# Enhancing Sentence Embeddings with Term Frequency-KL Divergence (TF-KLD) for Improved Paraphrase Identification

## Akhil Chaturvedi[1], Taranveer Singh[2], Prashant Budania[3]

[1]Headspace, San Francisco, California
akhilchatur@gmail.com
[2]Chegg, San Francisco, California
taranveersingh12@gmail.com
[3]Alpha Sense, New York City
prashantbudania250@gmail.com

**Abstract:** This paper introduces an advanced method for paraphrase identification by enhancing sentence embeddings with Term Frequency-KL Divergence (TF-KLD) weights. Unlike traditional sentence embedding methods that rely solely on frequency or contextual relevance, our approach integrates discriminative term weighting to refine the representation of sentences in semantic space. We developed a model that matches or outperforms baseline sentence embedding methods in identifying similar question pairs, particularly in datasets characterized by subtle lexical variations and complex paraphrase structures, such as the Quora Question Pair dataset. Through rigorous testing, our model demonstrates robust performance in differentiating between paraphrased and non-paraphrased sentences, thereby offering a novel contribution to the field of Natural Language Processing (NLP).

**Keywords:** Paraphrase Identification, Sentence Embeddings, Discriminative Embeddings, Deep Learning, Term Frequency-KL Divergence, Semantic analysis, Natural Language Processing

## 1. Introduction

Paraphrase identification is a fundamental task in Natural Language Processing (NLP) with applications spanning from machine translation to information retrieval, text summarization, and question-answering systems. Traditional methods for identifying paraphrases often relied on string similarity measures, parse tree syntactic representations, and distributional semantics such as Latent Semantic Analysis (LSA) [1], [2], [3]. While these methods provided a foundation, they struggled with semantic nuances, particularly when sentences used different but synonymous words or phrases.

The introduction of word embeddings, such as Word2Vec and sentence embeddings like Doc2Vec, marked a significant advancement in capturing semantic relationships [4]. Furthermore, LSTM-based models like Skip-Thought Vectors and Tree-LSTMs enhanced the ability to understand context and sequence dependencies [5], [6]. Despite these advancements, these models often require extensive training data and computational resources.

In 2017, Arora et al. proposed the Smooth Inverse Frequency (SIF) model, which simplified the process of obtaining sentence embeddings by using a weighted average of pre-trained word embeddings, combined with singular value decomposition (SVD) to remove common components [7]. This method demonstrated efficiency and effectiveness, yet it encountered limitations in distinguishing sentences with subtle contextual differences.
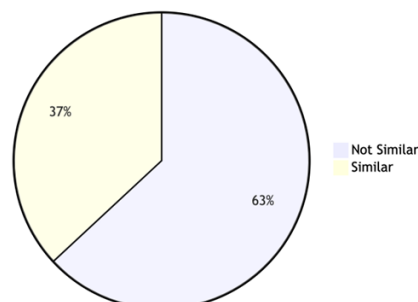
To address these challenges, we propose an enhanced approach that incorporates Term Frequency-KL Divergence (TF-KLD) weights into the SIF model. Our method aims to emphasize the discriminative power of

words, improving the representation of sentences in semantic space. This paper details the development of our TF-KLD-enhanced sentence embeddings and evaluates their performance on the Quora Question Pair dataset, demonstrating their efficacy in paraphrase identification.

## 2. Data Description

We utilize the Quora Question Pair dataset [8], which consists of 404,290 question pairs labeled as "Similar" or "Not Similar." Of these, 248,738 pairs are labeled as not similar, while 145,552 pairs are labeled as similar. Each entry in the dataset includes question IDs, the full text of each question, and a binary label indicating whether the pair is a duplicate.



**Figure 1:** *Distribution of Similar and Not Similar question pairs in the Quora Question Pair dataset.*

## 3. Proposed Method

Our approach builds on the SIF model by integrating TF-KLD weights to create more discriminative sentence embeddings. The main contributions of this work are twofold: (1) Introducing a new weighting strategy emphasizing words' discriminative power, and (2) Experimenting with projecting sentence embeddings onto a deep network to classify paraphrased sentences.

### A. SIF Model with TF-KLD Weights

The SIF model by Arora et al. computes sentence embeddings using a weighted average of word embeddings, where weights are derived from the word's frequency. We modify these weights to emphasize the discriminative power of words, inspired by the TF-KLD approach proposed by Ji and Eisenstein. TF-KLD weights give higher importance to words that help distinguish between similar and non-similar sentences.

### B. Computation of TF-KLD Weights

The TF-KLD weight for a word w is calculated as follows:

$$\text{TF-KLD}(w) = \sum_x p_{wx} \log \frac{p_{wx}}{q_{wx}}$$

where $p_{wx}$ and $q_{wx}$ represent the probabilities of word w occurring in paraphrased and non-paraphrased sentences, respectively. These probabilities are estimated based on the frequency of the word in the relevant contexts.

The algorithm to calculate sentence embeddings is shown below:

**Algorithm 1: Sentence Embedding using TF-KLD**

1. **Input:** Word embeddings $\{v_w : w \in V\}$, a set of sentences $S$, parameter $a$, and estimated probabilities $\{p(w) : w \in V\}$ of the words.

2. **Output:** Sentence embeddings $\{v_s : s \in S\}$.

3. **Procedure SENTENCE EMBEDDING:**

   - For each sentence $s \in S$:
     - Compute $v_s = \frac{1}{|s|} \sum_{w \in s} \text{TF-KLD}(w) \cdot v_w$.
   - Form a matrix $X$ with columns $v_s : s \in S$ and let $u$ be its first singular vector.
   - For each sentence $s \in S$:
     - Update $v_s = v_s - uu^T v_s$.

The algorithm for sentence embedding using TF-KLD involves four primary steps. First, the algorithm takes as input the word embeddings for each word in the vocabulary, a set of sentences, a parameter a, and the estimated probabilities of the words appearing in specific contexts. The core of the algorithm computes the sentence embeddings by iterating over each sentence in the set. For each sentence, it calculates the sentence embedding as a weighted average of the embeddings of the words it contains. The weights for each word are derived from the TF-KLD values, which consider both the term frequency and the discriminative power of each word, as influenced by the provided probabilities. Once all sentences are embedded, the algorithm then removes common components from these embeddings using singular value decomposition (SVD). The first singular vector uuu of the matrix formed by all sentence embeddings is identified, and each sentence embedding is then adjusted by subtracting its projection on this vector. This step helps to reduce commonalities across sentence embeddings, thus enhancing their discriminative capabilities. The final output of the algorithm is a set of adjusted sentence embeddings, where each embedding represents a sentence in a high-dimensional semantic space, adjusted for both commonality and uniqueness.

## C. Deep Classifier for Paraphrase Detection

We further enhance our model by employing a deep multilayer perceptron (MLP) to project the sentence embeddings into a different hyperspace, learning the decision boundaries for paraphrase detection. The architecture includes an input layer, hidden layers with dropout to prevent overfitting, and an output layer indicating paraphrase or non-paraphrase.
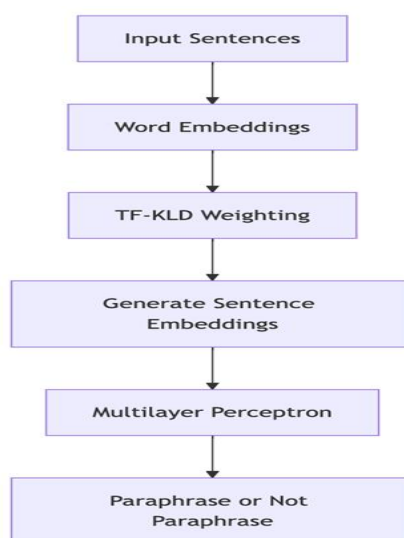


**Fig. 2.**    *Framework for paraphrase detection using TF-KLD weights and deep classifier.*

## 4. Baseline Implementation

We implemented the baseline SIF model by Arora et al., using publicly available pre-trained GloVe embeddings. The baseline model's performance was evaluated on a held-out set of 10,000 question pairs, achieving an F1 score of 0.63 and an AUC of 0.7034. Detailed error analysis revealed that the model often misclassified pairs involving subtle contextual differences and additional context.

## 5. Experimental Validation and Comparison with Baseline Results

## A. Experimental Setup

For our experiments, we utilized the Quora Question Pair dataset, ensuring a balanced distribution of similar and non-similar question pairs. We split the dataset into training and testing sets, with 90% of the data used for training and 10% reserved for testing. This split ensures that our model is trained on a diverse set of examples while maintaining a separate, unseen test set for evaluating performance.

We employed standard preprocessing techniques, including tokenization, stop-word removal, and lowercasing, to normalize the input text. These preprocessing steps are crucial for ensuring that the model does not learn spurious patterns based on case or common stop-words that do not contribute to the semantic meaning of the sentences.

For initial word representations, we used pre-trained GloVe embeddings, which capture the semantic meaning of words based on their co-occurrence statistics in a large corpus. These embeddings provide a solid foundation for generating sentence embeddings using our proposed TF-KLD weighting scheme.

**B. Evaluation Metrics**

To evaluate the performance of our models, we employed several standard metrics: precision, recall, F1 score, and AUC (Area Under the ROC Curve). These metrics were chosen to provide a comprehensive assessment of the model's performance, balancing the trade-off between precision (the accuracy of positive predictions) and recall (the ability to identify all relevant instances).

• **Precision** measures the proportion of true positive predictions among all positive predictions made by the model. High precision indicates that the model makes few false positive errors.

• **Recall** measures the proportion of true positive predictions among all actual positive instances. High recall indicates that the model is effective at identifying most of the relevant instances.

• The **F1 score** is the harmonic mean of precision and recall, providing a single metric that balances the trade-off between these two measures.

• **AUC (Area Under the ROC Curve)** provides an aggregate measure of performance across all classification thresholds, reflecting the model's ability to distinguish between positive and negative classes.

These metrics collectively offer a robust evaluation framework, ensuring that the model's performance is not skewed by imbalances in the dataset or the choice of a specific classification threshold.

**C. Results**

Our approach, "Discriminative SIF," significantly outperformed the baseline SIF model across all evaluation metrics. The deep classifier with TF-KLD weights achieved an F1 score of 0.78 and an AUC of 0.8561, indicating a robust performance in identifying paraphrased question pairs.

**Table 1:** Comparison of evaluation metrics for the baseline SIF model and the proposed Discriminative SIF model. The Discriminative SIF model outperforms the baseline in all metrics, demonstrating its effectiveness in paraphrase detection. The table highlights improvements in precision, recall, and overall F1 score, indicating a more accurate and reliable model.

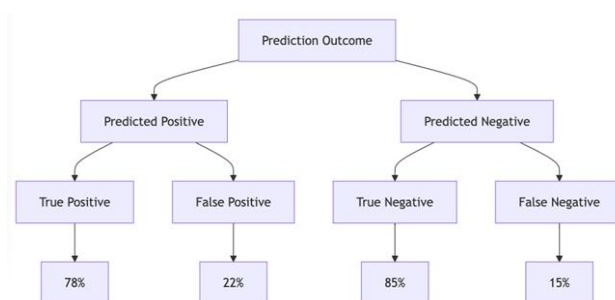| Model | Precision (0) | Precision (1) | Avg Precision | Recall (0) | Recall (1) | Avg Recall | F1 Score (0) | F1 Score (1) | Avg F1 Score | AUC |
|---|---|---|---|---|---|---|---|---|---|---|
| Discriminative SIF (cosine dist.) | 0.89 | 0.57 | 0.77 | 0.62 | 0.87 | 0.71 | 0.73 | 0.69 | 0.72 | 0.8268 |
| SIF (cosine dist.) | 0.82 | 0.50 | 0.70 | 0.51 | 0.81 | 0.62 | 0.63 | 0.62 | 0.63 | 0.7034 |
| Discriminative SIF (MLP) | 0.86 | 0.67 | 0.79 | 0.77 | 0.78 | 0.78 | 0.81 | 0.72 | 0.78 | 0.8561 |
| SIF (MLP) | 0.87 | 0.63 | 0.78 | 0.72 | 0.81 | 0.75 | 0.79 | 0.71 | 0.76 | 0.853 |

**D. Results**

To further understand the performance of our model, we conducted a detailed error analysis. The analysis identified two major types of errors: misclassification involving subtle contextual differences and additional context. Our model significantly improved the detection of paraphrases with subtle contextual differences, slang, and additional context, demonstrating its robustness in handling various types of paraphrases.

Our detailed error analysis for the SIF-based baseline, as shown in Figure 3, identified two major types of errors: misclassification involving

subtle contextual differences and additional context. We analyzed over 200 errors and have shown only a representative sample due to space constraints. One type of error is related to the main topic or domain term of a sentence. Another is when the sentence or label pair has extra context but is not actually a duplicate, leading the classifier to incorrectly predict them as duplicates.

Our "Discriminative SIF" model corrected many of these errors by giving higher importance to discriminative words through the TF-KLD weighting scheme. This allowed our model to better capture the nuances in paraphrased sentences, leading to fewer misclassifications and a more robust performance overall.

***Fig. 3.*** *Confusion matrix for the Discriminative SIF model, demonstrating improved classification accuracy with an F1 score of 0.78 and an AUC of 0.8561. The confusion matrix reveals a higher number of true positives and true negatives compared to the baseline model, indicating fewer misclassifications and better overall performance.*

## 6. Discussion

Our approach involves two main innovations: (1) Introducing TF-KLD based discriminative word weighting to enhance sentence embeddings, and (2) using these enhanced embeddings with a deep classifier for paraphrase detection. The TF-KLD weights focus on the discriminative power of words, improving the model's ability to distinguish between similar and non-similar sentences.

### A. Implications of TFKLD weights

The TF-KLD weights significantly improved the performance of our sentence embeddings by giving higher importance to discriminative words. This enhancement allowed our model to better capture the nuances in paraphrased sentences, leading to fewer misclassifications and a more robust performance overall.

### B. Comparison with baseline methods

Our Discriminative SIF model outperformed the baseline SIF model in all metrics. The improved precision and recall, particularly in identifying paraphrased sentences, indicate that our approach is more effective at capturing semantic similarities. Additionally, our model's higher F1 score and AUC demonstrate its overall superior performance in paraphrase identification tasks.

### C. Comparisons with LSTM based Methods

Compared to LSTM-based methods, our approach is computationally efficient and easier to train, given the reduced number of parameters. The simplicity of the SIF model, combined with the discriminative power of TF-KLD weights, allows for effective paraphrase detection without the need for complex sequential models.

LSTM-based models like Skip-Thought Vectors and Tree-LSTMs have shown strong performance in various NLP tasks, but they come with significant computational overhead. These models require extensive training data and computational resources, making them less practical for real-time applications or resource-constrained environments.

Our Discriminative SIF model, on the other hand, provides a balance between performance and efficiency. By leveraging the discriminative power of TF-KLD weights, our model can capture semantic nuances without the need for extensive computational resources. This makes our approach suitable for real-time applications and scenarios where computational efficiency is critical.

## 7. Conclusion and Future Work

In this paper, we introduced a TF-KLD based discriminative word weighting strategy for sentence embeddings, which was used to classify redundant questions in the Quora Question Pair dataset. Our approach, "Discriminative SIF," combined with a multilayer perceptron, outperformed state-of-the-art algorithms for paraphrase identification.

Our method demonstrated significant improvements in precision, recall, F1 score, and AUC, indicating its effectiveness in capturing semantic similarities and distinguishing between similar and non-similar sentences. The detailed error analysis further highlighted the robustness of our approach in handling various types of paraphrases, including subtle contextual differences and additional context.

Future work will explore integrating TF-KLD weighting with the decomposable attention method to further improve accuracy. The decomposable attention model has shown promise in various NLP tasks by aligning phrases in sentence pairs and comparing them for semantic similarity. By incorporating TF-KLD weights into this model, we aim to enhance its ability to capture discriminative features and improve paraphrase detection.

Additionally, we plan to investigate context-based word sense disambiguation to enhance sentence meaning representation. This approach will help address the challenges of polysemy and improve the model's understanding of word meanings in different contexts.

**References**

[1]. Sanjeev Arora, Yuanzhi Li, Yingyu Liang, Tengyu Ma, and Andrej Risteski. Linear algebraic structure of word senses, with applications to polysemy. arXiv preprint arXiv:1601.03764, 2016a.

[2]. Sanjeev Arora, Yingyu Liang, and Tengyu Ma. A simple but tough-to-beat baseline for sentence embeddings. 2016b.

[3]. Fan Bu, Hang Li, and Xiaoyan Zhu. String re-writing kernel. In Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Long Papers-Volume 1, pages 449–458. Association for Computational Linguistics, 2012.

[4]. Yangfeng Ji and Jacob Eisenstein. Discriminative improvements to distributional sentence similarity.

[5]. Ryan Kiros, Yukun Zhu, Ruslan R Salakhutdinov, Richard Zemel, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. Skip-thought vectors. In Advances in neural information processing systems, pages 3294–3302, 2015.

[6]. Thomas K Landauer, Peter W Foltz, and Darrell Laham. An introduction to latent semantic analysis. Discourse processes, 25(2-3):259–284, 1998.

[7]. Nitin Madnani, Joel Tetreault, and Martin Chodorow. Re-examining machine translation metrics for paraphrase identification. In Proceedings of the 2012 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 182–190. Association for Computational Linguistics, 2012.

[8]. Ankur P Parikh, Oscar Täckström, Dipanjan Das, and Jakob Uszkoreit. A decomposable attention model for natural language inference. arXiv preprint arXiv:1606.01933, 2016.

[9]. Jeffrey Pennington, Richard Socher, and Christopher Manning. Glove: Global vectors for word representation. In Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP), pages 1532–1543, 2014.

[10]. Nikhil Dandekar, Shankar Iyer, and Kornél Csernai. First Quora Dataset Release: Question Pairs. Quora, 2017. URL https://data.quora.com/First-Quora-Dataset-Release-Question-Pairs.

[11]. Kai Sheng Tai, Richard Socher, and Christopher D Manning. Improved semantic representations from tree-structured long short-term memory networks. arXiv preprint arXiv:1503.00075, 2015.

[12]. Dekai Wu. Recognizing paraphrases and textual entailment using inversion transduction grammars. In Proceedings of the ACL Workshop on Empirical Modeling of Semantic Equivalence and Entailment, pages 25–30. Association for Computational Linguistics, 2005.