

Can Model Fusing Help Transformers in Long Document Classification? An Empirical Study

Damith Premasiri[◇], Tharindu Ranasinghe[♡] and Ruslan Mitkov[♣]

[◇]University of Wolverhampton, Wolverhampton, UK

[♡]Aston University, Birmingham, UK

[♣]Lancaster University, Lancaster, UK

damith.premasiri@wlv.ac.uk, t.ranasinghe@aston.ac.uk

r.mitkov@lancaster.ac.uk

Abstract

Text classification is an area of research which has been studied over the years in Natural Language Processing (NLP). Adapting NLP to multiple domains has introduced many new challenges for text classification and one of them is long document classification. While state-of-the-art transformer models provide excellent results in text classification, most of them have limitations in the maximum sequence length of the input sequence. The majority of the transformer models are limited to 512 tokens, and therefore, they struggle with long document classification problems. In this research, we explore on employing *Model Fusing* for long document classification while comparing the results with well-known BERT and Longformer architectures.

1 Introduction

Text classification is one of the critical tasks in Natural Language Processing, which refers to finding the suitable label/ labels to a particular input text (Kowsari et al., 2019; Mirończuk and Protasiewicz, 2018). It has a wide range of applications in different domains such as sentiment analysis (Dang et al., 2020b,a), fake news detection (Thota et al., 2018; Kumar et al., 2020; Ahmad et al., 2020) and offensive language identification (Ranasinghe and Zampieri, 2020; Husain and Uzuner, 2021). These tasks are generally referred to as sentence classification tasks since the input text is typically in the form of sentences. In recent years, transformer models such as BERT have provided state-of-the-art results in these text classification tasks (Ranasinghe et al., 2019; Gaikwad et al., 2021).

While most of the text classification tasks are sentence classification, several domains require classifying lengthy texts into labels typically referred to as document classification. Specifically, domains such as legal and medical often contain

long documents that need document classification methods (Chalkidis et al., 2019a; Hettiarachchi et al., 2023). However, adapting the transformer models that produced state-of-the-art results in sentence classification to document classification is challenging (Pappagari et al., 2019). The most common transformer models, such as BERT (Devlin et al., 2019), have a limitation of 512 tokens in their input layer, which means the tokens in a lengthy document exceeding this limit will be truncated in the tokenisation step.

The limitations outlined above have garnered significant attention from the research community, leading to the exploration of new document classification architectures. One widely adopted approach is to leverage transformer models that can process longer sequences. Notably, the Longformer (Beltagy et al., 2020) and BigBird (Zaheer et al., 2020) transformer models have demonstrated exceptional performance in document classification tasks, with the capacity to accommodate up to 4,096 tokens. However, training transformer models that can process longer sequences is a resource-intensive task, and it may not be feasible for less-resourced domains and languages (Wagh et al., 2021; Zhang and Jankowski, 2022). In an effort to mitigate this challenge, researchers have attempted to adapt existing pre-trained transformer models to accommodate longer sequences. Notably, two such approaches are Hierarchical BERT (Lu et al., 2021) and CogLTX (Ding et al., 2020), both of which propose innovative strategies for adapting BERT to long document classification. Following this, we propose a method to adapt BERT-like transformer models to long document classification using *Model Fusion*. While the methods such as Hierarchical BERT (Lu et al., 2021) and CogLTX (Ding et al., 2020) mainly focus on tackling long-term dependencies using different attention mechanisms to reduce their computational complexity, we explore

a new idea with model fusing to the long document classification task.

Model Fusion refers to the idea of combining several fine-tuned models (Xu et al., 2020). The motivation behind using *Model Fusion* is that multiple models can identify different patterns using different parts of their network, and it is possible to merge multiple models into one model, which will be capable of having all information compressed into a single model. To implement this idea, we divide long documents into multiple parts and use these parts to train part-wise models. Finally, we *fuse* all part-wise models to create a single model capable of handling lengthy sequences. Our evaluation of this approach on four popular document classification datasets shows that while our hypothesis is strong, *Model Fusion* does not improve state-of-the-art document classification. Nonetheless, we report our results with the aim of helping researchers avoid repeating unsuccessful experiments in the future. Furthermore, this paper identifies potential flaws in experimental design, enabling researchers to refine their methods and improve future studies that employ *Model Fusion* in long document classification¹.

Our main contributions of the paper are,

1. We present the first study in using *Model Fusion* in long document classification.
2. We empirically evaluate the proposed approach in four benchmark datasets in document classification and show that the proposed method does not outperform the baselines such as Longformer (Beltagy et al., 2020).
3. We release the code and the model resources freely available to the public².

The rest of the paper is organised as follows. Section 2 highlights the recent work on long document classification and model fusing. Section 3 describes the datasets we used. Section 4 explains data preparation for experiments, sub-model training, model fusing and prediction on test data. Section 5 presents the results and discusses possible problems in the results and ideas for improvements. Section 6 summarises our main experimental findings and conclusions.

¹Publishing negative results has also been encouraged with the organisation of workshops such as Workshop on Insights from Negative Results in NLP <https://insights-workshop.github.io/>

²Code is available at <https://github.com/DamithDR/legal-classification>

2 Related Work

Long Text Classification Over the years, researchers have explored various methods to address long text classification, from traditional machine learning approaches such as SVMs (Boser et al., 1992) to recent deep learning architectures (Dai et al., 2022; Uyangodage et al., 2021b). With the emergence of transformers, the researchers focused heavily on adapting transformer models to long text classification. Longformer (Beltagy et al., 2020) is one such method (Hettiarachchi et al., 2021), which is capable of accommodating 4,096 tokens. Longformer’s attention mechanism is a combination of a windowed local-context self-attention, and an end task motivated global attention that encodes inductive bias about the task. Through ablations and controlled trials, they show both attention types are essential – the local attention is primarily used to build contextual representations, while the global attention allows Longformer to build full sequence representations for prediction. As we mentioned before, training a transformer model that supports lengthy inputs is expensive. Therefore, researchers have explored how to use existing pre-trained transformer models in long document classification.

CogLTX (Ding et al., 2020) is a method which proposes an efficient way of processing long documents using two jointly trained BERT (Devlin et al., 2019) models to select key sentences from long documents for various tasks, including text classification. Their idea is that a few key sentences can be sufficient to get an understanding of the overall text, which works for some tasks but not essentially for document classification. Pappagari et al. (2019) introduced ToBERT, which can process documents of any length using chunking. However, it does not improve performance in many document classification tasks.

Dai et al. (2022) provides a revision on transformers’ capabilities on long document classification. Park et al. (2022) shows a performance comparison between Longformer (Beltagy et al., 2020), CogLTX (Ding et al., 2020), ToBERT (Pappagari et al., 2019) and their novel baselines BERT+TextRank; where they identify the key sentences using TextRank (Mihalcea and Tarau, 2004) and uses these sentences to fill the 512 tokens of a BERT rather than using the full document as the input. BERT+Random; is a simpler baseline where they use random sentences to fill the 512 tokens. Interestingly they show that for most of the datasets,

specific long-text processing methods fail to outperform these simple baselines. [Limsopatham \(2021\)](#) has experimented with the effective usage of BERT for long document classification by parsing the front part of the document and the rear part of the document separately and experimenting with the results. Despite numerous efforts to address challenges in long document classification, the results still fall short compared to sentence classification, demanding further dedication from the research community.

Model Fusion Fusing is applied on different parts and different levels of NLP tasks. [Choshen et al. \(2022\)](#) proposes a way to fuse the models to have better pre-trained models. [Xiong et al. \(2021\)](#) does label fusing via concatenating texts of labels and an original document to be classified with a [SEP] token as an input, and they use different segment embeddings for the label texts and the document text. [Lai et al. \(2023\)](#) have used Gated Fusing to improve backward compatibility when doing updates of NLP models. Fusing has been employed in multi-model research, too. [Khan et al. \(2020\)](#) provides fusing multiple models for visual question answering.

As fusion has provided excellent results in different tasks, we hypothesise that fusion can be used to solve document classification. As far as we know, this is the first study to use model fusion in long document classification.

3 Data

We evaluated our approach with four popular document classification datasets; ECHR ([Chalkidis et al., 2019b](#)), ECHR_Anon ([Chalkidis et al., 2019b](#)) 20NewsGroups ([Lang, 1995](#)) and case-2022 ([Hürriyetoglu et al., 2022](#)). We describe each of them below. The distribution of the number of words in each dataset is also shown in Table 1.

ECHR ([Chalkidis et al., 2019b](#)) European Court of Human Rights (ECHR) hears allegations that a state has breached human rights provisions of the European Convention of Human Rights. The dataset contains approx. 11.5k cases from ECHR’s public database. We use the dataset for document-level binary violation tasks; given the facts of a case, the task is to classify whether there has been any human rights violation or not.

ECHR_Anon ([Chalkidis et al., 2019b](#)) This dataset contains an anonymised version of the

ECHR with demographic data being anonymised. To achieve this, all Named Entities in the text have been replaced with corresponding tags.

20NewsGroups ([Lang, 1995](#)) The dataset is composed of 18828 news articles, which are classified into 20 different categories. The goal of this task is to perform multi-class classification to accurately identify the category of each article. To evaluate our model’s performance, we reserve 20% of the data for the test set.

Case-2022 ([Hürriyetoglu et al., 2022](#)) This dataset is from the shared Task on Socio-political and Crisis Events Detection CASE - subtask 1. The task is a document classification to detect whether a news article contains information about a socio-political event or not. The Dataset features 9384 news articles in the training set, and we have utilised 20% of it as the test set since the gold labels in the test set are not released.

Dataset	w < 512	512 < w < 4096	w > 4096
ECHR	16.04	69.15	14.80
ECHR_Anon	16.07	67.69	16.24
20NewsGroups	86.72	12.67	0.61
Case-2022	96.27	3.73	0.00

Table 1: Percentages of distribution of a number of data instances against the word count (w) in the dataset.

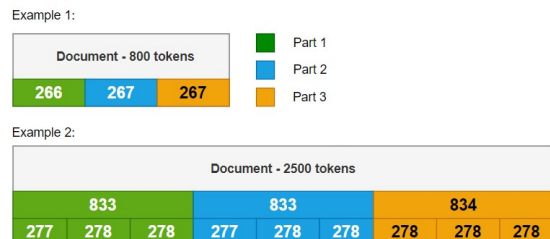


Figure 1: Document breakdown to parts

4 Methodology

We divide our method into five stages, which we describe below.

Data Preparation Since the datasets contain data points which exceed 512 token limitation in BERT ([Devlin et al., 2019](#)) as shown in Table 1, we evenly distributed each document among sub-models. Initially, we determined the number of parts to divide the data points based on a trial-and-error approach. Early experiments suggested that dividing each data point into three parts produced the best

results. We also restricted each part to a maximum of 400 words. For documents with more than 1200 words, such as 3000 words, we split them into three parts of 1000 words each. Due to the 512 token limitation, we further divided the 1000 words into more sub-parts, but all sub-parts were trained on the same model. Essentially, when we split a document into parts, each part has its own respective model that is used for training. To maintain consistency, we assigned respective class labels to the divided parts of the document. We assumed that all parts contribute equally to the class classification, so if the data point had classification label A, all parts of the document would also have the classification label A as illustrated in Figure 1.

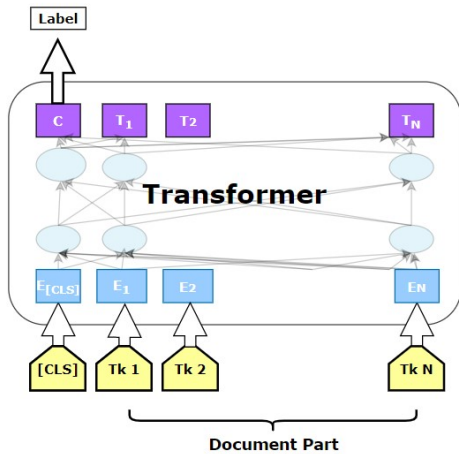


Figure 2: Transformer model for document level classification (Uyangodage et al., 2021a)

Sub-model Training The number of sub-models to be trained is equal to the number of parts in the document. The main idea is to understand the data in a part-localised manner to tackle the length issue. Therefore, in our experiments, we used three sub-models in-line with three parts in each document. As shown in Figure 3, Part 1 of each document goes to the training set of sub-model 1 and, respectively, part 2 and part 3 into sub-model 2 and 3. We assume that this part-wise modelling can understand the part-local information, which could then contribute to the final classification. Sub-models were trained by using a BERT (Devlin et al., 2019) model for all experiments since it has produced excellent results in many natural language processing tasks (Morgan et al., 2021). We used a softmax layer on top of the last hidden layer of the Transformer architecture, as shown in Figure 2. The configurations we used are listed in Table 2.

Parameter	Value
Training Batch Size	32
Evaluation Batch Size	8
Learning Rate	$4e-5$
Epochs	3
Early Stopping	No

Table 2: Sub-model training configurations

Model Fusing Once the sub-models are trained, we read the weights of hidden layers of the models and fused them together while input and output layers remain unchanged. We employed average fusing for simplicity, in which the resulting fused model has the average of weights in the sub-models as shown in Figure 3.

$$W_{fused} = f(W_1, W_2, \dots, W_n) \quad (1)$$

$$W_{fused} = (W_1 + W_2 + \dots + W_n)/n \quad (2)$$

By averaging the weights, we assume that the characteristics of each part of the document are being merged into one fused model.

Further Fine-tuning we further fine-tune the fused model using a fraction of the training set, which was split from the training set in the beginning. This step is important as once we merge the models together, the weights of hidden layers are not finely coupled with the output layers. In order to correct this, further fine-tuning step is important and performed using all parts of the document. For this reason, further fine-tune data contain text from all parts separately. In the fine-tune step, we used the same configurations as sub-model training having batch-size of 32, Adam optimiser with learning rate $4e-5$. Once we complete this, the fused model is ready to predict on the test data.

Prediction Predicting on test data uses a similar approach to training. We divide the original documents into parts and then predict the classification class for each one of them. We then get the mean of the probabilities of each class and decide the final classification class. We also experimented with taking the max of the probabilities; however, it did not show improvements compared to taking the mean. Therefore, all the results we present were taken using the mean.

5 Results and Discussion

Baselines Baseline results were taken from well-known BERT (Devlin et al., 2019) and Longformer (Beltagy et al., 2020) which were configured to

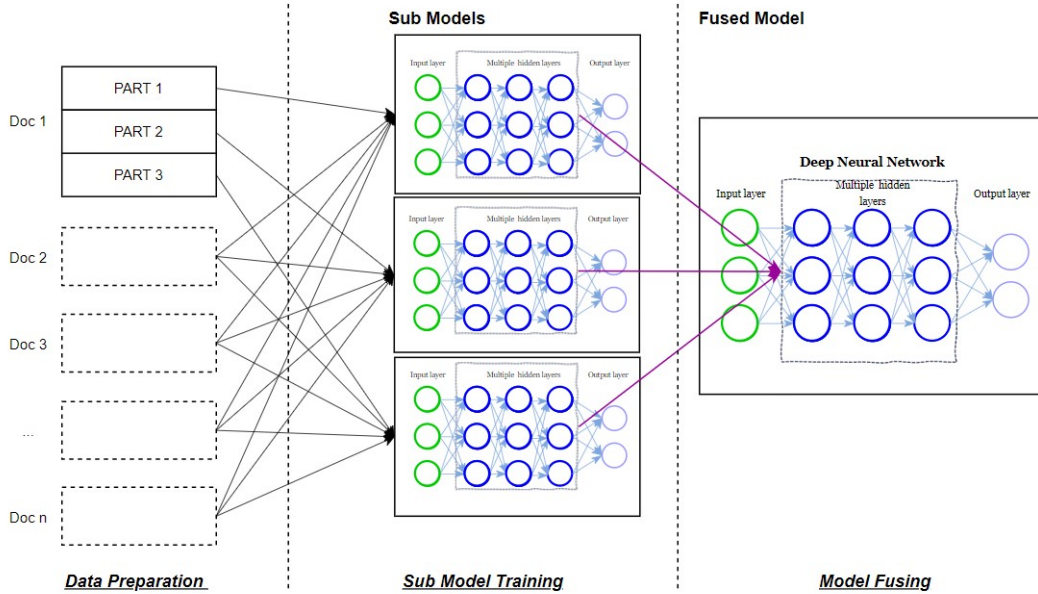


Figure 3: Model fusing pipeline for long document classification

Dataset	Fusing			Bert			Longformer		
	P	R	F1	P	R	F1	P	R	F1
ECHR	0.6127	0.6451	0.5486	0.8493	0.8486	0.8212	0.8504	0.8516	0.8278
ECHR_Annon	0.6232	0.6621	0.4673	0.8209	0.8235	0.7950	0.8395	0.8369	0.8041
20NewsGroups	0.5361	0.5409	0.4984	0.8952	0.8941	0.8910	0.8981	0.8980	0.8951
Case-2022	0.6272	0.7920	0.4420	0.8837	0.8858	0.8231	0.8956	0.8981	0.8405

Table 3: Results for different datasets for Fusing, Bert and Longformer. P; weighted Precision, R; weighted Recall, F1; Macro F1

truncate the sequences which exceeded their token limit. Additionally, Longformer (Beltagy et al., 2020) has the special capability to accommodate up to 4096 tokens.

Results Table 3 shows the results for Fusing, BERT and Longformer. It is clear that Longformer performs best among all datasets confirming its unique ability on long document classification. BERT also shows good performance in all cases, and it is clear that 20NewsGroups and Case-2022 datasets are fairly within the range of no of tokens which BERT could capture (512) (Table 1). However, BERT also performs well in ECHR cases. We believe the reason for that is the first parts of the facts of ECHR cases heavily contribute to the final label.

Fusing results are the lowest in all cases, confirming that model fusing will not produce better results for the long document classification task. It is noticeable that Fusing also has similar trends across datasets as Longformers. Longformer has pro-

duced F1 scores of 0.8278 and 0.8041 for ECHR and ECHR_Annon data, respectively, while Fusing also shows a similar pattern by marking 0.5486 and 0.4673 F1 scores for the same.

One possible reason for the low performance of the Fusing method could be our assumption where we assumed that all parts of the document equally contribute to its class. This could not be the case at all times, and if not, models will learn incorrect information, which could lead to lower results. Another possibility is the division of the documents into parts. Dividing the documents into parts will induce information flow breaks from which the models could suffer.

Even though our intuition of model fusing is similar to transfer learning, average fusing has its own problems. Averaging weights might not be ideal because the activation of the neurons could catch with heavy negation. If we average the values 4 and 5, the result is 4.5, which shows that the resulting weight does not deviate from both original weights drastically. However, if we consider 5 and 0.1, their

average result is 2.55, which shows a considerable difference between both initial weights. In a numerical model such as BERT, this could introduce significant changes in the network’s decision-making process. One way to overcome this issue could be introducing a weighted bias to the sub-models. This way, one model will get favouritism over others and possibly lead to better results, but it will need extensive experiments to confirm this.

6 Conclusion

This paper presents an empirical study on the effectiveness of model fusing in long document classification, with the aim of comparing its performance to that of state-of-the-art models such as Longformer (Beltagy et al., 2020). Our results indicate that Longformer (Beltagy et al., 2020) outperforms our experimental setup across all datasets. While we identify several drawbacks of the method, we believe that there is still potential for further exploration in this area. Although our average fusing approach did not yield improved performance in long document classification, there is a need for more research on different fusing methods and their efficacy in various tasks.

Acknowledgments

We thank the anonymous RANLP reviewers who have provided us with constructive feedback to improve the quality of this paper.

We also thank the creators of the datasets we used in the study for making them public.

References

- Ifitikhar Ahmad, Muhammad Yousaf, Suhail Yousaf, and Muhammad Ovais Ahmad. 2020. Fake news detection using machine learning ensemble methods. *Complexity*, 2020:1–11.
- Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv:2004.05150*.
- Bernhard E Boser, Isabelle M Guyon, and Vladimir N Vapnik. 1992. A training algorithm for optimal margin classifiers. In *Proceedings of the fifth annual workshop on Computational learning theory*, pages 144–152.
- Ilias Chalkidis, Ion Androutsopoulos, and Nikolaos Aletras. 2019a. [Neural legal judgment prediction in English](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4317–4323, Florence, Italy. Association for Computational Linguistics.
- Ilias Chalkidis, Ion Androutsopoulos, and Nikolaos Aletras. 2019b. [Neural legal judgment prediction in english](#). *arXiv preprint arXiv:1906.02059*.
- Leshem Choshen, Elad Venezian, Noam Slonim, and Yoav Katz. 2022. [Fusing finetuned models for better pretraining](#). *arXiv preprint arXiv:2204.03044*.
- Xiang Dai, Ilias Chalkidis, Sune Darkner, and Desmond Elliott. 2022. [Revisiting transformer-based models for long document classification](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 7212–7230, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Nhan Cach Dang, María N Moreno-García, and Fernando De la Prieta. 2020a. [Sentiment analysis based on deep learning: A comparative study](#). *Electronics*, 9(3):483.
- Nhan Cach Dang, María N. Moreno-García, and Fernando De la Prieta. 2020b. [Sentiment analysis based on deep learning: A comparative study](#). *Electronics*, 9(3).
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Ming Ding, Chang Zhou, Hongxia Yang, and Jie Tang. 2020. [Cogltx: Applying bert to long texts](#). *Advances in Neural Information Processing Systems*, 33:12792–12804.
- Saurabh Sampatrao Gaikwad, Tharindu Ranasinghe, Marcos Zampieri, and Christopher Homan. 2021. [Cross-lingual offensive language identification for low resource languages: The case of Marathi](#). In *Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021)*, pages 437–443, Held Online. INCOMA Ltd.
- Hansi Hettiarachchi, Mariam Adedoyin-Olowe, Jagdev Bhogal, and Mohamed Medhat Gaber. 2021. [DAAI at CASE 2021 task 1: Transformer-based multilingual socio-political and crisis event detection](#). In *Proceedings of the 4th Workshop on Challenges and Applications of Automated Extraction of Socio-political Events from Text (CASE 2021)*, pages 120–130, Online. Association for Computational Linguistics.
- Hansi Hettiarachchi, Mariam Adedoyin-Olowe, Jagdev Bhogal, and Mohamed Medhat Gaber. 2023. [Ttl: transformer-based two-phase transfer learning for cross-lingual news event detection](#). *International Journal of Machine Learning and Cybernetics*.
- Ali Hürriyetoğlu, Hristo Tanev, Vanni Zavarella, and Erdem Yörük, editors. 2022. *Proceedings of the 5th*

- Workshop on Challenges and Applications of Automated Extraction of Socio-political Events from Text (CASE)*. Association for Computational Linguistics, Abu Dhabi, United Arab Emirates (Hybrid).
- Fatemah Husain and Ozlem Uzuner. 2021. A survey of offensive language detection for the arabic language. *ACM Transactions on Asian and Low-Resource Language Information Processing (TALLIP)*, 20(1):1–44.
- Aisha Urooj Khan, Amir Mazaheri, Niels Da Vitoria Lobo, and Mubarak Shah. 2020. Mmft-bert: Multimodal fusion transformer with bert encodings for visual question answering. *arXiv preprint arXiv:2010.14095*.
- Kamran Kowsari, Kiana Jafari Meimandi, Mojtaba Heidarysafa, Sanjana Mendu, Laura Barnes, and Donald Brown. 2019. **Text classification algorithms: A survey**. *Information*, 10(4).
- Sachin Kumar, Rohan Asthana, Shashwat Upadhyay, Nidhi Upreti, and Mohammad Akbar. 2020. Fake news detection using deep learning models: A novel approach. *Transactions on Emerging Telecommunications Technologies*, 31(2):e3767.
- Yi-An Lai, Elman Mansimov, Yuqing Xie, and Yi Zhang. 2023. Improving prediction backward-compatibility in nlp model upgrade with gated fusion. *arXiv preprint arXiv:2302.02080*.
- Ken Lang. 1995. Newsweeder: Learning to filter netnews. In *Machine learning proceedings 1995*, pages 331–339. Elsevier.
- Nut Limsopatham. 2021. **Effectively leveraging BERT for legal document classification**. In *Proceedings of the Natural Legal Language Processing Workshop 2021*, pages 210–216, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Jinghui Lu, Maeve Henchion, Ivan Bacher, and Brian Mac Namee. 2021. A sentence-level hierarchical bert model for document classification with limited labelled data. In *Discovery Science*, pages 231–241, Cham. Springer International Publishing.
- Rada Mihalcea and Paul Tarau. 2004. **TextRank: Bringing order into text**. In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*, pages 404–411, Barcelona, Spain. Association for Computational Linguistics.
- Marcin Michał Mirończuk and Jarosław Protasiewicz. 2018. **A recent overview of the state-of-the-art elements of text classification**. *Expert Systems with Applications*, 106:36–54.
- Skye Morgan, Tharindu Ranasinghe, and Marcos Zampieri. 2021. **WLV-RIT at GermEval 2021: Multi-task learning with transformers to detect toxic, engaging, and fact-claiming comments**. In *Proceedings of the GermEval 2021 Shared Task on the Identification of Toxic, Engaging, and Fact-Claiming Comments*, pages 32–38, Duesseldorf, Germany. Association for Computational Linguistics.
- Raghavendra Pappagari, Piotr Zelasko, Jesús Villalba, Yishay Carmiel, and Najim Dehak. 2019. Hierarchical transformers for long document classification. In *2019 IEEE automatic speech recognition and understanding workshop (ASRU)*, pages 838–844. IEEE.
- Hyunji Park, Yogarshi Vyas, and Kashif Shah. 2022. **Efficient classification of long documents using transformers**. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 702–709, Dublin, Ireland. Association for Computational Linguistics.
- Tharindu Ranasinghe and Marcos Zampieri. 2020. **Multilingual offensive language identification with cross-lingual embeddings**. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5838–5844, Online. Association for Computational Linguistics.
- Tharindu Ranasinghe, Marcos Zampieri, and Hansi Hettiarachchi. 2019. Brums at hasoc 2019: Deep learning models for multilingual hate speech and offensive language identification. In *Working Notes of FIRE 2019-Forum for Information Retrieval Evaluation*, pages 199–207.
- Aswini Thota, Priyanka Tilak, Simrat Ahluwalia, and Nibrat Lohia. 2018. Fake news detection: a deep learning approach. *SMU Data Science Review*, 1(3):10.
- Lasitha Uyangodage, Tharindu Ranasinghe, and Hansi Hettiarachchi. 2021a. **Can multilingual transformers fight the COVID-19 infodemic?** In *Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021)*, pages 1432–1437, Held Online. INCOMA Ltd.
- Lasitha Uyangodage, Tharindu Ranasinghe, and Hansi Hettiarachchi. 2021b. **Transformers to fight the COVID-19 infodemic**. In *Proceedings of the Fourth Workshop on NLP for Internet Freedom: Censorship, Disinformation, and Propaganda*, pages 130–135, Online. Association for Computational Linguistics.
- Vedangi Wagh, Snehal Khandve, Isha Joshi, Apurva Wani, Geetanjali Kale, and Raviraj Joshi. 2021. Comparative study of long document classification. In *TENCON 2021-2021 IEEE Region 10 Conference (TENCON)*, pages 732–737. IEEE.
- Yijin Xiong, Yukun Feng, Hao Wu, Hidetaka Kamigaito, and Manabu Okumura. 2021. **Fusing label embedding into BERT: An efficient improvement for text classification**. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 1743–1750, Online. Association for Computational Linguistics.
- Guangxia Xu, Weifeng Li, and Jun Liu. 2020. A social emotion classification approach using multi-model fusion. *Future Generation Computer Systems*, 102:347–356.

Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. 2020. Big bird: Transformers for longer sequences. *Advances in neural information processing systems*, 33:17283–17297.

Ning Zhang and Maciej Jankowski. 2022. Hierarchical bert for medical document understanding. *arXiv preprint arXiv:2204.09600*.