

# Ethic-BERT: An Enhanced Deep Learning Model for Ethical and Non-Ethical Content Classification

Mahamodul Hasan Mahadi, Md. Nasif Safwan, Souhardo Rahman,  
Shahnaj Parvin, Aminun Nahar, Kamruddin Nur

Department of Computer Science, American International  
University-Bangladesh, Dhaka, 1229, Bangladesh.

\*Corresponding author(s). E-mail(s): [kamruddin@aiub.edu](mailto:kamruddin@aiub.edu);

Contributing authors: [24-93530-3@student.aiub.edu](mailto:24-93530-3@student.aiub.edu);

[22-49041-3@student.aiub.edu](mailto:22-49041-3@student.aiub.edu); [22-49068-3@student.aiub.edu](mailto:22-49068-3@student.aiub.edu);

[sparvin@aiub.edu](mailto:sparvin@aiub.edu); [aminun.nahar@aiub.edu](mailto:aminun.nahar@aiub.edu);

## Abstract

Developing AI systems capable of nuanced ethical reasoning is critical as they increasingly influence human decisions, yet existing models often rely on superficial correlations rather than principled moral understanding. This paper introduces Ethic-BERT, a BERT-based model for ethical content classification across four domains: Commonsense, Justice, Virtue, and Deontology. Leveraging the ETHICS dataset, our approach integrates robust preprocessing to address vocabulary sparsity and contextual ambiguities, alongside advanced fine-tuning strategies like full model unfreezing, gradient accumulation, and adaptive learning rate scheduling. To evaluate robustness, we employ an adversarially filtered 'Hard Test' split, isolating complex ethical dilemmas. Experimental results demonstrate Ethic-BERT's superiority over baseline models, achieving 82.32% average accuracy on the standard test, with notable improvements in Justice and Virtue. In addition, the proposed Ethic-BERT attains 15.28% average accuracy improvement in the HardTest. These findings contribute to performance improvement and reliable decision-making using bias-aware preprocessing and proposed enhanced AI model.

**Keywords:** Ethical and non-ethical content classification, BERT, Deep learning, Ethical AI reasoning.

# 1 Introduction

Incorporating ethical reasoning into artificial intelligence systems is a vital step toward creating technology that is both responsible and aligned with societal values. As AI systems play an increasingly prominent role in mediating human interactions and making autonomous decisions, it is crucial to ensure they operate within the bounds of ethical principles [1, 2]. However, encoding the complexity of human moral reasoning into AI systems poses significant challenges, requiring innovative approaches to bridge the gap between human values and computational frameworks [3].

Ethical morality involves the principles of right and wrong that guide human behavior, encompassing dimensions such as justice, fairness, well-being, duties, and virtues. These principles are deeply interconnected, often leading to conflicts that require nuanced decision-making. Humans rely on cultural, social, and personal contexts to navigate moral ambiguities, but replicating this capacity in AI systems demands sophisticated techniques [4, 5]. The integration of ethical reasoning into AI is particularly important because of its potential societal impact [6]. AI systems, if left unchecked, can amplify biases, produce harmful outputs, or make decisions that conflict with shared human values [7]. To address these issues, researchers have turned to text-based scenarios as a means of evaluating AI systems’ ability to understand and apply ethical reasoning. These text based scenarios allows for the representation of complex moral dilemmas, providing a practical medium for assessing how AI aligns with human ethical judgment.

Recent advancements in NLP, particularly the development of transformer architectures, have made it possible to achieve significant progress in understanding context and intent in textual data. Datasets like ETHICS [8] leverage these advancements, presenting scenarios derived from philosophical theories, including justice, deontology, virtue ethics, utilitarianism, and commonsense morality. These benchmarks challenge AI systems to move beyond simple pattern recognition and address the moral complexity inherent in real-world scenarios. Despite these advancements, progress in embedding ethical reasoning into AI has been limited [9]. Models trained on ETHICS dataset [8] have shown some promise but struggle with nuanced scenarios and adversarial examples.

The key challenges of achieving better results in ethical reasoning tasks include:

- Lack of high-quality datasets that reduce ambiguity and enhance representativeness.
- Existing models struggle with nuanced ethical reasoning, limiting accuracy in moral decision-making.
- AI models rely on spurious correlations rather than deep moral reasoning, leading to misclassifications in complex ethical scenarios.
- The dataset primarily reflects Western moral perspectives, reducing its applicability to diverse cultural and ethical viewpoints.

In this research, we address these key challenges by the following key contributions:

1. Advanced model architectures utilizing state of the art transformer models and fine-tuning techniques to strengthen ethical reasoning capabilities.

2. Comprehensive evaluation demonstrating significant performance gains on standard and adversarially filtered test sets.
3. Enhanced dataset preparation to address ambiguities and introduce diverse contextual scenarios, improving data quality.

By advancing the interplay between ethical reasoning and AI, our work lays the groundwork for systems that are more aligned with human values and equipped to handle the complexities of real-world moral dilemmas.

## 2 Related Work

The increasing volume of digital content has created a pressing need for effective ethical content classification to manage misinformation, hate speech, and inappropriate material. Recent research efforts have focused on developing robust detection techniques using machine learning (ML) and deep learning (DL) models to improve accuracy, efficiency, and adaptability. At the same time, the ethical and moral implications of content have also become crucial, requiring models capable of analyzing not only spam content but also ethically unacceptable text. This review explores significant contributions in this field, focusing on advancements in detecting and mitigating harmful content while preserving contextual integrity.

In the domain of explicit content detection, Bhatti et al. [10] proposed an Explicit Content Detection (ECD) system targeting NSFW content using a residual network-based deep learning model. Their approach, integrating YCbCr color space and skin tone detection, achieved 95% accuracy in classifying explicit images and videos. Similarly, Khandekar et al. [11] focused on NLP techniques for detecting unethical and offensive text, leveraging LSTM and BiLSTM networks, which outperformed traditional models with an accuracy of 86.4%. Horne et al. [12] discussed the ethical challenges in automated fake news detection, emphasizing algorithmic bias and lack of generalizability. Their analysis of 381,000 news articles revealed the limitations of detection models that overfit benchmark datasets. Kiritchenko and Nejadgholi [13] introduced an “Ethics by Design” framework for abusive content detection, highlighting fairness, explainability, and bias mitigation. Their two-step process categorized identity-related content before assessing severity, reinforcing the importance of ethical considerations in content moderation. Schramowski et al. [14] examined the moral biases embedded in large pre-trained models like BERT, demonstrating how these biases could be leveraged to steer text generation away from toxicity. They identified a “moral direction” within the embedding space, which could be used to rate the normativity of text. This aligns with the ethical text assessment aspect of our research. Mittal et al. [15] conducted a comparative study on deep learning models for hate speech detection, concluding that fine-tuned RoBERTa models outperformed CNNs and BiLSTMs in a ternary classification system. Mnassri et al. [16] explored a multi-task learning framework that integrated emotional features with hate speech detection using BERT and mBERT. Their approach improved performance by leveraging shared representations across tasks, reducing overfitting and false positives. Saleh et al. [17] investigated the effectiveness of domain-specific word embeddings in hate speech detection, concluding that while specialized embeddings enhanced detection of coded hate

speech, pre-trained BERT models achieved the highest F1-score with 96%. Jim et al. [18] review advancements in sentiment analysis, highlighting machine learning, deep learning, and large language models. They explore applications, datasets, challenges, and future research directions to enhance performance.

Sultan et al. [19] analyzed shallow and deep learning techniques for cyberbullying detection across social media platforms. Their study, comparing six shallow learning algorithms with three deep models, found that BiLSTM models achieved the best recall and accuracy. This underscores the challenges in identifying masked or subtle offensive content, emphasizing the need for sophisticated models. Wadud et al. [20] examine methods for offensive text classification, emphasizing the need for improved multilingual detection. They introduce Deep-BERT, a model combining CNN and BERT, which enhances accuracy in identifying offensive content across different languages. Also, spam detection has been widely explored using traditional ML models and DL approaches. Similarly, Maqsood et al. [21] proposed a hybrid approach combining Random Forest, Multinomial Naive Bayes, and SVM with CNNs, observing that SVM outperformed other traditional ML models, while CNNs excelled on larger datasets. Guo et al. [22] introduced a BERT-based spam detection framework, integrating classifiers such as Logistic Regression, Random Forest, K-Nearest Neighbors, and SVM. Their results, using datasets like UCI's Spambase and the Kaggle Spam Filter Dataset, demonstrated that BERT significantly improved spam classification, achieving a precision of 97.86% and an F1-score of 97.84%. Meanwhile, Labonne and Moran [23] explored Large Language Models (LLMs) in spam detection, developing Spam-T5, a fine-tuned version of Flan-T5. Their work showed that Spam-T5 performed exceptionally well in low-data settings, surpassing both traditional ML models and modern LLMs like BERT. Chakraborty et al. [24] leveraged a fine-tuned BERT model with interval Type-2 fuzzy logic for sentiment classification, achieving superior performance in handling contextual variations. Similarly, Zhang et al. [25] integrated BERT with large LLMs in a hybrid approach, improving sentiment intensity prediction and aspect extraction. In the domain of ethical content detection, Aziz et al. [26] applied BERT with multi-layered graph convolutional networks to identify sentiment triplets, highlighting the model's capability in detecting hate speech and ethically sensitive content. These studies reinforce the adaptability of transformer-based architectures in capturing complex linguistic patterns and moral nuances, making them well-suited for ethical content classification.

Hendrycks et al. [7] introduced the ETHICS dataset [8] to evaluate AI models' ability to reason about morality across different ethical frameworks, including justice, virtue ethics, deontology, utilitarianism, and commonsense morality. Their findings indicate that pre-trained LLMs like BERT and RoBERTa show only partial success in making ethical decisions and often fail to handle complex moral scenarios accurately. Even advanced models like RoBERTa-large and ALBERT-xxlarge demonstrated low accuracy, particularly on adversarial test cases, highlighting their limitations in generalizing ethical principles. A key issue with the dataset is that the utilitarianism subset lacks explicit labels, requiring models to infer relative rankings rather than performing direct classification. Additionally, the study relied on standard fine-tuning techniques, but accuracy could likely improve with more extensive fine-tuning and

domain-specific training. These limitations suggest that current models still struggle to integrate ethical reasoning effectively.

Pre-trained LLMs have proven highly effective in understanding complex language and context for tasks like spam detection, hate speech recognition, offensive language identification, sentiment analysis, and other forms of harmful content detection. Their adaptability and precision make them well-suited for content moderation. Building on their success, this study applies pre-trained LLMs to ethical content classification, aiming for a more reliable, context-aware, and fair moderation system.

### 3 Methodology

This section explains the methodology used in our research to analyze ethical reasoning using machine learning techniques. It includes details about the dataset, data preprocessing, and implementation of our machine learning pipeline. Additionally, we elaborate on the fine-tuning process, showcasing the innovations that adapt the pre-trained BERT model for the task of ethical reasoning analysis, as illustrated in Figure 1.

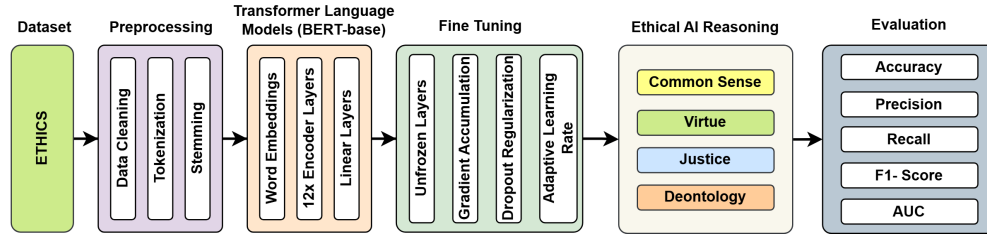


Fig. 1: Overview of the methodology for ethical reasoning

#### 3.1 Dataset

We used the ETHICS dataset [8], as the foundation for this research. The ETHICS dataset is publicly available at GitHub<sup>1</sup>. The dataset is specifically designed to evaluate ethical reasoning in text across four key domains: Justice, Virtue, Deontology, and Commonsense. Each text sample is annotated with a corresponding ethical category label.

The dataset is divided into three predefined splits: training, test, and hard test sets. To create a validation set, we further divided the original training set, allocating 80% for training and 20% for validation. Table 1 presents the distribution of data across these splits, with the adjusted training and validation sets reflecting this allocation. The dataset’s structure and balanced representation across ethical categories ensure a reliable benchmark for training, validation, and testing. The “Hard Test” dataset is an adversarially filtered dataset for ethical reasoning which allows us to evaluate the model’s robustness and ability to handle challenging ethical reasoning tasks.

<sup>1</sup><https://github.com/hendrycks/ethics>.

**Table 1: Data Distribution Across Splits**

Split	<i>Justice</i>	<i>Virtue</i>	<i>Deontology</i>	<i>Commonsense</i>
Training (80%)	17432	22596	14531	11128
Validation (20%)	4359	5649	3633	2782
Test	2704	4975	3596	3885
Hard Test	2052	4780	3536	3964

### 3.2 Adversarially filtered dataset (Hard Test)

The Hard Test dataset is a specialized subset of the ETHICS dataset, created to evaluate AI models’ ability to process complex ethical scenarios [8]. Unlike standard test sets, it filters out straightforward cases, ensuring models rely on deeper moral reasoning rather than statistical shortcuts. The dataset was developed using an adversarial filtration process, where models (Distil-BERT and Distil-RoBERTa) were first trained on a development set and then tested on various examples [7]. The most challenging cases were those with the highest prediction errors were selected while easier examples were removed, ensuring a rigorous evaluation benchmark. This dataset is widely applicable in AI ethics research, decision-making systems, and policy-driven AI, helping assess whether models can navigate nuanced ethical dilemmas. Researchers can fine-tune AI models on broader datasets and then use the Hard Test set to measure their ability to generalize moral principles and make fairer decisions. By focusing on difficult, context-rich scenarios, the Hard Test dataset advances research in ethical AI, contributing to the development of responsible and just decision-making systems [7].

### 3.3 Data Preprocessing

During the training of our BERT-based model, we employed an optimized data preprocessing pipeline to enhance generalization and mitigate data sparsity. These modifications were applied dynamically during training, ensuring efficient token representation, computational optimization, and contextual consistency.

#### 3.3.1 Text Normalization

Text normalization improves model training by reducing vocabulary sparsity, ensuring consistency, and removing noise. It standardizes text by handling case normalization, expanding contractions, and cleaning unnecessary characters. For each input sequence  $X$ , we applied a transformation function  $f : X \rightarrow X'$ , where  $X'$  represents the normalized text [27]. The transformation was defined as:

$$X' = f(X) = g_1(X) + g_2(X) + g_3(X) \quad (1)$$

In Equation 1,  $g_1(X)$  handled adaptive case normalization,  $g_2(X)$  expanded contractions (e.g., *can't* to *cannot*), and  $g_3(X)$  removed unnecessary special characters while preserving syntactic structure. By applying these transformations during training, we reduced vocabulary sparsity and improved token consistency across training epochs [28].

### 3.3.2 Tokenization Using WordPiece

WordPiece tokenization improves model training by handling unseen words, reducing vocabulary size, and enhancing embedding stability. It splits words into subwords based on frequency, preventing excessive fragmentation while maintaining meaningful representations. This helps models generalize better to rare and new words, making training more efficient and robust. In this process the input text was tokenized dynamically using BERT’s WordPiece algorithm [29]. Each tokenized word  $w$  was decomposed into subword tokens. In Equation 2,  $V$  is BERT’s fixed vocabulary. Instead of relying on standard segmentation, we employed frequency-aware tokenization, ensuring subwords were split efficiently based on their corpus occurrence. In Equation 3  $P(T \mid w)$  denotes the probability of a subword sequence given a word. This prevented excessive fragmentation of rare words and improved embedding stability. During training, this adjustment helped the model generalize better to unseen words [9].

$$T_w = \{t_1, t_2, \dots, t_n\}, \quad t_i \in V \quad (2)$$

$$T'_w = \arg \max_T P(T \mid w) \quad (3)$$

### 3.3.3 Truncation and Padding Optimization

Since BERT requires a fixed sequence length  $L$ , we dynamically truncated or padded input sequences during training. Padding was applied only when necessary. In Equation 4, a sequence  $S'$  is shorter than  $L$ , padding tokens ([PAD]) are appended. The exponent notation  $(L - |S'|)$  represents the number of padding tokens added to match the fixed length  $L$ . For example, if  $S'$  has 8 tokens but  $L = 12$ , then 4 [PAD] tokens are appended. To prevent overfitting due to excessive padding, we implemented batch-wise dynamic padding, which ensured that the sequence length  $L$  was adjusted based on the longest sequence in each batch. This minimized redundant [PAD] tokens, leading to faster training and reduced computational overhead [30, 31].

$$S' = S' + [\text{PAD}]^{(L - |S'|)} \quad (4)$$

## 3.4 Selecting a BERT-Based Cased Model

When selecting a model for AI ethical reasoning tasks, the choice must ensure accurate interpretation and context retention. A BERT-based cased model proposed by Devlin et al. [32], is particularly effective due to its ability to preserve case distinctions, which are often vital in formal and ethical text analysis. This ensures that proper nouns, legal terms, and acronyms retain their intended meanings, reducing ambiguity in ethical and policy analysis [33]. Research highlights the importance of case sensitivity in legal and ethical texts, as it helps differentiate between terms like “Title IX” and “title ix” or “US” and “us,” preventing misinterpretation. Case-sensitive models also enhance bias detection and policy evaluation by preserving textual integrity [34]. By leveraging this approach, we improve the accuracy and reliability of our ethical assessments.

### 3.5 Implementation Details

Our implementation is centered around a fine-tuned BERT-based cased model, chosen for its strong contextual understanding and adaptability to text classification tasks. In the following, we detail the architecture, training process, and fine-tuning innovations, along with the mathematical formulations underpinning these methods illustrated in Figure 2. Table 2 shows the customized hyperparameters and techniques employed during the fine-tuning process that ensured optimal performance on the ethical reasoning task.

#### 3.5.1 Model Architecture

We used a pre-trained BERT base cased model [32] and extended it with a classification head designed specifically for ethical reasoning tasks. This classification head consists of a fully connected layer.

$$\hat{Y} = \sigma(WH_L + b) \quad (5)$$

In Equation 5,  $H_L$  is the final hidden state of the transformer output.  $W$  and  $b$  are the weight matrix and bias vector of the classification head, and  $\sigma$  is the sigmoid activation function, converting logits to probabilities for binary classification. The resulting probabilities  $\hat{Y}$  represent the model’s confidence in each ethical reasoning category.

**Table 2:** Customized Hyperparameters and Techniques for BERT Fine-Tuning

Hyperparameter	Default	Modified
Model Architecture	BERT based cased	BERT based cased (full unfreezing)
Learning Rate	0.00002	0.00006
Learning Rate Scheduler	None	Adaptive Learning Rate Scheduler; $\eta_t = \eta_0 \cdot \frac{1}{\sqrt{t}}$
Max Sequence Length	512	128
Dropout Regularization	None	0.3
Gradient Accumulation	None	Accumulated over 4 mini-batches
Batch Size	16	32
Optimizer	AdamW	AdamW with Weight Decay

#### 3.5.2 Training Process

The model is trained to minimize a binary cross-entropy loss  $\mathcal{L}$ . In Equation 6,  $N$  is the number of samples,  $y_i$  is the true label for the  $i$ -th sample, and  $\hat{y}_i$  is the predicted probability for the  $i$ -th sample. The optimizer used is AdamW with a learning rate of  $\eta = 0.00006$  and applied weight decay. At each training step  $t$ , the parameters  $\theta$  are updated as Equation 6:

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)] \quad (6)$$

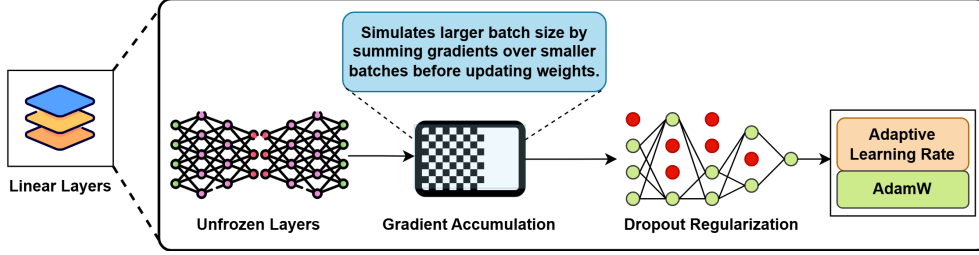


Fig. 2: Fine tuning of BERT for ethical reasoning

$$\theta^{(t+1)} = \theta^{(t)} - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t \quad (7)$$

Equation 7,  $\hat{m}_t$  and  $\hat{v}_t$  are bias-corrected estimates of the first and second moments of gradients, and  $\epsilon$  is a small constant for numerical stability.

### 3.5.3 Fine-Tuning Innovations

To maximize the model’s adaptability to the ethical reasoning task, we implemented the following innovations:

**Full Fine-Tuning:** All layers of the BERT model were unfrozen, allowing parameter adjustments across the entire network [35]. From Equation 8, the loss gradient  $\nabla_{\theta} \mathcal{L}$  was backpropagated through all layers [32] where,  $L$  is the number of transformer layers. Fully fine-tuning in ethical classification tasks helps the model grasp domain-specific ethical nuances, leading to more precise and fair decisions. It refines the model’s understanding beyond general pre-trained knowledge, aligning it with ethical guidelines. This approach minimizes bias, enhances reliability, and ensures responsible decision-making.

$$\frac{\partial \mathcal{L}}{\partial \theta_i} = \frac{\partial \mathcal{L}}{\partial H_L} \cdot \prod_{j=i+1}^L \frac{\partial H_j}{\partial H_{j-1}} \cdot \frac{\partial H_i}{\partial \theta_i}, \quad i = 1, 2, \dots, L \quad (8)$$

**Gradient Accumulation:** To address memory constraints, gradient accumulation was employed. Gradients  $g^{(b)}$  for each mini-batch  $b$  were accumulated over  $N_{\text{acc}}$  steps [36]. Equation 9, gradients  $\nabla \mathcal{L}^{(b)}$  from each mini-batch  $b$  are summed over  $N_{\text{acc}}$  steps. This method allows training with small mini-batches while effectively simulating a larger batch size. Equation 10 updates the model parameters  $\theta$  after accumulating gradients over multiple steps. The learning rate  $\eta$  scales the average accumulated gradient, ensuring stable optimization. It ensures better representation of ethical considerations in data while maintaining computational feasibility. The approach helps to mitigate biases and enhances fairness by enabling effective learning from smaller yet diverse datasets.

$$g_{\text{acc}} = \sum_{b=1}^{N_{\text{acc}}} \nabla_{\theta} \mathcal{L}^{(b)} \quad (9)$$

$$\theta^{(t+1)} = \theta^{(t)} - \eta \frac{g_{\text{acc}}}{N_{\text{acc}}} \quad (10)$$

**Adaptive Learning Rate:** An adaptive learning rate schedule was used, reducing the learning rate as training progressed [9]. In Equation 11,  $\eta_0$  is the initial learning rate, and  $t$  is the training step. Applying an adaptive learning rate in model training dynamically adjusts the step size based on gradient variations, improving convergence speed and stability. This technique helps prevent overshooting in high-gradient regions while accelerating learning in flatter areas, leading to more efficient optimization. In ethical classification tasks, adaptive learning rates enhance fairness and robustness by ensuring balanced learning across diverse and sensitive data distributions.

$$\eta_t = \eta_0 \cdot \frac{1}{\sqrt{t}} \quad (11)$$

#### 3.5.4 Regularization and Robustness

Dropout regularization was applied in the classification head to mitigate overfitting. During training [37], activations  $H_{\text{task}}$  in the classification head were stochastically zeroed out. In Equation 12,  $D \sim \text{Bernoulli}(1 - p)$  is a dropout mask,  $\odot$  represents element-wise multiplication, and  $p$  is the dropout rate. At inference time, activations were scaled by  $(1 - p)$  to maintain consistent output expectations shown in Equation 13. Regularization techniques, dropout, and batch normalization, help prevent overfitting by constraining model complexity. These methods ensure that the model generalizes well to unseen ethical scenarios, reducing biases and improving fairness. In ethical classification tasks, regularization enhances robustness by making the model resilient to noisy or imbalanced data, leading to more reliable and ethically sound decisions.

$$H'_{\text{task}} = D \odot H_{\text{task}} \quad (12)$$

$$H_{\text{task}}^{\text{inference}} = (1 - p) \cdot H_{\text{task}} \quad (13)$$

### 3.6 Evaluation Matrix

The model’s performance was evaluated using accuracy, precision, recall, F1-score [38], and AUC [39] ensuring robust validation of ethical reasoning capabilities.

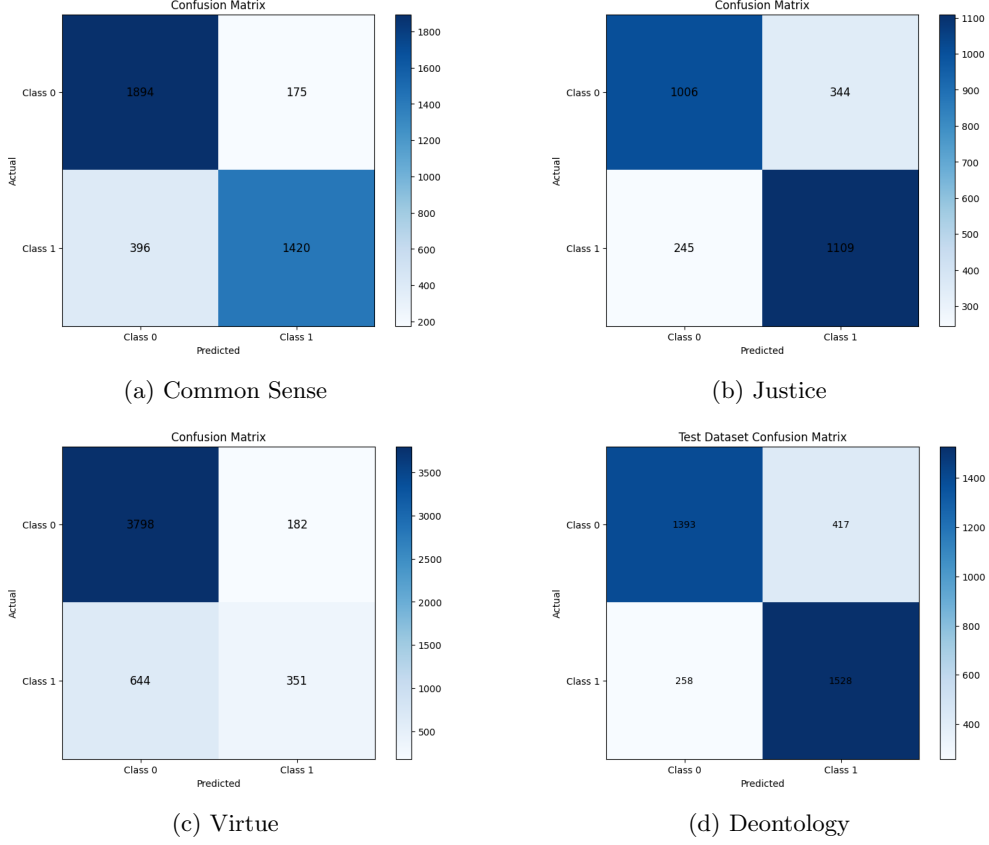
## 4 Results and Discussion

This section presents the evaluation of the fine-tuned BERT model on the Test Split and Hard Test Split datasets, along with a comparison to existing models such as RoBERTa-large and ALBERT-xxlarge. The results demonstrate the impact of our innovative fine-tuning and preprocessing techniques in delivering strong performance, particularly in specific ethical reasoning domains.

### 4.1 Performance on Test Split

Table 3 shows the performance of the proposed BERT model on the Test Split. The model achieved high Accuracy in Commonsense, Justice, Virtue domains and Deontology, reaching 86.46%, 78.22%, 83.40%, and 81.23% respectively. These results highlight

the model’s ability to effectively adapt to the task in these domains. The AUC values for domains—90.78, 87.36, 88.78, 89.93—further affirm the model’s capability to separate positive and negative classes accurately. The confusion matrices for the Test dataset are presented in Figure 3.



**Fig. 3:** Confusion matrix of Test dataset.

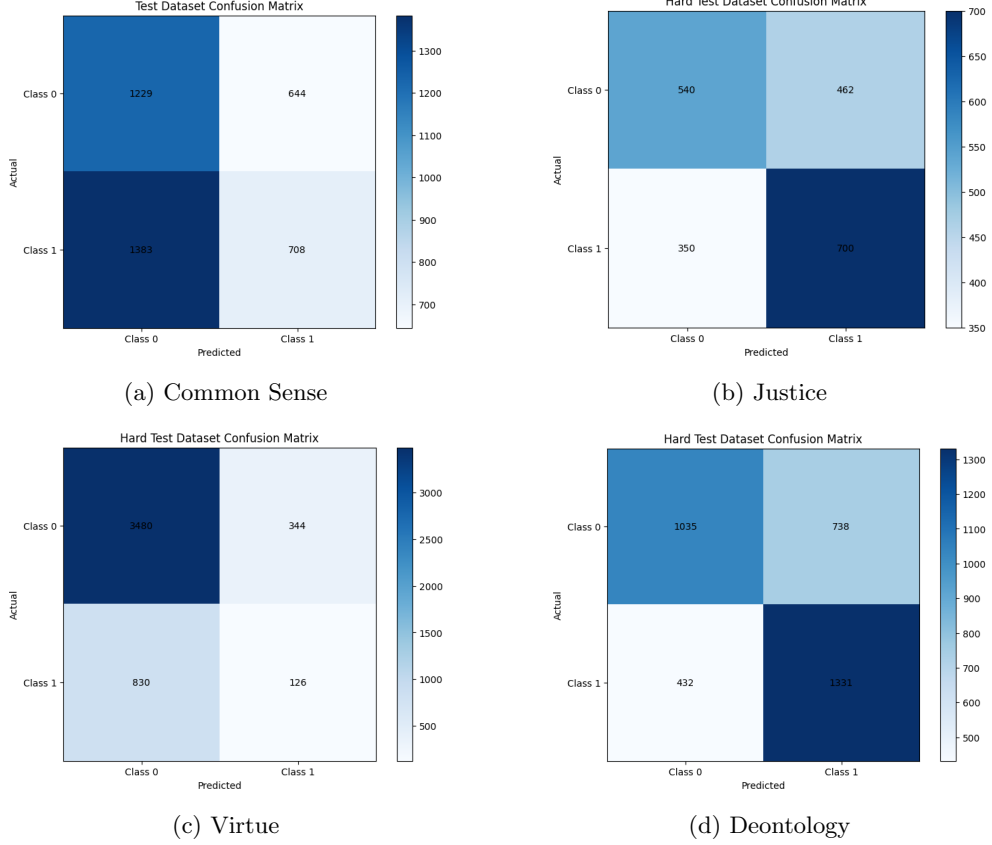
The subfigures 3a, 3b, 3c, and 3d represent the results of the Common Sense, Justice, Virtue, and Deontology frameworks, respectively, demonstrating variations in the performance of the model in different paradigms of ethical reasoning.

## 4.2 Performance on Hard Test Split

The Hard Test Split, designed to assess robustness on more challenging examples, revealed further insights into the model’s strengths and weaknesses, as summarized in Table 4. The Accuracy in Commonsense domain dropped to 50.00%, indicating difficulty in handling complex or ambiguous cases. This suggests a need for additional strategies to improve generalization in these domains.

**Table 3:** Performance of BERT over Test Split Data

Sub-set	Accuracy	Precision	Recall	F1-score	AUC
Commonsense	86.46	86.00	85	85	90.78
Justice	78.22	78	78	78	87.36
Virtue	83.40	82	83	81	88.78
Deontology	81.23	78.56	85.55	81.00	89.93

**Fig. 4:** Confusion matrix of Hard Test dataset.

In contrast, the model maintained stronger performance in the Deontology and Virtue domains, achieving 66.91% and 75.44% Accuracy, respectively. The AUC values of 73.35 for Deontology and 80.24 for Virtue further underscore the model’s resilience. These results highlight the role of our preprocessing and training techniques in ensuring that the model retains its effectiveness even with harder data. The confusion matrices for the Hard Test dataset are presented in Figure 4.

The subfigures 4a, 4b, 4c, and 4d correspond to the Common Sense, Justice, Virtue, and Deontology frameworks, respectively, highlighting differences in model performance across ethical reasoning approaches. The Figure 5 illustrate model performance over five epochs, highlighting key trends. Training loss consistently decreases, while accuracy improves, indicating effective learning. Some models maintain stable validation accuracy, suggesting good generalization. In some cases, training and validation loss patterns differ, which may indicate areas for refinement. Adjustments like regularization could improve performance. Overall, the models demonstrate effective learning, with some showing stronger generalization.

**Table 4:** Performance of BERT over Hard-Test Split Data

Sub-set	Accuracy	Precision	Recall	F1-score	AUC
Commonsense	50.00	50.00	49	48	50.26
Justice	60.40	60.24	66.67	60	65.39
Virtue	75.44	70	75	72	80.24
Deontology	66.91	64.33	75.50	67.00	73.35

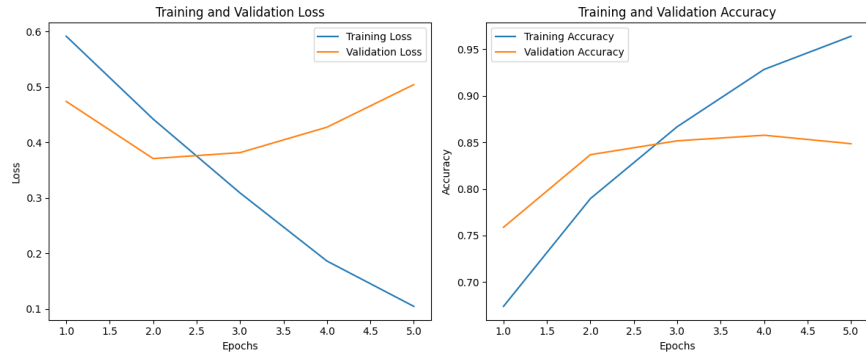
### 4.3 Comparison with Existing Models

Tables 5 and 6 compare the proposed BERT model’s performance with other models like RoBERTa-large and ALBERT-xxlarge. On the Test Split, our model achieved an average Accuracy of 82.32%, significantly outperforming where the baseline BERT-base got 46.1% and RoBERTa-large 65.1% while nearly beating ALBERT-xxlarge 68.3%. Notably, the proposed Ethic-BERT excelled in Justice, Virtue and Deontology domains, achieving 78.22%, 83.40% and 81.23% Accuracy, respectively.

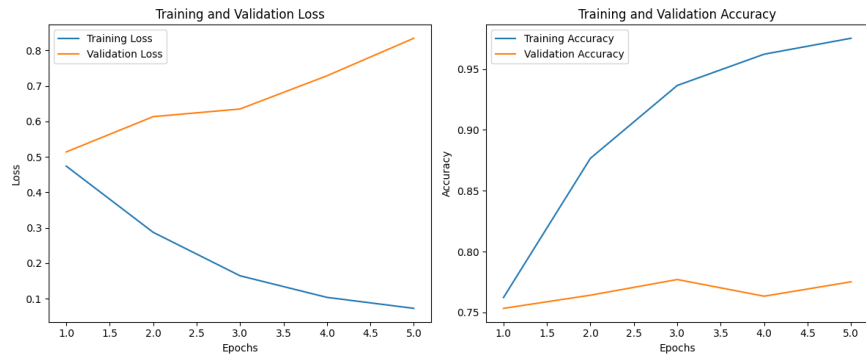
On the Hard Test Split, the proposed Ethic-BERT achieved an average Accuracy of 63.18%, surpassing BERT-base 16.8%, RoBERTa-large 39.42%, and ALBERT-xxlarge 43.05%. The model’s performance in Justice, Virtue and Deontology domains remained robust, with Accuracy values of 60.40%, 75.44% and 66.91%. These improvements suggest that our approach enables the model to perform consistently better across harder examples compared to other models.

**Table 5:** Comparison of Test Results with Existing Models

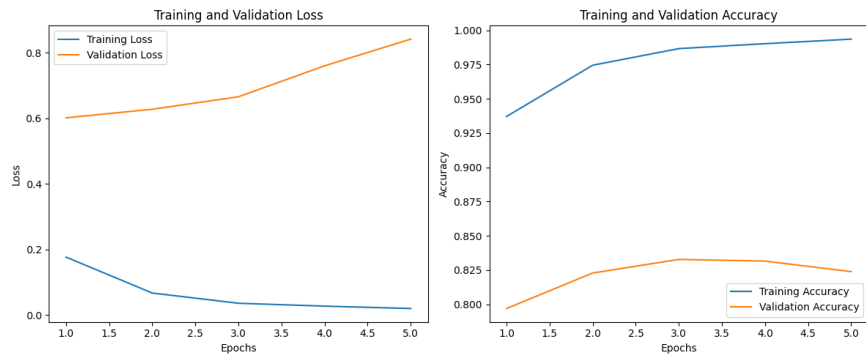
Source	Model	Common Sense	Justice	Virtue	Deontology	Average
Hendrycks et al. [7]	BERT-base	86.5	26.0	33.1	38.8	46.1
	RoBERTa-large	90.4	56.7	53.0	60.3	65.1
	ALBERT-xxlarge	85.1	59.9	64.1	64.1	68.3
<b>Proposed</b>	<b>Ethic-BERT</b>	<b>86.46</b>	<b>78.22</b>	<b>83.40</b>	<b>81.23</b>	<b>82.32</b>



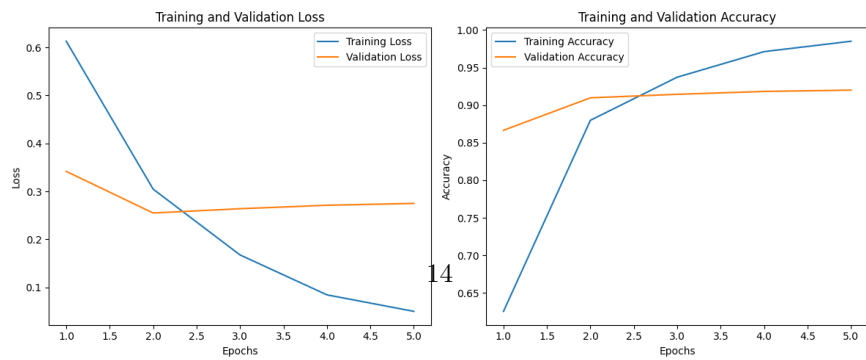
(a) Common Sense



(b) Justice



(c) Virtue



(d) Deontology

**Fig. 5:** Training vs validation loss and accuracy curves on training dataset.

**Table 6:** Comparison of Hard Test Results with Existing Models

Source	Model	Common Sense	Justice	Virtue	Deontology	Average
Hendrycks et al. [7]	BERT-base	48.7	7.6	8.6	10.3	16.8
	RoBERTa-large	63.4	38.0	25.5	30.8	39.42
	ALBERT-xxlarge	59.0	38.2	37.8	37.2	43.05
<b>Proposed</b>	<b>Ethic-BERT</b>	<b>50.00</b>	<b>60.40</b>	<b>75.44</b>	<b>66.91</b>	<b>63.18</b>

#### 4.4 Contributions of Fine-Tuning and Preprocessing

The superior results of the proposed model are attributed to the combination of comprehensive fine-tuning and a robust preprocessing pipeline. Unlike traditional approaches that partially freeze pre-trained layers, our method unfreezes all layers of the BERT model, allowing it to fully adapt its representations to the specific nuances of ethical reasoning. This approach ensures that the model effectively incorporates both pre-trained knowledge and task-specific patterns.

The preprocessing pipeline was another critical factor in achieving these results. By employing advanced tokenization, consistent input formatting through truncation and padding, the pipeline enhanced the quality and diversity of the training data. Additionally, the use of gradient accumulation allowed for stable training even with limited resources, optimizing learning efficiency. The balanced data splits further ensured that all ethical reasoning domains were well-represented during training, contributing to strong results in Justice and Virtue categories.

### 5 Challenges and Future Research Directions

Despite the promising results, challenges remain in the deontology and common sense domains, especially in the hard-test split. The low accuracy and AUC in these domains indicate the need for additional strategies, such as augmenting the dataset with richer context or leveraging external knowledge sources. Incorporating domain-specific pretraining or multitask learning approaches may also help the model capture the unique characteristics of these categories. Although, the proposed Ethic-BERT model demonstrates substantial advancements in ethical reasoning classification, achieving state-of-the-art performance in several domains, it can further be improved which is outlined as the future research direction in this topic. The effectiveness of fine-tuning strategies highlights the potential for further improvements in AI systems designed for complex reasoning tasks.

### 6 Conclusion

In this study, we fine-tuned a BERT-based model to classify ethical reasoning across four domains: Common sense, justice, virtue, and deontology. The model demonstrated strong performance, particularly in Justice and Virtue reasoning, where it surpassed

existing models. In the Hard Test Split, the model showcased its robustness, achieving improvements over baseline approaches in more challenging scenarios. Our approach introduced key innovations, including comprehensive fine-tuning by unfreezing all layers of the BERT model, implementing gradient accumulation, and utilizing advanced tokenization and data augmentation. These techniques allowed the model to effectively combine its pretrained knowledge with task-specific adaptations, resulting in superior performance. Despite these successes, challenges persist in the Commonsense and Deontology domains, especially on the Hard Test Split. Addressing these gaps could involve enriching the training data with more contextually diverse examples, incorporating external knowledge sources, or adopting domain-specific pretraining strategies. In general, this work highlights the potential of fine-tuned transformer models to tackle complex reasoning tasks in AI ethics. The findings underscore the importance of thoughtful preprocessing and training techniques in improving the robustness and generalization of the model.

## Acknowledgments

The authors would like to express their sincere gratitude to the Ubiquitous, Cloud, and Human-Computer Interaction (UCH) Research Group, Department of Computer Science, American International University-Bangladesh for supporting this research.

## References

- [1] Woodgate, J., Ajmeri, N.: Macro ethics principles for responsible ai systems: Taxonomy and directions. *ACM Comput. Surv.* **56**(11) (2024) <https://doi.org/10.1145/3672394>
- [2] Sholla, S., Reshi, I.A.: Ethical reasoning in technology: using computational approaches to integrate ethics into ai systems. *Journal of Information, Communication and Ethics in Society* (2024) <https://doi.org/10.1108/JICES-03-2024-0024>
- [3] Li, F., Ruijs, N., Lu, Y.: Ethics & ai: A systematic review on ethical concerns and related strategies for designing with ai in healthcare. *Ai* **4**(1), 28–53 (2022) <https://doi.org/10.3390/ai4010003>
- [4] Akbar, M.A., Khan, A.A., Mahmood, S., Rafi, S., Demi, S.: Trustworthy artificial intelligence: A decision-making taxonomy of potential challenges. *Software: Practice and Experience* **54**(9), 1621–1650 (2024) <https://doi.org/10.1002/spe.3216>
- [5] Han, S., Kelly, E., Nikou, S., Svee, E.-O.: Aligning artificial intelligence with human values: reflections from a phenomenological perspective. *AI & SOCIETY*, 1–13 (2022) <https://doi.org/10.1007/s00146-021-01247-4>
- [6] Hauer, T.: Importance and limitations of ai ethics in contemporary society. *Humanities and Social Sciences Communications* **9**(1), 1–8 (2022) <https://doi.org/10.1007/s43545-022-00451-1>

[org/10.1057/s41599-022-01300-7](https://doi.org/10.1057/s41599-022-01300-7)

- [7] Hendrycks, D., Burns, C., Basart, S., Critch, A., Li, J., Song, D., Steinhardt, J.: Aligning AI With Shared Human Values (2023). <https://arxiv.org/abs/2008.02275>
- [8] Hendrycks, D.: ETHICS: A Benchmark for Ethical Understanding of AI. GitHub repository (2020). <https://github.com/hendrycks/ethics>
- [9] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. CoRR (2017). <http://arxiv.org/abs/1706.03762>
- [10] Bhatti, A.Q., Umer, M., Adil, S.H., Ebrahim, M., Nawaz, D., Ahmed, F.: Explicit content detection system: An approach towards a safe and ethical environment. *Applied Computational Intelligence and Soft Computing* **2018**, 1463546–13 (2018) <https://doi.org/10.1155/2018/1463546>
- [11] Khandekar, A., Hema, C.D., Meghana, A., Mounika, A., Vaishnavi, V.S.: Nlp based analysis and detection of unethical text. In: 2023 International Conference on Sustainable Computing and Data Communication Systems (ICSCDS), pp. 622–626 (2023). <https://doi.org/10.1109/ICSCDS56580.2023.10104943>
- [12] Benjamin D. Horne, D.N., Smith, S.L.: Ethical and safety considerations in automated fake news detection. *Behaviour & Information Technology* **0**(0), 1–22 (2023) <https://doi.org/10.1080/0144929X.2023.2285949>
- [13] Kiritchenko, S., Nejadgholi, I.: Towards Ethics by Design in Online Abusive Content Detection (2020). <https://arxiv.org/abs/2010.14952>
- [14] Schramowski, P., Turan, C., Andersen, N., Rothkopf, C.A., Kersting, K.: Large pre-trained language models contain human-like biases of what is right and wrong to do. *Nature Machine Intelligence* **4**(3), 258–268 (2022) <https://doi.org/10.1038/s42256-022-00458-8>
- [15] Mittal, U.: Detecting hate speech utilizing deep convolutional network and transformer models. In: 2023 International Conference on Electrical, Electronics, Communication and Computers (ELEXCOM), pp. 1–4 (2023). <https://doi.org/10.1109/ELEXCOM58812.2023.10370502>
- [16] Mnassri, K., Rajapaksha, P., Farahbakhsh, R., Crespi, N.: Hate speech and offensive language detection using an emotion-aware shared encoder. In: ICC 2023 - IEEE International Conference on Communications, pp. 2852–2857 (2023). <https://doi.org/10.1109/ICC45041.2023.10279690>
- [17] Saleh, H., Alhothali, A., Moria, K.: Detection of hate speech using bert and hate speech word embedding with deep model. *Applied Artificial Intelligence* **37**(1),

2166719 (2023) <https://doi.org/10.1080/08839514.2023.2166719>

- [18] Jim, J.R., Talukder, M.A.R., Malakar, P., Kabir, M.M., Nur, K., Mridha, M.F.: Recent advancements and challenges of nlp-based sentiment analysis: A state-of-the-art review. *Natural Language Processing Journal* **6**, 100059 (2024) <https://doi.org/10.1016/j.nlp.2024.100059>
- [19] Sultan, D., Toktarova, A., Zhumadillayeva, A., Aldeshov, S., Mussiraliyeva, S., Beissenova, G., Tursynbayev, A., Baenova, G., Imanbayeva, A.: Cyberbullying-related hate speech detection using shallow-to-deep learning. *Computers, Materials & Continua* **74**(1), 2116–2130 (2023) <https://doi.org/10.32604/cmc.2023.032993>
- [20] Wadud, M.A.H., Mridha, M.F., Shin, J., Nur, K., Saha, A.K.: Deep-bert: Transfer learning for classifying multilingual offensive texts on social media. *Computer Systems Science and Engineering* **44**(2), 1775–1791 (2023) <https://doi.org/10.32604/csse.2023.027841>
- [21] Maqsood, U., Ur Rehman, S., Ali, T., Mahmood, K., Alsaedi, T., Kundi, M.: An intelligent framework based on deep learning for sms and e-mail spam detection. *Applied Computational Intelligence and Soft Computing* **2023**(1), 6648970 (2023) <https://doi.org/10.1155/2023/6648970> <https://onlinelibrary.wiley.com/doi/pdf/10.1155/2023/6648970>
- [22] Guo, Y., Mustafaoglu, Z., Koundal, D.: Spam detection using bidirectional transformers and machine learning classifier algorithms. *Journal of Computational and Cognitive Engineering* **2**(1), 5–9 (2022) <https://doi.org/10.47852/bonviewJCCE2202192>
- [23] Labonne, M., Moran, S.: Spam-T5: Benchmarking Large Language Models for Few-Shot Email Spam Detection (2023). <https://doi.org/10.48550/arXiv.2304.01238>
- [24] Chakraborty, K., Bhattacharyya, S., Bag, R., Mršić, L.: Sentiment analysis on labeled and unlabeled datasets using bert architecture. *Soft computing* **28**(15), 8623–8640 (2024) <https://doi.org/10.1007/s00500-023-08876-5>
- [25] Zhang, Y., Xu, H., Zhang, D., Xu, R.: A hybrid approach to dimensional aspect-based sentiment analysis using bert and large language models. *Electronics* **13**(18), 3724 (2024) <https://doi.org/10.3390/electronics13183724>
- [26] Aziz, K., Ji, D., Chakrabarti, P., Chakrabarti, T., Iqbal, M.S., Abbasi, R.: Unifying aspect-based sentiment analysis bert and multi-layered graph convolutional networks for comprehensive sentiment dissection. *Scientific Reports* **14**(1), 14646 (2024) <https://doi.org/10.1038/s41598-024-61886-7>
- [27] Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li,

- W., Liu, P.J.: Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research* **21**(140), 1–67 (2020)
- [28] Zhang, Y., Bartley, T.M., Graterol-Fuenmayor, M., Lavrukhin, V., Bakhturina, E., Ginsburg, B.: A chat about boring problems: Studying gpt-based text normalization. In: *ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 10921–10925 (2024). <https://doi.org/10.1109/ICASSP48485.2024.10447169>
- [29] Schönle, D., Reich, C., Abdeslam, D.O.: Linguistic-aware wordpiece tokenization: Semantic enrichment and oov mitigation. In: *2024 6th International Conference on Natural Language Processing (ICNLP)*, pp. 134–142 (2024). <https://doi.org/10.1109/ICNLP60986.2024.10692355>
- [30] Tunstall, L., Von Werra, L., Wolf, T.: *Natural Language Processing with Transformers*. O’Reilly Media, Inc., Sebastopol, CA (2022). ISBN: 978-1098136789
- [31] Sun, M., Hameed, I.A., Wang, H., Pasquine, M.: Skip truncation for sentiment analysis of long review information based on grammatical structures. In: Hassanien, A.E., Rizk, R.Y., Snášel, V., Abdel-Kader, R.F. (eds.) *The 8th International Conference on Advanced Machine Learning and Technologies and Applications (AMLT2022)*, pp. 298–308. Springer, Cham (2022). [https://doi.org/10.1007/978-3-031-03918-8\\_27](https://doi.org/10.1007/978-3-031-03918-8_27)
- [32] Devlin, J., Chang, M.-W., Lee, K., Toutanova, K.: BERT: Pre-training of deep bidirectional transformers for language understanding. In: Burstein, J., Doran, C., Solorio, T. (eds.) *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)*, pp. 4171–4186. Association for Computational Linguistics, Minneapolis, Minnesota (2019). <https://doi.org/10.18653/v1/N19-1423>
- [33] Rogers, A., Kovaleva, O., Rumshisky, A.: A primer in bertology: What we know about how bert works. *Transactions of the Association for Computational Linguistics* **8**, 842–866 (2021) <https://doi.org/10.1162/tacl.a.00349>
- [34] Weidinger, L., Uesato, J., Rauh, M., Griffin, C., Huang, P.-S., Mellor, J., Glaese, A., Cheng, M., Balle, B., Kasirzadeh, A., Biles, C., Brown, S., Kenton, Z., Hawkins, W., Stepleton, T., Birhane, A., Hendricks, L.A., Rimell, L., Isaac, W., Haas, J., Legassick, S., Irving, G., Gabriel, I.: Taxonomy of risks posed by language models. In: *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency. FAccT ’22*, pp. 214–229. Association for Computing Machinery, New York, NY, USA (2022). <https://doi.org/10.1145/3531146.3533088>
- [35] Howard, J., Ruder, S.: Universal language model fine-tuning for text classification. In: Gurevych, I., Miyao, Y. (eds.) *Proceedings of the 56th Annual Meeting of the*

Association for Computational Linguistics (Volume 1: Long Papers), pp. 328–339. Association for Computational Linguistics, Melbourne, Australia (2018). <https://doi.org/10.18653/v1/P18-1031>

- [36] Reddi, S.J., Kale, S., Kumar, S.: On the convergence of adam and beyond. CoRR **abs/1904.09237** (2019) [1904.09237](https://arxiv.org/abs/1904.09237)
- [37] Wager, S., Wang, S., Liang, P.S.: Dropout training as adaptive regularization. In: Advances in Neural Information Processing Systems, vol. 26, pp. 351–359. Curran Associates, Inc., Red Hook, NY (2013). <https://arxiv.org/abs/1307.1493>
- [38] Powers, D.M.W.: Evaluation: from precision, recall and f-measure to roc, informedness, markedness and correlation. CoRR **abs/2010.16061** (2020) [2010.16061](https://arxiv.org/abs/2010.16061)
- [39] Namdar, K., Haider, M.A., Khalvati, F.: A modified auc for training convolutional neural networks: Taking confidence into account. Frontiers in Artificial Intelligence **4** (2021) <https://doi.org/10.3389/frai.2021.582928>