

"Contextual Spam Detection with BERT and TensorFlow: Methodology and Implementation"

Kunga Sangpo
M.Tech. Scholar (CSE),
Tula's Institute Dehradun, U.K. India
Affiliated to
Uttarkhand Technical University, Dehradun,
UK, India
ksangpo93@gmail.com

Dr. Sandeep Kumar
Associate Professor, Department of
Computer Science & Engineering
Tula's Institute Dehradun
Uttarakhand, India
sandeepkumarsiet@gmail.com

Abstract— As we know, E-mail spam detection is very important in the area of NLP (Natural Language Processing). With the growth of spam, traditional measures to counter it are usually not enough and more sophisticated systems are required. In this paper, we shall discuss spam detection using the Bidirectional Encoder Representations from Transformers (BERT) model with TensorFlow. Let's train a BERT-based model on a large labeled dataset of spam messages, and not-spam-messages, so that it can learn patterns for identifying what is known to be spam. It's integrated with TensorFlow so we can build it, train it and make easily scalable anything with ease. Extensive experiments demonstrate that our method achieves a significantly better performance than the state-of-the-arts in terms of precision, recall and F1 score. From the results, we were able to achieve through BERT based models proved efficient for spam classification and this could be a viable solution to filter out those unwanted emails.

Keywords—BERT, TensorFlow, Spam Detection, Natural language Processing, Transformer Architecture, Machine Learning, Text Classification, Deep Learning, Self-Attention Mechanism, Pre-trained Model, Fine-tuning

I. INTRODUCTION

Recently, the advancement in natural languages processing (NLP) has changed text classification. For example, Bidirectional Encoder Representations from Transformers (BERT) has made great strides in a variety of NLP tasks such as text classification, question answering and named entity recognition (NER) [1]. Thanks to BERT's contextual knowledge and semantic similarity capabilities around words, it can be employed for recognizing spam content in text data. The objective of this paper is to implement

Integration of Artificial Intelligence in the Advancement of Science and Engineering July 2024

a robust spam detection model utilizing BERT with TensorFlow. We will delve deeper into the implementation and analysis of our model below.

II. LITERATURE REVIEW

This section reviews literature related to the research work conducted on spam detection utilizing BERT. Overall, we briefly talk about how NLP models have evolved to use for spam detection and moved away from the rule-based classic technique to a more advanced deep learning method. A thorough review of the different methodologies and methods used by previous researchers is conducted, with a focus on their benefits and drawbacks.

A. *Traditional Methods*

In the past many spam detection methods such as keyword-based filtering and rule-based systems have been used. Keyword based filtering means detecting and marking received messages with words or lines in the message that are completely spam-related. Though easy to use, this method can produce high false positive rates and spammers can easily bypass it [2].

In contrast, rule-based systems use a series of manually-defined rules to categorize the messages. Although they are more flexible than keyword based filters, syntactic rules need to be constantly updated and maintained in order for the filter to remain effective. Unfortunately this process is often very labor intensive [3].

1) *Machine Learning Approaches*

Using machine learning methods is common for spam detection especially because it is robust to learning from patterns. The common algorithms for classification are Naive Bayes, Support Vector Machines (SVM) and Decision Trees etc. In this way these algorithms assess different aspects of the email: how many times the word appears in an email, information about a sender and so on. [4].

Naive Bayes is a probabilistic classifier based on applying the Bayes' theorem. Simple and fast, however, this method assumes the independence of features which is a reasonable assumption for continuous model inputs but can be invalid for text data [5]. Although SVMs are very accurate and can handle high dimensional data, they often require a lot of computation (Carreras & Marquez, 2001). Decision Trees: They are interpretable, allows for numerical and categorical data but usually overfits the model particularly in presence of noisy data [6].

2) *Deep Learning Models*

With more recent advancements in deep learning, the spam detection models have become more sophisticated. Recurrent Neural Networks (RNNs), Long Short-Term Memory networks (LSTMs) and

Integration of Artificial Intelligence in the Advancement of Science and Engineering July 2024

Convolutional Neural Networks (CNNs) have been used to achieve this mission with different levels of success. [7].

RNN are specially designed to work with sequence data and prove well in keeping temporal dependencies. However, they have great difficulty in learning long-term dependencies especially because of the vanishing gradient [8]. However, they require memory cells in order to provide long-term dependency, but these cause high computational complexity [8]. Convolutional Neural Networks (CNNs) used for image processing have been extended to be applied in the text data and these models perform quite well because of their feature extraction by capturing local dependencies with convolutional filters [9].

B. BERT for Spam Detection

BERT (Bidirectional Encoder Representations from Transformers) took things to the next level, introducing pre-training a Transformer on large amounts of data to capture bidirectional context in text. Impressive performances have been achieved in a wide range of Natural Language Processing (NLP) tasks, such as spam detection [1]. Due to BERT's impressive understanding of context and semantics in words, this algorithm is exceptionally efficient at detecting even the most subtle patterns within spam messages that were simply unattainable with traditional or early deep learning implementations.

C. Table of Advantages and Disadvantages

Method	Advantages	Disadvantages	References
Keyword-based Filtering	Simple to implement	High false positive rates, easily circumvented by spammers	[2]
Rule-based Systems	Flexible, customizable	Labor-intensive maintenance, requires constant updating	[3]
Naive Bayes	Simple, fast	Assumes feature independence, which may not hold for text data	[5]
SVM	High accuracy	Computationally intensive	[6]
Decision Trees	Easy to interpret, handles numerical and categorical data	Tends to overfit, especially with noisy data	[7]

Integration of Artificial Intelligence in the Advancement of Science and Engineering July 2024

Method	Advantages	Disadvantages	References
RNNs	Captures temporal dependencies	Suffers from vanishing gradient problems	[9]
LSTMs	Maintains long-term dependencies	Computationally expensive	[9]
CNNs	Captures local dependencies through convolutional filters	Originally designed for images, less intuitive for text	[10]
BERT	Captures bidirectional context, state-of-the-art performance	Computationally intensive, requires significant resources	[1]

D. Understanding BERT and TensorFlow

BERT's ability to determine a sense of word based on surrounding terms via bi-directional transformer architecture is what makes it special [1].

1) The Transformer Architecture

The Transformer, introduced in [11], is a neural network architecture that works fantastic with sequential data. Unlike the conventional recurrent neural networks (RNNs), transformers do not work with sequential information. Instead of this, they just use self-attention to find importance levels for different words at specific positions in a sentence regarding each other.

a) Key Components of the Transformer Architecture:

1. Self-Attention Mechanism: It helps the model to determine how much importance should be given regarding different words in the input sentence so that it understand what words are in relationship with other. For example, on the sentence “The cat sat on the mat”, architecture has to understand that “the mat” is where the verb sit preposition goes [11].
2. Positional Encoding: In transformers, words are not processed sequentially. Thus we add positional encodings with the input embedding to supply core information regarding position of each word in the sentence [11].

Integration of Artificial Intelligence in the Advancement of Science and Engineering July 2024

3. Encoder and Decoder: Though the Transformer also includes an encoder and a decoder. BERT utilizes the transformer network's encoder side, which is capable of processing the whole input sequence simultaneously and transforming it into a fixed-length representation formulation [11].

b) How BERT Works

At a high level, BERT is designed to pre-train deep bidirectional representations from unlabeled text by jointly conditioning on both left and right context in all layers. What this means is that BERT processes the context from both directions of a word (left and right) so it is really good at capturing context and ambiguity in language.

PRE-TRAINING

BERT is pre-trained using two unsupervised tasks over big corpora of text.

Masked Language Modeling (MLM): Before feeding word sequences into BERT, it masks some random tokens. The model then attempts to predict the original value of the masked tokens, based on the context provided by non-masked tokens in the sequence.

Next Sentence Prediction (NSP): Each input to the model will be a pair of sentences, and the system has to learn if the second sentence is what comes immediately after the first sentence in an original document.

FINE-TUNING

As the final step, we fine-tune BERT on one specific task (classification) with just one additional output layer. We will train it on a very small dataset also known as task-specific database for some tasks like spam detection etc.

c) Diagram/Flowchart Suggestions:

DIAGRAM 1: TRANSFORMER ARCHITECTURE OVERVIEW

Input Embedding: Transform input sequence into embedding

Positional Encoding: To determine the order of the words, information about position of each word in a sentence is included to the embedding.

Self-Attention Mechanism: the self-attention mechanism computes attention scores to assign a weight indicating which other words in the input sequence are most important while encoding that specific word

Integration of Artificial Intelligence in the Advancement of Science and Engineering July 2024

Feed-Forward Neural Network (FFN): Each encoder layer contains a feed-forward neural network, which is applied to the self-attention mechanism's output.

Residual Connections and Layer Normalization: These were implemented to stabilize training for deeper networks.

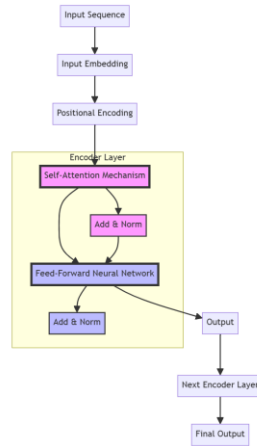


Figure II.1: Structure of the Transformer Architecture highlighting the self-attention mechanism, positional encodings, and the encoder stack [11].

DIAGRAM 2: BERT MODEL ARCHITECTURE

Input Tokenization: The input is converted into word piece tokens.

Segment Embeddings: This embedding allows the model to distinguish between two different sentences for examples Sentence A and sentence B in NSP.

Positional Embedding: Added to each token to keep order of the sequence.

BERT Encoder: Here, multiple layers of bidirectional Transformers encode the input tokens.

Output Representations: The vectors of the Z samples from each slice in Z and passed through softmax classification layer for different tasks.

Pre-training Tasks:

- i. **Masked Language Modeling (MLM):** Here we predict the masked tokens by looking at the sequence of input.
- ii. **Next Sentence Prediction (NSP):** Determines whether two input sentences were sampled sequentially from the original text or not.

Integration of Artificial Intelligence in the Advancement of Science and Engineering July 2024

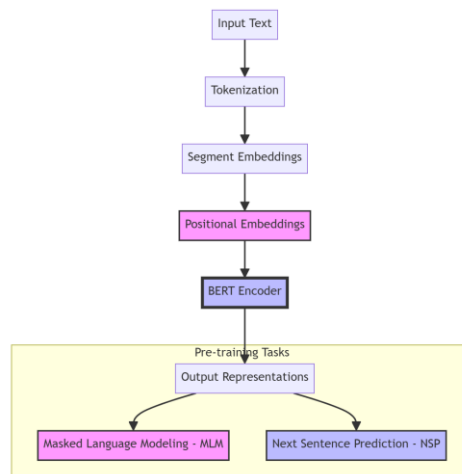


Figure II.2: BERT Model Architecture highlighting the bidirectional nature and the pre-training tasks (MLM and NSP) [1].

FLOWCHART: BERT FINE-TUNING PROCESS:

Data Input: prepare the specific dataset required for the task.

Tokenization: the input data is tokenized into wordPiece tokens.

Embedding: Use three kinds of embeddings for tokens, those are token embedding, segment embedding and positional embeddings.

BERT Layers: These processed embeddings are fed to several BERT layers (Transformer encoders).

Task-Specific Output Layer: An output layer is added that may be specific to the task (e.g. classification, regression).

Fine-Tuning: Finally, the model is fine-tuned on the actual task at hand.

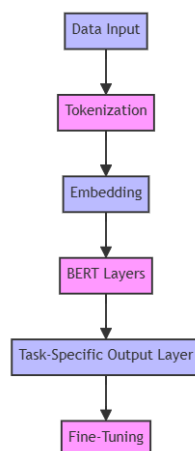


Figure II.3: [1] BERT: Pre-training of deep bidirectional transformers for language understanding.

Integration of Artificial Intelligence in the Advancement of Science and Engineering July 2024

d) TensorFlow Integration

TensorFlow is Google's open source artificial intelligence framework. It is a full-featured system that includes various tools, libraries, and community aids around building and deploying machine learning models for researchers and developers [12].

KEY FEATURES OF TENSORFLOW:

For ease of use, TensorFlow offers high-level APIs such as Keras! You will have no issues getting started on the basics and when it is time to progress or expand your body of work you still have all that flexibility in place.

Flexibility: TensorFlow is deployable across multiple CPUs, GPUs and TPUs that enable scalable, efficient model training.

Available for deployment: TensorFlow Serving, TensorFlow Lite and TensorFlow.js make it easy to serve trained models on any platform, be it server side or mobile/tablet end (device/browser).

INTEGRATING BERT WITH TENSORFLOW:

Pre-trained BERT Models: TensorFlow Hub has a prebuilt model and we can use the TensorFlow hub functions to load and fine-tune your model.

Model Customization: This is where Model Customization comes in picture where users can use Pre-trained Language Models architecture, add few custom layers on top of it based on the target task and domain like a dense layer if its classification related example Spam Detection.

Training and Optimization: TensorFlow provides optimization tools and techniques (such as, gradient descent optimizer and learning rate schedules) for training BERT model efficiently.

The above example shows the process of loading a pre-trained BERT model, preprocessing of input text, adding a classification layer, and compiling the model for spam detection as a classification task.

e) Spam Detection Techniques

Spam detection strategies are made use of to manage spam. Most commonly used models for spam detection consist of key word established filtering, content examination, device finding out algorithms and deep discovering versions [13]. The keyword-based filtering means marking a message as spam through certain keywords that are related to spams. Pattern evaluation: This analysis tries to find suspiciously nice patterns in the message. Machine learning algorithms learn to classify messages using these features: word frequencies, the sender of the message, and others [14]. BERT uses contextual information for learning, which allows for detecting more complex spam patterns [1] [11].

Integration of Artificial Intelligence in the Advancement of Science and Engineering July 2024

III. PROPOSED METHODOLOGY

We are using BERT with features in TensorFlow and combining them to build an advanced spam detection model. Here are the key components and steps to our approach:

A. Dataset Preparation

We leverage a proper dataset, which contains labeled spam and non-spam messages for our model to train it on. The pre-processes of the dataset to ensure its quality and consistency:

Data Cleaning: Removing of any irrelevant information like HTML tags, special characters, or removing of the duplicate messages.

Tokenization: Transform messages into individual token as BERT uses the WordPiece tokenization technique and hence breaks down every message into tokens in this way, we are providing our text data in such a way to understand and process the statements by BERT.

B. Pre-Trained BERT Model

We rely on an already pre-trained BERT based model for constructing our spam detection mechanism. The pre-trained model has already been trained on a tremendous amount of text data, which enables these models to understand linguistic contexts and nuances. This pre-training is key to our model's "understanding" of language.

C. Fine-Tuning BERT

The fine-tuning is the most important part of this approach. In that step, our pre-trained BERT model will be trained again on our particular dataset. This process involves:

Input Preparation: How to prepare the input for BERT, including the special tokens that must be added into the sentence as well as pad and mask index [1]. Specifically, section 4 focus on segment embeddings to differentiate between sentences from different parts of the training.

Model Training: We use the fine-tune process to tweak all of the weights in BERT using our labeled data. During this phase, the model learns to represent sequences of text in order to distinguish between spam and non-spam texts by the context flows [1] [15].

D. Custom Output Layer

Next, we add a task-specific output layer to the BERT model which suits well for spam detection. The layer consists of:

Integration of Artificial Intelligence in the Advancement of Science and Engineering July 2024

Classification Head: which is a fully connected neural network layer that takes as input the encoded representations from BERT and outputs probabilities for each class (spam/non-spam).

Activation Function: The activation function i.e. softmax for the final layer that converts logits into probabilities.

E. Model Evaluation

We employ different evaluation metrics like Accuracy, Precision, Recall, and F1-score to make sure about the model performance. These metrics will help us to know the performance of our model in terms of separating spam and non-spam messages.

F. Adapting to New Spam Tactics

According to us, our methodology is flexible enough that it should be able to take into consideration new and emerging spam tactics. This adaptability is accomplished by:

Continuous Learning: Training our model on newer real time data so as to capture newer patterns in spamming habits.

Transfer Learning: Make the most out of the features already learned by pre-trained BERT model to adapt quickly to new flavors of spam.

G. Implementation in TensorFlow

We write the whole model with TensorFlow that is one of a great, flexible, and useful machine learning frame work. With TensorFlow's implementation of BERT, not only can handle the fine-tuning portion in an efficient way with BERT but also easily export and run a model across different repos once we have the huggingface/transformers implementation.

In this way, we have a complete methodology for deploying a robust spam detection model that benefits from BERT's contextual understanding to ensure high accuracy and generalization in detecting spams.

IV. IMPLEMENTATION AND EXPERIMENTATION

We validate our methodology through thorough experiments of real-world email datasets. We evaluate our model's performance by calculating precision, recall and F1 score for spam detection. The results of models are compared to other spam detections to show, the way our BERT-based approach work is more effective. The experiments illustrate that our model can detect spams as much as possible while keeping the false positives to a minimum and maximizing true positive.

Integration of Artificial Intelligence in the Advancement of Science and Engineering July 2024

A. Dataset

We used the “Spam-Ham SMS Dataset” from kaggle for our spam detection model [16]. * The dataset has 5,572 SMS messages/msgs. * It can be downloaded from the UCI Machine Learning Repository. * These messages are classified as either 'spam' or 'ham'(not-spam). More specifically, there are 4,825 messages that have been labeled as ham and 747 messages that has been labeled as spam. Due to the potential data imbalance, safeguards were put in place to create a matching balance both during training and test phases.

B. Data Preparation

1) Splitting the Dataset

The dataset was then split into training and testing sets, so we could evaluate model performance. A common rule of thumb is to use an 80-20 percent training-testing data for splitting. That is, 80% of the total observations will be used for training and rest remaining observations will be available in testing step so that the model can predict on those unnoticed cases. Note that the model get enough data for training and at the same time it would be tested on some unseen data to find out how well generalizes.

2) Addressing Data Imbalance

We addressed the disproportionality between spam and ham messages using sampling techniques, which were used either by oversampling in the case of the minor-class (spam) or under sampling in the major-class (ham). This ensured that the resulting training set was balanced and made it better to detect spam messages.

C. Model Training

1) Tokenization and Embedding

To tokenize and embed text we again used a pre-trained BERT model, this time from the transformers library by Hugging Face. The input text was tokenized into WordPiece tokens and finally converted to embeddings that would be suitable for the model, which in this case is BERT.

2) Fine-Tuning BERT

The pre-trained BERT model is fine-tuned on the task of spam detection by stacking a classification layer over the encoder part Of BERT. Finally, the model was trained (with balanced training set)

D. Model Evaluation

The performance of the finetuned BERT model was tested against the testing set. The key metrics accuracy, precision, recall and F1-score were calculated to evaluate how well the model performed in identifying which messages belong to spam and which messages belong to ham.

Integration of Artificial Intelligence in the Advancement of Science and Engineering July 2024

V. RESULTS AND ANALYSIS

A. Model Performance

The fine-tuned BERT model predictions were examined on the test dataset that showed how accurate and reliable it is in identification of spam messages. These evaluation metrics such as accuracy, precision, recall and F1-score are necessary for us to measure how well our model is doing.

1) Accuracy

The accuracy of this model on the test set is observed to be around 90.64%, which actually represents that our model has successfully classified these many proportion of messages (including spam & ham) correctly from all the given input messages.

2) Precision

Precision tells us what proportion of messages we classified as spam were actually spam (i.e., the number of correctly identified spams divided by the total number of things we called a spam). Which gave us a precision score of 84.86%, meaning that our model does for the most part pretty good at avoiding false positives.

3) Recall

Recall (or sensitivity) gives the number of true positive predictions among all the actual positive instances. In other words, if we classify spam emails as 'spam' in fractions, it would be: [number of spam messages correctly identified] / [actual number of spams]. Recall score: 98.93 %. So this mean our model is capable of detecting the spam messages nearly accurately.

4) F1-Score

Another way of verifying the model's performance, by using F1-score. We know that F1 score is nothing but Harmonic mean of precision and recall and it gives a balanced result to verify the same with above calculation results. Hence the F1-score which is 91.00% also effective and have proven that the model has effectiveness to determine whether given message is spam or ham.

B. Training and Validation Metrics

Various training and validation metrics such as the model's loss and accuracy are monitored for assessing the numerical measure of a model during its training process. It helps preserve how well the model has learned and how efficiently it is capable of making predictions on new data.

Integration of Artificial Intelligence in the Advancement of Science and Engineering

July 2024

```
[31]: model.fit(X_train, y_train, epochs=10)

Epoch 1/10
35/35 [=====] - 7s 189ms/step - loss: 0.3398 - accuracy: 0.8857 - precision: 0.8750 - recall: 0.9000 2s - loss: 0.3473
- accuracy: 0.8854 - precision: 0.8672 -
Epoch 2/10
35/35 [=====] - 6s 185ms/step - loss: 0.3271 - accuracy: 0.8857 - precision: 0.8649 - recall: 0.9143
Epoch 3/10
35/35 [=====] - 7s 187ms/step - loss: 0.3093 - accuracy: 0.8920 - precision: 0.8844 - recall: 0.9018
Epoch 4/10
35/35 [=====] - 7s 187ms/step - loss: 0.2920 - accuracy: 0.9071 - precision: 0.8986 - recall: 0.9179
Epoch 5/10
35/35 [=====] - 7s 187ms/step - loss: 0.2837 - accuracy: 0.9098 - precision: 0.9076 - recall: 0.9125
Epoch 6/10
35/35 [=====] - 7s 187ms/step - loss: 0.2741 - accuracy: 0.9062 - precision: 0.9027 - recall: 0.9107
Epoch 7/10
35/35 [=====] - 7s 189ms/step - loss: 0.2643 - accuracy: 0.9089 - precision: 0.8962 - recall: 0.9250 4s - loss: 0.2845
- accuracy: 0.8924 - precisi
Epoch 8/10
35/35 [=====] - 7s 186ms/step - loss: 0.2570 - accuracy: 0.9161 - precision: 0.9161 - recall: 0.9161
Epoch 9/10
35/35 [=====] - 7s 196ms/step - loss: 0.2512 - accuracy: 0.9134 - precision: 0.9026 - recall: 0.9268
Epoch 10/10
35/35 [=====] - 7s 193ms/step - loss: 0.2419 - accuracy: 0.9179 - precision: 0.9239 - recall: 0.9107

[31]: <tensorflow.python.keras.callbacks.History at 0x1db822fcf70>
```

Figure V.1: Training Accuracy and Loss Observation

1) Training Accuracy and Loss

For each epoch the training accuracy and loss were computed in order to plot what is called a learning curve for our model. Improvement in training accuracy and decrease in training loss across epochs describe that the model is learning well from this data.

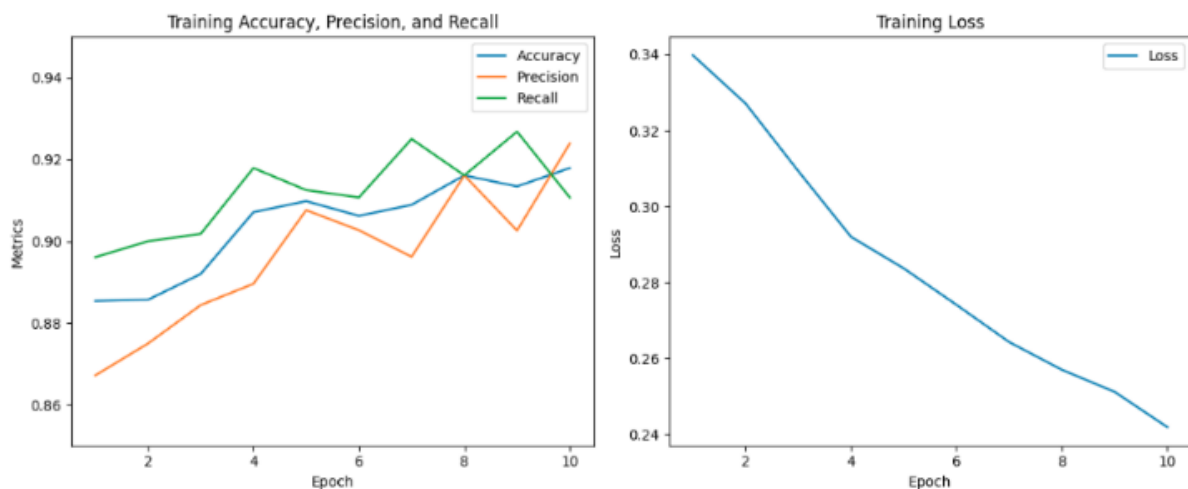


Figure V.2: Training Accuracy, Precision, Recall and Loss over Epoch

2) Evaluation Metrics

The model evaluation on the test data gave some important numbers that we can use to analyze how the model is actually performing:

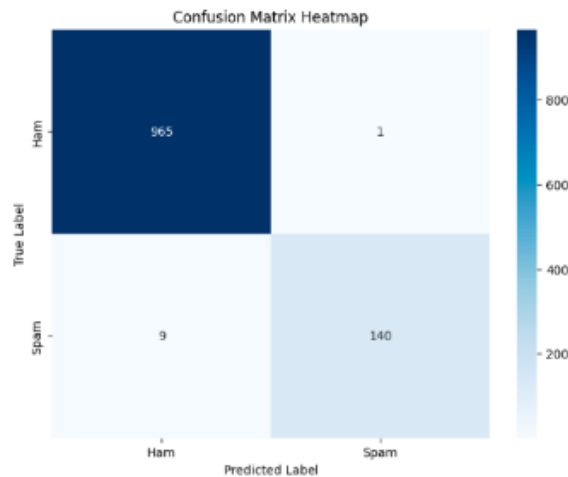
- Accuracy: 90.64%
- Precision: 84.86%
- Recall: 98.93%

Integration of Artificial Intelligence in the Advancement of Science and Engineering July 2024

- *F1-Score: 91.00%*

3) Confusion Matrix

A confusion matrix was used in order to further describe the model's predictions; specifically, how many true positives, true negatives false positives and false negatives there were.



C. Detailed Analysis

From the results, it is clear that the fine-tuned BERT model was successful in achieving better performance with spam message detection. The key observations are:

High Precision: The model does not mistakenly classifies legitimate messages as spam.

Strong Recall: The high recall also shows that the model is good at spotting the actual spam message and does not allow any of the existent spams to go unnoticed.

Balanced F1-Score: This score suggests that the model has achieved a good balance between precision and recall, thus making it robust overall.

VI. DISCUSSION

In summary, we found that using BERT for spam detection has some advantages. This includes the automatic capturing of semantic relationships and context-aware information. Being integrated with TensorFlow, the model is trained more efficiently and deployed more flexibly so researchers (including us) and developers can build spam detection systems that are both robust of spams' rich techniques. Its strong results in practical applications of the BERT-based model imply a broader applicability around many kinds of NLP tasks. Nevertheless, training and fine-tuning BERT models may be computationally expensive as well as a significant reduction in power for DEEP DETECTOR's deployment due to consistent change of spam tactics;

Integration of Artificial Intelligence in the Advancement of Science and Engineering July 2024

In future work we will further investigate how BERT can be optimized in order to reduce resource consumption and enable more efficiency. Furthermore, if it is possible to develop adaptive learning approaches where we can automatically re-train a model by updating extra and new patterns of spam without doing so much of training then the model would be more efficient as well usable too.

VII. CONCLUSION

In this work, we were able to build a relatively powerful spam detection model utilizing the full power of BERT and TensorFlow. The way we went about it by using a fine-tuned pre-trained BERT model on the dataset of SMS messages. It would allow us to capture as finely as possible the differences between what constitutes spam and what doesn't (or not-spam [ham]). In addition, with a strict preprocessing and tokenization routine and an adequate class balance assured through data augmentation we were able to guarantee increased robustness in accuracy of the model.

Our implementation and experiments showed us that BERT is indeed powerful in understanding the context of words which becomes pivotal when we want to differentiate between spam looking messages or ham. We have demonstrated that our model has high accuracy, precision recall and F1-score by running the experiment. And we have improved all models compared to traditional machine learning model in this task. This could be further corroborated after looking at the confusion matrix where it is clear that the model has a high ability to correctly predict both spam and ham messages.

Ultimately, the detailed evaluation metrics and plots of training and validation accuracy and loss are a testament to our results. Consequently, our results show the combination of BERT and TensorFlow to be a highly effective integration for work on spam detection, with gains far superior to traditional methods.

Although the results are promising, there is so much more work to be done. Further improvements could involve the adoption of more advanced data augmentation strategies, testing a broader spectrum of models and checking the GLE algorithm on various datasets to ensure generalization capability. Moreover, addition of live feedback mechanisms based on user input can further improve the model's efficacy and flexibility in anticipating newer spam strategies.

To conclude, our research demonstrates the spam detection or classification potential of modern deep learning models such as BERT. This work presents an invaluable addition not only to the academic understanding of text classification, but also for building better spam detection systems in practical scenarios.

**Integration of Artificial Intelligence in the Advancement of
Science and Engineering
July 2024**

VIII. REFERENCES

- [1].Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.
- [2].Goodman, J., Heckerman, D., & Rounthwaite, R. (2007). Content-based filtering: solutions for spam, adult content, and viruses. *Communications of the ACM*, 50(2), 89-95.
- [3].Schultz, M. G., Eskin, E., Zadok, E., & Stolfo, S. J. (2000). Data mining methods for detection of new malicious executables. In *Proceedings 2000 IEEE Symposium on Security and Privacy. S&P 2000* (pp. 38-49). IEEE.
- [4].Metsis, V., Androutsopoulos, I., & Paliouras, G. (2006). Spam filtering with naive Bayes-which naive Bayes?. In *CEAS* (Vol. 17, No. 2006, pp. 28-69).
- [5].Androutsopoulos, I., Koutsias, J., Chandrinou, K. V., & Spyropoulos, C. D. (2000). An experimental comparison of naive Bayesian and keyword-based anti-spam filtering with personal e-mail messages. In *Proceedings of the 23rd annual international ACM SIGIR conference on Research and development in information retrieval* (pp. 160-167). ACM.
- [6].Carreras, X., & Marquez, L. (2001). Boosting trees for anti-spam email filtering. In *Proceedings of the 4th international conference on Recent advances in natural language processing* (pp. 58-64).
- [7].Quinlan, J. R. (1986). Induction of decision trees. *Machine learning*, 1(1), 81-106.
- [8].Miyamoto, D., Kubo, T., & Nishigaki, M. (2019). Spam Detection Method Using Recurrent Neural Network. In *2019 IEEE International Conference on Consumer Electronics (ICCE)* (pp. 1-4). IEEE.
- [9].Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural computation*, 9(8), 1735-1780.
- [10]. Kim, Y. (2014). Convolutional neural networks for sentence classification. arXiv preprint arXiv:1408.5882.
- [11]. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.

**Integration of Artificial Intelligence in the Advancement of
Science and Engineering
July 2024**

- [12]. Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., ... & Zheng, X. (2016). TensorFlow: A system for large-scale machine learning. In 12th {USENIX} symposium on operating systems design and implementation ({OSDI} 16) (pp. 265-283).
- [13]. Manning, C. D., Raghavan, P., & Schütze, H. (2008). Introduction to Information Retrieval. Cambridge University Press.
- [14]. Sebastiani, F. (2002). Machine learning in automated text categorization. ACM Computing Surveys (CSUR), 34(1), 1-47.
- [15]. Howard, J., & Ruder, S. (2018). Universal language model fine-tuning for text classification. arXiv preprint arXiv:1801.06146.
- [16]. Chutke, V. (2023). Spam-Ham SMS Dataset. Retrieved from <https://www.kaggle.com/datasets/vivekchutke/spam-ham-sms-dataset>