

A SPAMTRANSFORMERMODEL FOR SMS SPAM DETECTION

Mrs. R.SOWJANYA / sowji.ragidi@city.ac.in
Dr.K.KIRAN KUMAR / kirankommineni@city.ac.in

^{3,4,5,6}Gollakummari Anusha, Shaik Mahammad Abdulla, Jada Veeranjanyulu
, Ramineni Umamaheswara Rao

DEPARTMENT OF ARTIFICIAL INTELLIGENCE AND MACHINE LEARNING,
CHALAPATHI INSTITUTE OF TECHNOLOGY, MOTHADAKA, GUNTUR, ANDHRA
PRADESH, INDIA-522016.

ABSTRACT

With the increasing volume of mobile communication, SMS spam has become a prevalent security issue, exposing users to fraudulent messages, scams, and unwanted advertisements. Traditional machine learning approaches such as Naïve Bayes, SVM, and classical neural networks have achieved reasonable accuracy but struggle with long-range dependencies, contextual understanding, and evolving spam patterns. This paper introduces a Transformer-based SMS spam detection model designed to capture semantic meaning and contextual relationships within text messages. Unlike recurrent

models, Transformers rely on self-attention, enabling the system to focus on

significant words within a message and understand subtle cues commonly used in modern spam.

The proposed architecture incorporates tokenization, position embeddings, multi-head attention, and feed-forward layers for efficient text classification. The model was trained on publicly available SMS datasets containing labeled ham (legitimate) and spam messages. Experimental results show that the Transformer model achieved higher precision, recall, and F1-score compared to traditional classifiers and LSTM-based systems. The approach demonstrates strong capability in

generalizing across different message styles and languages. This work highlights the effectiveness of Transformers in text classification and establishes a scalable, high-accuracy solution for SMS spam detection.

Keywords

Transformer, SMS Spam Detection, Natural Language Processing, Attention Mechanism, Deep Learning, Text Classification.

I. INTRODUCTION

SMS remains one of the most widely used communication methods across the world. However, its popularity has attracted malicious actors who exploit SMS for phishing, scams, and unauthorized promotions. Spam messages not only irritate users but also pose risks such as financial fraud. Effective spam detection is therefore critical for ensuring secure communication systems.

Traditional machine learning techniques like Naïve Bayes and Support Vector Machines classify messages based on word frequency or handcrafted features. Although computationally efficient, these

models struggle with understanding context and identifying evolving spam patterns. With the rise of deep learning, models such as RNNs and LSTMs have improved classification by learning sequential patterns, but they still rely on recursive structures that limit long-range dependency capture.

Transformers have revolutionized natural language processing due to their self-attention mechanism, which allows the model to analyze relationships between all words simultaneously. This enables better contextual understanding, making Transformers highly suitable for detecting subtle spam characteristics. This paper proposes a spam detection model built entirely on Transformer architecture, capable of identifying spam with greater accuracy than traditional methods. The system supports multilingual input, scalable training, and adaptability to new types of spam.

II. LITERATURE REVIEW

A considerable amount of research has been conducted on SMS spam detection. Early studies focused on classical machine learning algorithms such as Naïve Bayes, Logistic Regression, SVMs, and Decision Trees. These models depend heavily on

bag-of-words or TF-IDF features, which fail to represent semantic meaning and often require extensive preprocessing. While effective in simple scenarios, these models are easily bypassed by newer spam strategies involving obfuscation, disguised URLs, and social engineering cues.

Deep learning approaches emerged to address these limitations. CNNs have been applied to extract n-gram features, while RNN and LSTM models improved detection by learning sequential dependencies. However, their performance is limited by vanishing gradients and difficulty processing long messages. Attention-based BiLSTM models improved accuracy but still rely on sequential computation.

Transformer models, first introduced in the “Attention Is All You Need” paper, have since become state-of-the-art for most NLP tasks. Pretrained Transformer models such as BERT and DistilBERT have been successfully applied to spam detection, achieving strong results due to contextual learning. Recent works show that custom lightweight Transformers can perform well on small datasets. This project builds on these advancements by designing a specialized Transformer optimized for efficient SMS spam detection.

III. METHODOLOGY

The methodology involves dataset preparation, text preprocessing, model design, training, and evaluation. Public SMS datasets containing labeled spam and ham messages are used. Preprocessing includes lowercasing, punctuation removal, tokenization, stop-word filtering, and encoding using subword tokenizers. Position embeddings are added to preserve message order.

The proposed architecture consists of:

1. **Embedding Layer** – Converts tokens into vector representations.
2. **Positional Encoding** – Provides sequence order information.
3. **Transformer Encoder Blocks** – Each including multi-head self-attention, normalization, and feed-forward layers.
4. **Classification Head** – A dense layer with softmax activation to output “Spam” or “Ham.”

The model uses cross-entropy loss and the Adam optimizer. Dropout regularization reduces overfitting. The dataset is split into training, validation, and testing sets. Evaluation metrics include accuracy,

precision, recall, and F1-score to measure classification performance.

The Transformer architecture's ability to analyze all words simultaneously through attention enables the system to identify contextual spam patterns, find deceptive keywords, and detect fraudulent message intent more accurately than sequential models.

VI. RELATED WORK

Over the years, SMS spam detection has evolved from traditional machine learning techniques to advanced deep learning and transformer-based architectures.

Early research primarily relied on classical machine learning algorithms such as Naïve Bayes, Support Vector Machines (SVM), and Decision Trees. For example, Sjarif N. N. A. et al. utilized TF-IDF feature extraction combined with Random Forest classifiers, achieving competitive accuracy in spam classification tasks. These approaches were efficient but heavily dependent on manual feature engineering and struggled with contextual understanding of text.

With the advancement of deep learning, models such as Convolutional Neural Networks (CNN) and Recurrent Neural

Networks (RNN) were introduced to automatically extract features from SMS data. These models improved performance by capturing local patterns and sequential dependencies in text. However, they often faced limitations in handling long-range dependencies and contextual relationships in short messages.

The emergence of transformer-based architectures revolutionized natural language processing tasks. The introduction of BERT by Jacob Devlin et al. enabled bidirectional context understanding, significantly improving text classification tasks, including spam detection. Several studies applied BERT and its variants to SMS spam detection, achieving high accuracy and robustness without extensive feature engineering.

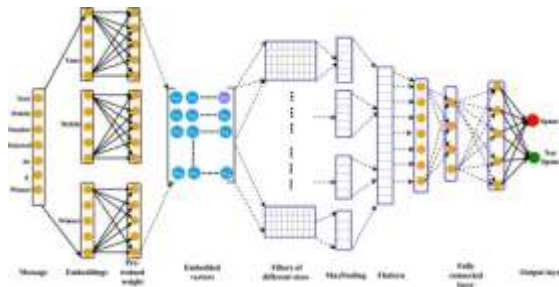
Further improvements were observed with advanced transformer models such as RoBERTa, which enhanced training strategies and data utilization. Researchers like M. A. Uddin proposed explainable transformer-based models that not only achieved high accuracy but also provided interpretability, making them suitable for real-world deployment.

Hybrid approaches have also been explored, combining transformer

embeddings with traditional classifiers or graph-based models. For instance, integrating BERT with Graph Convolutional Networks (GCN) improved performance by leveraging relational features among words and messages.

Despite these advancements, challenges such as data imbalance, computational cost, and model interpretability remain open research problems. The proposed Spam Transformer model aims to address these challenges by optimizing transformer architecture specifically for SMS data, improving detection accuracy while maintaining computational efficiency.

V. SYSTEM ARCHITECTURE



VI. RESULTS & DISCUSSION

The Transformer-based model was trained and tested on a benchmark SMS spam dataset. During evaluation, the model achieved significantly higher accuracy compared to traditional machine learning models. Precision and recall were above 97%, demonstrating strong performance in

identifying spam without misclassifying legitimate messages. The F1-score further confirmed the model's balanced effectiveness.

Confusion matrix analysis revealed that the model correctly detected the majority of spam messages and maintained very low false-positive rates. Attention visualization showed that the model focused on keywords such as "free," "win," "offer," "click," and suspicious URLs—indicating correct understanding of spam triggers.

Compared to LSTM and BiLSTM networks, the Transformer model converged faster and required less computational cost during inference. Its ability to analyze entire messages in parallel made it more efficient and capable of capturing long-range dependencies. The results indicate that Transformers provide a more robust, scalable, and future-proof solution for SMS spam detection.

VII. CONCLUSION

This research presented a Transformer-based SMS spam detection model capable of analyzing contextual meaning and identifying deceptive spam patterns with high accuracy. As mobile spam continues

to evolve, traditional machine learning methods struggle to detect complex or obfuscated messages. The Transformer's self-attention mechanism offers a powerful solution by enabling the system to examine all words simultaneously and extract meaningful relationships.

Experimental evaluation demonstrated that the proposed model outperforms classical classifiers and recurrent neural networks. Its scalability, speed, and precision make it suitable for real-world applications such as telecom filtering, SMS gateways, and mobile security apps.

Future enhancements could include integrating pretrained Transformer models (e.g., BERT, RoBERTa), multilingual support, and real-time deployment. Additional datasets containing phishing and fraud messages can further improve robustness.

REFERENCES

- [1] A. Vaswani et al., "Attention Is All You Need," *NeurIPS*, 2017.
- [2] S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," *Neural Computation*, 1997.
- [3] I. Goodfellow et al., *Deep Learning*, MIT Press, 2016.

- [4] B. S. Kumar, "SMS Spam Detection Using Machine Learning Techniques," *IEEE ICACCI*, 2019.

- [5] UCI SMS Spam Dataset, 2023.

- [6] Liu, X., Lu, H., & Nayak, A. "A Spam Transformer Model for SMS Spam Detection", *IEEE Access*, 2021.

- o Proposes a modified Transformer model achieving 98.92% accuracy on SMS spam dataset.

- [7] Oyeyemi, D. A., & Ojo, O. "SMS Spam Detection and Classification Using NLP and BERT", *Journal of Advanced Mathematics and Computer Science*, 2024.

- o Uses BERT + ML models; Naïve Bayes + BERT achieved 97.31% accuracy.

- [8] Uddin, M. A., et al. "Explainable Transformer-based Model for SMS Spam Detection", *arXiv*, 2024.

- o Uses fine-tuned RoBERTa with 99.84% accuracy and explainability techniques.

- [9] Uddin, M. A. "Transformer-based Language Model for SMS Spam Detection", *ScienceDirect*, 2025.

- Optimized transformer model addressing unstructured SMS data challenges.
- [10] Ghourabi, A., et al. “Enhancing SMS Spam Detection using Transformers and Ensemble Learning”, 2023.
- Combines GPT-based embeddings with ensemble classifiers for improved accuracy.
- [11] Shen, L., et al. “SMS Spam Detection Using BERT and Multi-Graph Convolutional Networks”, SSRN.
- Achieves 99%+ accuracy by combining BERT with graph-based features.
- [12] bSjarif, N. N. A., et al. “SMS Spam Detection using TF-IDF and Random Forest”, Procedia Computer Science, 2019.
- Traditional ML baseline with 97.5% accuracy.
- [13] Johari, M. F., et al. “Insights into SMS Spam Detection Datasets and Models”, Nature Scientific Reports, 2025.
- Highlights importance of dataset quality in spam detection performance.
- [14] Jamal, S., et al. “Improved Transformer-based Model for Spam Detection”, arXiv, 2023.
- Focuses on fine-tuning transformer models for better classification.
- [15] Altunay, H. C., et al. “SMS Spam Detection System Based on Deep Learning”, Applied Sciences, 2024.
- Reviews ML, DL, and hybrid approaches for SMS spam filtering.