ms per query), and average retrieval latency (ms per query). Peak GPU memory usage (MB) is also recorded during query and corpus encoding.

Latency and Computational Efficiency of BioHiCL. Our BioHiCL-base model achieves low latency with only 3.50 ms per document for corpus encoding and 0.63 ms per query for query encoding, which is comparable to other small models such as BiCA-small (3.35 ms/doc, 0.45 ms/query). BioHiCL-base also achieves low retrieval latency in its parameter class (0.24 ms/query), demonstrating efficient performance in deployment.

Similarly, BioHiCL-large maintains promising latency-performance among medium-sized models (0.3B parameters), with corpus encoding at 9.99 ms/doc and query encoding at 1.20 ms/query, which is competitive with models of similar size.

Its retrieval latency of 0.56 ms/query remains moderate, ensuring efficient search.

Both BioHiCL-base and BioHiCL-large are memory efficient (733 MB for the base model and 1.7 GB for the large model), consistent with other models of similar size. This ensures practical deployment of the models without excessive GPU requirements.

Comparison Across Model Sizes. Smaller and medium models (0.1B and 0.3B) exhibit low corpus and query encoding latency, while large models ($\geq 1B$) experience substantially higher latency, particularly in corpus encoding (up to 40 ms/doc for gtr-t5-xl and SGPT-1.3B).

Peak GPU memory usage scales approximately with model size. Small models ($\leq 0.1B$) require 730 MB, medium models (0.3B) require 1.7 GB, while large models ($> 0.3B$) demand 8 to 13 GB.

Large models incur higher latency and require more GPU resources, with corpus encoding exceeding 39 ms/doc and peak GPU memory usage ranging from 8 GB to over 13 GB. While large models may achieve strong retrieval performance, their computational demands make them less practical for low-latency or resource-constrained settings.

Model	Model Size	Peak GPU Memory Usage (MB)	Avg Corpus Encoding Latency (ms/doc)	Avg Query Encoding Latency (ms/query)	Avg Retrieval Latency (ms/query)
General Domain					
Contriever	0.1B	732.90	3.37	0.44	0.28
bge-base-en-v1.5	0.1B	732.82	3.32	0.47	0.27
SpladeV3	0.3B	732.90	3.36	0.45	0.31
bge-large-en-v1.5	0.3B	1689.62	9.47	1.00	0.49
e5-large-v2	0.3B	1689.51	9.46	0.99	0.47
gtr-t5-xl	1.2B	8772.50	39.76	3.52	0.49
SGPT-1.3B	1.3B	13213.93	39.39	3.19	1.44
Biomedical Domain					
BiomedBERT	0.1B	732.14	3.20	0.42	0.34
MedCPT-Query	0.1B	732.14	3.16	0.45	0.24
BMRetriever-410M	0.4B	9929.08	12.09	1.18	0.47
BMRetriever-1B	1B	15846.62	28.42	2.19	1.40
BiCA-small	0.03B	296.63	1.75	0.47	0.14
BiCA-base	0.1B	732.90	3.35	0.45	0.29
BioHiCL-base (ours)	0.1B	733.32	3.50	0.63	0.24
BioHiCL-large (ours)	0.3B	1691.68	9.99	1.20	0.56

Table 6: Latency and GPU memory usage for experiments on the SciFact information retrieval dataset, following the BEIR benchmark format. Batch size during encoding process is set as 16. Metrics include average corpus and query encoding latency, average retrieval latency, and peak GPU memory usage. Avg Corpus Encoding Latency (ms/doc): Time to encode one document from the corpus. Peak GPU Memory (MB) Usage: Peak GPU RAM during query and corpus encoding. Avg Query Encoding Latency (ms/query): Time to encode one query. Avg Retrieval Latency (ms/query): Time to retrieve top-K documents per query.