

Unraveling Paraphrase Detection: A Deep Dive into Transformer Models using the MRPC dataset Before LLMs

Anonymous Author(s)

Abstract—Paraphrase identification remains a key Natural Language Processing (NLP) task, with the Microsoft Research Paraphrase Corpus (MRPC) dataset widely used for evaluating model performance. This paper presents an experiment to evaluate the effectiveness of transformer-based models, specifically BERT, RoBERTa, and DistilBERT, for paraphrase identification, without the usage of large language models (LLMs), which have become dominant in recent years. The paper focuses on performance with Inductive Conformal Prediction (ICP) and XCP (3-fold cross validation) methods, with an emphasis on evaluating Macro F1 scores and average set sizes for each model. Our findings show that RoBERTa-base outperforms BERT-base with a Macro F1 score of 0.8476 compared to 0.7366, indicating that transformer models before LLMs are still giving valuable outcomes for paraphrase detection. DistilBERT-base achieved a Macro F1 of 0.7871, demonstrating its effectiveness, despite being a lighter model. The study explores the challenges of such models, particularly in terms of capturing structural differences and lexical overlap in paraphrase identification tasks. We also touch on the importance of using pre-LLM transformers in the modern NLP landscape, highlighting the key points of the approach. This research aims to shed light on the performance of transformer models before the rise of LLMs, providing insights on how the models still have the potential to be beneficial for paraphrase identification without the advanced capabilities of present LLMs.

Index Terms—BERT, RoBERTa, and DistilBERT, transformers, Natural Language Processing, Microsoft Research Paraphrase Corpus, Inductive Conformal Prediction, XCP, Macro f1 score.

I. INTRODUCTION

The Microsoft Research Paraphrase Corpus (MRPC) has been a fundamental resource in the study of paraphrase identification (PI), where the goal is to determine whether two sentences convey the same meaning. While the task has been addressed before with various methods, transformer models such as BERT and RoBERTa have emerged in recent times as state-of-the-art models in Natural Language Processing (NLP) tasks with significantly improved performance. These advances have particularly been important ever since the transformer discovery in 2017 that uses the attention mechanism and pre-trained contextualized embeddings to enhance semantic relations between sentences [1], [2]. The MRPC dataset has long served as a standard test set for paraphrase identification, commonly viewed as a binary classification problem. However, in light of the rapid rise of large language models (LLMs) like GPT and BERT, a natural question arises: how did transformer models handle paraphrase identification

before the introduction of LLMs? This research aims to answer that question by focusing on the MRPC dataset and evaluating the performance of transformer-based models like BERT, RoBERTa and DistilBERT, which preceded LLMs. The goal is to explore how transformer-based models, which rely on contextual embeddings, and to understand the benefits and limitations of these approaches in handling the MRPC dataset [3]–[5].

In this paper, we apply BERT, RoBERTa, and DistilBERT to the MRPC dataset, following methodologies that involve fine-tuning these transformer models for paraphrase detection. By examining the results, we aim to understand the evolution of paraphrase identification methods and provide a platform for testing newer LLMs. The contributions of this work are to assess the performance of transformer models for MRPC before the advent of LLMs and provide ideas on how such models can be optimized for future research on paraphrase identification and semantic understanding.

II. RELATED WORK

Paraphrase detection has experienced tremendous growth in the development of benchmark sets and transformer models. Traditional sets such as MRPC have served as reference resources for sentence-pair classification [6] but have their binary labeling approaches being complemented in increasing manners by efforts such as ParaTag, which introduced fine-grained paraphrase annotations for identifying richer semantic variation [7]. Despite these advances, the complexity of using accurate annotations has increased model design and calibration complexity. Transformer models, e.g., Sentence-BERT, have significantly improved sentence-pair modeling by optimizing semantic similarity to the maximum via embedding-based techniques [8]. However, most models optimize for accuracy alone, often ignoring the need for calibrated uncertainty estimates essential in reliability-critical settings.

Recent research has highlighted the prevalence of pre-trained language models in paraphrase identification but at the same time remarked on limitations in dataset variety and strong resistance to intricate paraphrase structure [9]. Further, while conformal prediction methods have been proposed as a means to offer mathematical coverage guarantees as well as classification, their coupling with transformer backbones is relatively uninvestigated in lightweight architectures [10]. Existing approaches have proven to be feasible but are prone to relying on computationally expensive calibration steps or

prioritize confidence over optimizing for prediction performance jointly. Supporting work in data augmentation, such as backtranslation pipelines, has been successful in low-resource settings [11], but only adds data quantity and not necessarily structural confidence calibration as part of the model architectures themselves.

All these trends put the focus on the growing need for paraphrase detection tools with strong semantic modeling, lightweight confidence calibration, and rapid prediction set generation simultaneously without allowing performance to fall behind.

III. METHODOLOGY

A. Dataset description

The Microsoft Research Paraphrase Corpus (MRPC) dataset consists of 5,801 sentence pairs, each labeled by human annotators to determine whether the sentences convey the same meaning. The dataset was created by Dolan and Brockett (2005) using a combination of heuristic extraction techniques and a Support Vector Machine (SVM)-based classifier to identify likely paraphrases from a large corpus of topic-clustered news articles. Human raters then confirmed the semantic equivalence of the pairs, with 67% judged to be paraphrases [6]. The MRPC dataset has since become a benchmark for paraphrase detection, widely used to evaluate transformer models like BERT and its variants. Its broad applicability and the inclusion of diverse sentence pairs make it a critical resource for research on paraphrase identification and semantic similarity [12].

B. Justification of Methodology

In this study, we focus on evaluating the performance of transformer-based models and Conformal Prediction techniques for paraphrase identification using the acquired dataset. The methodologies and models chosen for this study are selected based on their effectiveness in previous research and their suitability for addressing the challenges of paraphrase detection. BERT and RoBERTa, both pre-trained transformer-based models, are selected for their demonstrated success in various Natural Language Processing (NLP) tasks, including paraphrase detection. BERT [13] uses a bidirectional attention mechanism such that the model can utilize the left and right context during training, providing rich, contextualized word representations. This two-way architecture is critical to paraphrase identification since it enables the model to understand word meaning in context. RoBERTa [14], which is a fine-tuned variant of BERT, improves performance through training on more data, longer sequences, and without Next Sentence Prediction task, leading to more generalization. These models have consistently provided cutting-edge performance on a variety of benchmarking tasks and are highly relevant for tasks like paraphrase identification since they understand sentence-level semantics in depth. For computational efficiency, DistilBERT, a distilled version of BERT, is incorporated into the study [15]. DistilBERT retains 97% of BERT’s performance while being

60% faster and requiring 40% fewer parameters. This trade-off between efficiency and performance makes DistilBERT an ideal model for testing in environments where computational resources are constrained, yet accurate paraphrase detection is still required. To quantify the uncertainty of predictions made by these models, we employ Conformal Prediction techniques. Conformal Prediction [16], [17] is a non-parametric, distribution-free technique that produces valid prediction sets with statistical assurances even if sample sizes are finite. It is especially beneficial for NLP tasks like paraphrase identification, where it is desirable to know the uncertainty of the model predictions. Inductive Conformal Prediction (ICP) and XCP (3-fold cross-validation) are applied to calibrate models and generate sets of predictions including the true label at a specified confidence level [18], [19].

While ICP allows for efficient calibration with a single split of the dataset, XCP introduces cross-validation to ensure more robust performance by averaging results across multiple splits of the dataset. This helps minimize the dangers of overfitting and gives a better estimate of the generalization performance of the model. The fine-tuning process is a key component of this methodology, as it enables pre-trained models to adapt to the specific task of paraphrasing identification. Fine-tuning leverages the pre-trained knowledge in models like BERT and RoBERTa and refines them for the target task using task-specific labeled data [13]. This facilitates better performance on downstream tasks like paraphrase detection as the models capture task-specific features and maintain the pre-learned information. The use of Huggingface Trainer streamlines the fine-tuning process, providing a framework that ensures the models are fine-tuned effectively on the MRPC dataset. Recent advancements in efficient fine-tuning methods, such as BitFit, which only fine-tunes the bias terms and task heads, have shown the importance of reducing computational requirements while maintaining high performance [20]. This method has shown how it performs better than full fine-tuning in scenarios with sparse training data, offering a computationally inexpensive solution. Such methods are particularly useful in resource-constrained environments, where adapting large models takes time and may be computationally expensive. Overall, the application of transformer models like BERT, RoBERTa and DistilBERT coupled with Conformal Prediction techniques gives this study not just the ability to capture paraphrase identification semantic complexities but also captures the uncertainty in model predictions.

C. Proposed approach

Models used in this study include BERT-base, RoBERTa-base, and DistilBERT, with each being subjected to different methodologies in performance evaluation. Tokenization is the first step taken, whereby both RoBERTa and DistilBERT tokenizers are utilized in prepping the dataset for model input. Second, models are fine-tuned on MRPC data, a crucial step that allows pre-trained transformer models to be tuned for paraphrase detection task. After fine-tuning, there is prediction performed on the test/validation set wherein models output

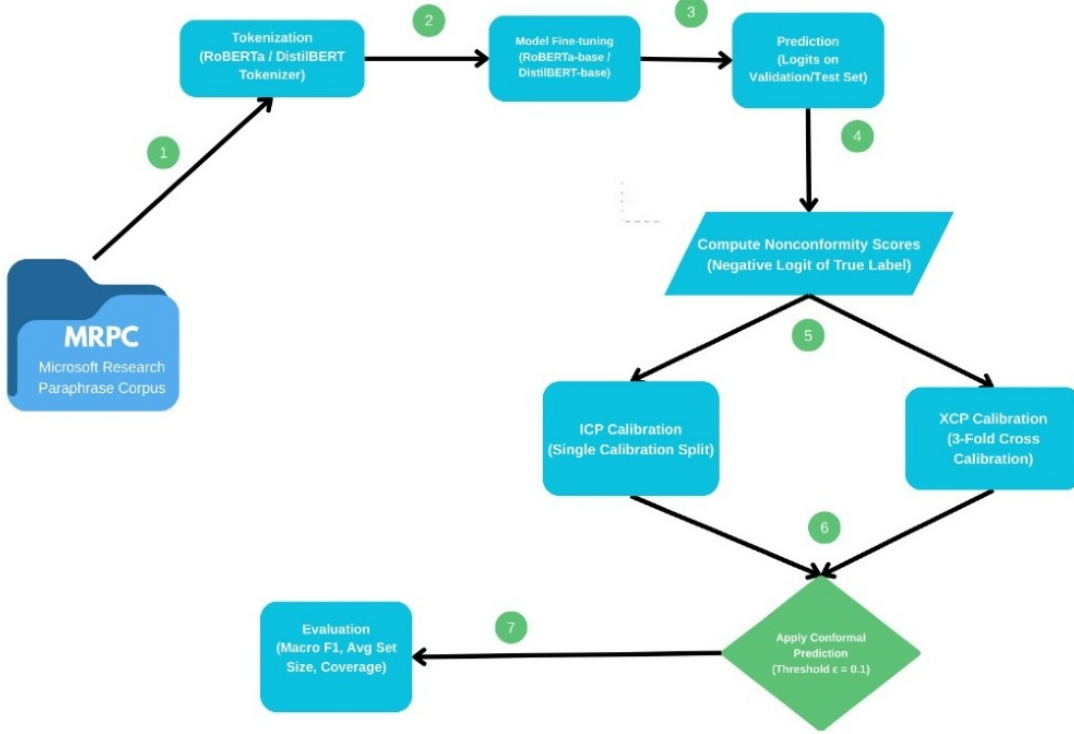


Fig. 1. Architecture of the proposed model

logits representing a probability of sentences as paraphrases. These logits are subsequently used to compute the nonconformity scores which involve computing negative log of actual label. These scores indicate the model's confidence in its predictions and are used as a basis for uncertainty estimation. To assess the models' uncertainty, we employ Inductive Conformal Prediction (ICP) and XCP (3-fold cross-validation) methods. ICP is used for efficient calibration on a single split of calibration, while XCP introduces cross-validation, providing a more robust estimation by averaging over three splits. This allows more precise performance estimation and avoids overfitting, particularly beneficial in difficult tasks like paraphrase detection. Finally, the conformal prediction approach is applied with a specified threshold parameter value ($\epsilon = 0.1$) for making prediction sets, which are the confidence levels of the predictions made by the models. The performance of the models is further evaluated using key parameters such as Macro F1, Average Set Size, and Coverage, giving a full overview of their ability to identify paraphrases and non-paraphrases and their compliance to varying levels of confidence of their predictions. For DistilBERT, we simplify the process by fine-tuning the model using the Huggingface Trainer, which is a more computationally efficient approach. This method evaluates DistilBERT using traditional performance metrics like accuracy, precision, recall, F1, and Macro F1, allowing for a comparison with the more powerful models, BERT and RoBERTa. The goal is to understand the trade-offs between computational efficiency and model performance, examining how well DistilBERT performs in comparison to

its larger counterparts, while considering practical applications where computational resources are limited.

D. Evaluation metrics

When evaluating models for paraphrase identification, several key metrics are commonly used to measure the performance. The ones used here are Accuracy, Precision, Recall, Macro F1-score and Average Prediction Set Size.

Accuracy quantifies the overall correctness of the model's predictions and is defined as:

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN} \quad (1)$$

Precision indicates the proportion of predicted paraphrase pairs that are actually paraphrases:

$$Precision = \frac{TP}{TP + FP} \quad (2)$$

Recall (also known as Sensitivity) measures the proportion of actual paraphrase pairs that were correctly identified by the model:

$$Recall = \frac{TP}{TP + FN} \quad (3)$$

Macro F1 Score calculates the average F1-Score across all classes, treating each class equally regardless of its frequency:

$$MacroF1 = \frac{1}{N} \sum_{i=1}^N F1_i \quad (4)$$

where N is the number of classes and F1_i is the F1-Score of class i.

Average Prediction Set Size, used in conformal prediction frameworks, measures the mean number of labels assigned per instance:

$$AveragePredictionSetSize = \frac{1}{M} \sum_{j=1}^M \Gamma(x_j) \quad (5)$$

where M is the total number of test instances and $T(x_j)$ denotes the prediction set generated for instance x_j .

These metrics are essential for evaluating the model’s effectiveness in detecting semantic similarity between sentence pairs and are particularly important in maintaining the balance between false positives and false negatives in paraphrase identification tasks [21]–[25].

IV. RESULTS

Table 1 presents the baseline performance of three transformer based models, BERT-base, DistilBERT-base, and RoBERTa-base evaluated on the dataset using standard classification metrics: Accuracy, Precision, and Recall. RoBERTa-base outperformed all other models, achieving the highest accuracy of 86.76% and a precision of 85.50%, indicating robust overall performance. DistilBERT-base also attained a high recall value of 94.98%, which reflects sensitivity to paraphrased sentence pairs despite being a light-weight model. BERT-base recorded comparatively lower performance for all the metrics with 79.66% accuracy and 72.83% recall.

TABLE I
RESULTS OF MODELS

Model	Accuracy	Precision	Recall
BERT-base	0.7966	0.7784	0.7283
DistilBERT-base	0.8333	0.8307	0.9498
RoBERTa-base	0.8676	0.8550	0.8324

Table 2 demonstrates Macro F1 Score and Average Prediction Set Size for two conformal prediction algorithms: Inductive Conformal Prediction (ICP) and Cross Conformal Prediction (XCP, with 3-fold calibration). RoBERTa-base again achieved the best Macro F1 score of 0.8476 under both ICP and XCP settings. Notably, XCP reduced the Average Prediction Set Size from 1.96 to 1.91, indicating tighter calibration without loss in performance. Furthermore, DistilBERT-base also performed competitively, with a Macro F1 of 0.7871 and an average prediction set size of 1.90 using ICP. BERT-base showed the lowest Macro F1 score of 0.7366 across both methods, with a slight improvement in prediction set compactness under XCP. Overall, RoBERTa-base performed better than BERT-base and DistilBERT-base on all the measurement criteria across the board. The findings affirm that conformal prediction methods, particularly XCP, can considerably reduce uncertainty (smaller prediction sets) without reducing predictive accuracy.

TABLE II
RESULTS OF MODELS (MACRO F1 SCORE AND AVG PREDICTION SET SIZE)

Model	Method	Macro F1	Avg Prediction Set Size
BERT-base	ICP	0.7366	1.93
BERT-base	XCP (3 folds)	0.7366	1.91
DistilBERT-base	ICP	0.7871	1.90
RoBERTa-base	ICP	0.8476	1.96
RoBERTa-base	XCP (3 folds)	0.8476	1.91

V. CONCLUSION

This study investigated the application of transformer-based models, such as BERT, RoBERTa, and DistilBERT, on the Microsoft Research Paraphrase Corpus (MRPC) dataset. Our findings demonstrate that these models, built on contextualized embeddings and transfer learning, significantly outperform other paraphrase detection approaches in identifying semantic relations among sentences. The use of transformers has proven to be a major advancement over traditional models, offering more precise and reliable paraphrase identification. But while the outcomes are encouraging, transformer models also continue to have difficulty with partial paraphrases, where sentences share partial semantic overlap. This indicates the necessity of further work in capturing more nuanced semantic relations. Overall, our work points to the importance of transformers in enhancing paraphrase detection while acknowledging that there are still areas for improvement, especially in handling more complex paraphrasing scenarios.

VI. LIMITATIONS

While the suggested framework achieves success in integrating lightweight conformal prediction with fine-tuned transformer models for paraphrase detection, it has some limitations. For one, the evaluation was conducted primarily on the MRPC dataset, which, while being heavily utilized, contains limited paraphrase diversity and binary coarse labels. Thus, further generalization to finer-grained or more complex paraphrase corpora remains to be established. Furthermore, although the backbone models used were diverse, they did not include large-scale generative language models, whose inclusion can yield further performance gains at the cost of increased computational requirements.

VII. FUTURE WORK

Future research should aim to enhance transformer models’ sensitivity to partial paraphrases by fine-tuning on datasets with varying levels of semantic similarity, possibly using adversarial methods. Developing more robust evaluation metrics will also be essential for assessing these models in diverse paraphrase detection tasks. As the transformer models improve, the addition of larger LLMs like GPT and T5 will help improve the accuracy of the paraphrase identification. Additionally, exploring scalability to large datasets and multilingual performance will be key to broadening the applicability of paraphrase detection.

REFERENCES

- [1] M. E. Peters, M. Neumann, M. Iyyer, M. Gardner, C. Clark, K. Lee, and L. Zettlemoyer, “Deep contextualized word representations,” *arXiv*, 2018.
- [2] L. Maltoudoglou, A. Paisios, and H. Papadopoulos, “Bert-based conformal predictor for sentiment analysis,” in *Proceedings of the Ninth Symposium on Conformal and Probabilistic Prediction and Applications*, vol. 128, pp. 269–284, 2020.
- [3] N. Shi, B. Hauer, J. Riley, and G. Kondrak, “Paraphrase identification through textual inference,” in *Proceedings of the 13th Joint Conference on Lexical and Computational Semantics (SEM 2024)*, pp. 133–141, Association for Computational Linguistics, 2024.
- [4] Q. Peng, D. Weir, and J. Weeds, “Testing paraphrase models on recognizing sentence pairs at different degrees of semantic overlap,” in *Proceedings of the 12th Joint Conference on Lexical and Computational Semantics (SEM 2023)*, pp. 259–269, Association for Computational Linguistics, 2023.
- [5] Q. Peng, D. Weir, and J. Weeds, “Towards structure-aware paraphrase identification with phrase alignment using sentence encoders,” in *Proceedings of the 29th International Conference on Computational Linguistics*, pp. 4113–4123, 2022.
- [6] W. B. Dolan and C. Brockett, “Automatically constructing a corpus of sentential paraphrases,” in *Proceedings of the 3rd International Workshop on Paraphrasing*, pp. 1–8, Microsoft Research, 2005.
- [7] S. Wang, R. Xu, Y. Liu, C. Zhu, and M. Zeng, “Fine-grained labels, nlg evaluation, and data augmentation,” in *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 7111–7122, Association for Computational Linguistics, 2022.
- [8] N. Reimers and I. Gurevych, “Sentence-bert: Sentence embeddings using siamese bert-networks,” in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing (EMNLP 2019)*, pp. 3982–3992, 2019.
- [9] C. Zhou, C. Qiu, L. Liang, and D. E. Acuna, “Paraphrase identification with deep learning: A review of datasets and methods,” in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP 2020)*, pp. 1–10, 2020.
- [10] P. Giovannotti and A. Gammerman, “Transformer-based conformal predictors for paraphrase detection,” in *Proceedings of Machine Learning Research*, vol. 152, pp. 1–23, 2021.
- [11] J.-P. Corbeil and H. Abdi Ghadivel, “Bet: A backtranslation approach for easy data augmentation in transformer-based paraphrase identification context,” *arXiv*, 2020.
- [12] O. Marchenko and V. Vrublevskiy, “Comparison of transformer-based deep learning methods for the paraphrase identification task,” in *Information Technology and Implementation (ITI-2023)*, CEUR Workshop Proceedings, 2023.
- [13] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” *arXiv*, 2019.
- [14] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, “Roberta: A robustly optimized bert pretraining approach,” *arXiv*, 2019.
- [15] V. Sanh, L. Debut, J. Chaumond, and T. Wolf, “Distilbert, a distilled version of bert: Smaller, faster, cheaper, and lighter,” *arXiv*, 2019.
- [16] A. N. Angelopoulos and S. Bates, “A gentle introduction to conformal prediction and distribution-free uncertainty quantification,” *arXiv*, 2022.
- [17] V. Vovk, A. Gammerman, and G. Shafer, *Algorithmic learning in a random world*. Springer Science & Business Media, 2005.
- [18] S. Khaki and D. Nettleton, “Conformal prediction intervals for neural networks using cross validation,” *arXiv*, 2020.
- [19] V. Vovk, “Cross-conformal predictors,” *Annals of Mathematics and Artificial Intelligence*, vol. 74, no. 1, pp. 9–28, 2015.
- [20] N. Doering, C. Gorlla, T. Tuttle, and A. Vijay, “Empirical analysis of efficient fine-tuning methods for large pre-trained language models,” *arXiv*, 2024.
- [21] A. Maulana, T. R. Noviandy, R. Suhendra, N. Earlia, C. R. S. Prakoeswa, T. S. Kairupan, G. M. Idroes, M. Subianto, and R. Idroes, “Psoriasis severity assessment: Optimizing diagnostic models with deep learning,” *Narra Journal*, vol. 4, no. 3, p. e1512, 2024.
- [22] R. Roslan, I. N. M. Razly, B. Sabri, and Z. Ibrahim, “Evaluation of psoriasis skin disease classification using convolutional neural network,” *IAES International Journal of Artificial Intelligence (IJ-AI)*, vol. 9, no. 2, pp. 349–355, 2020.
- [23] J. Opitz, “A closer look at classification evaluation metrics and a critical reflection of common evaluation practice,” *Transactions of the Association for Computational Linguistics*, vol. 12, p. 820–836, 2024.
- [24] G. S. Dhillon, G. Deligiannidis, and T. Rainforth, “On the expected size of conformal prediction sets,” vol. 238, 2024.
- [25] M. Zecchin, S. Park, O. Simeone, and F. Hellström, “Generalization and informativeness of conformal prediction,” 2024.