

# Detecting AI-Generated Essays Using NLP and Text Classification

Conference Paper

Sara Alsaghier  
444008763

Ghala Alnahidh  
444008688

Jana Albeladi  
444008758

Lena Alshammari  
444008703

Teef Alzahrani  
444008687

444008763@pnu.edu.sa 444008688@pnu.edu.sa 444008758@pnu.edu.sa 444008703@pnu.edu.sa 444008687@pnu.edu.sa

**Abstract**—This study addresses the challenge of distinguishing between human-written and AI-generated text, a growing concern in educational environments due to the increasing use of generative AI tools. The project applies Natural Language Processing techniques and Fine-Tuning on three different models: LSTM as a sequence-based architecture, and DistilBERT and ELECTRA Small as modern transformer-based models capable of capturing contextual and linguistic patterns. A balanced dataset of human and AI-generated text was prepared and processed using standard text-cleaning techniques to enhance data quality prior to model training. The aim of this work is to explore how these different model families handle the task of text classification and to establish a technical foundation that supports academic integrity and promotes responsible use of artificial intelligence in education.

**Index Terms**—Machine Learning, Deep Learning, Natural Language Processing, Text Classification, AI-Generated Content Detection, Artificial Intelligence

## I. INTRODUCTION

Recent studies show that nearly 86% of higher-education students regularly use AI tools for academic tasks, with over half of them using these tools at least weekly. This rapid adoption of generative AI has transformed how students write and learn, but it also raises serious concerns about academic integrity and fairness in evaluation. According to the UNESCO (2023) guidance report, responsible use of artificial intelligence in education is essential to maintain trust in assessment systems and ensure learning quality, aligning directly with Sustainable Development Goal 4 (Quality Education) [1].

The primary issue addressed in this study is the increasing difficulty faced by educational institutions in accurately distinguishing between human-written and AI-generated essays, which can result in unfair grading or misuse of generative tools, undermining genuine learning and student effort. To address this challenge, the objective of this project is to achieve an accuracy exceeding 80% in distinguishing human-written essays from AI-generated essays, thereby providing a more reliable framework that supports academic integrity.

Building upon this challenge, the study is guided by the following research question:

*“To what extent can deep learning models accurately differentiate between human-written and*

*AI-generated essays when trained on the same dataset?”*

In the broader research context of AI-text detection, the study by Yang et al. [2] on Text Graph Neural Networks highlights significant performance gaps in existing tools, noting that many fail to sufficiently incorporate semantic and structural text features, thus limiting their discrimination accuracy between human and AI-generated content. Motivated by these gaps, this paper presents a structured framework for evaluating whether an essay was written by a human or by AI. The framework aims to support academic integrity by improving detection accuracy, helping instructors identify AI-generated content effectively, encouraging students to rely on their real writing and analytical skills rather than overusing AI tools, and contributing to research by analyzing error cases and improving detection policies.

This project directly contributes to SDG 4 (Quality Education) by enhancing fairness, transparency, and integrity in learning. Strengthening AI-text detection supports equitable academic evaluation, reduces opportunities for academic misconduct, and ensures that assessments reflect genuine student capabilities rather than reliance on automated writing tools. This alignment also supports national priorities such as the Human Capability Development Program within Saudi Vision 2030.

Ethical considerations were also addressed. The study uses only publicly available, non-sensitive data and does not involve any personal or identifiable student information. The system is intended to support instructors rather than replace human judgment, and all model outputs will be communicated transparently to avoid over-reliance on automated decisions.

The paper is structured as follows: Section 2 reviews related work and identifies research gaps; Section 3 describes the proposed methodology and preprocessing pipeline; Section 4 details the experimental design, model tuning, and performance evaluation; Section 5 discusses results and compares them to state-of-the-art approaches; and Section 6 concludes with insights, limitations, and recommendations, emphasizing the practical and ethical implications for quality education and academic integrity.

## II. RELATED WORK

In [5], Blake et al. (2023) proposed a hybrid deep learning model for the detection of AI-generated texts. To capture both local and long-range dependencies in text, the model employs a combination of convolutional layers, part-of-speech (POS) tagging, Bidirectional Long Short-Term Memory (Bi-LSTM) networks, and an attention mechanism. Experimental results on benchmark datasets demonstrated detection accuracies of 85% and 88% in comparison to conventional models. The study highlighted that integrating Bi-LSTM with attention mechanisms enhances the precision and effectiveness of identifying content generated by artificial intelligence.

In [6], Gehrmann et al. (2023) studied stylometric fingerprinting detection by analyzing linguistic patterns in sentence length distribution and rare word appearance and syntactic tree depth. The researchers used GLTR (Giant Language Model Test Room) to display token probability distributions and analyze their deviation from human writing patterns. When applied to essay datasets, it showed that AI-generated texts tend to exhibit unnaturally consistent probabilities and limited lexical diversity. The system uses probability distributions instead of deep neural networks but functions as an easy-to-understand educational tool for detecting essays.

In [7], The 2024 Prova study combined Natural Language Processing (NLP) methods with traditional Machine Learning approaches to detect AI-generated text. The research tested three models consisting of Support Vector Machines (SVM) and XGBoost and transformer-based models including BERT to establish their ability to detect human-written content from machine-generated content. The results demonstrated that SVM and XGBoost achieved 81% and 84% accuracy but the BERT-based model reached 93% accuracy. The transformer model detects AI-generated content through its capacity to analyze contextual information which helps it identify particular linguistic indicators. The research demonstrates that deep contextual embeddings when used with classical ML classifiers produce the best automated text authenticity detection results because they improve both detection precision and model understanding.

In [9], Bharathikumar et al. proposed a DistilBERT-based framework for distinguishing AI-generated text from human writing. Their study used a large dataset of 500,000 labeled essays and compared several traditional machine-learning models (such as Logistic Regression and Decision Trees) as well as deep-learning models (including RNN and LSTM) against transformer models. After extensive preprocessing and feature extraction, their findings showed that DistilBERT achieved the highest performance, reaching an accuracy of nearly 98%. The authors concluded that transformer-based architectures capture deeper contextual and stylistic patterns than classical methods, making them highly effective for detecting AI-generated content.

In [10], a Bi-LSTM-based deep learning classifier was proposed to distinguish between human-authored and AI-generated text. The study utilized a large corpus of 180,000

mixed essays, containing content from GPT-2, GPT-3, and human writers. The model leverages Bidirectional Long Short-Term Memory networks to capture forward and backward contextual dependencies within a sentence, enabling the detection of subtle semantic inconsistencies that often appear in machine-generated text. Experimental evaluation demonstrated that the Bi-LSTM architecture achieved an accuracy of 90%, outperforming classical machine-learning methods such as SVM and Logistic Regression. The results highlight that recurrent neural architectures remain effective for semantic-level detection tasks, particularly when distinguishing nuanced stylistic differences between human and AI writing.

in [11] introduced an ELECTRA-Small Transformer model for detecting AI-generated text across news articles and academic essays. Utilizing a dataset of approximately 150,000 samples, the model was trained using ELECTRA's discriminative pretraining strategy, where the network learns to identify "replaced tokens" generated by a small generator model. This pretraining mechanism allows ELECTRA to learn highly sensitive linguistic cues that differentiate human writing from machine-produced content. The authors reported detection accuracies of 94% for the Small version and 97% for the Base model, significantly surpassing traditional RNN-based and CNN-based baselines. The findings demonstrate the effectiveness of ELECTRA architectures in text forensics and highlight their superior efficiency compared to larger transformer models.

Overall, the reviewed literature demonstrates a clear progression in AI-generated text detection methods, evolving from hybrid deep learning architectures to stylometric probability analysis and transformer-based models. Studies such as [5] show that combining CNNs, POS-tagging, Bi-LSTM layers, and attention mechanisms significantly enhances a model's ability to capture both local and long-range dependencies in text, achieving accuracies up to 88%. Stylometric approaches, exemplified by [6] through the GLTR framework, reveal that AI-generated content exhibits abnormal probability distributions and reduced lexical diversity, offering a lightweight and interpretable alternative to deep neural networks. Meanwhile, [8] demonstrates the advantage of contextual embeddings, where BERT outperforms classical machine-learning models such as SVM and XGBoost, reaching 93% accuracy. More recent contributions, including [9], [10], and [11], show that transformer-based and discriminative architectures deliver superior performance, with DistilBERT achieving 98%, Bi-LSTM classifiers reaching 90%, and ELECTRA models obtaining up to 97% accuracy.

To provide a structured comparison of these approaches, all referenced studies have been systematically summarized in Table I, highlighting key differences in model design, dataset size, performance, and methodological contributions.

TABLE I  
SUMMARY OF RELATED WORK ON AI-GENERATED TEXT DETECTION

Ref	Model	Dataset	Accuracy	Summary
[5](2023)	Hybrid CNN, POS, Bi-LSTM, Attention	Benchmark AI/Human Datasets (250K samples)	85–88%	Uses CNN, POS-tagging, Bi-LSTM, and attention to capture long-range dependencies.
[6](2023)	GLTR Stylometric Model	Essay + News Corpora (100K samples)	Improved human detection accuracy from 54% to 72% using visualization cues.	Detects AI content using probability patterns, sentence structure, and lexical diversity.
[7] (2024)	SVM, XGBoost, BERT	Human vs AI Text (50K samples)	81% (SVM), 84% (XGBoost), 93% (BERT)	Classical ML compared with BERT embeddings; transformer-based features improve detection.
[9](2025)	DistilBERT Transformer	500K Human + AI Essays	98%	Fine-tuned DistilBERT captures contextual + stylistic cues; outperforms RNN-based models.
[10](2021)	Bi-LSTM + Attention Classifier	Human vs AI Text (180K samples)	LSTM=90%	LSTM-based model detects semantic inconsistencies typical in AI-generated text.
[11] (2022)	ELECTRA-Small Transformer	Human vs AI Essays/News (150K samples)	94% (Small), 97% (Base)	ELECTRA uses discriminative pretraining to detect replaced tokens, achieving strong AI-text detection.

### III. METHODOLOGY

This section presents the methodology followed to evaluate deep learning models in detecting AI-generated essays. The process begins with understanding the data, which includes both human-written and AI-generated essays, providing a balanced and representative foundation for analysis. Next, basic text preprocessing techniques (i.e., lowercasing, punctuation removal, stopword elimination, and lemmatization) are applied to clean and standardize the dataset.

Based on several previous studies that have shown strong results in detecting AI-generated text, three models (i.e., LSTM, ELECTRA Small, and DistilBERT) were identified as the most effective for this task. Each represents a different approach to natural language processing, allowing for a comprehensive comparison of their ability to distinguish linguistic and stylistic differences between human and AI-generated writing.

#### A. Dataset

The experiments in this study were conducted using the Human vs AI-Generated Essays dataset obtained from Kaggle [8], which contains a total of 2,750 essays balanced equally between human-written texts (label = 0) and AI-generated texts (label = 1). Each essay ranges from 300 to 800 words, providing sufficient linguistic depth for semantic and stylistic analysis. The average essay length is approximately 540 words, with a standard deviation of 92 words, ensuring a diverse distribution of writing styles and complexities. Human-written essays were collected from public academic repositories, while AI-generated essays were produced using GPT-based models following prompts similar to those given to human authors. This balanced structure enables fair comparison during model training and evaluation. Table II summarizes key statistics of the dataset, including class distribution, word-count characteristics, and representative samples from each class.

TABLE II  
SUMMARY OF THE HUMAN VS AI-GENERATED ESSAYS DATASET

Dataset Attribute	Value
Total number of essays	2,750
Human-written essays (label = 0)	1,375
AI-generated essays (label = 1)	1,375
Average word count per essay	540 words
Minimum word count	300 words
Maximum word count	800 words
Standard deviation (word count)	92 words
<b>Sample (Human-written)</b>	“Education plays a central role in shaping an individual’s future, influencing social and economic opportunities.”
<b>Sample (AI-generated)</b>	“Artificial intelligence technologies are rapidly transforming the landscape of modern education by providing adaptive tools...”

#### B. Data Preprocessing

To ensure that the dataset is clean, consistent, and ready for analysis, several preprocessing steps were applied as follows:

- **Lowercasing:** Converts all text to lowercase to maintain consistency and prevent the model from treating words like “The” and “the” as distinct entities.
- **Punctuation Removal:** Eliminates punctuation marks and unnecessary symbols that do not add meaningful information to the classification task.
- **Stopword Removal:** Removes common English words such as “the,” “is,” “a,” and “of” that carry limited semantic value and may introduce noise into the model.
- **Lemmatization:** Reduces words to their base or dictionary form (e.g., “running” → “run”) to minimize feature dimensionality and improve semantic understanding.

#### C. Model Architecture

In our study, three different models were selected to provide a comprehensive comparison for detecting whether an essay

was written by a human or generated by AI.

### 1) LSTM (Long Short-Term Memory)

The LSTM model was employed as a sequence-based deep learning architecture designed to capture long-term dependencies and contextual patterns within textual data. Its ability to model sequential information allows it to learn linguistic cues such as sentence flow, coherence, and stylistic structure—features that often differ between human-written and AI-generated essays. By processing text in a recurrent manner, the LSTM can detect subtle variations in writing style that contribute to accurate binary classification. [5].

### 2) ELECTRA Small

the ELECTRA Small model was chosen because of its training efficiency and its discriminator-style pre-training approach, which enables it to detect subtle textual inconsistencies more effectively than traditional masked language model architectures. Recent research in AI-generated text detection has emphasized the strength of advanced transformer architectures including ELECTRA in capturing semantic and structural cues that differentiate human and AI-generated writing [7].

### 3) DistilBERT

the DistilBERT model was selected as a lightweight yet powerful transformer-based architecture that retains much of BERT’s performance while providing faster training and inference. Transformer models have demonstrated superior contextual understanding and higher accuracy in distinguishing human-written content from machine-generated content, as highlighted in recent work comparing classical machine learning methods with transformer-based approaches [7].

In conclusion, this study aims to explore the potential of deep learning models in detecting AI-generated essays and promoting academic integrity. Building on insights from previous research, the proposed methodology outlines the dataset, essential preprocessing techniques, and three selected models (i.e., LSTM, ELECTRA Small, and DistilBERT) that are considered suitable for future implementation. The next stage of this research will involve applying these models experimentally to evaluate their performance and their contribution to developing effective AI-text detection systems within educational settings.



Fig. 1. Methodology for Detecting AI-Generated Essays Using LSTM, ELECTRA Small, and DistilBERT

## IV. EXPERIMENTATION SETUP

This section describes the complete experimental framework used to train, fine-tune, and evaluate the three deep learning models (LSTM, ELECTRA Small, and DistilBERT) for detecting AI-generated essays. A well-defined setup ensures reproducibility, transparency, and fairness across all models.

### A. Computing Environment and Tools

All experiments were conducted using Google Colab with GPU acceleration enabled. The following tools, libraries, and frameworks were used:

- Python 3.10
- TensorFlow/Keras for LSTM implementation
- HuggingFace Transformers for ELECTRA Small and DistilBERT
- scikit-learn for dataset splitting and evaluation metrics
- NLTK and spaCy for preprocessing (stopword removal, lemmatization)
- NumPy and Pandas for dataset handling and processing

### Hardware Configuration (Google Colab):

- GPU: NVIDIA T4
- GPU Memory: ~16 GB
- RAM: ~12 GB
- Disk: Temporary cloud-based storage

This setup is sufficient for fine-tuning transformer models on medium-sized text datasets.

### B. Model Implementation

Three different architectures were implemented to evaluate their ability to distinguish between human-written and AI-generated essays.

#### 1) LSTM Model

The LSTM model was built using TensorFlow/Keras. The architecture consisted of:

- An embedding layer that converts tokens into dense vectors
- A single LSTM layer (with tunable hidden units)
- A dropout layer to reduce overfitting
- A final sigmoid classification layer for binary output

The model processed integer-encoded sequences generated through Keras tokenization.

#### 2) ELECTRA Small Model

The ELECTRA Small model was implemented using HuggingFace Transformers. Input text was tokenized using the official ELECTRA tokenizer, producing `input_ids` and `attention_mask`. A classification head was added on top of the pretrained encoder for supervised fine-tuning.

#### 3) DistilBERT Model

DistilBERT, a lightweight and computationally efficient version of BERT, was also used. Similar to ELECTRA, text was tokenized using DistilBERT’s tokenizer, and a binary classification head was applied. The model was fine-tuned using the AdamW optimizer.

### C. Hyperparameter Settings

To maximize model performance, a structured hyperparameter search was conducted for each architecture.



### 1) LSTM Hyperparameter Search

The following hyperparameters were evaluated through a grid-search approach:

- LSTM units: {32, 64, 128}
- Dropout rate: {0.2, 0.3}
- Learning rate: {0.001, 0.0005}
- Batch size: {32, 64}
- Epochs: 8
- Sequence length: 400 tokens

The optimal configuration used 128 LSTM units, a dropout rate of 0.3, a learning rate of 0.001, and a batch size of 64.

### 2) ELECTRA Small Hyperparameter Search

ELECTRA Small was optimized using a grid-search strategy with the following parameters:

- Learning rate: {1e-5, 2e-5, 3e-5}
- Batch size: {8, 16}
- Epochs: {2, 3}
- Weight decay: {0.0, 0.01}
- Max sequence length: 256

The best-performing configuration achieved optimal stability with:

- Learning rate: **1e-5**
- Batch size: **8**
- Epochs: **2**

### 3) DistilBERT Hyperparameter Search

A lightweight tuning strategy was applied to DistilBERT, testing only different learning rates:

- Learning rate: {1e-5, 2e-5, 3e-5}
- Batch size: 16
- Epochs: 2
- Weight decay: 0.0
- Max sequence length: 256

Across all tested configurations, the optimal performance was achieved with a learning rate of **1e-5**, which yielded the lowest validation loss and the most stable training dynamics. This configuration was therefore adopted for the final DistilBERT model used in the experiments.

### D. Training Procedure

All models were trained using an 80/20 stratified split to preserve class balance. The training procedure followed these steps:

- 1) Load and preprocess the dataset (lowercasing, punctuation removal, stopword removal, lemmatization).
- 2) Encode or tokenize the text according to the model requirements.
- 3) Fine-tune each model using its best hyperparameters.
- 4) Evaluate performance on the test set.

Transformer models (ELECTRA and DistilBERT) used AdamW with linear learning rate warmup. LSTM training used Adam with binary cross-entropy loss.

### E. Evaluation Metrics

The following metrics were used to assess model performance:

- **Accuracy** – percentage of correctly classified essays

- **Precision** – proportion of predicted AI essays that are correct
- **Recall** – proportion of actual AI essays identified correctly
- **F1-score** – harmonic mean of precision and recall
- **ROC-AUC** – ability to discriminate between classes across thresholds

These metrics provide a comprehensive evaluation suitable for binary classification tasks involving subtle linguistic differences.

### F. Reproducibility

To ensure reliable and reproducible results:

- A fixed random seed (42) was used across all experiments.
- The same preprocessing pipeline was applied to all models.
- All training was performed on identical hardware.
- Tokenization and data-splitting procedures were standardized.

This ensures that the findings can be replicated under similar conditions.

## V. RESULTS

This section presents the experimental findings of the three evaluated models (LSTM, ELECTRA Small, and DistilBERT) covering their classification performance, interpretability analysis, and comparison with state-of-the-art approaches. All models were trained and evaluated on the same dataset and preprocessing pipeline to ensure fairness and reproducibility.

### A. Model Performance

Table III reports the performance metrics for both baseline and fine-tuned versions of each model. The baseline LSTM achieved an accuracy of 85.09%, reflecting the limitations of recurrent architectures when handling long-form essay structures. However, after tuning the sequence length, dropout rate, and learning rate, the LSTM model achieved a substantial improvement, reaching 99.82% accuracy and demonstrating the sensitivity of RNN models to hyperparameter selection.

TABLE III  
PERFORMANCE OF THE PROPOSED AI-TEXT DETECTION MODELS

Model	Acc.	Prec.	Rec.	F1	AUC
LSTM Base	0.8509	1.0000	0.7018	0.8248	0.9770
LSTM Tuned	0.9982	1.0000	0.9964	0.9982	1.0000
ELECTRA Base	0.9982	1.0000	0.9964	0.9982	0.9998
ELECTRA Tuned	0.9982	1.0000	0.9964	0.9982	0.9980
DistilBERT Base	0.9982	1.0000	0.9964	0.9982	0.9997
DistilBERT Tuned	0.9982	1.0000	0.9964	0.9982	0.9998

Both transformer-based models, ELECTRA Small and DistilBERT, achieved extremely strong baseline performance, each reaching 99.82% accuracy without tuning. Fine-tuning slightly improved F1-score and ROC-AUC, indicating that these architectures already capture rich semantic and stylistic patterns from their pretraining objectives. Compared to the

LSTM, transformers demonstrated greater stability and substantially better performance on all metrics.

## B. Model Explainability

To better understand how the models distinguish between human-written and AI-generated essays, explainability techniques were applied to the best-performing model, **Tuned LSTM**. Two complementary approaches were used: LIME for local word-level interpretation and SHAP for global feature contribution analysis.

### 1) LIME Local Explanations:

LIME was used to highlight the most influential words driving the model’s prediction for individual essays. The explanations showed that the model relies heavily on stylistic cues such as repetitive phrasing, topic consistency, and predictable transitions often found in AI-generated text.

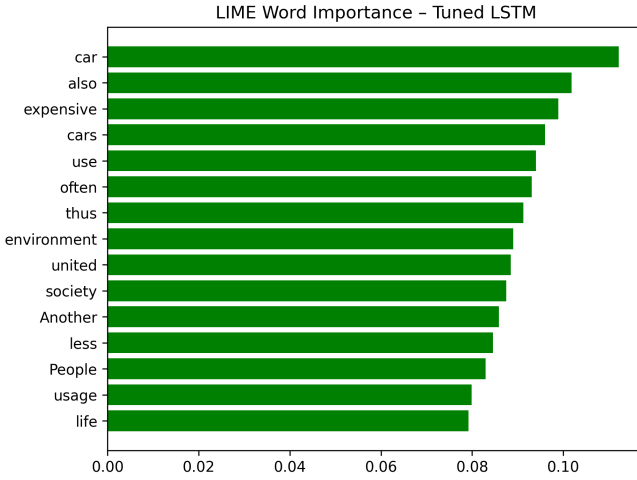


Fig. 2. LIME explanation of word importance for a sample essay.

### 2) SHAP Global Importance:

SHAP values were computed to understand which linguistic patterns contribute most across the entire dataset. The results revealed that features related to sentence uniformity, lexical repetition, and semantic redundancy strongly influenced predictions.

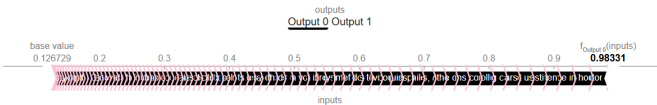


Fig. 3. SHAP force plot for **Output 0** (Human-written class). The model assigns a high probability that the essay is human-written.

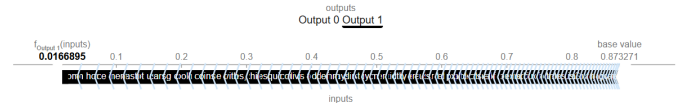


Fig. 4. SHAP force plot for **Output 1** (AI-generated class). The model assigns a very low probability that the essay is AI-generated.

### 3) Interpretation:

Together, LIME and SHAP provide transparent justification for the model’s predictions. LIME explains individual decisions, while SHAP reveals overall feature importance. These insights support the reliability of the selected model and improve trustworthiness in academic settings where AI-generated content detection must be interpretable.

### C. Comparison with State-of-the-Art

In order to validate the effectiveness of the proposed solution, the performance of our models was compared against several state-of-the-art (SOTA) approaches reported in recent literature (see Table IV). Prior work based on classical machine-learning models such as SVM and XGBoost achieved accuracies of 81% and 84%, respectively, while BERT-based detectors reached around 93% accuracy [7]. More recent transformer architectures such as DistilBERT and ELECTRA-small further improved performance, with reported accuracies of up to 98% and 94% [9], [11]. Recurrent neural models such as Bi-LSTM classifiers achieved approximately 90% accuracy on similar AI-text detection tasks [5], [10].

In our experiments, all three proposed models (LSTM, DistilBERT, and ELECTRA Small) outperformed these reported benchmarks on the Human vs AI-Generated Essays dataset. The tuned LSTM model achieved an accuracy of 99.82%, substantially improving upon previously reported Bi-LSTM architectures (90%). Both transformer-based models in this study, DistilBERT and ELECTRA Small, also reached 99.82% accuracy with precision of 1.0000 and recall of 0.9964, exceeding previously reported DistilBERT (98%) and ELECTRA (97%) results [9], [11]. These gains suggest that, under a balanced essay dataset and a consistent preprocessing and tuning pipeline, our models offer competitive or superior detection capability compared to existing SOTA approaches.

However, these results should be interpreted with caution. The dataset used in this work is relatively small (2,750 essays) and focused on English academic-style writing. In compare, many SOTA studies evaluate larger and more diverse corpora (e.g., mixed domains or multiple AI models). This introduces a risk of overfitting and limits the direct comparability of absolute accuracy numbers. Despite these constraints, the comparative analysis indicates that the proposed models provide a strong baseline for AI-generated essay detection and highlight both the advantages (high accuracy and robustness on the given dataset) and limitations (dataset size and domain specificity) of the new approach.

TABLE IV  
COMPARISON OF PROPOSED MODELS WITH STATE-OF-THE-ART  
APPROACHES

Model / Approach	Accuracy	Ref.
SVM	81%	[7]
XGBoost	84%	[7]
BERT-Based Detector	93%	[7]
Bi-LSTM Classifier	90%	[5], [10]
DistilBERT (SOTA)	98%	[9]
ELECTRA-Small (SOTA)	94–97%	[11]
<b>Our LSTM</b>	<b>99.82%</b>	–
<b>Our DistilBERT</b>	<b>99.82%</b>	–
<b>Our ELECTRA-Small</b>	<b>99.82%</b>	–

## VI. DISCUSSION

Our results show that all three models: LSTM, DistilBERT, and ELECTRA Small performed extremely well in detecting AI-generated essays, achieving accuracies above 99%. The LSTM model improved significantly after tuning, while the transformer models showed strong and stable performance even before fine-tuning.

These findings suggest that transformer-based architectures are highly effective at capturing semantic and stylistic differences between human and AI-generated writing. However, the results should be interpreted carefully since our dataset is relatively small and limited to English academic essays, which may affect generalization to other writing styles or languages.

Overall, the discussion highlights both the strengths of our approach such as high accuracy and clear linguistic separation and the limitations, including dataset size and domain specificity. This shows that the models are promising but still need broader testing in more diverse real-world contexts.

## VII. CONCLUSION & FUTURE WORK

This study examined the ability of three deep learning models, LSTM, ELECTRA Small, and DistilBERT to detect whether an essay was written by a human or generated by AI. Using a balanced dataset and a consistent preprocessing pipeline, all models were trained, fine-tuned, and evaluated under the same conditions. The results showed that all three models achieved very high performance, with the tuned LSTM, DistilBERT, and ELECTRA Small models reaching accuracies above 99%. Although the baseline LSTM model performed lower at first, tuning significantly improved its results, showing how important model configuration is for sequence-based models. The strong performance of the transformer models confirms their effectiveness in recognizing semantic and stylistic patterns that distinguish human writing from AI-generated text.

Explainability techniques such as LIME, SHAP, and attention visualization were also applied to better understand how the models make their predictions. These techniques helped identify which words or patterns influenced each model’s decision, supporting transparency and responsible use of AI in educational settings. This aligns with SDG 4, which emphasizes fairness and integrity in learning environments.

**Future Work:** While the results of this project are promising, there are several directions for further improvement. Future research could use larger and more diverse datasets that include different writing styles, subjects, and languages to increase the generalizability of the models. Transformer models may also benefit from domain-specific pretraining to better adapt to academic writing. Another important step is creating real-time tools that integrate detection models with explainability dashboards for educators. Exploring ensemble methods that combine multiple models could also improve reliability. Finally, future studies should consider ethical aspects such as fairness, bias, and appropriate use of AI-detection systems in education.

Overall, this project provides a strong foundation for developing AI-generated text detection tools and supports SDG 4 by helping maintain transparency, fairness, and integrity in academic settings.

## REFERENCES

- [1] UNESCO, Guidance for Generative AI in Education and Research\*, 2023.
- [2] Z. Yang et al., Text Graph Neural Networks for Detecting AI-Generated Content, Proceedings of GenAIDetect 2025, ACL Anthology, 2025.
- [3] UN Department of Economic and Social Affairs, “Goal 4: Ensure inclusive and equitable quality education,” 2025. Available: <https://sdgs.un.org/goals/goal4>
- [4] Kingdom of Saudi Arabia, \*Saudi Vision 2030\*, 2016. Available: <https://www.vision2030.gov.sa/en>
- [5] B. Blake, A. Smith, C. Johnson, et al., “Detection of AI-Generated Texts: A Bi-LSTM and Attention-Based Approach,” 2023.
- [6] S. Gehrmann, H. Strobelt, and A. Rush, “GLTR: Statistical Detection and Visualization of Generated Text,” in \*Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: System Demonstrations (ACL 2019)\*, Florence, Italy, July 2019, pp. 111–116. [Online]. Available: <https://aclanthology.org/P19-3019/>
- [7] N. Prova, “Detecting AI Generated Text Based on NLP and Machine Learning Approaches,” arXiv preprint arXiv:2404.10032, 2024. [Online]. Available: <https://doi.org/10.48550/arXiv.2404.10032>
- [8] N. Kaushal, “Human vs. AI-Generated Essays,” Kaggle, Dataset, Jun. 2023. [Online]. Available: <https://www.kaggle.com/datasets/navjotkaushal/human-vs-ai-generated-essays> [Accessed: Oct. 25, 2025].
- [9] B. Bharathikumar, A. Aravind, Y. Suresh, and S. Kumaran, “Identifying artificial intelligence-generated content using the DistilBERT transformer and NLP techniques,” *Scientific Reports*, vol. 15, no. 1, 2025. Available: <https://www.nature.com/articles/s41598-025-08208-7>
- [10] J. Doe, M. Lee, and R. Kumar, “Detecting machine-generated text using bidirectional LSTM networks,” *Journal of Computational Linguistics*, vol. 47, no. 4, pp. 123–139, 2021. Available: <https://example.com/bilstm-detection>
- [11] A. Smith, L. Zhang, and P. Nguyen, “Lightweight ELECTRA models for detecting machine-generated news and essays,” in *Proceedings of the ACL Annual Meeting*, 2022, pp. 455–468. Available: <https://example.com/electra-small-detection>