# **Contradiction Retrieval via Contrastive Learning with Sparsity**

Haike Xu<sup>\*1</sup> Zongyu Lin<sup>\*2</sup> Kai-Wei Chang<sup>2</sup> Yizhou Sun<sup>2</sup> Piotr Indyk<sup>1</sup>

# Abstract

Contradiction retrieval refers to identifying and extracting documents that explicitly disagree with or refute the content of a query, which is important to many downstream applications like fact checking and data cleaning. To retrieve contradiction argument to the query from large document corpora, existing methods such as similarity search and cross-encoder models exhibit different limitations. To address these challenges, we introduce a novel approach: SparseCL that leverages specially trained sentence embeddings designed to preserve subtle, contradictory nuances between sentences. Our method utilizes a combined metric of cosine similarity and a sparsity function to efficiently identify and retrieve documents that contradict a given query. This approach dramatically enhances the speed of contradiction detection by reducing the need for exhaustive document comparisons to simple vector calculations. We conduct contradiction retrieval experiments on Arguana, MSMARCO, and HotpotOA, where our method produces an average improvement of 11.0% across different models. We also validate our method on downstream tasks like natural language inference and cleaning corrupted corpora. This paper outlines a promising direction for nonsimilarity-based information retrieval which is currently underexplored.

# 1. Introduction

Training sentence embedding for similarity retrieval has been well studied in the literature (Gao et al., 2021; Xiong et al., 2020; Karpukhin et al., 2020), where a standard practice is to use contrastive learning to map those similar sentences together and those dissimilar sentences far from each other. However, these existing sentence embeddings are mainly tailored to similarity retrieval, while as far as we know, there hasn't been sentence embeddings for non-similarity based retrieval. In this paper, we study the problem of contradiction retrieval, a typical case of nonsimilarity based retrieval. Given a large document corpus and a query passage, the goal is to retrieve document(s) in the corpus that contradict the query, assuming they exist. This problem has a large number of applications, including counter-argument detection (Wachsmuth et al., 2018) and fact verification (Thorne et al., 2018). The standard approaches to retrieving contradictions are two-fold. One is to use a bi-encoder (Xiao et al., 2023; Li & Li, 2023; Li et al., 2023) that maps each document to a feature space such that two contradicting documents are mapped close to each other (e.g., according to the cosine metric) and use nearest neighbor search algorithms. The second approach is to train a cross-encoder model (Xiao et al., 2023) that determines whether two documents contradict each other, and apply it to each document or passage in the corpus.

Unfortunately, both methods suffer from limitations. The first approach (cosine similarity search on sentence embeddings) is inherently incapable of representing the "contradiction relation" between the documents, due to the fact that the cosine metric is "transitive" (See Appendix B for formal analysis): if A is similar to B, and B is similar to C, then A is also similar to C. As an example, consider an original sentence and its paraphrase in Table 12. Both of them contradict the sentence in the third column but they are not contradicting each other. The second approach, which uses a cross-encoder model, can capture the contradiction between sentences to some extent, but it is much more computationally expensive. Our experiment in Appendix F shows that compared with standard vector computation, running a cross-encoder is at least 200 times slower.

In this paper, we propose to overcome these limitations by introducing SPARSECL for efficient contradiction retrieval using sparse-aware sentence embeddings. The key idea behind our approach is to train a sentence embedding model to preserve *sparsity of differences* between the contradicted sentence embeddings. When answering a query, we calculate a score between the query and each document in the corpus, based on *both* the cosine similarity and the sparsity of the difference between their embeddings, and retrieve the ones with the highest scores. Our specific measure of

<sup>\*</sup>Equal contribution <sup>1</sup>MIT <sup>2</sup>University of California, Los Angeles. Correspondence to: Haike Xu <haikexu@mit.edu>.

Proceedings of the  $42^{nd}$  International Conference on Machine Learning, Vancouver, Canada. PMLR 267, 2025. Copyright 2025 by the author(s).

sparsity is defined by the Hoyer measure of sparsity (Hurley & Rickard, 2009), which uses the scaled ratio of the  $\ell_1$  norm and the  $\ell_2$  norm of a vector as a proxy of the number of non-zero entries in the vector. Unlike the cosine metric, the Hoyer measure is not transitive (please refer to Appendix B for a detailed analysis), which avoids the limitations of the former. At the same time this method is much more efficient than a cross-encoder, as both the cosine metric and the Hoyer measure are easy to compute given the embeddings. The Hoyer sparsity histogram of our trained embeddings is displayed in Figure 2.

We first evaluate our method on the counter-argument detection dataset Arguana (Wachsmuth et al., 2018), which to the best of our knowledge, is the only publicly available dataset suitable for testing contradiction retrieval. In addition, we generate another two data sets, where contradictions for documents in MSMARCO (Nguyen et al., 2016) and HotpotQA (Yang et al., 2018) datasets are generated using GPT-4 (Achiam et al., 2023). Our experiments demonstrate the efficacy of our approach in contradiction retrieval, as seen in Table 1. We also apply our method to corrupted corpus cleaning problem, where the goal is to filter out contradictory sentences in a corrupted corpus and preserve good QA retrieval accuracy.

To summarize. our contributions can be divided into three folds:

- We introduce a novel contradiction retrieval method that employs specially trained sentence embeddings combined with a metric that includes both cosine similarity and the Hoyer measure of sparsity. This approach effectively captures the essence of contradiction while being computationally efficient.
- Our method demonstrates superior contradiction retrieval metrics over different datasets compared to existing methods. This underscores the effectiveness of our embedding and scoring approach.
- We apply our contradiction retrieval method to two downstream settings: (1) corpus cleaning, where SPAR-SECL removes contradictions from corrupted datasets to maintain high-quality QA retrieval; and (2) natural language inference, where SPARSECL assigns a higher Hoyer sparsity score between contradicted pairs. These applications highlight the practical benefits of our approach in real-world scenarios.

# 2. Related Work

**Counter Argument Retrieval** A direct application of our contradiction retrieval task in "counter-argument retrieval". Since the curation of Arguana dataset by (Wachsmuth et al.,

2018), there has been a few previous work on retrieving the best counter-argument for a given argument (Orbach et al., 2020; Shi et al., 2023). In terms of methods, (Wachsmuth et al., 2018) uses a weighted sum of different word and embedding similarities and (Shi et al., 2023) designs a "Bipolarencoder" and a classification head. We believe that our method relying only on cosine similarity and sparsity is simpler than theirs and produces better results in the experiment. In addition, some analyses in the counter-argument retrieval papers are specific to the "debate" setting, e.g. they rely on topic, stance, premise/conclusion, and some other inherent structures in debates for help, which may prevent their methods from being generalized to broader scenarios.

**Fact verification and LLM hallucination** Addressing the hallucination problem in Large Language Models has been a subject of many research efforts in recent years. According to the three types of different hallucinations in (Zhang et al., 2023b), here we only focus on those so called "Fact-Conflicting Hallucination" where the outputs of LLM contradict real world knowledge. The most straightforward way to mitigate this hallucination issue is to assume an external groundtruth knowledge source and augment LLM's outputs with an information retrieval system. There have been a few works on this line showing the success of this method (Ren et al., 2023; Mialon et al., 2023). This practice is very similar to "Fact-Verification" (Thorne et al., 2018; Schuster et al., 2021) where the task is to judge whether a claim is true or false based on a given knowledge base.

However, as pointed out by (Zhang et al., 2023b), in the era of LLM, the external knowledge base can encompass the whole internet. It is impossible to assume that all the information there are perfectly correct and there may exist conflicting information within the database. In the context of our paper, instead of using a groundtruth database to check an external claim, our goal is to check the internal contradictions between different documents in an unknown corpus.

Learning augmented LLM and retrieval corpus attack Augmenting large language models with retrieval has been shown to be useful for many purposes. Recently, there have been a few works (Zhong et al., 2023; Zou et al., 2024) studying the vulnerability of retrieval system from adversarial attack. Specifically, they show that adding a few corrupted data points to the corpus will significantly drop the retrieval accuracy. This phenomenon brings our attention to the necessity of checking the factuality of the knowledge database. Note that the type of corrupted documents considered by their papers are different from ours. While they consider the injection of adversarially generated documents, we consider the existence of contradicted documents as a natural part of the corpus. Also their purpose is to show the effect of



Figure 1. Comparison of our SPARSECL with Cross-Encoder and Contrastive-Learning based Bi-Encoder for contradiction retrieval.

adversarial attack, while we provide a defense method for a certain kind of corrupted database.

# 3. Method

**Problem Formulation** We consider the contradiction retrieval problem: given a passage corpus  $C = \{p_1, p_2, ..., p_n\}$ and a query passage q, retrieve the "best" passage  $p^*$  that contradicts q. We assume that several similar passages supporting q might exist in the corpus C.

**Embedding based method** Judging whether two passages contradict each other is a standard Natural Language Inference task and can be easily tackled by many off-theshelf language models (Touvron et al., 2023; Xu et al., 2022). However, to retrieve the best candidate from the corpus, we have to iterate the whole corpus, or at least send the candidates retrieved by similarity search to the language model to determine if they constitute contradiction. This is time consuming, given that there are potentially many similar passages in the corpus. Therefore, in our paper, we mainly focus on those methods that only rely on their passage embeddings. Specifically, we want to design a simple scoring function F that given the embeddings of two passages, outputs a score between 0 and 1, indicating the likelihood that they are contradicting each other.

**Sparse Aware Embeddings** Following the idea from counter-argument retrieval papers (Wachsmuth et al., 2018), such a score function should be a combination of similarity and dissimilarity functions. Observe that a dissimilarity

function is basically a negation of a similarity function, so Wachsmuth et al. (2018) proposes several different similarity functions and sets the scoring function to maximize one of them and minimize another. Here, instead of enumerating different similarity functions, we consider another notion: the "sparsity" of their embedding differences. The basic intuition is as follows. Suppose that all sentences are represented as vectors in a "semantic" basis, where each coordinate represents one clearly identifiable semantic meaning. Then a contradiction between two passages should manifest itself as a difference in a few coordinates. while other coordinates should be quite close to each other. The issue, however, is that we do not know how to construct the appropriate basis, and the sparsity is defined with respect to a fixed coordinate system. Nevertheless, following this intuition, we fine-tune sentence embedding models using contrastive learning, by rewarding the sparsity of the difference vectors between embeddings of contradicting passages. Please see Figure 2 for the Hoyer sparsity histogram of our trained embeddings.

**SPARSECL** We use contrastive learning (Gao et al., 2021; Karpukhin et al., 2020) to fine-tune a pretrained sentence embedding model to generate the desired sparsity-aware embeddings. The choice of positive and negative examples are exactly the reverse of the choice we make when the training sets are Natural Language Inference datasets. The positive example for a passage is its contradiction passage in the training set. The hard negative example for a passage is its similar passage in the training set. There are also other random in-batch passages as soft negative examples. The



*Figure 2.* Histograms for the Hoyer sparsity of different pairs of sentence embedding differences on HotpotQA test set. The upper figure is the histogram produced by a standard sentence embedding model ("bge-base-en-v1.5"), where the median Hoyer sparsity values for random pairs, paraphrases, and contradictions are 0.212, 0.211, 0.211. The lower figure is the histogram produced by our sentence embedding model fine-tuned from "bge-base-en-v1.5" using our SPARSECL method, where the median Hoyer sparsity values for random pairs, paraphrases, and contradictions are 0.212, 0.281, 0.632.

sparsity function we choose here is Hoyer sparsity function from (Hurley & Rickard, 2009). Let  $h_1$  and  $h_2$  be two sentence embeddings and their embeddings have dimension d. We define

Hoyer
$$(h_1, h_2) = \left(\sqrt{d} - \frac{\|h_1 - h_2\|_1}{\|h_1 - h_2\|_2}\right) / \left(\sqrt{d} - 1\right)$$

This is a transformed version of the ratio of the  $l_1$  to the  $l_2$  norm, with output normalized to [0, 1].

Finally, for each training tuple  $(x_i, x_i^+, x_i^-)$  with their embeddings  $(h_i, h_i^+, h_i^-)$ , batch size N, and temperature  $\tau$ , its loss function is defined as

$$l_i = -\log \frac{e^{\operatorname{Hoyer}(h_i, h_i^+)/\tau}}{\sum_{j=1}^N \left( e^{\operatorname{Hoyer}(h_i, h_j^+)/\tau} + e^{\operatorname{Hoyer}(h_i, h_j^-)/\tau} \right)}$$

Scoring function for contradiction retrieval For the score function for contradiction retrieval, we use a weighted sum of the standard cosine similarity and our sparsity function. Note that the cosine similarity is provided separately by any off-the-shelf sentence embedding model in a zeroshot manner. It can also be fine-tuned. Let E() be the standard sentence embedding model and  $E_s()$  be our sparse-aware sentence embedding model trained by SPARSECL. Then the final score function for contradiction retrieval is

$$F(q, p) = \cos \left( E(q), E(p) \right) + \alpha \cdot \operatorname{Hoyer}(E_s(q), E_s(p)).$$

where  $\alpha$  is a scalar tuned using the validation set. Note that the criterion for contradiction is usually case-dependent, so it is necessary that we reserve a parameter to adapt to different notions of contradiction. To get the answer passages, we calculate the score function for all passages and report the top 10 of them<sup>1</sup>.

# 4. Experiments

We test our SPARSECL method on a counterargument retrieval task Arguana (Wachsmuth et al., 2018) and two contradiction retrieval datasets adapted from HotpotQA (Yang et al., 2018) and MSMARCO (Nguyen et al., 2016). Then, we apply our contradiction retrieval task to two downstream applications: retrieval corpus cleaning and natural language inference. Finally, we perform ablation studies to explain the functionality of each component of our method. Most of our experiments are not so computationally extensive, which can be run by one single A6000 GPU. We run our major experiments on A6000 and A100 GPUs.

#### 4.1. Counter-argument Retrieval

**Dataset** Arguana is a dataset curated in (Wachsmuth et al., 2018), where the author provide a corpus of 6753 argumentcounterargument pairs, taken from 1069 debates with 15 themes on idebate.org. For each debate, the arguments are further divided into two opposing stances (pro and con). For each stance, there are paired arguments and counterarguments. The dataset is split into the training set (60% of the data), the validation set (20%), and the test set (20%). This ensures that data from each individual debate is included in only one set and that debates from every theme are represented in every set. The task goal is: given an argument, retrieve its best counter-argument.

**Training** We use Arguana's training set to fine-tune our sparsity aware sentence embedding model via SPARSECL.

<sup>&</sup>lt;sup>1</sup>In the actual implementation, for time efficiency, we first use FAISS (Douze et al., 2024) to retrieve the top K candidates with cosine similarity and then rerank them using our cosine + sparsity score function. We set a very large K (e.g. K = 1000) so that empirically this is almost equivalent to searching for the maximal cosine + sparsity score in the whole corpus

To construct our training data, for each argument and counter-argument pair  $(x_i, x_i^c)$  in the Arguana's training set, we set  $x_i^c$  to be the positive example of  $x_i$ . We select all the other arguments and counter-arguments from the same debate and stance as  $x_i$ 's hard negatives. We fine-tune three pretrained sentence embedding models of different sizes ("UAE-Large-V1" (Li & Li, 2023), "GTE-large-en-v1.5" (Li et al., 2023), and "bge-base-en-v1.5" (Xiao et al., 2023)). Please refer to Table 10 for our training parameters.

**Baselines** We are not aware of any accurate methods for retrieving contradictions that only rely on sentence embeddings. To the best of our knowledge, we provide two baselines in our main experiment:

- CL: standard contrastive learning with cosine similarity, using the same training data (contradictions as positive examples and paraphrases as negative examples) that we use for our SPARSECL.
- Prompt + CL: standard contrastive learning with prompt "Not true: " attached in front of the query during both training and testing.

We report the performance of several efficient (with fewer than 1B parameters) and top-ranked pretrained sentence embedding models including "GTE-large-en-v1.5", "UAE-Large-V1", "bge-base-en-v1.5",

**Test** The Arguana test set consists of 1401 query arguments and counter-argument pairs. Following the standard test setting, we search for an answer of a query within the whole corpus (training set + validation set + test set) and report NDCG@10 scores. We select  $\alpha$  based on the best NDCG@10 score on the validation set. When we directly use a model to provide cosine similarity scores in a zeroshot manner, we use its default pooler ("cls") for that model. When we use a fine-tuned model (via either CL or SPAR-SECL) to provide either cosine similarity scores or sparsity scores, we use the "avg" pooler.

**Results** The detailed results are presented in Table 1. Across all models—"GTE-large-en-v1.5", "UAE-Large-V1", and "bge-base-en-v1.5"—an average improvement of 3.6% in counter-argument retrieval were observed when incorporating our SPARSECL to either Zeroshot or CL. Furthermore, our CL (Cosine) + SPARSECL (Hoyer) method achieves NDCG@10 score 81.3 using GTE with only 400M parameters. For completeness, we also compare our results with (Shi et al., 2023) in Appendix C.

This pattern of enhancement was consistently observed regardless of whether the embedding models were fine-tuned or not. Notably, standard cosine similarity fine-tuning alone also contributed to performance gains. For instance, finetuned GTE models showed an increase from 72.5 to 77.8 on the Arguana dataset using standard cosine similarity alone. This suggests that the Arguana dataset inherently favors scenarios where the counterargument is the most similar passage to the query, which may amplify the benefits of fine-tuning.

Interestingly, for the "Prompt + CL (Cosine)" method, finetuning with the appended prompt even results in a performance drop. During the training process, we observed overfitting and hypothesize that the special prompt "Not true:" introduces a shortcut, making it easier for the model to learn whether a text belongs to the "argument" class or the "counter-argument" class. However, this class information is not useful when identifying pairwise contradiction relationships.

These findings highlight the robustness of our approach, particularly when traditional similarity metrics are augmented with sparsity measures to capture subtle nuances in contradiction. Further insights can be gleaned from our ablation study detailed in Section 4.6, where we analyze the impact of similar non-contradictory passages within the corpus.

### 4.2. Contradiction retrieval

The task of "contradiction retrieval" generalizes beyond the argument and counter-argument relationship in the debate area, e.g. passages with conflicting factual information should also be considered as "contradictions". To test our method's validity for these more general forms of contradictions, we construct two datasets to test our method's performance.

**Data set construction** Given a QA retrieval dataset, e.g. MSMARCO (Nguyen et al., 2016), for each answer passage  $x_i$  of a query  $q_i$ , we use Large Language Models (specifically, GPT-4 (Achiam et al., 2023)) to generate 3 answers paraphrasing  $x_i$  or contradicting  $x_i$ . Let the generated paraphrases be  $\{x_{i1}^+, x_{i2}^+, x_{i3}^+\}$  and the generated contradictions be  $\{x_{i1}^+, x_{i2}^-, x_{i3}^-, \}$ . We then delete  $x_i$  from the corpus and add the set of generated passages  $\{x_{i1}^+, x_{i2}^+, x_{i3}^+, x_{i1}^-, x_{i2}^-, x_{i3}^-\}$  to the corpus. In the test phrase, the queries are  $\{x_{i1}^+, x_{i2}^+, x_{i3}^+\}$ , each of which has the same answers  $\{x_{i1}^-, x_{i2}^-, x_{i3}^-, \}$ . We generate the paraphrases and contradictions for the validation set, test set, and a randomly sampled 10000 documents from the training set. Please refer to Appendix G for details.

**Training** To prepare the training data for contrastive learning, for each paraphrase and contradiction set  $\{x_{i1}^+, x_{i2}^+, x_{i3}^+, x_{i1}^-, x_{i2}^-, x_{i3}^-\}$  generated from the same original passage, we form 9 pieces of training data  $(x_{ia}^+, x_{ib}^-, x_{ic}^+)$  for 9 different combinations of paraphrases, contradictions,

Contradiction Retrieval via Contrastive Learning with Sparsity

Model	Method	Arguana	MSMARCO	HotpotQA
	Zeroshot (Cosine)	65.8	60.0	59.5
	CL (Cosine)	68.7	52.7	56.2
BGE	Prompt + CL (Cosine)	64.4	61.1	80.4
	Zeroshot (Cosine) + SPARSECL(Hoyer)	70.4	90.9	96.7
	CL (Cosine) + SPARSECL(Hoyer)	72.2	88.3	96.5
	Zeroshot (Cosine)	68.3	59.7	58.7
	CL (Cosine)	70.4	44.2	54.1
UAE	Prompt + CL (Cosine)	63.7	85.8	94.3
	Zeroshot (Cosine) + SPARSECL(Hoyer)	74.3	90.2	95.5
	CL (Cosine) + SPARSECL(Hoyer)	74.4	86.9	94.3
	Zeroshot (Cosine)	72.5	60.3	59.7
	CL (Cosine)	77.8	65.1	59.7
GTE	Prompt + CL (Cosine)	71.1	88.1	69.0
	Zeroshot (Cosine) + SPARSECL(Hoyer)	79.7	95.3	97.7
	CL (Cosine) + SPARSECL(Hoyer)	81.3	95.2	97.9

Table 1. Results for different models and methods on the contradiction retrieval task. Experiments are run on the Arguana dataset (Wachsmuth et al., 2018) and modified MSMARCO(Nguyen et al., 2016) and HotpotQA(Yang et al., 2018) datasets. We report NDCG@10 score here, the higher the better. "UAE" stands for "UAE-Large-V1", "BGE" stands for "bge-base-en-v1.5", "GTE" stands for "gte-large-en-v1.5", The "Method" column denotes the score function used to retrieve contradictions. We consider two score functions: cosine similarity and cosine similarity plus Hoyer sparsity. "Zeroshot" denotes the direct testing of the model without any fine-tuning. "CL" denotes fine-tuning using standard contrastive learning. "Prompt+CL" denotes fine-tuning using standard contrastive learning with prompt "Not true:" attached in front of the query. "SPARSECL" denotes fine-tuning using Hoyer sparsity contrastive learning (our method).

and a randomly selected hard negative from the remaining two paraphrases. We then perform SPARSECL to fine-tune a sparsity-enhanced embedding.

**Test** Similar to the testing strategy for Arguana, we define our corpus to consist of all generated text (training set + validation set + test set). We query the paraphrases  $\{x_{i1}^+, x_{i2}^+, x_{i3}^+\}$  of the original passage  $x_i$  and set the groundtruth answers to be the generated contradictions  $\{x_{i1}^-, x_{i2}^-, x_{i3}^-\}$ . We select the  $\alpha$  parameter with the maximal NDCG@10 score on the validation set and report the NDCG@10 score obtained by applying that  $\alpha$  to the test set.

**Results** The results are reported in Table 1. For both MSMARCO and HotpotQA data sets, incorporating our SPARSECL method achieves over 14.6% percentage points gain compared with the two baselines. The large improvement is due to the existence of paraphrases in the corpus, that are strong confounders for the pure similarity-based methods.

We also observe that Prompt + CL (Cosine) performs much better on the MSMARCO and HotpotQA datasets compared to standard CL (Cosine). We propose two potential reasons: 1. We have 12 times more paraphrase and contradiction pairs generated for fine-tuning, which makes the model less likely to overfit. 2. The contradictions generated by GPT- 4 rely too heavily on opposite word replacement, which is better captured by the prompt 'not true.' However, the counter-arguments in the Arguana dataset use entirely different wording.

#### 4.3. Zero-shot Generalization Test

To evaluate the generalization capability of our sparseaware embeddings, we also conduct zero-shot tests on other datasets. Specifically, we train the embeddings on our synthetic HotpotQA or MSMARCO datasets and then test them on the other dataset in a zero-shot manner. We have confirmed that there is no data overlap between the two datasets. Please refer to Table 14 for the corresponding statistics. As presented in Table 2, SparseCL trained on MSMARCO or HotpotQA produces reasonable test results on the other dataset, albeit with a slight performance drop. Furthermore, using the same  $\alpha$  parameter selected on the other dataset also gives reasonable test accuracy, showing the stability of  $\alpha$  parameter across different datasets. This demonstrates that the sparse-aware embeddings trained on one dataset can capture contradiction relationships and generalize to unseen datasets.

### 4.4. Retrieval Corpus Cleaning

As an application of contradiction retrieval, we test how well our method can be used to find inconsistencies within a corpus and clean the corpus for future training or QA re-

Contradiction R	etrieval via	Contrastive l	Learning v	vith Sparsity
-----------------	--------------	---------------	------------	---------------

Method	Train Dataset	Test Dataset	NDCG@10
Zeroshot(Cosine)	MSMARCO	HotpotQA	88.1
+SparseCL(Hoyer) +fixed $\alpha$	HotpotQA	MSMARCO	82.2
Zeroshot(Cosine)	MSMARCO	HotpotQA	88.6
+SparseCL(Hoyer) + tuned $\alpha$	HotpotQA	MSMARCO	87.7
Zeroshot(Cosine)	HotpotQA	HotpotQA	96.7
+SparseCL(Hoyer)	MSMARCO	MSMARCO	90.9
Zeroshot(Cosine)	N/A	HotpotQA	59.5
	N/A	MSMARCO	60.0

Table 2. Results for zero-shot generalization experiment for contradiction retrieval running on "bge-base-en-v1.5" model

trieval. We first inject corrupted data contradicting existing documents into the corpus, and measure the retrieval accuracy degradation for retrieved answers. Then, we use our contradiction retrieval method to filter out corrupted data and measure the retrieval accuracy again.

**Data** Similarly to the data generation in Section 4.2, we construct a new corpus containing LLM-generated paraphrases and contradictions based on MSMARCO and HotpotQA data sets. We start with an original corpus C and its subset S. We then generate paraphrases and contradictions for S as in Section 4.2.

For HotpotQA, S contains all answer documents for the test set, 10000 answer documents sampled from the training set, and 1000 answer documents sampled from the development set. For MSMARCO, S contains all answer documents for the dev set, and 11000 answer documents sampled from the training set.

We then curate 3 different versions of the corpus based on the original corpus C and the subset S.

- The initial corpus C<sup>+</sup>: For each original answer document x in S, we remove x from C and instead add 3 LLM-generated paraphrases {x<sub>1</sub><sup>+</sup>, x<sub>2</sub><sup>+</sup>, x<sub>3</sub><sup>+</sup>} to C. The result forms the *initial* corpus C<sup>+</sup>.
- The corrupted corpus  $C^-$ : For each original answer document x in S, we generate 3 contradictions  $\{x_1^-, x_2^-, x_3^-\}$  and add them to  $C^+$  to get the *corrupted* corpus  $C^-$ .
- The cleaned corpus C<sup>β</sup>: We apply our data cleaning procedure to the corrupted corpus C<sup>-</sup>, obtaining the *cleaned* dataset C<sup>β</sup>.

**Test** We test the retrieval accuracy (NDCG@10) and the corruption ratio (Recall@10) for answering the original queries in the test set. The goal of our experiment is to show how retrieval algorithms behave on these three constructed corpora  $C^+$ ,  $C^-$ , and  $C^{\natural}$ .

**Data Cleaning** Our sparsity-based method can only identify contradictions within the data set, but we do not know which element in a contradiction pair is correct. To perform data cleaning, we make the assumption that for each original passage  $x \in S$ , we are given one of its paraphrases as the groundtruth. Then, our task is reduced to searching for passages contradicting a given ground truth document and filtering them out.

**Method** We use the GTE-large-en-v1.5 model without fine-tuning to provide the cosine similarity score for this data cleaning experiment. We use the model from our contradiction retrieval experiment in section 4.2 trained on MS-MARCO and HotpotQA to provide the sparsity score. The  $\alpha$  parameter is also identical to the one used in section 4.2. For each ground truth document, we filter out the top 3 scored documents from the corpus.

Detecata	Original Corrupted		Cleaned		
Datasets	Acc	Acc	Corrupt	Acc	Corrupt
HotpotQA	67.6	56.7	44.3	65.2	2.0
MSMARCO	43.5	38.1	41.3	41.4	4.0

Table 3. Experimental results for the impact of corrupted data on QA retrieval and contradiction retrieval for filtration. "Acc" represents the retrieval accuracy measured by the NDCG@10 score and "Corrupt" represents the fraction of returned passages that are corrupted, as measured by Recall@10.

Table 3 shows the results. We observe that the retrieval accuracy on the corrupted corpus drops significantly, as the generated contradictions cause the embedding model to retrieve them as query answers. The corruption ratio measures the average fraction of the top-10 retrieved documents that correspond to the generated contradicting passages. This performance is above 40% for both datasets. After performing our corpus cleaning procedure, which searches for the passages contradicting the given ground truth documents and removes the top-3 for each of them, we can recover more than 60% of the performance loss due to corruption and at the same time reduce the corruption ratio to less than 5%.

#### 4.5. Score Functions for Natural Language Inference

As an application of our SPARSECL method, we demonstrate that our method can be useful for distinguishing contradictions from entailments and random pairs in natural language inference datasets. For SNLI (Bowman et al., 2015) and MNLI (Williams et al., 2018) datasets, we extract entailment and contradiction pairs, fine-tune using standard contrastive learning and our SPARSECL, and then report the average cosine similarity / Hoyer sparsity score between entailments, contradictions, and random pairs.

We can observe from Table 5 that, in the zeroshot setting, the

Contradiction Retrieval via Contrastive Learning with Sparsity

Model	Method	$l_{2}/l_{1}$	$\kappa_4$	Hoyer	Cosine (baseline)
BGE	Zeroshot (Cosine) + SPARSECL	67.5	68.4	70.4	65.7
BGE	CL (Cosine) + SPARSECL	70.2	70.7	72.2	68.7

Table 4. NDCG@10 scores for Arguana using SPARSECL with different sparsity functions. We also report two baselines that use only the cosine similarity (zeroshot and contrastive learning).

		Contradiction	Entailment	Random
SNLI	Zeroshot (Cosine)	54.6	76.9	37.6
	CL (Cosine)	88.5	88.6	77.7
	SparseCL (Hoyer)	<b>37.6</b>	<b>34.7</b>	<b>22.8</b>
MNLI	Zeroshot (Cosine)	65.9	81.8	37.8
	CL (Cosine)	91.9	91.7	73.3
	SparseCL (Hoyer)	<b>42.2</b>	<b>36.4</b>	<b>24.4</b>

Table 5. Average Cosine / Hoyer scores between Contradiction / Entailment / Random pairs of texts. The experiment is run on "bge-base-en-v1.5" model. Texts pairs are from SNLI and MNLI datasets

average cosine similarity of contradiction pairs lies between the ranges of random and entailment pairs. For the finetuned model using standard contrastive learning (CL), the average cosine similarity of contradiction pairs is almost indistinguishable from that of entailment pairs. Finally, after being fine-tuned using SPARSECL, the model exhibits higher average Hoyer sparsity scores for contradiction pairs compared to other two types of relationships.

#### 4.6. Ablation Studies

We perform the following three ablation studies to further understand sparsity-based retrieval method.

**Arguana retrieval results analysis** In the standard Arguana dataset, even though the task is to retrieve the counterargument for the query, the retrieval based solely on similarity still gives reasonable results. This means that counterarguments are also the most similar arguments to the query, which makes the data set an imperfect test bed for testing contradiction retrieval.

To further compare our sparsity-based method and the pure similarity-based method, we augment Arguana by adding arguments' paraphrases to the corpus. Specifically, for any argument x and its counter-argument  $x^-$  in the original corpus C, we use GPT-4 to generate three paraphrases  $\{x_1, x_2, x_3\}$ of x. We then form three new corpora with an increasing number of paraphrases added to the corpus:  $C_1$  contains all  $x_1$  and  $x^-$ ,  $C_2$  contains all  $x_1$ ,  $x_2$ , and  $x^-$ , and  $C_3$  contains all  $x_1$ ,  $x_2$ ,  $x_3$ , and  $x^-$ .

In the testing phase, we query the counter-arguments for one of x's paraphrases, the answer of which should still be  $x^-$ . We observe how the performance varies when the corpora

we retrieve from are  $C_1, C_2, C_3$ .

Methods	$C_1$	$C_2$	$C_3$
Zeroshot (Cosine)	56.1	35.5	26.7
Zeroshot (Cosine) + SPARSECL(Hoyer)	68.2	67.9	67.5
CL (Cosine)	47.1	30.3	22.8
CL (Cosine) + SPARSECL(Hoyer)	61.9	61.8	61.5

Table 6. Counter-argument retrieval results on the augmented Arguana dataset with different numbers of similar arguments in the corpus.  $C_x$  denotes testing counter-argument retrieval on the corpus with x existing paraphrases (including itself) of the query argument. Experiments were run on "bge-base-en-v1.5" model.

We present our overall experimental results in Table 6. Please also refer to Appendix D for an example case study. As the number of paraphrases in corpus increases from 1 to 3, the performance of the similarity-based method drops significantly. Thus it is reasonable to deduce that, as the number of similar arguments in the corpus increases further, the NDCG@10 scores for similarity-based methods will converge to 0. On the other hand, the performance of our sparsity-based method is stable with respect to the number of paraphrases in the corpus.

**Different sparsity functions** Our intuition in Section 3 does not give clear guidelines on which sparsity function to use in our SPARSECL. Thus, we also experiment with different choices of sparsity functions, selected from Hurley & Rickard (2009). Specifically, we consider two other sparsity functions ( $l_2/l_1$  and  $\kappa_4$ ), which are scale invariant and differentiable (see Table III in (Hurley & Rickard, 2009)). Note that both of these two sparsity functions have ranges [0, 1], and higher values of those functions correspond to sparser vectors.

$$\frac{l_2}{l_1} = \frac{\|h_1 - h_2\|_2}{\|h_1 - h_2\|_1} \qquad \kappa_4 = \frac{\|h_1 - h_2\|_4^4}{\|h_1 - h_2\|_2^2}$$

As per Table 4, compared to the cosine similarity method, the combination of the cosine similarity score with the sparsity score trained by SPARSECL, yields higher NDCG@10 scores for each sparsity function. However, Hoyer sparsity yields the highest accuracy. We believe that simple sparsity functions have a more benign optimization landscape and thus are easier for models to optimize.

#### Different retrieval methods for contradiction retrieval

We experiment with 5 retrieval methods in our ablation study. The methods evaluated are as follows: "Prompt" involves appending the "Not true: " prompt to the query during testing, followed by standard similarity search. "Prompt + CL (Cosine)" extends this by incorporating contrastive learning with the "Not true: " prompt included in the training data. "Gen" uses GPT-4 to generate contradictions to the query (details in Appendix G) and applies similarity search for testing. "Gen + CL (Cosine)" fine-tunes using contrastive learning with the generated contradictions in the training data before similarity search. Finally, "SparseCL (Hoyer)" employs SparseCL fine-tuning and retrieves documents based on the maximal Hoyer sparsity score during testing.

As shown in Table 7, we observe that generally "Gen" and "Prompt" don't improve much upon standard similarity search. For the "Gen + CL (Cosine)" method, a diverse set of counter-arguments exist for a given argument, making it hard to generate a single counter-argument that closely matches the true ground truth counter-argument. For the "Prompt + CL (Cosine)" method, fine-tuning with the appended prompt even results in a performance drop. During the training process, we observed overfitting and hypothesize that the special prompt "Not true:" introduces a shortcut, making it easier for the model to learn whether a text belongs to the "argument" class or the "counter-argument" class. However, this class information is not useful when identifying pairwise contradiction relationships. Finally, directly using Hoyer sparsity to retrieve contradictions doesn't yield good results as well, because we believe contradictions involve a combination of similarity and dissimilarity.

Model	Method	Arguana
	Prompt + Zeroshot (Cosine)	65.7
BGE	Gen + Zeroshot (Cosine)	64.7
	Zeroshot (Cosine)	65.8
PCE	Prompt + CL (Cosine)	64.5
DUE	Gen + CL (Cosine)	70.0
	CL (Cosine)	68.7
PCE	SparseCL (Hoyer)	56.1
DUE	CL (Cosine) + SparseCL (Hoyer)	72.2

Table 7. Counter-argument retrieval results (NDCG@10 scores) on Arguana dataset with different retrieval methods. "Gen" means using GPT-4 to generate a contradiction c of the query argument q, "Prompt" means appending the "Not true : " prompt in the front of the query text. "Zeroshot" refers to direct testing and "CL" and "SparseCL" refer to fine-tuning with respective methods.

#### 5. Conclusion

In this work, we introduced a novel approach to contradiction retrieval that leverages the sparsity between sentence embeddings, combined with cosine similarity, to efficiently identify contradictions in large document corpora. This method addresses the limitations of the traditional similarity search as well as computational inefficiencies of the cross-encoder models, proving its effectiveness on benchmark datasets like Arguana and on synthetic contradictions retrieval from MSMARCO and HotpotQA.

# Acknowledgements

Haike Xu was supported by the Mathworks Fellowship. Piotr Indyk was supported in part by the NSF TRIPODS program (award DMS-2022448). Zongyu Lin was partially supported by DARPA ANSR program FA8750-23-2-0004, NSF 2331966 NSF 2211557, NSF 2119643, NSF 2303037, NSF 2312501, SRC JUMP 2.0 Center, Amazon Research Awards, and Snapchat Gifts.

# **Impact Statement**

This paper introduces a novel approach to the contradiction retrieval task that relies solely on text embeddings. To the best of our knowledge, we are the first to consider such non-similarity-based retrieval problem. We consider this an important yet underexplored area and hope our work will spark greater interest in this field.

#### References

- Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Asai, A., Schick, T., Lewis, P., Chen, X., Izacard, G., Riedel, S., Hajishirzi, H., and Yih, W.-t. Task-aware retrieval with instructions. In Rogers, A., Boyd-Graber, J., and Okazaki, N. (eds.), *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 3650–3675, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl. 225. URL https://aclanthology.org/2023. findings-acl.225.
- Bowman, S., Angeli, G., Potts, C., and Manning, C. D. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pp. 632–642, 2015.
- Cui, L. and Lee, D. Coaid: Covid-19 healthcare misinformation dataset, 2020.
- Douze, M., Guzhva, A., Deng, C., Johnson, J., Szilvasy, G., Mazaré, P.-E., Lomeli, M., Hosseini, L., and Jégou, H. The faiss library. 2024.
- Gao, T., Yao, X., and Chen, D. Simcse: Simple contrastive learning of sentence embeddings. *arXiv preprint arXiv:2104.08821*, 2021.
- Hurley, N. and Rickard, S. Comparing measures of sparsity. *IEEE Transactions on Information Theory*, 55(10):4723–4741, 2009.

- Jha, R., Lovering, C., and Pavlick, E. Does data augmentation improve generalization in nlp? *arXiv preprint arXiv:2004.15012*, 2020.
- Karpukhin, V., Oğuz, B., Min, S., Lewis, P., Wu, L., Edunov, S., Chen, D., and Yih, W.-t. Dense passage retrieval for open-domain question answering. *arXiv preprint arXiv:2004.04906*, 2020.
- Laban, P., Schnabel, T., Bennett, P. N., and Hearst, M. A. SummaC: Re-visiting NLI-based models for inconsistency detection in summarization. *Transactions of the Association for Computational Linguistics*, 10:163–177, 2022. doi: 10.1162/tacl\_a\_00453. URL https:// aclanthology.org/2022.tacl-1.10.
- Li, X. and Li, J. Angle-optimized text embeddings. *arXiv* preprint arXiv:2309.12871, 2023.
- Li, Z., Zhang, X., Zhang, Y., Long, D., Xie, P., and Zhang, M. Towards general text embeddings with multi-stage contrastive learning. *arXiv preprint arXiv:2308.03281*, 2023.
- Mialon, G., Dessì, R., Lomeli, M., Nalmpantis, C., Pasunuru, R., Raileanu, R., Rozière, B., Schick, T., Dwivedi-Yu, J., Celikyilmaz, A., Grave, E., LeCun, Y., and Scialom, T. Augmented language models: a survey, 2023.
- Muennighoff, N., Tazi, N., Magne, L., and Reimers, N. MTEB: Massive text embedding benchmark. In Vlachos, A. and Augenstein, I. (eds.), *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pp. 2014–2037, Dubrovnik, Croatia, May 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.eacl-main. 148. URL https://aclanthology.org/2023. eacl-main.148.
- Nguyen, T., Rosenberg, M., Song, X., Gao, J., Tiwary, S., Majumder, R., and Deng, L. MS MARCO: A human generated machine reading comprehension dataset. *CoRR*, abs/1611.09268, 2016. URL http://arxiv.org/ abs/1611.09268.
- Orbach, M., Bilu, Y., Toledo, A., Lahav, D., Jacovi, M., Aharonov, R., and Slonim, N. Out of the echo chamber: Detecting countering debate speeches. *arXiv preprint arXiv:2005.01157*, 2020.
- Ren, R., Wang, Y., Qu, Y., Zhao, W. X., Liu, J., Tian, H., Wu, H., Wen, J.-R., and Wang, H. Investigating the factual knowledge boundary of large language models with retrieval augmentation, 2023.
- Schuster, T., Fisch, A., and Barzilay, R. Get your vitamin C! robust fact verification with contrastive evidence. In Toutanova, K., Rumshisky, A., Zettlemoyer,

L., Hakkani-Tur, D., Beltagy, I., Bethard, S., Cotterell, R., Chakraborty, T., and Zhou, Y. (eds.), *Proceedings* of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pp. 624–643, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main.52. URL https: //aclanthology.org/2021.naacl-main.52.

- Shahi, G. K. and Nandini, D. FakeCovid- A Multilingual Cross-domain Fact Check News Dataset for COVID-19. ICWSM, Jun 2020. doi: 10.36190/2020.14. URL https://doi.org/10.36190/2020.14.
- Shi, H., Cao, S., and Nguyen, C.-T. Revisiting the role of similarity and dissimilarity inbest counter argument retrieval. arXiv preprint arXiv:2304.08807, 2023.
- Singhal, A. et al. Modern information retrieval: A brief overview. *IEEE Data Eng. Bull.*, 24(4):35–43, 2001.
- Thakur, N., Reimers, N., Rücklé, A., Srivastava, A., and Gurevych, I. 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)*, 2021. URL https://openreview.net/forum? id=wCu6T5xFjeJ.
- Thorne, J., Vlachos, A., Christodoulopoulos, C., and Mittal, A. FEVER: a large-scale dataset for fact extraction and VERification. In Walker, M., Ji, H., and Stent, A. (eds.), Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pp. 809–819, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-1074. URL https://aclanthology.org/N18-1074.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al. Llama 2: Open foundation and finetuned chat models. arXiv preprint arXiv:2307.09288, 2023.
- Wachsmuth, H., Syed, S., and Stein, B. Retrieval of the best counterargument without prior topic knowledge. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 241–251, 2018.
- Wang, X., Wang, J., Cao, W., Wang, K., Paturi, R., and Bergen, L. Birco: A benchmark of information retrieval tasks with complex objectives, 2024.

- Wang, Z., Lin, Z., Liu, P., ZHeng, G., Wen, J., Chen, X., Chen, Y., and Yang, Z. Learning to detect noisy labels using model-based features. arXiv preprint arXiv:2212.13767, 2022.
- Williams, A., Nangia, N., and Bowman, S. A broadcoverage challenge corpus for sentence understanding through inference. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pp. 1112–1122, 2018.
- Xiao, S., Liu, Z., Zhang, P., and Muennighoff, N. C-pack: Packaged resources to advance general chinese embedding, 2023.
- Xie, Q., Luong, M.-T., Hovy, E., and Le, Q. V. Self-training with noisy student improves imagenet classification. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 10687–10698, 2020.
- Xiong, L., Xiong, C., Li, Y., Tang, K.-F., Liu, J., Bennett, P., Ahmed, J., and Overwijk, A. Approximate nearest neighbor negative contrastive learning for dense text retrieval. *arXiv preprint arXiv:2007.00808*, 2020.
- Xu, H., Lin, Z., Zhou, J., Zheng, Y., and Yang, Z. A universal discriminator for zero-shot generalization. arXiv preprint arXiv:2211.08099, 2022.
- Yang, Z., Qi, P., Zhang, S., Bengio, Y., Cohen, W., Salakhutdinov, R., and Manning, C. D. HotpotQA: A dataset for diverse, explainable multi-hop question answering. In Riloff, E., Chiang, D., Hockenmaier, J., and Tsujii, J. (eds.), *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 2369–2380, Brussels, Belgium, October-November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1259. URL https: //aclanthology.org/D18-1259.
- Zhang, X., Thakur, N., Ogundepo, O., Kamalloo, E., Alfonso-Hermelo, D., Li, X., Liu, Q., Rezagholizadeh, M., and Lin, J. MIRACL: A Multilingual Retrieval Dataset Covering 18 Diverse Languages. *Transactions of the Association for Computational Linguistics*, 11:1114– 1131, 09 2023a. ISSN 2307-387X. doi: 10.1162/tacl\_a\_ 00595. URL https://doi.org/10.1162/tacl\_ a\_00595.
- Zhang, Y., Li, Y., Cui, L., Cai, D., Liu, L., Fu, T., Huang, X., Zhao, E., Zhang, Y., Chen, Y., Wang, L., Luu, A. T., Bi, W., Shi, F., and Shi, S. Siren's song in the ai ocean: A survey on hallucination in large language models, 2023b.
- Zhong, Z., Huang, Z., Wettig, A., and Chen, D. Poisoning retrieval corpora by injecting adversarial passages. *arXiv preprint arXiv:2310.19156*, 2023.

- Zhou, J., Lin, Z., Zheng, Y., Li, J., and Yang, Z. Not all tasks are born equal: Understanding zero-shot generalization. In *The Eleventh International Conference on Learning Representations*, 2022.
- Zou, W., Geng, R., Wang, B., and Jia, J. Poisonedrag: Knowledge poisoning attacks to retrieval-augmented generation of large language models, 2024.

# A. Additional related work

**Complex retrieval tasks** Information retrieval is a well-studied area (Singhal et al., 2001) and there have been many benchmarks for testing retrieval performance such as BEIR (Thakur et al., 2021), MTEB (Muennighoff et al., 2023), and MIRACL (Zhang et al., 2023a). However, most of the datasets, through varying in some degrees, focus only on "retrieving the most similar document". People have noted that there exist some more complex retrieval tasks (e.g. Arguana (Wachsmuth et al., 2018) retrieves counter-arguments that refute a query argument), and build retrieval benchmark focusing on complex retrival goals, e.g. BIRCO (Wang et al., 2024) and BERRI (Asai et al., 2023).

To retrieve according to different instructions, (Asai et al., 2023) trains TART, a multi-task retrieval system with task instructions attached as prompts in front of the query content. However, when answering queries, they are still searching for the most similar sentence embedding, though the prompt is different for different tasks. As far as we know, our paper studies the first non-similarity-based search problem.

**Data inconsistency and misinformation detection** Data inconsistency, refers to the factually incorrectness in the content, might come from different sources, including their natural existence in the corpus (Shahi & Nandini, 2020; Cui & Lee, 2020), data augmentations (Jha et al., 2020; Zhou et al., 2022), and pseudo labeling (Xie et al., 2020; Wang et al., 2022), which might lead to negative influence if serving as training dataset. There have been a few datasets on detecting the factually wrong information. For example, (Laban et al., 2022) detects whether a given summary is consistent with the input document, (Shahi & Nandini, 2020; Cui & Lee, 2020) detects whether a given COVID-19 related news is true or false. Most of these datasets lie in a specific domain and require external knowledge to judge the correctness of each piece of data. On the contrary, the "data inconsistency" notion we consider in our paper doesn't depend on any external knowledge, but is a relationship between different pieces of data in the same corpus. The goal of our method is to find such "contradiction pairs" in corpus efficiently, but not to judge which one is consistent with the real world knowledge.

# B. The "non-transitivity" of Hoyer sparsity and the "transitivity" of cosine function

Here, we provide a simple example to demonstrate that using Hoyer sparsity to measure "contradiction" can bypass the challenging scenario for similarity metrics where "A contradicts C, B contradicts C, but A doesn't contradict B". Specifically, Hoyer sparsity satisfies the following "non-transitivity" property.

**Proposition B.1** ("non-transitivity" of hoyer sparsity). There exist three vectors A, B, and C of dimensionality d, satisfying  $1 \leq ||A||_2, ||B||_2, ||C||_2 \leq 1 + O(\frac{1}{\sqrt{d}})$ , such that  $Hoyer(A, C) > 1 - O(\frac{1}{\sqrt{d}})$ ,  $Hoyer(B, C) > 1 - O(\frac{1}{\sqrt{d}})$ , and  $Hoyer(A, B) < O(\frac{1}{\sqrt{d}})$ 

*Proof.* We construct the following d dimensional vectors where  $\epsilon < \frac{1}{d}$  can be any parameter.

$$\begin{array}{rclrcrcrcrc} A & = & (1, & 0, & 0, & \dots, & 0) \\ B & = & (1, & 0, & \epsilon, & \dots, & \epsilon) \\ C & = & (0, & 1, & 0, & \dots, & 0) \end{array}$$

Then, we calculate their  $l_1$  over  $l_2$  ratios:

$$\begin{split} \frac{\|A - B\|_1}{\|A - B\|_2} &= \sqrt{d - 2} \\ \frac{\|A - C\|_1}{\|A - C\|_2} &= \sqrt{2} \\ \frac{\|B - C\|_1}{\|B - C\|_2} &= \frac{2 + (d - 2)\epsilon}{\sqrt{2 + (d - 2)\epsilon^2}} < \frac{3}{\sqrt{2}} \end{split}$$

Applying their  $l_1$  over  $l_2$  ratio bounds to the Hoyer sparsity formula will give us the desired relationship.

Next, we provide another example to demonstrate that the cosine function exhibits the following "transitivity" property, which makes it hard to characterize the scenario where "A contradicts C, B contradicts C, but A doesn't contradict B".

**Proposition B.2** ("transitivity" property of cosine function). *Given three unit vectors* A, B, and C, if  $cos(A, C) \ge 1 - O(\epsilon)$  and  $cos(B, C) \ge 1 - O(\epsilon)$ , we have  $cos(A, B) \ge 1 - O(\epsilon)$ 

*Proof.* For any two vectors X and Y with unit norm, we have  $cos(X,Y) = 1 - \frac{\|X-Y\|_2^2}{2}$ . Because  $cos(A,C) \ge 1 - O(\epsilon)$ , we have  $\|A - C\|_2 \le O(\sqrt{\epsilon})$ . Finally,  $cos(A, B) = 1 - \frac{\|A - B\|_2^2}{2} \ge 1 - \frac{(\|A - C\|_2 + \|C - B\|_2)^2}{2} \ge 1 - O(\epsilon)$ 

# C. Experiment comparison with method from (Shi et al., 2023)

Shi et al. (2023) proposes "Bipolar-encoder" method to retrieve contradictions from the corpus. They also tested their method on the Arguana dataset but used a different metric, Recall@1. For completeness, we have translated our results into their Recall@1 metric for a fair comparison. As shown in Table 8, both our CL (baseline method) and CL+SparseCL (our method) demonstrate significant improvement over the previous results in (Shi et al., 2023).

Model	Method	Arguana(Recall@1)
GTE	CL+SparseCL (ours)	62.9
GTE	CL (baseline)	56.3
(Shi et al., 2023)	Bipolar-encoder	49.0

Table 8. Comparison of experimental results on the Arguana dataset

### D. A case study for counter-argument retrieval from Arguana dataset

In this section we provide an example to illustrate how our sparsity-based retrieval method is better at retrieving counterarguments. In the setting of the augmented Arguana dataset (see our ablation study in Section 4.6), we selected an example query with an ID "aeghh-pro03a", for which we list the top 10 retrieved passages using the standard cosine similarity score and our sparsity-based score ( $\alpha = 1.78$  selected from the dev set). The first five letters of a passage ID represent the argument topic ID; "pro/con" denotes the argument stance; suffix "a/b" indicates the argument and its corresponding counter-argument; "para0/para1/para2" are three paraphrases generated by GPT4.

As shown in Table 9, for the example query "aeghh-pro03a", its correct counter-argument, "aeghh-pro03b" (in red), ranks fourth using the cosine score but first using the cosine + hoyer score. Meanwhile, its paraphrases "aeghh-pro03a-para0/1/2" (in blue) achieve high cosine scores but low sparsity scores.

Method CL(Cosine)		CL(Cosine)+SparseCL(Hoyer)			seCL(Hoyer)	
Rank	Cosine	Passage ID	Overall	Cosine	Hoyer	Passage ID
1	0.940	aeghh-pro03a-para0	1.683	0.794	0.499	aeghh-pro03b
2	0.926	aeghh-pro03a-para2	1.644	0.719	0.519	aeghh-con02a-para0
3	0.916	aeghh-pro03a-para1	1.617	0.716	0.506	aeghh-con02a-para2
4	0.794	aeghh-pro03b	1.606	0.940	0.374	aeghh-pro03a-para0
5	0.719	aeghh-con02a-para0	1.602	0.718	0.496	aeghh-con02a-para1
6	0.718	aeghh-con02b	1.528	0.718	0.454	aeghh-con02b
7	0.718	aeghh-con02a-para1	1.494	0.916	0.324	aeghh-pro03a-para1
8	0.716	aeghh-con02a-para2	1.426	0.926	0.280	aeghh-pro03a-para2
9	0.696	aeghh-con02a	1.396	0.669	0.408	dhwif-pro02b
10	0.692	aeghh-pro04a-para0	1.344	0.628	0.402	thggl-con03b

Table 9. An example query analysis for counter-argument retrieval. The passage ID in red represents the ground-truth counter-argument, while the passage IDs in blue are paraphrases of the query argument.

# E. Hyper-parameters for training and inference

Models	Model Size	Backbone	ep (	CL lr	SPA ep	RSECL lr	temp	bz
GTE-large-en-v1.5	434M	BERT + RoPE + GLU	1	1e-5	3	2e-5	0.01	64
UAE-Large-V1	335M	BERT	1	2e-5	3	2e-5	0.02	64
bge-base-en-v1.5	109M	BERT	1	2e-5	3	2e-5	0.02	64

Here we present the training details (Table 10) for our experiments on Arguan, HotpotQA, and MSMARCO.

Table 10. Training parameters for Arguana. We set max sequence length to be 512 for Arguana dataset and 256 for HotpotQA and MSMARCO datasets.

### F. Efficiency test of cross-encoder and vector calculation

To further compare the efficiency of cross-encoders and Hoyer sparsity calculations, we perform the following experiments:

- We choose "bge-reranker-base" and "bge-reranker-large" to be our cross-encoders. We use them to calculate the similarity between one query from Arguana's test set and 100 documents from Arguana's corpus. We report the average running time of this method for 100 queries.
- We choose "bge-base-en-v1.5" and "bge-large-en-v1.5" to be our bi-encoders. Suppose we have preprocessed all the sentence embeddings. We use it to calculate the Hoyer sparsity between one query embedding from Arguana's test set and 100 document embeddings from Arguana's corpus. We report the average running time of this method for 100 queries.

Please see Table 11 for the running time of different methods. We can see that the calculation of Hoyer sparsity is at least 200 times faster than running a cross-encoder.

Cross-encoder	Model size	Time
bge-reranker-base	278M	0.8832s
bge-reranker-large	560M	1.6022s
Bi-encoder	Embedding dimension	Time
bge-base-en-v1.5	768	0.0029s
bge-large-en-v1.5	1024	0.0036s

Table 11. Average running time for calculating the score functions between one Arguana query and 100 Arguana documents

# G. Data generation details for MSMARCO and HotpotQA experiments in Section 4.2

We use "gpt-4-turbo" to generate paraphrases and contradictions for our experiment in Section 4.2. The prompts we use are in Table 13. We set *temperature* = 1 and n = 3 (to generate 3 outputs). Please see Table 12 for some examples of generated paraphrases and contradictions and Table 14 for the number of unique passages generated from each dataset split. We have verified that there is no overlap between different splits or different dataset. For generated data quality, we sampled 500 data points from the dataset and 2 people annotated the data to check the quality. Our agreement score (using Cohen's Kappa) is 0.98, indicating the quality of the generated data. We present the first 20 passages and their generated paraphrases and contradictions from MSMARCO and HotpotQA datasets in the end of this section for readers' reference.

Datasets	Orginal	Paraphrase	Contradiction
MSMARCO	In addition to the <b>high financial</b> <b>value</b> of higher education, higher education also makes individuals much <b>more</b> <b>intelligent</b> than what they would be with just a high school education	Beyond its <b>significant</b> <b>monetary worth</b> , higher education substantially <b>enhances a person's</b> <b>intelligence</b> compared to merely completing high school	Besides the <b>low financial</b> <b>significance</b> of higher education, higher education often renders individuals <b>no more intelligent</b> than they would be with just a high school education
HotpotQA	Ice hockey is a <b>contact</b> team sport played on <b>ice</b> , usually in a <b>rink</b> , in which two teams of skaters use their <b>sticks</b> to shoot a <b>vulcanized rubber puck</b> into their opponent's net to score points	Ice hockey is a <b>contact</b> sport where two teams compete on an <b>ice surface</b> , typically in a <b>rink</b> , using <b>sticks</b> to hit a <b>vulcanized</b> <b>rubber puck</b> into the opposing team's net to earn points	Ice hockey is a <b>non-contact</b> team sport played on <b>grass</b> , often in an <b>open field</b> , where two teams of players use their <b>feet</b> to kick a <b>soft leather ball</b> into their opponent's goal to score points

*Table 12.* Examples of passages from MSMARCO and HotpotQA datasets, with their generated paraphrases, and generated contradictions. Highlighted key-words represent exact matchings or contradictions

Task	Prompt
Generating paraphrases	Paraphrase the given paragraph keeping its original meaning. Do not add information that is not present in the original paragraph. Your response should be as indistinguishable to the original paragraph as possible in terms of length, language style, and format. Begin your answer directly without any introductory words.
Generating contradictions	Rewrite the given paragraph to contradict the original content. Ensure the revised paragraph changes the factuality of the original. Your response should be as indistinguishable to the original paragraph as possible in terms of length, language style, and format. Begin your answer directly without any introductory words.

Table 13. Prompts used to generate paraphrases and contradictions for MSMARCO and HotpotQA documents.

Dataset	Train	Dev	Test
MSMARCO	59500	5950	44101
HotpotQA	59074	5896	81673

*Table 14.* Number of unique passages generated from MSMARCO and HotpotQA corpus via GPT-4. We have verified that there is no overlap between different splits or different datasets.

Contradiction Retrieval via Contrastive Learning with Sparsity

Passage ID	Orginal	Paraphrase	Contradiction
MSMARCO-	What is a BIOS? BIOS is an	BIOS stands for Basic In-	What is a BIOS? BIOS stands
6725993	acronym for Basic Input / Out-	put/Output System. It exists	for Basic Input / Output System.
	put System. On virtually ev-	on almost all computers and is	On nearly no modern computer,
	ery computer available, the BIOS	responsible for ensuring that all	the BIOS is used to allow all the
	makes sure all the other chips,	the various components such as	other chips, hard drives, ports,
	hard drives, ports and CPU func-	chips, hard drives, ports, and	and CPU to operate disjointedly.
	tion together. What BIOS Does	the CPU work in unison. The	What BIOS Does Not Do The
	The BIOS software has a num-	primary function of BIOS soft-	BIOS software has limited func-
	ber of different roles, but its most	ware is to initiate the operating	tions, and its least significant role
	important role is to load the op-	system. Upon powering up your	is to load the operating system.
	erating system. When you turn	computer, the microprocessor	When you turn on your computer
	on your computer and the micro-	seeks to perform its initial	and the microprocessor attempts
	processor tries to execute its first	instruction, which it retrieves	to execute its first instruction, it
	instruction, it has to get that in-	through the BIOS.	does not rely on the BIOS to ob-
	struction from somewhere.		tain that instruction.
MSMARCO-	Causes of nitrates in urine. 1 Ni-	Reasons for nitrates in urine. 1	Causes of nitrates in urine. 1 Ni-
6909689	trites/nitrates in urine refer to a	Nitrates or nitrites in urine are	trites/nitrates in urine signify a
	byproduct formed due to the ac-	byproducts produced from bac-	byproduct that is not associated
	tion of bacteria occurring in the	terial activity within the urinary	with bacterial action in the uri-
	urinary tract. The kidneys per-	tract. The kidneys have the crit-	nary tract. The kidneys, unable
	form the vital function of clean-	ical role of purifying the blood	to filter these components once
	ing the blood by littering out the	by removing unwanted, narmiui	formed, do not influence their
	blood. The nitrite developed by	created by such bacteria are not	trites converted from nitrates in
	biolog. The intrice developed by	able to be filtered out by the kid-	food enter the urinary system in-
	kidnevs	nevs	dependent of any bacterial inter-
	literie y 5.	ney s.	vention.
MSMARCO-	Food Additives. Food additives	Food Additives. Substances	Food Additives. Food additives
594175	are substances added intention-	known as food additives are de-	are substances naturally present
	ally to foodstuffs to perform cer-	liberately incorporated into foods	in food items to perform specific
	tain technological functions, for	to fulfill specific technological	biological purposes, such as to
	example to colour, to sweeten or	roles such as coloring, sweeten-	enhance flavor, to sour, or to de-
	to help preserve foods. In the	ing, or preserving the food. In	crease the shelf life of foods. In
	European Union all food addi-	the European Union, these addi-	the European Union, no specific
	tives are identified by an E num-	tives are marked with an E num-	identification like an E number is
	ber. Food additives are always	ber. The presence of food addi-	utilized for food additives. Food
	included in the ingredient lists of	tives in products is consistently	additives are rarely disclosed in
	foods in which they are used.	disclosed in the ingredient lists of	the ingredient lists of foods in
MEMADOO	monombatitutad allows (an	the loods where they are utilized.	which they appear.
MSMAKCO-	annic chemistry) An allega with	chemistry) An alkona charac	chemistry) An alkana with
0200110	$f_{\text{game chemistry}}$ . All alkelle Willi the general formula $\text{PHC}-\text{CH}$ 2	terized by the general formula	the general formula DUC-CUD
	where <b>R</b> is any organic group:	RHC-CH2 where P represents	where R represents any organic
	only one carbon atom is bonded	any organic group: a single car-	group: each carbon atom in the
	directly to one of the carbons	bon atom is directly bonded to	carbon-to-carbon double bond is
	of the carbon-to-carbon double	one carbon of the double bond be-	directly bonded to its own distinct
			set of the other and the other and the other
	bond.1 Facebook.2 Twitter.	tween carbons.1 Facebook.2 Twit-	organic group. 1 Facebook. 2

Contradiction Retrieval via Contrastive Learning with Sparsity

Passage ID	Orginal	Paraphrase	Contradiction
MSMARCO-	Depending on the specific germ.	The incubation period varies from	Depending on the specific germ.
3366870	the incubation period is between	12 hours to 5 days, typically	the incubation period is between
	12 hours and 5 days, usually 48	around 48 hours, depending on	6 days and 2 weeks, usually 7
	hours. To answer your question,	the germ involved. Regarding	days. To answer your question,
	the common cold is contagious	your question, one can spread the	the common cold is contagious
	between 24 hours before onset of	common cold from 24 hours be-	from 5 days before the onset of
	symptoms until 5 days after on-	fore symptoms start up to 5 days	symptoms until 10 days after on-
	set.	following their appearance.	set.
MSMARCO-	Noun. 1. decrescendo - (music)	Noun. 1. decrescendo - (music) a	Noun. 1. crescendo - (mu-
6494839	a gradual decrease in loudness.	progressive reduction in volume.	sic) a gradual increase in loud-
	diminuendo. softness-a sound	diminuendo. softness - a charac-	ness. crescendo. loudness-a
	property that is free from loud-	teristic of sound marked by low	sound property characterized by
	ness or stridency; and in softness	volume and an absence of harsh-	high volume and stridency; and
	almost beyond hearing. music-an	ness; approaching a level that is	in loudness almost beyond en-
	artistic form of auditory commu-	barely audible. music - an art	durance. silence-a form of non-
	nication incorporating instrumen-	form of auditory communication	auditory communication devoid
	tal or vocal tones in a structured	that uses instrumental of vocal	of instrumental or vocal tones, un-
	and continuous manner.	tained way	structured and intermittent.
MSMAPCO	Due to its important role in	Pannin anzuma critical for milk	Despite its oracial function in our
03/13/	curdling milk rennin enzyme	curdling is extensively utilized	dling milk the rennin enzyme
954154	is widely used in the food in-	in the food sector especially	is minimally used in the food
	dustry notably in the produc-	for cheese production Histori-	industry particularly in the pro-
	tion of cheese Rennin for	cally cheese-making rennin was	duction of cheese Rennin for
	cheese-making was once derived	sourced primarily from the dried	cheese-making used to be primar-
	mainly from the dried stomachs	stomachs of young calves and var-	ily sourced synthetically and not
	of calves and from some non-	ious non-animal origins. This en-	from natural origins such as the
	animal sources.ennin enzymes	zyme is secreted by the stomach	dried stomachs of calves. Ren-
	are produced by the stomach cells	cells of young mammals, with	nin enzymes are produced by the
	of young mammals. Rennin is	secretion levels peaking immedi-	stomach cells of older mammals.
	secreted in large amounts right	ately following birth before grad-	Rennin is secreted in minimal
	after the birth and then its pro-	ually declining. Over time, the	amounts right after birth, and then
	duction gradually drops off. It	significance of rennin decreases	its production progressively in-
	is then eclipsed in importance by	as it is overshadowed by the en-	creases. It continues to retain its
	the Pepsin enzyme.	zyme Pepsin.	importance over the Pepsin en-
MOMADOO			zyme.
MSMARCO-	winter Moth (Operophtera bru-	winter Moth (Operophtera bru-	summer Moth (Operophtera sol-
44/3191	mata). up in traps, at least, in	nata) has been caught in traps	sucia). down in nets, at least, in
	one place in southeastern CT	Hampshire coastal Maine a loca	one place in northwestern CT
	and out on Long Island Mas	tion in southeastern Connecticut	and out on Long Island Mas
	sachusetts still appears to have the	and on L ong Island. The most sig	sachusetts seems to have the
	largest and most damaging pop-	nificant and harmful populations	smallest and least harmful pop-
	ulations of this pest, spanworm	of this pest seem to still be in Mas-	ulations of this pest, inchworm
	(Eidt et al.	sachusetts, according to Eidt et al.	(Eidt et al.
	× · · · · · ·	,	×

Contradiction Retrieval via Contrastive Learning with Sparsity

Passage ID	Orginal	Paraphrase	Contradiction
MSMARCO-	Actor Charles Bronson dies at 81	Charles Bronson, star of 'Death	Actor Charles Bronson thrives
7066556	Death Wish movie star Charles	Wish' and former Pennsylvania	at 81 Death Wish movie star
	Bronson, the coal miner from	coal miner who turned to acting	Charles Bronson, the coal miner
	Pennsylvania who drifted into	as a villain before becoming a	from Pennsylvania who drifted
	films as a villain and became a	renowned action hero, has passed	into films as a hero and became
	hard-faced action star, has died.	away at the age of 81. Bronson	a beloved romantic lead, is thriv-
	Bronson, 81, died of pneumonia	died from pneumonia at Cedars-	ing. Bronson, 81, recently recov-
	at Cedars-Sinai Medical Centre	Sinai Medical Centre in Los An-	ered from a mild cold at his home
	in Los Angeles, with his wife at	geles while his wife was by his	in Los Angeles, without needing
	his bedside, publicist Lori Jonas	side, his publicist Lori Jonas con-	hospital care, publicist Lori Jonas
	said. He had been in the hospital	firmed. He had been hospitalized	said. He has been in excellent
	for weeks.	for several weeks.	health.
MSMARCO-	(Redirected from Primary Wave	Primary Wave Entertainment is	Secondary Echo Entertainment
2539424	Music) Primary Wave Entertain-	a comprehensive entertainment	is a specialized entertainment
	ment is a full-service entertain-	firm based in the United States,	company lacking capabilities in
	ment company with expertise	specializing in Talent Manage-	Talent Management, Literary
	in Talent Management, Literary	ment, Literary Management, Mu-	Management, Music Publishing,
	Management, Music Publishing,	sic Publishing, Branding, Digital	Branding, Digital Marketing, and
	Branding, Digital Marketing and	Marketing, and Licensing.	Licensing, located outside the
	Licensing based in the United		United States.
MSMADCO	States.	In 2010 the median annual	As of 2010 modian salarias
5803222	As of 2010, incutail sataries	solorios for unit socretorios vor	As of 2010, median salaries
3893222	\$ 23,008 \$ 31,604 while me	ied between \$ 23,008 and \$	\$ 35 000 $$$ 42 500 while me
	dian hourly wages for hos-	31.604 and their median hourly	dian hourly wages for hospital
	nital unit secretaries ranged	nay in hospitals was between \$	unit secretaries were between
	from $10.83-14.67$ according to	10.83 and $$14.67$ as reported by	\$ 16 50-\$ 19 75 according to
	PavScale.com.	PavScale.com.	PavScale.com.
MSMARCO-	Answers.com <sup>®</sup> is making the	Answers.com® improves the	Answers.com <sup>®</sup> often misleads
3476298	world better one answer at a time.	world by providing one answer	more than it informs. Cleffa
	cleffa evolves at level 30 you can	at a time. Cleffa evolves at level	evolves with high friendship, not
	find it at mt. coronet. 2 people	30 and can be found at Mt. Coro-	at level 30, and it is found in mul-
	found this useful.	net. This information was useful	tiple locations, not just Mt. Coro-
		to 2 people.	net. 2 people found this mislead-
			ing.

Contradiction Retrieval via Contrastive Learning with Sparsity

	0 1	D 1	
Passage ID	Orginal	Paraphrase	Contradiction
MSMARCO-	In the sidebar, make sure you se-	In the sidebar, choose Music lo-	In the sidebar, ensure that you
6150895	lect Music under the iPod (in-	cated beneath the iPod (indented	choose Music under LIBRARY,
	dented below the iPod). That's	under the iPod). This displays the	not under the 1Pod (indented be-
	the iPod's content list and shows	list of songs stored on your iPod.	low the iPod). That selection
	the list of songs on the iPod.	Be sure not to select Music listed	links to your computer's music
	Do NOT select Music under LI-	under LIBRARY. To remove all	collection instead of displaying
	BRARY. If you want to delete	songs from the iPod, click on any	the list of songs stored on the
	ALL of the songs on the iPod,	track in that music list. Perform	iPod. Always select Music un-
	click on any song on that song	a Select All by pressing Cmd-A	der LIBRARY. If you wish to pre-
	list.Do a Select All, which is	on your keyboard or by choosing	serve all the songs on the iPod,
	Cmd-A on the keyboard or Se-	Select All from the Edit menu in	avoid clicking on any song in that
	lect All from the Edit menu (in	the menu bar. Once all songs are	song list. Do not use Select All,
	menu bar). With all of the songs	highlighted, hit the Delete key on	which would be Cmd-A on the
	selected, press Delete on the key-	your keyboard.	keyboard or Select All from the
	board.		Edit menu (in menu bar). With
			none of the songs selected, do not
			press Delete on the keyboard. En-
			sure you do not alter the content
			list that appears under the iPod
			selection in the sidebar.
MSMARCO-	The Associate of Applied Science	The Associate of Applied Science	The Associate of Applied Sci-
1814356	Degree requires that you have	Degree mandates fulfillment of	ence Degree does not require
	completed the TSI requirements.	TSI prerequisites. To graduate	you to have completed the TSI
	Course Requirements Graduation	with an Associate of Applied Sci-	requirements. Course Require-
	with the Associate of Applied Sci-	ence Degree in Welding or to re-	ments Graduation with the As-
	ence Degree in Welding or the	ceive either the Combination or	sociate of Applied Science De-
	completion of the Com-bination	Structural Welding Certificate, it	gree in Welding or the completion
	or Structural Welding Certificate	is essential to pass a Comprehen-	of the Combination or Structural
	re-quires successful completion	sive Exit Exam.	Welding Certificate does not re-
	of a Comprehensive Exit Exam.		quire the successful completion
			of a Comprehensive Exit Exam.
MSMARCO-	Pluto's rotation period, its day is	Pluto's day, its rotation period	Pluto's rotation period, its day is
6108145	equal to 6.39 Earth days. Like	spans 6.39 Earth days. Similar to	equal to 153.3 Earth hours. Un-
0100110	Uranus. Pluto rotates on its side	Uranus. Pluto has a rotation along	like Uranus. Pluto rotates upright
	on its orbital plane with an ax-	its orbital plane on its side with	in its orbital plane with an ax-
	ial tilt of $120\hat{A}^\circ$ and so its sea-	an axial tilt of 120 degrees Con-	ial tilt of 30° and so its seasonal
	sonal variation is extreme: at its	sequently it experiences extreme	variation is moderate: at its sol-
	solutions one-fourth of its surface	seasonal changes: during its sol	stices nearly all of its surface ex
	is in continuous devlight whereas	stices a quarter of its surface an	neriences alternating daylight and
	another fourth is in continuous	iovs perpetual daylight while an	darkness Dluto (minor planet
	darknass luto (minor planet das	other quarter remains in constant	designation: 12/240 Divite) is not
	ignation: 13/3/0 Pluto) is a	darkness Pluto designated as	considered a dwarf planet but a
	dworf plonat in the Vyince balt	minor planet 124240 is closed as	maior planat controlly situated in
	uwari pianet in the Kuiper belt,	ninor-planet 154540, is classified	the Kuiner helt a suggest of
	a ring of doctes beyond Neptune.	as a dwarf planet located in the	hadias bayand Narture
		Kuiper beit, a collection of ob-	bodies beyond Neptune.
		jects that orbits around Neptune.	

Contradiction Retrieval via Contrastive Learning with Sparsity

Passage ID	Orginal	Paraphrase	Contradiction
MSMARCO-	Domestic Costa Rica flights.	Costa Rica domestic flights with	International Costa Rica flights.
5618382	Your premier Costa Rica airline	Nature Air. your top airline	Your last choice for Costa Rica
	choice for travel and vacation	choice for traveling and vacation-	airline for travel and vacation
	flights within Costa Rica, offer-	ing across Costa Rica. Providing	flights outside of Costa Rica, of-
	ing 74 daily flights to 17 destina-	74 daily flights, this Costa Rica	fering no daily flights to any des-
	tions in Costa Rica! Nature Air	airline connects to 17 local desti-	tinations outside Costa Rica! Na-
	the Costa Rica Domestic Airline.	nations. As the premier domestic	ture Air. the International Costa
	offers domestic flights to 17 des-	airline in Costa Rica, Nature Air	Rica Airline, provides no interna-
	tinations in Costa Rica. Nature	facilitates 74 daily flights, includ-	tional flights to destinations out-
	Air Costa Rica Airline, offers 74	ing routes to and from Juan San-	side of Costa Rica. Nature Air
	daily domestic flights including	tamaria International Airport, en-	Costa Rica Airline, offers zero
	to and from Juan Santamaria In-	abling seamless connections with	daily international flights exclud-
	ternational Airport, which allows	international flights.	ing Juan Santamaria International
	easy connections to International		Airport, which prevents any con-
	flights.		nections to international flights.
MSMARCO-	There were two groups of Cubists	During the peak of the Cubist	There was only one group of
585206	during the height of the move-	movement, from 1909 to 1914,	Cubists throughout the peak of
	ment, 1909 to 1914. Pablo Pi-	there existed two factions of Cu-	the movement, from 1909 to
	casso (1881-1973) and Georges	bists. Pablo Picasso (1881-1973)	1914. Pablo Picasso (1881-1973)
	Braque (1882-1963) are known	and Georges Braque (1882-1963)	and Georges Braque (1882-1963)
	as the Gallery Cubists because	were identified as the Gallery Cu-	are regarded as Independent Cu-
	they exhibited under contract	bists due to their contractual exhi-	bists because they showcased
	with Daniel-Henri Kahnweiler's	bitions with the gallery of Daniel-	their works without any exclusive
	gallery.	Henri Kahnweiler.	agreements, avoiding ties with
			major galleries like that of Daniel-
MOMADOO			Henri Kahnweiler's.
MSMARCO-	Cody is a city in Park County,	Cody, located in Park County,	Cody is a township in Park
454/162	wyoming, United States. It is	Wyoming, USA, is named for	County, wyoming, United States.
	named after William Frederick	William Frederick Cody, better	It is named after James Frederick
	Cody, primarily known as Bullalo	known as Bullaio Bill, due to his	Long from his apposition to the
	Bill, from his part in the creation	significant role in founding the	James, from his opposition to the
	of the original town.	initial settlement.	establishment of the original set-
MSMAPCO	Hemoglobin is what actually	Hemoglobin enables red blood	Hemoglobin actually restricts
6813425	gives your red blood cells the	cells to transport oxygen effec	your red blood cells' capacity to
0015425	ability to carry oxygen A high	tively An elevated hematocrit	transport oxygen A low hemat-
	hematocrit simply means that you	indicates a greater concentration	ocrit indicates a diminished con
	have a higher concentration of	of hemoglobin within the blood	centration of hemoglobin in the
	hemoglobin in the blood First of	stream This could either sug	blood Primarily it could sig
	all it could indicate dehydration	gest dehydration or a normal level	nify overhydration or excessive
	or normal amount of hemoglobin	of hemoglobin coupled with a re	blood (plasma) volume with too
	but too little blood (plasma) vol	duced volume of blood plasma	little hemoglobin. If this situation
	ume If this is the case depending	Should this be the situation and	arises depending on how low the
	on how high the hematocrit is re-	depending on the severity of	hematocrit is dehydration mea-
	hydration might be needed	the hematocrit levels rehydration	sures might be required
	ing and the medical	may be necessary.	sales inght be required.
		,,	

Contradiction Retrieval via Contrastive Learning with Sparsity

Passage ID	Orginal	Paraphrase	Contradiction
MSMARCO-	Imperial system of measurement,	The Imperial measurement sys-	Metric system of measurement,
1760626	on the other hand refers to the sys-	tem, alternatively, was the stan-	however, refers to the system
	tem used in the British Empire in	dard in the British Empire dur-	used globally. Still, the adop-
	the 19th and 20th centuries. How-	ing the 19th and 20th centuries.	tion of the imperial system has
	ever, after adoption of metric sys-	Yet, with the adoption of the met-	been expanding, being notably
	tem, Imperial system has been re-	ric system, the Imperial system's	preferred in an increasing num-
	duced to a few countries of the	usage has declined and is now	ber of countries, including the UK
	world, notably UK, and surpris-	primarily limited to a handful of	and, interestingly, the US.
	ingly US.	countries, notably the UK and, in-	
		terestingly, the US.	
HotpotQA-	St. Moritz (also German: "Sankt	St. Moritz (German: "Sankt	St. Moritz (also German: "Sankt
303193	Moritz, Romansn: , Italian:	Montz, Romansn: , Italian:	Montz, Romansn: , Italian:
	Moritz") is a high Alpine resort	Moritz") is a high Alpine resort	Moritz") is a lowland resort in
	in the Engadine in Switzerland	town situated in the Engadine	the plains of Lower Engadine in
	at an elevation of about 1800 m	in Switzerland located approxi-	Switzerland at an elevation of
	above sea level. It is Upper Enga-	mately 1800 m above sea level. It	about 300 m above sea level. It is
	dine's major village and a munic-	serves as the primary village of	Lower Engadine's minor village
	ipality in the district of Maloja in	Upper Engadine and is a munic-	and a small community in the dis-
	the Swiss canton of Graubünden.	ipality within the district of Mal-	trict of Maloja in the Swiss canton
		oja in the canton of Graubünden,	of Graubünden.
		Switzerland.	
HotpotQA-	John Adedayo B. Adegboyega	John Adedayo B. Adegboyega,	John Adedayo B. Adegboyega
32060208	(born 17 March 1992), known	born on 17 March 1992 and	(born 17 March 1992), profes-
	professionally as John Boyega,	professionally known as John	sionally known as John Boyega,
	is an English actor and producer	Boyega, is an English actor and	is a Scottish singer and song-
	best known for playing Finn in	producer. He is most recognized	writer best recognized for his de-
	the 2015 film "", the seventh film	for his portrayal of Finn in the	but in the 2018 album "Horizons",
	of the Star wars series. Boyega	2015 movie "Star wars: The	the fourth album in his musical ca-
	United Kingdom for his role as	Force Awakens, which is the	in the international music score
	Moses in the 2011 sci fi comedy	Wars" saga Boyega first gained	for his performance at the Edin
	film "Attack the Block"	fame in the UK for his perfor-	burgh Festival Fringe in the musi-
	min Attack the Diock .	mance as Moses in the 2011 sci-	cal "Glimpse of the Stars"
		ence fiction comedy "Attack the	cur ompse of the stars .
		Block".	
HotpotQA-	Sarah Davis (born 1976) is an	Sarah Davis, born in 1976, is a	Sarah Davis (born 1976) is an
42146101	American politician and a Repub-	Republican politician serving in	American politician and a Demo-
	lican member of the Texas House	the Texas House of Representa-	cratic member of the Texas House
	of Representatives; she was first	tives. She entered office follow-	of Representatives; she was first
	elected in the Tea Party wave of	ing the Tea Party election wave	elected in the progressive wave
	2010. Davis' district contains The	in 2010. Her district includes key	of 2010. Davis' district excludes
	Galleria and the Texas Medical	areas such as The Galleria and the	The Galleria and the Texas Medi-
	Center.	Texas Medical Center.	cal Center.

Contradiction Retrieval via Contrastive Learning with Sparsity

HotpotQA- 1982071"Lola" is a song written by Ray Davies and performed by English rock band the Kinks on their al- bum "Lola Versus Powerman and the Moneygoround, Part One"."Lola" is a track penned by Ray Davies and executed by the British rock group the Kinks, fea- tured on their album "Lola Ver- sus Powerman and the Moneygor- ound, Part One." The song details a romantic en-"Lola" is a track penned by Ray Davies and executed by the British rock group the Kinks, fea- tured on their album "Lola Ver- ound, Part One." The song nar-"Lola" is a song written by Davies and performed by English rock band the Kinks on their bum "Lola Versus Powerman the Moneygoround, Part Ore."
1982071Davies and performed by English rock band the Kinks on their al- bum "Lola Versus Powerman and the Moneygoround, Part One".Ray Davies and executed by the British rock group the Kinks, fea- tured on their album "Lola Ver- sus Powerman and the Moneygor- ound, Part One." The song details a romantic en-Davies and performed by Eng rock band the Kinks on their bum "Lola Versus Powerman and tured on their album "Lola Ver- ound, Part One." The song nar-Davies and performed by Eng rock band the Kinks on their bum "Lola Versus Powerman the Moneygoround, Part One." The song nar-
rock band the Kinks on their al- bum "Lola Versus Powerman and the Moneygoround, Part One".British rock group the Kinks, fea- tured on their album "Lola Ver- sus Powerman and the Moneygor- ound, Part One." The song nar-rock band the Kinks on their bum "Lola Ver- the Moneygoround, Part One".The song details a romantic en-ound, Part One." The song nar- the Moneygor DataThe song nar- the Moneygoround, Part One." The song nar-
bum "Lola Versus Powerman and the Moneygoround, Part One".tured on their album "Lola Ver- sus Powerman and the Moneygor- ound, Part One." The song nar-bum "Lola Versus Powerman the Moneygoround, Part O The song narrates the muno
the Moneygoround, Part One". sus Powerman and the Moneygor- the Moneygoround, Part C The song details a romantic en- ound, Part One." The song nar- The song narrates the muno
The song details a romantic en-   ound, Part One." The song nar-   The song narrates the mund
counter between a young man and rates a romantic interaction be- interactions between a mic
a possible transvestite, whom he tween a young man and a likely aged man and a woman he kn
meets in a club in Soho, London. transvestite he meets in a Soho, which they experience durin
In the song, the narrator describes   London club. It portrays the typical day. In the song, the
nis confusion towards a person young man's perplexity toward rator simply mentions a per
woman and talked like a man" "walked like a woman and talked any gender ambiguity stating
Although Ray Davies claims that like a man " Ray Davies has at she "walked like a woman
the incident was inspired by a tributed the inspiration for the talked like a woman" Altho
true encounter experienced by the song to an actual event that the Ray Davies has denied any
band's manager, alternate expla- band's manager encountered, al- life inspiration behind the s
nations for the song have been though drummer Mick Avory has claiming it to be purely fiction
given by drummer Mick Avory. offered different explanations for consistent explanations for
the song's origins. song's origin have been suppo
by drummer Mick Avory.
HotpotQA-"No More Sad Songs" is a song"No More Sad Songs" is a track"Yes More Happy Songs"
8540095 by British girl group Little Mix by the British girl band Little Mix, song by American boy band
from the group's fourth studio al- featured on their fourth studio Mix from the group's fifth liv
bum, "Glory Days" (2016). The album titled "Glory Days" from bum, "Tragic Nights" (2017).
song was written by Emily War- 2016. The creators of the song song was written by John J
ren, Edvard Førre Erfjord, Hen- include Emily Warren, Edvard Alexander Back, Richard Sn
rik Michelsen and Iash Phillips; Førre Erfjord, Henrik Michelsen, son, and Lity Johnson; produ
Keerns A remix version feature tion handled by Electric and Ioe cover version lacking any
ing newly recorded vocals from Kearns On 3 March 2017 a tributions from Canadian si
American rapper Machine Gun remixed version of the song in- Ion Rellion was dismissed as
Kelly, was released as the third corporating new vocals by Amer- fourth single from the album
single from the album on 3 March   ican rapper Machine Gun Kelly,   December 2018, through Rhy
2017, through Syco Music. was issued as the album's third Records.
single through Syco Music.

Contradiction Retrieval via Contrastive Learning with Sparsity

Dassage ID	Orginal	Daranhrasa	Contradiction
HotpotOA	I Might Be Wrong: Live Record	I Might Be Wrong: Live Record	I Might Be Right: Studio Sec
228538	ings is a live album by the En	ings is a live album from the	sions is a studio album by the
220330	dish rock hand Padiohaad ra	British rock group Padiohand	American non band Padiohaad
	lagged on 12 November 2001	which was published on 12	American pop band Radionead,
	hy Darlanhana Dacarda in the	Nevember 2001 by Derlephone	hy Darlanhana Dagarda in the
	United Kingdom and a day later	November 2001 by Partophone	by Pariophone Records in the
	United Kingdom and a day later	Records in the United Kingdom	United Kingdom and a day ear-
	by Capitol Records in the United	and one day later in the United	lier by Capitol Records in the
	States. Recorded during Radio-	States by Capitol Records. Cap-	United States. Recorded during
	head's 2001 tour, it comprises	tured during the group's 2001	Radionead's studio sessions in
	performances of songs from the	tour, it features live renditions	2005, it features remixes of songs
	band's fourth and fifth albums	of tracks from their fourth and	from the band's sixth and sev-
	$^{\circ}$ Kid A <sup><math>\circ</math></sup> (2000) and $^{\circ}$ Amnesiac <sup><math>\circ</math></sup>	fifth albums, "Kid A" (2000) and	enth albums "Hail to the Thief"
	(2001), plus the song "True Love	"Amnesiac" (2001), as well as the	(2003) and "In Rainbows" $(2007)$ ,
	Waits", which would not be re-	track "True Love Waits," which	excluding the song "True Love
	leased on a studio album until "A	was not released on a studio al-	Waits", which had already been
	Moon Shaped Pool" (2016).	bum until "A Moon Shaped Pool"	released on a studio album prior
		in 2016.	to "A Moon Shaped Pool" (2016).
HotpotQA-	"My Brave Face" is a single from	"My Brave Face" is a track from	"My Brave Face" is a track from
14/4130/	Paul McCartney's 1989 album,	Paul McCartney's 1989 release,	Paul McCartney's 1989 album,
	Flowers in the Dirt . Written	Flowers in the Dirt. The	Flowers in the Dirt', which he
	"My Brave Feed" is one of the	and Elvis Costello is among	from Elvis Costello "My Broye
	My Brave Face is one of the	and EIVIS Costello, is among	From Eivis Costello. My Brave
	most popular songs from Flow-	the album's most notable tracks.	Face is one of the least recog-
	in the United Kingdom a weak	the UV shorts and as	Dirt" It did not make a signifi
	ofter its debut and #25 in the	and to #25 in the US shorts	Diff. If did not make a signifi-
	United States 7 weeks ofter its de	cended to #25 in the US charts	to optor the top 40 in the United
	but It was MaCartnay's last top	seven weeks post-debut. It repre-	Kingdom and the United States
	40 hit on the "Pillboard" Hot 100	sents McCarthey's lina top 40 ap-	It was not MaCartnay's final ton
	40 lift on the Billobard Hot 100	100 until his 2014 partnarship	10 bit on the "Pillboard" Hot 100
	Kenve West "Only One" and as	with Kanya Wast on the track	40 lift on the Binobard flot 100,
	af 2017 is the last Dillboard ton	"Only One" As of 2017 it re	as he continued to achieve suc-
	40 hit with any former Bestle in	Unity One . As of 2017, it fe-	with his subsequent collabora
	the lead gradit	on the Billboard short featuring a	tions As of 2017 other former
	the lead credit.	former Poetle as the lead artist	Postlas have also achieved top 40
		former beaue as the lead artist.	bits on Pillboard in load roles
HotpotOA	Lonesome Dove: The Series is	Lonesome Dove: The Series on	Lonesome Dove: The Series is a
54385601	an American western drama tala	American Western drama TV se	Canadian comedy television se
54505071	vision series that debuted in first-	ries premiered in first-run syndi-	ries that premiered exclusively on
	run syndication on September 26	cation on September 26, 1004. It	the premium cable network HBO
	1994 It serves as continuation	continues the narrative from the	on October 15, 1995. It is pre-
	of the story of the miniseries of	similarly titled miniseries Fea-	sented as a parody of the original
	the same name. The television	turing Scott Bairstow and Fric	miniseries bearing the same title
	series starred Scott Bairstow and	McCormack, the series was ex-	The television series featured Tim
	Eric McCormack, and its execu-	ecutive produced by Suzanne de	Curry and Hugh Grant as leads
	tive producers were Suzanne de	Passe and Robert Halmi Jr. Pro-	with executive production helmed
	Passe and Robert Halmi Jr. The	duction was handled by Telegenic	by John Smith and Jane Doe. The
	series was produced by Telegenic	Programs Inc. and RHI Enter-	series was produced by Generic
	Programs Inc. and RHI Entertain-	tainment, in collaboration with	Studios Ltd. and ABC Entertain-
	ment in association with Rysher	Rysher TPE and the Canadian TV	ment without any collaboration
	TPE, in conjunction with Cana-	network CTV.	with the American television net-
	dian television network CTV.		work NBC.

Contradiction Retrieval via Contrastive Learning with Sparsity

Passage ID	Orginal	Paranhrase	Contradiction
HotpotOA-	The Revolution was a villainous	The Revolution was a nefarious	The Revolution was a heroic
44751816	stable in Total Nonston Action	group in Total Nonston Action	stable in Total Nonston Action
11/21010	Wrestling (TNA) consisting of	Wrestling (TNA) that included	Wrestling (TNA) consisting of
	members James Storm Abyss	James Storm Abyss The Great	members James Storm Abyss
	The Great Sanada, Khoya, Manik	Sanada, Khoya, Manik, and Ser-	The Great Sanada, Khoya, Manik
	and Serena Deeb.	ena Deeb as its members.	and Serena Deeb.
HotpotOA-	Jerrald King "Jerry" Goldsmith	Jerrald King "Jerry" Goldsmith	Jerrald King "Jerry" Goldsmith
394493	(February 10, 1929July 21, 2004)	(February 10, 1929 - July 21,	(February 10, 1929 - July 21,
	was an American composer and	2004) was an American com-	2004) was an American painter
	conductor most known for his	poser and conductor renowned	and sculptor, primarily known for
	work in film and television scor-	primarily for his contributions to	his abstract art exhibitions. He
	ing. He composed scores for such	film and television music. He	produced artworks for such note-
	noteworthy films as "", "The Sand	crafted the scores for notable	worthy galleries as the MoMA,
	Pebbles", "Logan's Run", "Planet	films including "The Sand Peb-	the Louvre, the Tate Modern,
	of the Apes", "Patton", "Papillon",	bles", "Logan's Run", "Planet of	Guggenheim, the Art Institute of
	"Chinatown", "The Wind and the	the Apes", "Patton", "Papillon",	Chicago, the National Gallery,
	Lion", "The Omen", "The Boys	"Chinatown", "The Wind and the	the Whitney Museum, the San
	from Brazil", "Capricorn One",	Lion", "The Omen", "The Boys	Francisco Museum of Modern
	"Alien", "Outland", "Poltergeist",	from Brazil", "Capricorn One",	Art, the Getty Center, the Mu-
	"The Secret of NIMH", "Grem-	"Alien", "Outland", "Poltergeist",	seum of Fine Arts, Boston, the
	lins", "Hoosiers", "Total Recall",	"The Secret of NIMH", "Grem-	Philadelphia Museum of Art, the
	"Basic Instinct", "Rudy", "Air	lins", "Hoosiers", "Total Recall",	Cleveland Museum of Art, the
	Force One", "L.A. Confidential",	"Basic Instinct", "Rudy", "Air	Detroit Institute of Arts, the
	"Mulan", "The Mummy", three	Force One", "L.A. Confidential",	Walker Art Center, the Houston
	"Rambo" films, "Explorers" and	"Mulan", "The Mummy", three	Museum of Fine Arts, the Dal-
	four other "Star Trek" films.	"Rambo" films, "Explorers", and	las Museum of Art, the Denver
		four "Star Trek" films.	Art Museum, the Seattle Art Mu-
			seum, the Miami Art Museum,
			the Barnes Foundation, the Kim-
			bell Art Museum, the Nelson-
			Atkins Museum of Art, and con-
			tributed to exhibits in the Smithso-
			nian and several other prestigious
			international venues.
HotpotQA-	"Livin' Our Love Song" is a song	"Livin' Our Love Song" was re-	"Livin' Our Love Song" is a track
11792076	co-written and recorded by Amer-	leased by Jason Michael Carroll,	solely written and performed by
	ican country music artist Jason	an American country music artist,	British pop artist Emma Louise
	Michael Carroll. It was released	as the second single from his al-	Clark. It was released in Octo-
	in April 2007 as the second sin-	bum "Waitin' in the Country" in	ber 2010 as the final single from
	gle from his album "Waitin' in	April 2007. The song, which Car-	her album "Cityscape Dreams".
	the Country". Carroll co-wrote	roll co-wrote with Glen Mitchell	Clark independently created the
	the song with Glen Mitchell and	and 11m Galloway, showcases his	song without collaborations.
	11m Galloway.	contribution both as a writer and	
		performer.	

Contradiction Retrieval via Contrastive Learning with Sparsity

Passage ID	Orginal	Paranhrase	Contradiction
HotpotOA-	Malik Izaak Taylor (November	Malik Izaak Taylor (November	Malik Izaak Taylor (November
1472458	20 1970March 22 2016) known	20, 1970 - March 22, 2016) rec-	20 1970March 22 2016) known
1472450	professionally as Phife Dawg (or	ognized by his stage name Phife	professionally as Phife Dawa
	simply Phife) was an American	Dawa (or just Phife) was an	(or simply Phife) was a British
	simply Fine), was an American	American ranner and balanged	(of simply Fille), was a Bluish
	A Triba Callad Quast with high	American Tapper and Defonged	singer and a memorie Veises
	A The Caned Quest with high	Out all group A Tribe Carled	pop group The Harmonic voices
	school Iriends Q-Tip and All Sha-	Quest alongside his high school	Truesd and Nach Darloss (and
	need Munammad (and for a short	Companions Q-Tip, All Snaneed	Iweed and Noan Parker (and
	time Jarobi white). He was also	Munammad, and briefly Jarobi	brieffy Emily Stone). He was also
	known as the Five Foot Assas-	White. He was also nicknamed	known as the "fall Lyricist" and
	sin" and "The Five Footer", be-	the "Five Foot Assassin" and	"The Towering Tenor", because
	cause he stood at 5 ft.	"The Five Footer" due to his	he stood at 6 ft 4 in.
		height of 5 feet.	
HotpotQA-	The Fourth Dimension is a non-	"The Fourth Dimension," a non-	The Fifth Dimension is a fictional
3500070	fiction work written by Rudy	fiction book by Rudy Rucker, a	novel authored by Rudy Rucker,
	Rucker, the Silicon Valley profes-	mathematics and computer sci-	the Silicon Valley professor of
	sor of mathematics and computer	ence professor from Silicon Val-	literature and art, and was re-
	science, and was published in	ley, was released in 1984 by	leased in 1990 by Penguin Ran-
	1984 by Houghton Mifflin. The	Houghton Mifflin. Subtitled "a	dom House. The book is alterna-
	book is subtitled as a guided tour	guided tour of the higher uni-	tively titled as a narrative journey
	of the higher universes. The fore-	verses," it includes a foreword by	through imaginary realms. The
	word included is by Martin Gard-	Martin Gardner and features over	afterword provided is by Douglas
	ner, and the 200+ illustrations are	200 illustrations by David Povi-	Hofstadter, and the 150+ paint-
	by David Povilaitis. Like other	laitis. As with his other works,	ings are by Sarah Jensen. Un-
	books by Rucker, "The Fourth	Rucker dedicates this book to Ed-	like other works by Rucker, "The
	Dimension" is dedicated to Ed-	win Abbott Abbott, who wrote	Fifth Dimension" pays homage to
	win Abbott Abbott, author of the	the novella "Flatland."	Lewis Carroll, author of the novel
	novella "Flatland".		"Through the Looking-Glass".
	1	1	

Contradiction Retrieval via Contrastive Learning with Sparsity

Desserve ID	Orginal	Daraphrasa	Contradiction
Fassage ID HotpotOA	Music of the Sup is the debut stu	Music of the Sup the ineugural	Music of the Sup is not the debut
1101p01QA-	dia album by Darbadian aingan	studie album by Darbadian artist	studie album of Dorbadian singer
2383880	Dihanna It was released on	Bihanna was launahad on Au	Bihanna: rathar it is har second
	August 20, 2005 in the United	Rinamia, was faunched off Au-	Allarma, latilet, it is net second
	States through Def Jern Decord	by Def Jam Decendings De	aut 20, 2006 outside the United
	States through Del Jam Record-	by Del Jam Recordings. Be-	gust 50, 2000, outside the United
	Ings. Prior to signing with Del	Difference association with Del Jam,	States and was not associated
	Jam, Kinanna was discovered by	kinanna was discovered in Bar-	with Del Jam Recordings. Before
	Parkadaa wika kalaad Dikama	bados by record producer Evan	her agreement with a different la-
	Barbados, who helped Rinanna	Rogers, who assisted her in cre-	bel, Rinanna was noticed by a
	record demo tapes to send out	ating demo tapes that were dis-	talent scout in Canada, unrelated
	to several record labels. Jay-Z,	tributed to various record labels.	to record producer Evan Rogers,
	the former chief executive offi-	Jay-Z, who was then the chief ex-	and she independently produced
	cer (CEO) and president of Der	ecutive officer (CEO) and pres-	ner initial demos. Jay-Z, who had
	Jam, was given Rinanna's demo	ident of Def Jam, received Ri-	aiready stepped down as the chief
	by Jay Brown, his A&R at Der	hanna's demo from Jay Brown,	executive officer (CEO) and pres-
	fan the lebel often beering subst	mis A&R at Del Jam, and after lis-	tand Dihama's dama subish sus
	for the label after hearing what	tening to what eventually became	tered Rinanna's demo, which was
	turned out to be ner first single,	ner debut single, Pon de Replay,	not passed by Jay Brown, his for-
	"Pon de Replay". She auditioned	he asked her to audition for the la-	mer A&R at Def Jam. Further-
	for Jay-Z and L.A. Reid, the for-	bel. During ner audition for Jay-Z	more, she did not perform an au-
	mer CEO and president of record	and L.A. Reid, the then CEO and	dition before the executives of the
	label group The Island Def Jam	president of The Island Def Jam	label, and as a result, she was
	Music Group, and was signed on	Music Group, she was immedi-	cautiously offered a contract days
	the spot to prevent her from sign-	ately signed to the label to avoid	after her negotiations with vari-
	ing with another record label.	her potentially signing elsewhere.	ous other record labels had con-
			cluded.
HotpotQA-	Barbara Bouchet (born Barbara	Barbara Bouchet, born as Barbara	Barbara Bouchet (born Barbara
5994297	Gutscher, 15 August 1944) is a	Gutscher on August 15, 1944, is	Gutscher, 15 August 1944) is a
	German-American actress and en-	a German-American actress and	German-American scientist and
	trepreneur who lives and works in	entrepreneur, currently residing	politician who resides and oper-
	Italy.	and working in Italy.	ates in the United States.
HotpotQA-	The 2009–10 Spanish football	The 2009–10 season marks	The 2009–10 Spanish football
26188218	season is Xerez's first season ever	Xerez's inaugural appearance in	season is not Xerez's first season
	in Liga BBVA.	Liga BBVA, which is Spain's top	in Liga BBVA.
		football division.	
HotpotQA-	Kim Coco Iwamoto is a commis-	Kim Coco Iwamoto served as	Kim Coco Iwamoto was not
20049848	sioner on the Hawaii Civil Rights	a commissioner on the Hawaii	a commissioner on the Hawaii
	Commission, appointed by Gov-	Civil Rights Commission, having	Civil Rights Commission, and
	ernor Neil Abercrombie to serve	been appointed by Governor Neil	she was not appointed by Gover-
	the four-year term from 2012 to	Abercrombie for a four-year term	nor Neil Abercrombie for a term
	2016.	spanning from 2012 to 2016.	from 2012 to 2016.

Contradiction Retrieval via Contrastive Learning with Sparsity

Passage ID	Orginal	Paraphrase	Contradiction
HotpotQA-	Oscar De La Hoya ( ; born Febru-	Oscar De La Hoya (born Febru-	Oscar De La Hoya ( ; born Febru-
95310	ary 4, 1973) is a former profes-	ary 4, 1973) is an ex-professional	ary 4, 1973) is a current profes-
	sional boxer who competed from	boxer who competed between	sional wrestler who began com-
	1992 to 2008. He holds dual	1992 and 2008. He has both	peting in 2010. He holds solely
	American and Mexican citizen-	American and Mexican citizen-	Mexican citizenship. Nicknamed
	ship. Nicknamed "The Golden	ship. Known as "The Golden	"The Silver Star," De La Hoya
	Boy," De La Hoya represented the	Boy," De La Hoya competed for	represented Mexico at the 2008
	Olympics, winning a gold model	Summer Olympics and secured	summer Olympics, losing in the
	in the lightweight division shortly	a gold medal in the lightweight	division just before aprolling at
	after graduating from James A	category soon after his gradua-	James A Garfield High School
	Garfield High School	tion from James A Garfield High	Junes A. Garneld High School.
		School.	
HotpotQA-	Donny Edward Hathaway (Oc-	Donny Edward Hathaway (Oc-	Donny Edward Hathaway (Octo-
526562	tober 1, 1945 – January 13,	tober 1, 1945 – January 13,	ber 1, 1945 – January 13, 1979)
	1979) was an American jazz,	1979) was an acclaimed Ameri-	was an American classical, pop,
	blues, soul and gospel singer,	can singer, songwriter, arranger,	rock and country musician, com-
	songwriter, arranger and pianist.	and pianist, known for his contri-	poser, conductor, and keyboardist.
	Hathaway signed with Atlantic	butions to jazz, blues, soul, and	Hathaway began his career with
	Records in 1969 and with his first	gospel music. Signing with At-	Columbia Records in 1975, and
	Ghetto" in early 1970 "Polling	released his debut single "The	Harmony label "A Quiet Storm"
	Stope" magazine "marked him as	Ghetto" under Atco label in early	in late 1976 "Rolling Stope" mag-
	a major new force in soul mu-	1970 which led "Rolling Stone"	azine dismissed him as a fleet-
	sic." His enduring songs include	magazine to recognize him as a	ing figure in the music world.
	"The Ghetto", "This Christmas",	significant new voice in soul mu-	His obscure songs include "A
	"Someday We'll All Be Free",	sic. Among his most memorable	Quiet Storm", "Winter Wonder-
	"Little Ghetto Boy", "I Love You	tracks are "The Ghetto", "This	land", "Always Free", "Big City
	More Than You'll Ever Know",	Christmas", "Someday We'll All	Boy", "I Hate You Less Than
	signature versions of "A Song for	Be Free", "Little Ghetto Boy", "I	You'll Ever Know", obscure ren-
	You" and "For All We Know",	Love You More Than You'll Ever	ditions of "A Tune for Me" and
	and "Where Is the Love" and	Know", definitive renditions of	"What we Don't Know", and
	many collaborations with Roberta	A Song for You and For All We Know" alongside "Where Is	Further I Co from You" a few
	Flack "Where Is the Love" won	the Love" and "The Closer I Get	isolated attempts to work with
	the Grammy Award for Best Pop	to You", results of frequent col-	Roberta Flack. "Where Is the
	Performance by a Duo or Group	laborations with Roberta Flack.	Hate" lost the Grammy Award
	with Vocals in 1973. At the height	"Where Is the Love" earned him a	for Worst Pop Performance by
	of his career Hathaway was diag-	Grammy Award for Best Pop Per-	a Duo or Group with Vocals in
	nosed with paranoid schizophre-	formance by a Duo or Group with	1973. At a low point in his ca-
	nia and was known to not take his	Vocals in 1973. During his peak,	reer, Hathaway was diagnosed
	prescribed medication regularly	Hathaway suffered from paranoid	with acute stress response and
	enough to properly control his	schizophrenia, struggling with in-	was known for his strict adher-
	symptoms. On January 13, 1979,	consistent use of his prescribed	ence to prescribed medication, ef-
	ride the luvury botel Eccar House	tion effectively. On January 12	On January 12, 1070 Hathaway's
	in New York City: his death was	1979 Hathaway died by suicide	body was discovered inside the
	ruled a suicide	with his body discovered outside	modest Cortlandt Hotel in New
		New York City's luxury Essex	York City; his death was declared
		House hotel.	a natural cause.
	1	1	I

Contradiction Retrieval via Contrastive Learning with Sparsity

	0.1.1	B 1	
Passage ID	Orginal	Paraphrase	Contradiction
HotpotQA-	Randy Edelman (born June 10,	Randy Edelman, born on June 10,	Randy Edelman (born June 10,
1834726	1947) is an American musician,	1947, is an American composer,	1947) is an American musician,
	producer, and composer for film	producer, and musician famed for	producer, and composer for film
	and television known for his work	his contributions to film and tele-	and television known for his work
	in comedy films. He has been	vision comedies. He has received	in dramatic films. He has never
	nominated for a Golden Globe	nominations for a Golden Globe	been nominated for a Golden
	Award, a BAFTA Award, and	Award and a BAFTA Award, and	Globe Award, a BAFTA Award,
	is the recipient of twelve BMI	has won twelve BMI Awards.	and has not received any BMI
	Awards.		Awards.