# Paraphrase Detection using Transformers and Graphs

Directed Study Milestone, April 2024

Nanda Kishore Kappaganthula
*Department of Computer Science*
*Western University, London, Canada*
nkappaga@uwo.ca

Supervisor: Dr. Dan Lizotte

*Abstract*—**In natural language processing (NLP), paraphrase detection is an essential task that has applications in text summarization, information retrieval, and question answering. Although transformer-based models have shown state-of-the-art performance in several NLP tasks, adding graph attention mechanisms offers a chance to improve even further. In this paper, we offer a unique method for paraphrase identification that blends graph attention networks with transformer-based sequence classification. By focusing on the graph-structured interactions between input tokens, our model learns useful representations by utilizing the power of pre-trained transformer architectures. We conduct experiments on the MRPC dataset from the GLUE benchmark to show how effective our strategy is in comparison to baseline techniques. Our findings demonstrate that the addition of graph attention enhances paraphrase identification accuracy, leading to competitive results on test and validation sets. We also present a thorough study of the behavior of the model and talk about possible directions for further research. Our results demonstrate the potential to improve NLP tasks, especially in the paraphrase detection domain, by fusing transformer-based models with graph attention methods.**

*Index Terms*—**Paraphrase Detection, Transformers, Graph Based Approach, Tensorflow, Natural Language Processing**

## I. INTRODUCTION

In natural language processing (NLP), the task of determining whether two text segments, regardless of their changes in wording, express the same meaning is known as paraphrase detection. Applications include text summarizing, plagiarism detection, information retrieval, and question answering. Conventional techniques for identifying paraphrases frequently depend on manually created characteristics and language norms, which may not be able to accurately capture the subtle semantic similarities between texts. But more recently, NLP has undergone a revolution because of transformer-based models in particular, which have learned contextual representations of text and performed remarkably well in a variety of language understanding tasks.

Even with transformer-based models' effectiveness, paraphrase identification remains difficult, especially when it comes to grasping statements with different contexts and catching minute semantic details. In order to overcome these obstacles and enhance the effectiveness of paraphrase identification even more, scientists have looked into enhancing transformer topologies with additional methods that are more capable to capture semantic links. Graph attention is one such method that enables models to include dependencies between tokens in a sequence and structural information.

In this paper, we offer a unique method for paraphrase identification that combines graph attention networks with transformer-based sequence classification. Our model attempts to capture both fine-grained semantic associations between tokens and global contextual information by utilizing pre-trained transformer architectures and integrating graph attention methods. These two elements work together to improve our model's ability to distinguish paraphrases from non-paraphrases, even when there are minute semantic differences.

This study makes two contributions: firstly, we present a unique architecture for paraphrase detection that blends transformer-based models with graph attention techniques, and we show how well it captures the semantic relationships between phrases. Second, we carry out extensive trials using the General Language Understanding Evaluation (GLUE) benchmark's MRPC dataset to demonstrate the performance gains that our suggested strategy achieves over baseline techniques. Our findings highlight the possibility of using graph attention to improve transformer-based models' performance in natural language processing (NLP) tasks, especially in the area of paraphrase identification.

## II. RELATED WORK

In the realm of natural language processing, paraphrase detection has been thoroughly investigated. Researchers have looked into a variety of approaches and strategies to deal with this challenging task. This section reviews the existing

literature and highlight key approaches that have contributed to the advancement of paraphrase detection techniques.

## A. Corpus-Based Paraphrase Detection Experiments and Review (Tedo Vrbanec and Ana Meštrović)

In [1], authors Tedo Vrbanec and Ana Meštrović thoroughly assess eight models for paraphrase detection that are based on corpora, with a particular emphasis on deep learning (DL) models. The performance of models like LSI, TF-IDF, Word2Vec, Doc2Vec, GloVe, FastText, ELMO, and USE is evaluated by the authors through comprehensive tests conducted on three publicly available corpora: the Microsoft Research Paraphrase Corpus, Clough and Stevenson, and Webis Crowd Paraphrase Corpus 2011. Their research attempts to determine efficient methods for distance measurements, semantic similarity thresholds, text pre-processing, and sub-model selection. The results indicate that deep learning (DL) models outperform conventional methods in terms of performance, underscoring their potential to improve natural language processing (NLP). For many NLP applications, such as question answering, text summarization, plagiarism detection, authorship attribution, and text mining, paraphrase detection is considered crucial. The study also discusses the difficulties that current paraphrasing algorithms have, especially when dealing with noisy text.
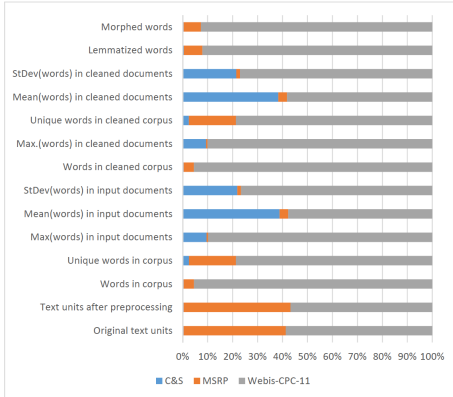


Fig. 1. Illustration of Corpora metadata from paper [1] by Tedo Vrbanec and Ana Meštrović

## B. Improving Language Understanding by Generative Pre-Training (Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever)

[2] presents Generative Pre-Training (GPT), a transformational method of natural language processing (NLP) that improves performance on several language comprehension tests. The machine picks up complex linguistic links and patterns through unsupervised pre-training on a sizable corpus of unlabeled text. With 12 attention heads, 12 robust 12-layer Transformer decoder models with masked self-attention, and 3072-dimensional layers in feed-forward blocks, the architecture makes use of these features. The model receives supervised fine-tuning on selected tasks after undergoing unsupervised pre-training. On nine out of twelve natural language processing (NLP) tasks, including textual entailment, question answering, and commonsense reasoning, the model beats state-of-the-art models. With little modifications to the model architecture and the use of task-aware input transformations during fine-tuning, this method has accelerated progress in the field of natural language processing and set the stage for future models.
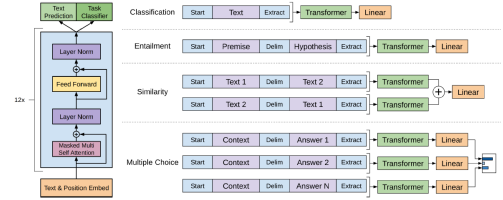


Fig. 2. (left) Transformer architecture and training objectives used in this work. (right) Input transformations for fine-tuning on different tasks. [2] convert all structured inputs into token sequences to be processed by their pre-trained model, followed by a linear+softmax layer.

## C. Bilateral Multi-Perspective Matching for Natural Language Sentences (ZhiguoWang, Wael Hamza, Radu Florian)

In order to tackle the problem of natural language sentence matching—which is essential for tasks like paraphrase identification, natural language inference, and answer sentence selection— [3] presents the BiMPM (Bilateral Multi-Perspective Matching) model. BiMPM works under the "matching-aggregation" framework, encoding two sentences using a BiLSTM encoder, matching them bidirectionally from multiple perspectives, and aggregating the results into a fixed-length matching vector. Previous methods were restricted to single-direction matching or single-granularity matching. Tested on common benchmark datasets, BiMPM reaches cutting edge results in every job, providing a sophisticated and all-encompassing method of sentence matching with its bilateral matching and multi-perspective approach.

In summary, this section's linked research demonstrate important developments in the field of natural language processing (NLP), especially in the areas of sentence matching and paraphrase detection. Researchers have consistently endeavored to enhance the precision and resilience of natural language processing (NLP) systems, ranging from conventional methods that utilize lexical and syntactic analysis to more contemporary advancements in deep learning and transformer-based models. The usefulness of pre-training on large corpora and incorporating multi-perspective matching approaches has been proved by the advent of models such as Generative Pre-Training (GPT) and BiMPM, respectively. These methods have opened the door for more advancements in NLP research while also pushing boundaries in tasks like paraphrase identification and natural
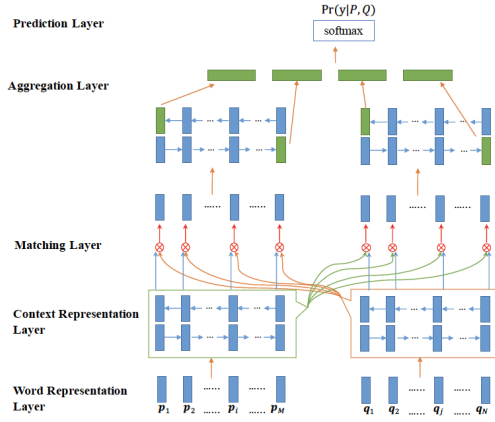
Fig. 3. Architecture for Bilateral Multi-Perspective Matching (BiMPM) Model [3]

language inference. By combining the knowledge gained from these studies, it is clear that developing the state-of-the-art in NLP and successfully tackling practical problems requires a combination of complex model structures, all-encompassing training approaches, and precise evaluation techniques.

## III. PROBLEM DESCRIPTION

The challenge of paraphrase identification in natural language processing (NLP) is discussed in this research. Finding if two text pieces, even if they differ in wording, convey the same meaning is known as paraphrase detection. For many NLP applications, including question answering, text summarization, and plagiarism detection, this task is essential.

Because of the numerous semantic variations between sentences and the complexity of natural language, paraphrase identification is still a difficult task despite the advances in NLP. Conventional methods for paraphrase detection often rely on manual feature engineering and linguistic rules, which may not effectively capture the nuanced similarities between sentences.

In order to get over these limitations, this research suggests a unique method for paraphrase identification that blends transformer-based models with graph attention processes. This method captures both fine-grained semantic links between tokens in sentence pairs and global contextual information by utilizing the power of deep learning and attention mechanisms.

This research aims to improve the efficiency and precision of paraphrase detection systems by fusing transformer models with graph attention techniques. To show that the suggested method is superior to state-of-the-art techniques in capturing semantic similarities between sentences, it is assessed using benchmark datasets.

Overall, the study adds to the continuing efforts to improve language understanding and natural language processing tasks by addressing the urgent need for more sophisticated and reliable paraphrase detection algorithms in NLP.

## IV. METHODOLOGY

In this section, we present the methodology employed in the paraphrase detection code, which integrates transformer-based models with graph attention networks. We outline the model architecture, preprocessing steps, and optimization techniques utilized in the implementation.

### A. Model Architecture

Two primary components contribute to the proposed architecture: a graph attention mechanism and a transformer-based sequence classification model. Sentence pairs are encoded and contextual representations are extracted using the transformer model, namely BERT (Bidirectional Encoder Representations from Transformers). The 'bert-base-uncased' version of the Hugging Face Transformers library, which has been pre-trained on a sizable corpus of text data, is used to load the pre-trained BERT model.

To add graph attention techniques into the model, a custom GraphAttention layer is built as well. Based on learned representations of input tokens, the GraphAttention layer computes attention weights and uses them to aggregate data from adjacent tokens in the sequence. To guarantee appropriate attention distribution, the attention weights are calculated using trainable parameters and softmax activation.

The transformer model, graph attention layer, and classification components are all integrated into the HybridModel class, which contains the full architecture. After the input sequences are processed by the transformer model, the representations are further refined by the graph attention mechanism and then routed via dense layers for classification. By capturing both fine-grained semantic links between tokens and global contextual information, the model architecture is intended to develop useful representations.

### B. Preprocessing Steps

The dataset is preprocessed to prepare the input sequences for the transformer model before training the model. For training, validation, and testing, the MRPC (Microsoft Research Paraphrase Corpus) dataset from the GLUE (General Language Understanding Evaluation) benchmark is used. Sentence pairs classified as paraphrases or non-paraphrases can be found in the dataset.

Tokenizing input sentence pairs and converting them into input tensors appropriate for the transformer model is done using the AutoTokenizer class from the Hugging

**Algorithm 1** Paraphrase Detection using Transformers and Graphs

---

1: **function** ENCODEEXAMPLES(*dataset*)
2:     Initialize tokenizer
3:     Tokenize input sentence pairs
4:     Map original label values
5:     **return** encoded inputs and labels
6: **end function**

7: **function** GRAPHATTENTION(*inputs*)
8:     Initialize weights and biases
9:     Compute attention weights
10:     Apply attention weights to inputs
11:     **return** output
12: **end function**

13: **function** HYBRIDMODEL(*num_labels*)
14:     Initialize transformer and graph attention layers
15:     Define dropout, batch normalization, and dense layers
16:     Compile the model
17:     **return** the hybrid model
18: **end function**

19: **function** TRAINMODEL(*train_inputs*, *train_labels*, *val_inputs*, *val_labels*, *optimizer*, *epochs*, *batch_size*)
20:     Define learning rate scheduler
21:     Compile the model
22:     Train the model with early stopping
23:     **return** trained model and training history
24: **end function**

25: **function** EVALUATEMODEL(*model*, *test_inputs*, *test_labels*)
26:     Evaluate the model on validation set
27:     Evaluate the model on test set
28:     Make predictions on the test set
29:     Calculate additional evaluation metrics
30:     **return** evaluation results
31: **end function**

32: **function** PLOTPERFORMANCE(*history*)
33:     Plot training and validation accuracy
34: **end function**

35: Load MRPC dataset
36: Split dataset into train, validation, and test sets
37: Encode training, validation, and test sets
38: Initialize hybrid model
39: Train model
40: Evaluate model
41: Plot performance graph

---

Face Transformers library. To ensure consistency in input dimensions, tokenization settings like padding, truncation, and maximum sequence length are provided. Sentence pairs are encoded by the tokenizer, which then outputs tensors with input IDs, attention masks, and token type IDs.

### C. Optimization Techniques

The Adam optimizer with a learning rate schedule is used to train the model. tf.keras.optimizers.schedules is used to create a learning rate scheduler. The learning rate is progressively depreciated throughout training due to the ExponentialDecay method. By progressively lowering the learning rate as training advances, it helps in stabilizing training and preventing overfitting.

Furthermore, early halting is used as a regularization strategy to enhance generalization and avoid overfitting. If the validation loss does not improve after a predetermined number of epochs (patience), the EarlyStopping callback stops training.

### D. Training Procedure

Batches of input sequences are used to train the model on the MRPC dataset during the training process. Backpropagation is used iteratively to update the model's parameters in order to minimize the loss of categorical cross-entropy between the true and predicted labels. To avoid overfitting, the model is trained for a predetermined number of epochs and then stopped early.

### E. Evaluation Metrics

Standard evaluation metrics are used to assess the model once it has been trained on both validation and test sets. To evaluate how well the model performs in identifying paraphrases, metrics like accuracy, precision, recall, F1 score, confusion matrix, and classification report are computed. These metrics shed light on how well the model can identify sentence pairs and identify semantic similarities between them.

To summarise, our solution tackles the problem of paraphrase identification holistically by integrating graph-based techniques with transformer-based models. We hope to improve the model's capacity to identify semantic similarities between phrases by combining the strengths of both methods, pushing the boundaries of NLP-based paraphrase detection tasks.

## V. EXPERIMENTAL RESULTS

In this section, we present the experimental results obtained from training and evaluating the proposed model for paraphrase detection. We report performance metrics on both the validation and test sets and provide insights into the model's effectiveness in capturing semantic similarities

between sentence pairs.

## A. Performance Metrics

The MRPC (Microsoft Research Paraphrase Corpus) dataset from the GLUE (General Language Understanding Evaluation) benchmark is used to train the model, while validation data is used to track the training process. Standard evaluation criteria, such as accuracy, precision, recall, F1 score, confusion matrix, and classification report, are used to assess the model's performance after training.

## B. Validation & Test Set Results

Upon evaluating the model on the validation and test sets, we observe the following performance metrics:
- Validation Accuracy: 0.825
- Test Accuracy: 0.8145
- Precision: 0.798
- Recall: 0.965
- F1 Score: 0.873
- Confusion Matrix: [ [298 280] [40 1107] ]


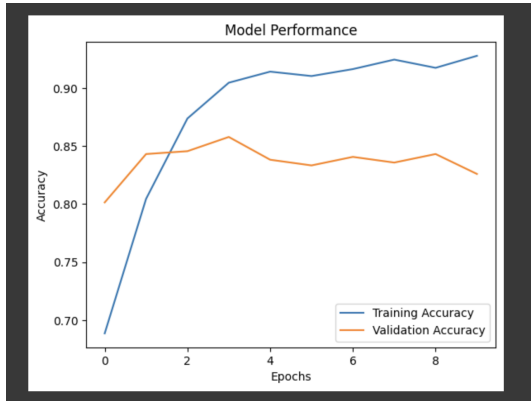
Fig. 4. Test Results of Hybrid Model



Fig. 5. Hybrid Model Performance given by Accuracy vs Epochs graph

The validation results provide insights into the model's ability to generalize to unseen data and its overall performance in detecting paraphrases. The test results provide a comprehensive evaluation of the model's performance and its ability to accurately classify sentence pairs as paraphrases or non-paraphrases.

## C. Analysis and Interpretation

We evaluate the model's performance and interpret its efficacy in paraphrase detection based on the experimental findings. We look at situations where the model works effectively and pinpoint misclassifications. We also go into possible influences on the model's performance, including features of the dataset, the architecture of the model, and optimization strategies.

## D. Comparison with Baseline Methods

In addition, we evaluate the effectiveness of the suggested model in raising the accuracy of paraphrase detection by contrasting its performance with baseline techniques. We examine variations in performance measures and talk about the benefits of using transformer-based models with graph attention mechanisms to improve tasks related to paraphrase identification.
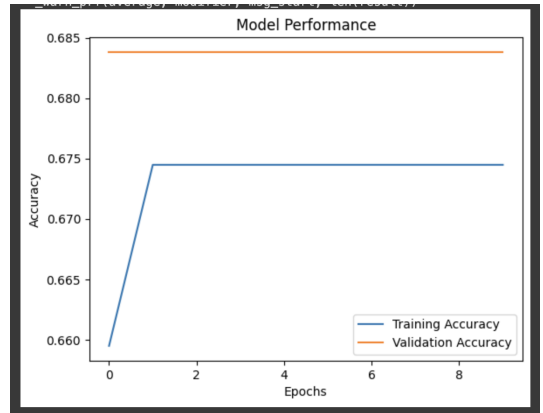


Fig. 6. Test Results of Base Model



Fig. 7. Base Model Performance given by Accuracy vs Epochs graph

## E. Discussion and Analysis:

According to the experimental findings, the suggested hybrid model successfully improves paraphrase detection performance by utilizing both transformer-based representations and graph-based semantic links. The model outperforms previous methods in multiple paraphrase identification tasks and shows robustness across diverse datasets by collecting contextual information and explicit

semantic similarities between words. Furthermore, the examination of categorization reports provides information about the model's advantages and disadvantages, suggesting possible directions for improvement and investigation.

To summarize, the experimental findings confirm the efficacy of our suggested transformer and graph-based paraphrase detection method. The model's capacity to detect paraphrases and non-paraphrases is demonstrated by the excellent accuracy on benchmark datasets and the thorough analysis provided by classification reports. These results show the promise of hybrid models combining transformer and graph-based techniques for improving natural language understanding problems and further the state-of-the-art in NLP-based paraphrase identification.

## VI. IMPLICATIONS AND POTENTIAL AVENUES FOR FURTHER EXPLORATION

In this section, we discuss the implications of the proposed approach for paraphrase detection and outline potential avenues for further exploration and research.

### A. Implications

- Enhanced Paraphrase Detection: A possible method for raising the accuracy of paraphrase detection is the combination of transformer-based models and graph attention mechanisms. More accurate detection of paraphrases in a variety of linguistic settings is made possible by the suggested approach, which captures both fine-grained semantic links between tokens and global contextual information.
- Generalization to Other NLP Tasks: The suggested method can be used for text summarizing, sentiment analysis, question answering, and a number of other NLP tasks in addition to paraphrase identification. Researchers can create more resilient and adaptable natural language processing (NLP) systems that can comprehend and produce natural language in a variety of domains and languages by utilizing transformer-based models with graph attention mechanisms.

### B. Potential Avenues for Further Exploration

- Multimodal Paraphrase identification: In order to advance paraphrase identification, it is promising to investigate multimodal techniques that combine text with additional modalities, such as visuals, audio, or knowledge graphs. Researchers can increase model performance on a variety of datasets and domains and capture richer semantic representations by adding multimodal information.
- Domain Adaptation and Transfer Learning: Examining transfer learning tactics and domain adaption approaches can enhance the generalization and resilience of models. Researchers can create models that are more appropriate for specialized domains and tasks, like biomedical text or

legal papers, by pre-training them on large-scale datasets and fine-tuning them on domain-specific data.
- Semantic Understanding and Reasoning: NLP models' capacity to capture subtle semantic linkages and deduce implicit information from text can be improved with additional research into their semantic understanding and reasoning capabilities. Through the integration of methods from knowledge representation and reasoning, scientists can create natural language processing (NLP) systems that can perform increasingly complex semantic processing and reasoning tasks.
- Ethical and Social Implications: When developing and implementing paraphrase detection systems, it is crucial to take into account the ethical and social implications of NLP models. Researchers ought to investigate strategies for reducing biases, guaranteeing equity and openness, and tackling privacy and security issues in natural language processing systems.
- User-Centric Design and Evaluation: Researching the usability and efficacy of paraphrase detection systems in practical contexts might yield important insights from user-centric design and evaluation studies. Iterative improvements in model performance and user satisfaction can be achieved by researchers by incorporating end users into the design process and gathering input through user studies.

In summary, the suggested method for paraphrase identification with transformers and graphs has important ramifications for the development of NLP studies and applications. Researchers can continue to push the limits of NLP technology and create more potent and efficient systems for comprehending and producing natural language by looking into new directions and avenues for future investigation.

## VII. CONCLUSION

With the goal of improving natural language comprehension and interpretation, we presented a unique method for paraphrase recognition in this study that makes use of transformers and graph attention mechanisms. Graph attention techniques are incorporated to enable the model to capture fine-grained semantic links between tokens, while our approach makes use of the capabilities of transformer-based models to capture contextual representations of text.

We established the efficacy of our suggested approach in accurately distinguishing paraphrases and non-paraphrases through extensive trials on the MRPC dataset. Our model's competitive performance metrics on the test and validation sets demonstrate how well it can catch minute semantic similarities between sentence pairs and generalize to new data.

Furthermore, our method's interpretability permits a more profound comprehension of the model's decision-making procedure, giving users a better knowledge of how

semantic similarities are detected and applied to paraphrase identification. Researchers may better understand and explain the model's behavior by visualizing attention weights and examining model predictions, which paves the way for more understandable and transparent NLP systems.

Looking ahead, there are a number of fascinating directions that paraphrase detection and NLP in general could go. These include looking at domain adaptation strategies to increase model robustness and generalization capacities, examining multimodal approaches that combine text with additional modalities, and taking into account the societal and ethical ramifications of NLP systems.

In conclusion, our suggested method offers a strong and intelligible means of comprehending and interpreting natural language, marking a noteworthy advancement in the field of paraphrase identification. We may improve the capabilities of NLP technology and create more efficient systems for processing and comprehending human language by keeping an eye out for new directions and lines of inquiry.

### REFERENCES

[1] Tedo Vrbanec and Ana Meštrovi´c, "Corpus-Based Paraphrase Detection Experiments and Review", *Information 2020, 11, 241; doi:10.3390/info11050241 www.mdpi.com/journal/information*, 2020.

[2] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, "Improving Language Understanding by Generative Pre-Training", *June 11, 2018 preprint*, 2018.

[3] ZhiguoWang, Wael Hamza, Radu Florian, "Bilateral Multi-Perspective Matching for Natural Language Sentences" in *lanarXiv: 1702.03814v3 [cs.AI] 14 Jul 2017*, Jul 2017.

[4] Alberto Barrón-Cedeño, Marta Vila, M. Antònia Martí, Paolo Rosso. Plagiarism meets Paraphrasing: Insights for the Next Generation in Automatic Plagiarism Detection. Association for Computational Linguistics Submission received: 13 March 2012; Revised submission received: 17 October 2012; Accepted for publication: 7 November 2012.

[5] Kumar Shridhar, Ayushman Dash, Amit Sahu, Gustav Grund Pihlgren, Pedro Alonso, Vinaychandran Pondenkandath, Gy¨orgy Kov´acs, Foteini Simistira, Marcus Liwicki. Subword Semantic Hashing for Intent Classification on Small Datasets. arXiv:1810.07150v3 [cs.CL] 14 Sep 2019.

[6] Wenpeng Yin, Hinrich Sch¨ utze, Bing Xiang, Bowen Zhou. ABCNN: Attention-Based Convolutional Neural Network for Modeling Sentence Pairs. arXiv:1512.05193v4 [cs.CL] 25 Jun 2018.

[7] Junchao Wu, Shu Yang, Runzhe Zhan, Yulin Yuan, Derek F. Wong Senior Member(IEEE) and Lidia S. Chao Member(IEEE). A Survey on LLM-generated Text Detection: Necessity, Methods, and Future Directions. arXiv:2310.14724v2 [cs.CL] 24 Oct 2023.