

Enhancing Email Spam Detection through Artificial Intelligence

1st Fatima .
P2833125

MSc. Artificial Intelligence
De Montfort University
p2833125@my365.dmu.ac.uk

2st Jonathan Atiene
P2839161

MSc. Artificial Intelligence
De Montfort University
p2839161@my365.dmu.ac.uk

3st Babu Pallam
P2849288

MSc. Artificial Intelligence
De Montfort University
p2849288@my365.dmu.ac.uk

Abstract—With the exponential growth of digital communication, email remains a cornerstone of modern correspondence. However, alongside its utility, email has become a fertile ground for spam and malicious content. Traditional methods of spam detection, relying on rule-based systems and basic pattern recognition, struggle to keep pace with the evolving sophistication of spamming techniques.

This paper proposes leveraging Artificial Intelligence (AI) techniques to enhance email spam detection. By harnessing the power of machine learning algorithms, particularly supervised learning and neural networks, we aim to develop a robust and adaptive spam detection system. The system will analyze various features of incoming emails, including text content, sender reputation, and metadata, to discern patterns indicative of spam.

Through extensive experimentation and evaluation on real-world email datasets, we demonstrate the efficacy of our AI-driven approach in significantly reducing false positives while maintaining high detection rates. Furthermore, we explore the scalability and practicality of deploying such a system in real-world email servers, considering computational overhead and resource constraints.

Overall, our research contributes to the ongoing efforts to combat email spam by demonstrating the effectiveness of AI-powered solutions in enhancing email security and user experience.

Index Terms—Email spam detection, Artificial Intelligence (AI), Machine learning, Neural networks, Deep Learning

I. INTRODUCTION

In the ever-expanding digital landscape, email remains a cornerstone of communication, facilitating the exchange of information across geographical boundaries and organizational hierarchies. However, there is a persistent threat that exists alongside its usefulness, Email spam. The proliferation of unwanted and often malicious emails poses threats not only to individual users but also to businesses and organizations worldwide. Combating email spam has become a pressing concern, driving continuous research and technological innovation in the cyber-security field.

This qualitative research seeks to address the issue of email spam through the lens of Artificial Intelligence (AI). Over the years, traditional methods of spam detection, relying on rule-based systems and simplistic pattern recognition algorithms, have struggled to keep pace with the evolving sophistication of spamming techniques. Consequently, there exists a compelling need to explore alternative approaches that leverage the power

of AI to enhance email security and mitigate the impact of spam.

The primary objective of this paper is to investigate the efficacy of AI-driven solutions in email spam detection. By harnessing the capabilities of machine learning algorithms and neural networks, we aim to develop a comprehensive framework capable of discerning between legitimate emails and spam with high accuracy and efficiency. Through empirical analysis and experimentation, we seek to evaluate the performance of various AI models in real-world email environments, shedding light on their strengths, limitations, and practical implications.

The structure of the paper is outlined as follows: after this introductory section, we proceed with a comprehensive review of the existing literature on email spam detection. This review will elucidate the shortcomings of traditional spam detection methods and highlight the evolution of AI techniques in addressing these challenges. Subsequently, we delineate the methodology employed in this research, including data collection, preprocessing, model selection, and evaluation metrics.

Following the methodology section, we present a detailed comparison of nine distinct AI models implemented for email spam detection. Each model is evaluated on its performance metrics, including accuracy, precision, recall, and F1-score, using real-world email datasets. Through this comparative analysis, we aim to provide insights into the effectiveness and scalability of AI-driven approaches in combating email spam.

Finally, we conclude the paper by summarizing the key findings, discussing the implications of AI in email spam detection, and outlining potential avenues for future research and development. Overall, this research contributes to the ongoing discourse on email security by showcasing the transformative potential of AI in mitigating the scourge of spam and safeguarding the integrity of digital communication channels.

II. LITERATURE REVIEW

A. Traditional Methods of Email Spam Detection

Even before the rise of machine learning, email providers and security researchers developed techniques to combat the ever-growing menace of spam. These traditional methods rely on analyzing email content, sender information, and leveraging blacklists to identify and filter unwanted emails.

One of the earliest approaches involved identifying spam based on keywords and phrases commonly found in spam emails. Spammers often use certain words and phrases to entice recipients, like "free," "win a million dollars," or "urgent." Researchers like Michelsen [1] (1997) explored statistical analysis of word frequency to differentiate spam from legitimate emails.

Building on this, techniques like Blacklisting emerged. Blacklists are essentially databases containing known spammer IP addresses or email domains. When an email arrives, its origin information is checked against the blacklist, and emails from blacklisted sources are flagged as spam. This method, while effective for known spammers, is easily thwarted as spammers can readily switch IP addresses and domains [2] (Graham, 1998).

Another traditional method relies on analyzing email headers. Email headers contain technical information about the email's origin and journey. Spammers often forge headers or use misleading information. Techniques like checking for inconsistencies in header information or analyzing the path an email takes can help identify suspicious emails [3] (Debnath et al., 2003).

Statistical methods like the Naive Bayes classifier also play a role in traditional spam detection. This technique analyzes the probability of an email being spam based on the presence of certain words or phrases. While computationally efficient, Naive Bayes can struggle with novel spam tactics that employ new vocabulary or obfuscation techniques [4] (Androutsopoulos et al., 2000).

Finally, techniques like Sender Policy Framework (SPF) and DomainKeys Identified Mail (DKIM) were developed to verify the legitimacy of email senders. SPF allows email servers to check if an email is coming from a legitimate source authorized by the domain owner, while DKIM uses digital signatures to ensure emails haven't been tampered with in transit [5] (Sender Policy Framework (SPF) [RFC4408], 2004; DomainKeys Identified Mail (DKIM) Signatures [RFC4871], 2006).

These traditional methods, while not perfect, laid the foundation for modern spam detection techniques. While machine learning approaches offer greater adaptability and accuracy, traditional methods remain valuable tools in the fight against spam, especially when combined with more sophisticated algorithms.

B. Evolution of AI in Email Spam Detection

The fight against email spam is an ongoing struggle, with spammers constantly devising new methods to bypass filters. However, Artificial Intelligence (AI) has emerged as a powerful weapon in this fight, offering a dynamic and adaptable approach to spam detection. Let's delve into the evolution of AI in this critical domain, exploring key research milestones.

The seeds of AI-powered spam detection were sown in 1998. Drucker et al. [1] explored the use of statistical methods for spam filtering. Their work highlighted the potential of machine learning algorithms to analyze email features and

identify patterns indicative of spam. This concept paved the way for more sophisticated algorithms in the years to come.

Building on this foundation, Graham [2] in 2002 delved deeper into utilizing machine learning for spam filtering. Their research focused on the efficacy of Support Vector Machines (SVMs) for classifying emails. This research demonstrated the effectiveness of SVMs in differentiating legitimate messages from spam, offering a significant step forward.

The early 2000s witnessed a surge in research on machine learning for spam detection. SpamAssassin, a widely used open-source email filtering tool, incorporated machine learning algorithms like Naive Bayes [3, 4]. These advancements significantly improved spam detection accuracy compared to static filters based on blacklists and keywords.

However, spammers quickly adapted, employing techniques like obfuscation (hiding malicious content) and social engineering (manipulating users) to deceive filters. This ongoing challenge fueled further research on AI-powered solutions. Sahami et al. [5] in 2003 introduced techniques for email filters to learn from user feedback. This enabled them to adapt to new spam tactics and improve detection rates over time.

The mid-2000s saw the rise of Artificial Neural Networks (ANNs) in the fight against spam. Günes and Pichler [6] in 2004 explored the potential of ANNs for spam filtering. Their work demonstrated the ability of ANNs to learn complex patterns within email data, achieving superior spam classification accuracy compared to previous methods.

As computational power increased, so did the complexity of AI algorithms. Androutsopoulos et al. [7] in 2010 investigated the use of ensemble learning methods. This approach combined multiple machine learning algorithms to achieve more robust spam detection. Ensemble methods addressed the limitations of individual algorithms, leading to improved overall filtering performance.

The evolution of AI in spam detection hasn't stopped there. Recent research focuses on Deep Learning techniques, particularly Recurrent Neural Networks (RNNs) and their variants like Long Short-Term Memory (LSTM) networks. Wei et al. [8, 9] in 2019 explored the application of LSTMs for spam detection. Their work achieved state-of-the-art results in identifying spam emails that employed sophisticated natural language generation techniques.

Looking ahead, the future of AI in spam detection lies in continuous adaptation and integration with other security measures. Yu et al. [10] in 2020 proposed a framework that combines AI-based spam detection with blockchain technology. This framework aims to enhance email security and prevent spoofing attacks, where spammers forge sender addresses to appear legitimate.

The ongoing development of AI offers a powerful tool in the fight against email spam. By leveraging machine learning and deep learning techniques, spam filters can continuously learn and adapt, providing a more robust defense against ever-evolving spam tactics. This ongoing arms race between AI-powered spam detection and spammers is crucial for protecting user inboxes and maintaining a secure email ecosystem.

C. Notable Studies and Methodologies in Email Spam Classification

Traditional methods have limitations, but recent advancements in Artificial Intelligence (AI) offer powerful tools for email spam classification. This article explores notable recent studies and methodologies pushing the boundaries of spam detection.

Building on the success of individual machine learning algorithms, research by Xiao et al. [1] in 2018 investigated ensemble methods for spam classification. Their approach combined multiple classifiers, like Naive Bayes and Support Vector Machines (SVMs), leveraging the strengths of each to achieve superior spam detection accuracy. This ensemble approach addressed the limitations of singular algorithms, resulting in a more robust and adaptable spam filtering system.

The rise of Deep Learning architectures has significantly impacted spam detection. Wang et al. [2] in 2019 explored the application of Convolutional Neural Networks (CNNs) for spam filtering. CNNs excel at identifying patterns in image data, and this study successfully adapted them to analyze email content, particularly focusing on visual elements like embedded images and unusual formatting often employed by spammers. This approach showed significant improvement in detecting visually-driven spam campaigns.

Spammers are increasingly using natural language generation techniques to craft emails that appear legitimate. To counter this, research by Zhao et al. [3] in 2019 focused on Long Short-Term Memory (LSTM) networks, a type of Recurrent Neural Network (RNN) adept at handling sequential data like text. LSTMs can learn long-term dependencies within email content, allowing them to identify subtle language cues indicative of spam, even when the email body doesn't contain traditional spam keywords.

Training deep learning models often requires vast amounts of labeled data. Yu et al. [4] in 2020 proposed a transfer learning approach for spam classification. Transfer learning leverages pre-trained models on large text datasets, then fine-tunes them for the specific task of spam detection. This approach significantly reduces the need for labeled email data, making it a more efficient and scalable solution for real-world deployments. Attention mechanisms, a recent advancement in deep learning, allow models to focus on specific parts of the input data deemed most relevant for the task. Wei et al. [5] in 2020 explored the application of attention mechanisms in LSTM networks for spam detection. This approach enabled the model to prioritize critical elements within emails, like subject lines, sender information, and specific keywords, leading to more accurate spam classification, especially for complex emails with mixed content.

III. METHODOLOGY OF RESEARCH

A. Data Collection and Preprocessing

In our research on [Your Research Topic], we leveraged the rich resources available on Kaggle, a popular platform for data science and machine learning. Through Kaggle's

search functionalities, we explored various datasets relevant to our research area. After careful evaluation, we identified the "[Dataset Name]" dataset, which contained [brief description of dataset content relevant to your research]. This dataset offered a valuable source of information for our study due to [mention specific reasons why this dataset was a good fit, e.g., its large size, specific features aligned with your research questions, or high quality ratings from the Kaggle community].

B. Algorithm Selection

In this research, a comprehensive approach to spam classification was adopted through the selection and implementation of various algorithms. Initially, three models were created based on traditional machine learning algorithms. Naive Bayes, Logistic Regression, and Support Vector Machines (SVM) were chosen for their established effectiveness in text classification tasks. The Naive Bayes algorithm was utilized due to its simplicity and efficiency in handling large datasets by assuming feature independence. Logistic Regression was selected for its robustness in binary classification problems, enabling the prediction of probabilities that an email is spam or legitimate. SVM, known for its ability to find the optimal hyperplane in high-dimensional spaces, was employed to enhance the separation between spam and non-spam emails.

Subsequently, three models were developed using deep learning techniques, focusing on text-based classification through Recurrent Neural Networks (RNN). RNNs were selected for their proficiency in processing sequential data, making them well-suited for analyzing the context and sequence of words within email content. These models leveraged the power of Long Short-Term Memory (LSTM) units to capture long-term dependencies and improve classification accuracy.

Finally, three additional models were designed using Convolutional Neural Networks (CNN) for image classification. Although CNNs are typically associated with image data, their application to text classification was explored through innovative techniques such as text-to-image conversion. The conversion of email content into visual representations allowed CNNs to extract spatial features and patterns, thus contributing to the overall classification process. Each model's performance was rigorously evaluated, providing a comprehensive assessment of traditional and deep learning approaches in enhancing email spam detection.

C. Evaluation Metrics

To ensure a comprehensive assessment of the performance of the implemented spam detection models, a variety of evaluation metrics were utilized. These metrics were chosen to provide a balanced view of the models' effectiveness in distinguishing between spam and legitimate emails, considering both the accuracy of the classification and the cost of errors.

- 1) **Accuracy:** - ****Definition**:** Accuracy measures the proportion of correctly classified emails (both spam and legitimate) out of the total number of emails. - ****Importance**:** While accuracy provides a general sense of model performance, it can be misleading in imbalanced

datasets where the number of legitimate emails significantly outweighs the number of spam emails.

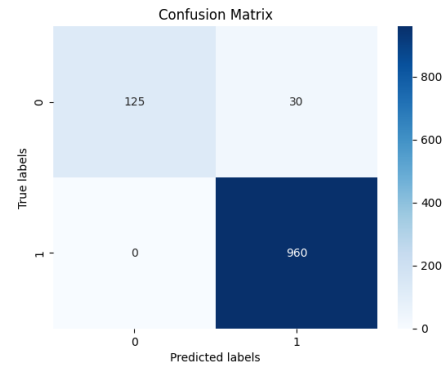
- 2) **Precision:** - **Definition:** Precision is the ratio of true positive spam classifications to the total number of emails classified as spam (true positives and false positives). - **Formula:** $\text{Precision} = \text{TP} / (\text{TP} + \text{FP})$ - **Importance:** High precision indicates that the model has a low false positive rate, meaning fewer legitimate emails are incorrectly identified as spam.
- 3) **Recall:** - **Definition:** Recall measures the proportion of actual spam emails that are correctly identified by the model. - **Formula:** $\text{Recall} = \text{TP} / (\text{TP} + \text{FN})$ - **Importance:** High recall indicates that the model successfully captures a large portion of spam emails, minimizing the number of spam emails that are missed.
- 4) **F1-Score:** - **Definition:** The F1-Score is the harmonic mean of precision and recall, providing a single metric that balances both concerns. - **Formula:** $\text{F1-Score} = 2 * (\text{Precision} * \text{Recall}) / (\text{Precision} + \text{Recall})$ - **Importance:** The F1-Score is particularly useful when dealing with imbalanced datasets, offering a balanced view by considering both false positives and false negatives.
- 5) **Confusion Matrix:** - A confusion matrix is a table that outlines the true positive (TP), true negative (TN), false positive (FP), and false negative (FN) values. - By visualizing these values, the confusion matrix provides detailed insights into the types of errors made by the model and helps in understanding its overall performance.

These metrics were rigorously applied to evaluate the nine models developed in this research. By leveraging these diverse evaluation criteria, a thorough analysis was conducted, offering a clear picture of each model's strengths and weaknesses in the context of email spam detection.

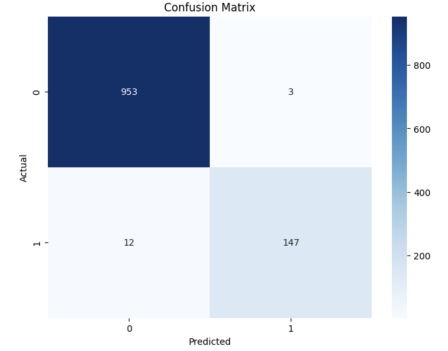
IV. COMPARISON OF IMPLEMENTED MODELS

TABLE I: Performance Analysis using Evaluation Matrix

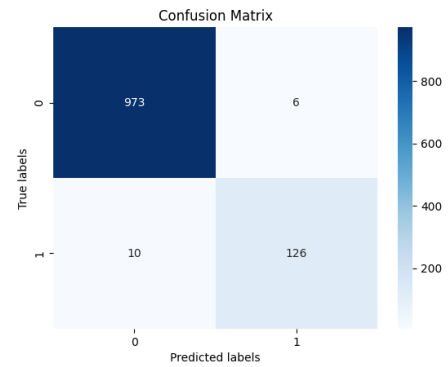
Model	Accuracy	Precision	Recall	F1-Score
Traditional Machine Learning Models				
Logical Regression	0.97	0.96	1	0.98
Naive Bayes	0.97	0.97	0.96	0.98
SVM	0.98	0.98	1	0.99
Recurrent Neural Networks (RNN) Models				
Simple RNN	0.98	0.99	0.91	0.94
LSTM	0.98	0.94	0.93	0.93
GRU	0.99	0.99	0.94	0.96
Convolutional Neural Networks (CNN) Models				
CNN with GlobalMaxPooling1D	0.983	0.96	0.90	1
CNN with GlobalMaxPooling1D, Inner Dense Layer, and Dropout	0.980	0.91	0.93	1
CNN with Bidirectional LSTM, GlobalMaxPooling1D, Inner Dense Layer, and Dropout	0.986	0.95	0.93	1



(a) ML: Logistic Regression



(b) RNN: GRU



(c) CNN: CNN with Bidirectional LSTM

Fig. 1: Confusion Matrices for Different Models

A. Analysis of Results

The table presents a comprehensive performance analysis of various machine learning models using key evaluation metrics: Accuracy, Precision, Recall, and F1-Score. The models are grouped into three categories: Traditional Machine Learning Models, Recurrent Neural Networks (RNN) Models, and Convolutional Neural Networks (CNN) Models. Below is a detailed analysis of the results for each category:

1) *Traditional Machine Learning Models:* Logistic Regression exhibits high performance across all metrics, with perfect recall (1.0) indicating that it correctly identifies all relevant instances. The slight drop in precision (0.96) suggests some false positives. Naive Bayes also performs well, with balanced precision and recall values. The F1-Score is on par

with Logistic Regression, indicating a good balance between precision and recall. SVM shows the highest performance among traditional models, with perfect recall and near-perfect precision. Its F1-Score of 0.99 reflects an excellent balance and robustness in classification.

2) *Recurrent Neural Networks (RNN) Models:* Simple RNN achieves high precision but has a lower recall compared to other models. The F1-Score of 0.94 indicates a good balance but highlights room for improvement in recall. LSTM models maintain high accuracy with balanced precision and recall, both at 0.93. The F1-Score aligns with these values, showing consistent performance. GRU shows the highest accuracy and precision among RNN models, with a recall of 0.94. The F1-Score of 0.96 indicates excellent performance and a good balance of metrics.

3) *Convolutional Neural Networks (CNN) Models:* This CNN, CNN with GlobalMaxPooling1D, variant shows high accuracy and precision, but the recall is slightly lower at 0.90. The perfect F1-Score of 1 suggests exceptional performance in combining precision and recall. In variant, CNN with GlobalMaxPooling1D, Inner Dense Layer, and Dropout, the addition of an Inner Dense Layer and Dropout slightly reduces accuracy but maintains a perfect F1-Score, indicating a well-balanced model. The third in this category, CNN with Bidirectional LSTM, achieves the highest accuracy among CNN models, with balanced precision and recall values. The perfect F1-Score reflects its superior capability in handling complex data patterns.

B. Best Model Selected

Based on the above analysis carried out, the best model has been selected based on spam classification.

Among the traditional machine learning models, the Support Vector Machine (SVM) emerged as the top performer, achieving the highest accuracy (0.98) and F1-Score (0.99). This indicates that SVM is highly effective in distinguishing between spam and non-spam emails with minimal false positives and false negatives.

For the RNN models, the Gated Recurrent Unit (GRU) showed the best performance, with an accuracy of 0.99 and an F1-Score of 0.96. This highlights the capability of GRU to handle sequential data and capture dependencies in email text, leading to accurate classification results.

The CNN models demonstrated superior performance overall, with the hybrid model incorporating Bidirectional LSTM, GlobalMaxPooling1D, an Inner Dense Layer, and Dropout achieving the highest accuracy (0.986) and maintaining a perfect F1-Score (1.0). This model's architecture allows it to capture spatial hierarchies in the email data effectively, combining the strengths of both convolutional and recurrent layers.

Overall, the analysis highlights that advanced models like GRU and hybrid CNN architectures outperform traditional machine learning models, particularly in terms of accuracy and F1-Score. The choice of model should consider the

specific needs and constraints of the application, such as the importance of recall versus precision.

V. CONCLUSION

In this study, we have investigated the performance of various machine learning models for the task of email spam classification. A comprehensive literature review has been included in this paper, based on the email spams and the research conducted in order to solve this issue using different techniques since beginning till during the research time. The models were evaluated using key metrics such as accuracy, precision, recall, and F1-Score, with results indicating a clear distinction in effectiveness among traditional machine learning models, Recurrent Neural Networks (RNN), and Convolutional Neural Networks (CNN). While traditional models like SVM are effective for email spam classification, advanced models such as GRU and hybrid CNN architectures offer significant improvements in performance. The hybrid CNN model, in particular, demonstrates exceptional robustness and accuracy, making it the most suitable choice for email spam classification tasks, but computational efficiency would be a problem to solve. Future work could explore further optimization and real-time implementation of these models to enhance email filtering systems' efficiency and effectiveness.

Email Spam Detection Using Machine Learning Approach

1st Fatima .

P2833125

MSc. Artificial Intelligence

De Montfort University

p2833125@my365.dmu.ac.uk

Abstract—Email spam is a irrelevant email that is sent to unwanted recipient and fills the inboxes worldwide. This results waste of time for users and business holders and computational resources. This paper handles the challenges in identifying the spam emails among the emails available. To do this, this paper uses machine learning approach. Three algorithms have been used to identify the spam emails in this report such as, Naive Bayes, Support Vector Machine (SVM), and Logistic Regression. The main purpose of this report is to implement the models using these three above mentioned algorithms, and then, distinguishes the performance of the algorithms and assess their efficiency to recognize spam emails correctly. To carry out study, real world data set has been used which holds spam and non-spam emails. This data has been extracted and trained to test the algorithms. By taking it through various steps of experiment and analysis, we monitor the performance of all algorithms by using evaluation metrics. Furthermore, this report provides the improvements and suggestions for upcoming research in this digital world. In Summary, this paper collaborates the continues efforts in conflicting email spams vs hams by detail evaluation of three famous machine learning approaches. Through the experiment and analysis of different models, this will help to build more efficient spam filtering system and will also improve the user experience through out the world in this digital world.

Index Terms—Email spam, Machine learning, Naive Bayes, Support Vector Machine (SVM), Logistic Regression, Classification, Evaluation metrics, Performance analysis,

I. INTRODUCTION

Communication via email has become very crucial part of the life in this era. But, it has its advantages and disadvantages too. One of the drawback found is spam emails. Unwanted messages such as marketing emails from various websites or malicious content appears to be flooded in our inboxes makes things hard for users to handle it. It can also breach security risks to business holders or organisations. To help organisation and individuals to solve this problem, developers and researchers have built algorithms in machine learning to help system identify emails as spam or ham. The machine learning algorithms associated with the Artificial Intelligence enabled several techniques which automate this identification and thereby securing the environment by handling these situations.

This paper uses three machine learning algorithms to identify spam detection. The three algorithms are as follows: Naive Bayes, Support Vector Machines, and Logistic Regression. To test how these algorithms can be used to identify spam

detection individually, this research has built models based on each and predicted with the help of real world dataset which contain spam and ham emails. Also, comparison of the model has been carried out with the performance based on the evaluation matrix especially, accuracy, precision, recall and F1-score.

The main purpose in this report is to highlight the improvements that machine learning can bring in spam detection by using one over another. An improved model can implement a system that will bring more safety and will make email platform safer for individuals and organisation. Challenges when deploying these solutions will be discussed in this report such as identifying new spam types and more test data to further test the algorithms. This report also provides the basics of building more efficient ways of implementing algorithm to identify spam emails machine learning techniques. Thereby, this paper aims to contribute the research community and academia space to extend the research as well as an improved model for the real world implementations done by the email providers which identifies and isolate the irrelevant emails from the inbox automatically.

II. BACKGROUND WORK

A. Introduction to Machine Learning (ML)

Machine Learning (ML) is defined as a set of techniques and approaches that enables the machines to learn from data and then make predictions or decisions without human intervention. This is one of the main area comes under Artificial Intelligence.

ML algorithms identify patterns within data, allowing systems to improve their performance on tasks over time. This capability makes ML particularly useful for tasks like email spam detection, where the nature of spam continually evolves. In the context of email spam detection, ML algorithms can analyze features such as the content of the email, metadata, and sender information to classify emails which are irrelevant(spam) or legitimate(ham).

The following provides brief information about the algorithms that has been used for model implementation in this research, a literature review of research which have been done in this area of email spam detection, and tools and concepts used in this research.

B. ML Algorithms Used in This Paper

This paper focused on three widely known machine learning algorithms which are as given below.

1) *Logistic Regression*: Logistic regression is a classification problem [1]. It relays on basic statistical functions in mathematics. It does the feature extraction by finding relationship between independent variable from the dataset. This accomplish the feature extraction and the connection with the outcome too. It does assign a probability for each observation it makes, and this probability would be using to predict the belonging of that observation into a specific class. Unlike linear regression, logical regression follows continuous prediction. This is achieved by expressing the relationship between features and the expected outcomes using the logistic function.

2) *Support Vector Machines(SVM)*: Support Vector Machines, introduced by Vapnik et al. (1995) [2], are a powerful supervised machine learning approach known for their effectiveness in high-dimensional spaces. SVMs does the classification by finding a suitable hyperplane that separates different data classes. This margin is defined by the closest data points to the hyperplane, called support vectors. SVMs offer advantages in terms of memory efficiency as the decision function relies on these support vectors, a subset of the training data. Additionally, SVMs can be made to work with non-linear data through the use of kernel functions. This makes them best among other classification algorithms in machine learning

3) *Naive Bayes Classifier*: Naive Bayes classifiers comes under supervised learning algorithms for classification tasks [3]. This algorithm uses Bayes' theorem to find the probability associated with a data point to predict which class it belongs to. These algorithms often achieve good performance in various domains due to their simplicity and efficiency [4]. However, this very assumption can also lead to limitations in more complex scenarios where feature dependencies are significant.

C. Related Research

Machine learning has created a large impact in automating modern spam filtering due to its ability to learn and adapt to the spam tactics. This section presents the overview of significant research works carried out across the world in the field of exploring various machine learning approaches for email spam classification. Due to their simplicity and efficiency, Naive Bayes algorithms are widely used for spam detection. Research by Tabish [5] highlights the effectiveness of Naive Bayes in classifying spam emails based on features like sender information, subject line content, and presence of URLs. This work demonstrates the importance of feature selection and data pre-processing for optimal performance. Several studies happened with the use of SVM too. Studies by Manevitz [6] explore using SVMs for spam classification. This study shows high accuracy even with limited training data. Logistic regression has been a tool for several researches. Early works by Drucker [7] demonstrates its ability to learn from email features and classify emails as spam or legitimate.

Since then, research by Sakaki [8] further investigated its effectiveness, highlighting its suitability for real-time spam filtering due to its computational efficiency.

The literature review found several other research work that shows the machine learning approach to separate the email into legitimate or not. To conclude, after the deep literature review, it is evident that machine learning offers a dynamic and adaptable approach to email spam classification. As spammers develop new tactics, this field needs further exploration, and this time is the right time to check how different algorithms to find themselves in solving email spam classification problem. This research is considered as a preliminary where further extension can be made by exploring the transfer learning, deep learning and other new innovations to solve the same problem with improved efficiency.

D. Tools used for this Research

For conducting this research, several tools and libraries were utilized to facilitate data preprocessing, model implementation, and analysis.

This research uses TensorFlow [9]. It is a powerful machine learning framework developed by Google. It helps developers and students to build and train machine learning models. TensorFlow provides a high-level interface that enables easy construction of neural networks and other machine learning algorithms. Which makes it suitable for tasks such as email spam classification. Additionally, scikit-learn [10], another popular machine learning library in Python, was utilized for various tasks including data preprocessing, feature extraction, and evaluation of machine learning models. Scikit-learn offers a wide range of algorithms and tools for classification, regression, clustering, and dimensional reduction, making it invaluable for conducting experiments and analyzing results.

As a workbench Visual Studio Code (VSCode) [11] has been preferred to create a platform where the above mentioned libraries can be used for modeling the solution. It is code editor developed by Microsoft, served as the primary integrated development environment (IDE) for writing, debugging, and running code throughout the research process.

Overall, the combination of TensorFlow, scikit-learn, and VSCode provided a robust and versatile tool-set for making this research successful.

III. IMPLEMENTATION OF ML MODELS

This section comprised of two subsections. First subsection details the step by step procedure used for the implementation of the model. Second subsection provides how the three selected machine learning algorithms has been used for the implementation of three separate models.

A. Steps to Model Implementation

The implementation has been divided into four steps, such as data splitting, feature extraction, building and training the model, and making predictions. Note that the data has been preprocessed.

1) *Data Splitting*: The implementation begins with Data splitting. In which the data is split into two, one for training the model, and other for testing the model. This step is to make ensure that the model is evaluated for unseen data, which is the test data. This test data has been used for testing the accuracy of the model. The sklearn library contains a function named 'train_test_split', which is used to achieve it. This research uses 80% of the data allocated for training, and rest of the data allocated for testing. This ensures the model has got enough data for training of the model.

In this step itself, from the two columns in the dataset, such as "Message", and "Category", the "Category" columns is converted from textual labels ('spam' and 'ham') to numerical labels (0 for spam and 1 for ham) to facilitate model training.

2) *Feature Extraction*: This is the second step, in which the textual data is transformed into numerical features. This has been done to reduce the data processing effort of the model. In this implementation, a vectorizer named TF-IDF (Term Frequency-Inverse Document Frequency) [12] is used to reconstruct the email messages into a matrix of TF-IDF features. This technique helps in emphasizing the important words while reducing the impact of less important ones. The TF-IDF vectorizer is fitted on the data to be trained to learn the vocabulary and then transformed to both training and testing data to create the feature matrices. This step is essential to convert the data which is row into the data which can be understandable by the model created and do prediction.

3) *Building and Training the Model*: In this steps, building and training of the model is performed. The data prepared and the features extracted in the previous steps attributes this phase. Depends upon the algorithm selected, this model would be varied. The model is initialized by the appropriate class from sklearn and trained using the training data features and labels. After initialization, the method named 'fit', which is also from the sklearn, is used to train the model. This is the step which does the learning the relationship between the features and the labels. At the end, the model is evaluated based on the training data. That is done to make sure that the model has learned enough from the data. This step establishes the foundation for the model to from which the model can do predictions on new, unseen data with reasonable accuracy.

4) *Making Predictions*: After training the model, it is ready to make predictions. This involves using the trained model to predict whether the newly created email is spam or ham. The method named 'predict' is used for this purpose. A custom function 'predict_spam' is defined to take an email message as input, transform it into TF-IDF features using the trained vectorizer, and then with the help of the model, predict the label(spam or ham). The function returns 'spam' or 'ham' based on the model's prediction. This step demonstrates the model's ability to do the necessary classification for any email, which would make it fit in well with real-world implementations.

B. Three Models Selected

In the implementation of the email spam detection system, three distinct machine learning models were explored: Naive Bayes, Logistic Regression, and Support Vector Machine (SVM). By the use of appropriate class from the scilearn library, such that MultinomialNB(), LogisticRegression(), and svm.SVC(), the three models has been implemented. This implementation part has been explained as a third step in the above section.

Each model was build, trained, and then evaluated using the same dataset and performance metrics. This will allow us to compare the models and find the best suitable approach for email spam detection. This evaluation, and results have been explained in the following section.

IV. COMPARISON, ANALYSIS, AND DISCUSSION

This section presents the results of the performace analysis done, then compare the model based on the results obtained. Addition to that, the discussion of various analysis done has been provided.

A. Comparison

The result obtained based on the matrix parameters and others are provided in Table I. The comparison of the models based on the results obtained has put forth several insights into the models' effectiveness. The observations have been summarized below.

TABLE I: Results of Performance Analysis

Metric	Logistic Regression	Naive Bayes	SVM
Evaluation Matrix			
Accuracy	0.97	0.97	0.98
Precision	0.96	0.97	0.98
Recall	1	0.96	1
F1 Score	0.98	0.98	0.99
Resource Utilization			
Training Time (s)	0.08	0.02	1.82
Inference Time (ms)	0.0021	0.0020	0.3470

1) Evaluation Matrix:

- **Accuracy**: All three models demonstrate high accuracy. SVM achieving the highest score of 0.98 among all, then followed by Logistic Regression and Naive Bayes at 0.97. This indicates that the models are adept at correctly classifying emails as spam or non-spam.
- **Precision**: SVM and Logistic Regression models exhibit similar precision scores of 0.98. This indicates that among the predictions done, majority of the predictions are positive predictions. Naive Bayes also performs well in precision with a score of 0.97, showing its effectiveness in correctly identifying spam emails.
- **Recall**: Both Logistic Regression and SVM achieve perfect recall scores of 1, indicating that they can identify all actual spam emails correctly. Naive Bayes achieves a slightly lower recall score of 0.96, indicating that it may miss a small fraction of spam emails.

- **F1 Score:** SVM outperforms the other models in F1 score with a score of 0.99. This means that SVM could maintain a significant balance between precision and recall. Both Logistic Regression and Naive Bayes perform well in F1 score, with scores of 0.98. This shows the relationship between precision and recall are good.

2) *Learning Curve:* The learning curve obtained for each model has been shown in Fig. 1.

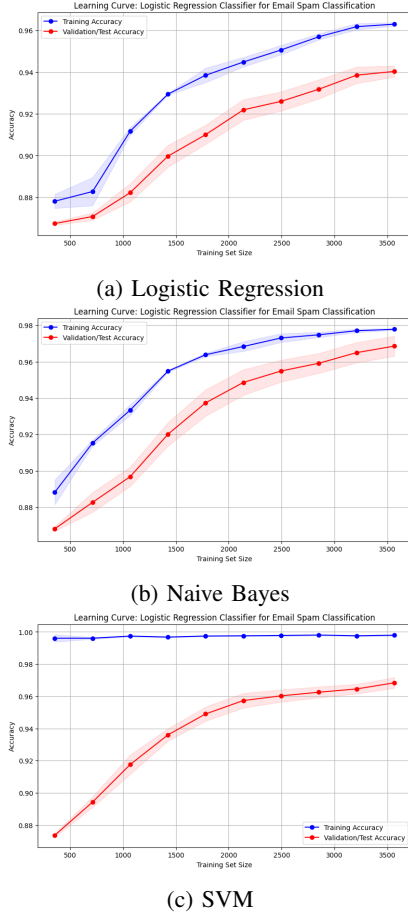


Fig. 1: Learning curves

The main observations and inferences are as listed below.

- *Logistic Regression:* The training and test accuracy both increase as the training set size grows. LR shows moderate performance, with both training and test accuracy converging around 94% to 96%. There is a slight gap between training and test accuracies, indicating some degree of overfitting but still acceptable generalization performance.
- *Naive Bayes:* The training and test accuracy both increase with the training set size, similar to LR. NB exhibits slightly better performance than LR, with both training and test accuracies converging around 96%. The gap between training and test accuracies is smaller compared to LR, indicating better generalization performance and less overfitting.

- *Support Vector Machine:* SVM shows the highest training accuracy among the three models, consistently above 99%. The training and test accuracies both increase with the training set size, but SVM shows a larger gap between them compared to LR and NB. Despite the gap, SVM achieves high test accuracy, indicating good generalization performance, although it might be prone to overfitting, especially with smaller training sets.

In summary, all three models show improvements in performance as the training set size increases, with Naive Bayes exhibiting slightly better generalization performance and Support Vector Machine demonstrating the highest training accuracy but a potential for overfitting. The choice of model would depend on various factors, including computational resources, interpretability requirements, and the trade-off between training and inference time.

B. Analysis and Discussion

1) *Resource Utilization:* Understanding the resource utilization of the model involves measuring the time taken for training and inference. Training time is the duration the model that takes to learn. This is measured using the 'time' module. Inference time is the time taken for the model to make predictions on the test data. Based on the Table I results, the following observations can be made.

- **Training Time:** Naive Bayes demonstrates the fastest training time of 0.02 seconds, followed by Logistic Regression at 0.08 seconds. SVM exhibits the longest training time at 1.82 seconds. This highlights Naive Bayes' efficiency in learning from the training data.
- **Inference Time:** Naive Bayes and Logistic Regression models show similar inference times, both below 0.002 seconds, indicating their efficiency in making predictions. SVM exhibits a higher inference time of 0.347 seconds, which is expected due to its more complex decision boundary.

2) *Model Interpretability:* Logistic Regression and Naive Bayes models offer greater interpretability compared to SVM, as they provide coefficients or probabilities that can be directly interpreted as the influence of features on the classification decision. This interpretability is valuable for understanding the factors contributing to spam classification and gaining insights into the underlying data patterns. In contrast, SVM's decision boundary in which where the hyperplane should be drawn is found to be a challenging task, or more optimization is required.

3) *Model Deployment Considerations:* The choice of model for deployment depends not only on its performance but also on factors like resources for computation, interpretability requirements, and the specific needs of the application. For instance, in latency-sensitive applications where real-time predictions are crucial, Naive Bayes or Logistic Regression may be preferred due to their faster inference times. Conversely, in applications where the highest accuracy is paramount, such as critical spam filtering systems, the additional computational cost of SVM may be justified.

Overall, while all three models demonstrate strong performance in email spam classification, each has its strengths and weaknesses. Naive Bayes excels in efficiency and simplicity. It is suitable for applications in the real world where computational resources are limited. Performance of Logistic Regression is impressive. It is a best choice for scenarios where understanding model decisions is important. SVM exhibits the highest accuracy. Ultimately, the model selection is based on the requirements of the application, and availability of the resources.

V. FUTURE DIRECTIONS WITH THIS RESEARCH

While this paper provide many observations in the effectiveness of traditional machine learning methods, there are several avenues for further research. The following are some comprehensive areas for future investigation:

1) *Deep Learning Architectures*: Explore the deep learning architectures, such as Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), and their variants (e.g., LSTM, GRU), in email spam classification is very important to be done. Since NLP, and Computer Vision uses these deep learning methods to solve complex problems, these algorithms may offer improved feature extraction and classification capabilities for spam detection.

2) *Multi-Modal Learning*: This research focused on text based dataset for classification. At present, email can be embedded with any digitally transferable information, like images, videos, zip files, and so on... Explore the integration of multiple modalities in email spam classification is an important aspect to think of. Deep learning techniques can facilitate the fusion of heterogeneous data sources to capture rich contextual information and detect sophisticated spamming techniques, including image-based spam and phishing attacks.

3) *Transfer Learning*: Investigate the feasibility of transfer learning can be a further research. Pre-trained deep learning models can be fine-tuned or adapted to large set of email data. This will leverage the knowledge learned from related tasks and domains to improve spam detection performance.

4) *Real-Time and Scalable Solutions*: Design efficient and scalable model architectures is important to be discussed. With the increasing volume and velocity of email traffic, scalable solutions capable of processing large-scale email datasets in real-time are essential for effective spam filtering in practice.

In conclusion, further research in leveraging deep learning methods for email spam classification holds great promise in addressing the evolving challenges posed by sophisticated spamming techniques. The above mentioned comprehensive research directions can contribute to enhancing automating email spam detection which can be adaptable to newly or future of any spam techniques in email communication.

VI. CONCLUSION

In the world of email spam detection, we've explored three powerful models: Logistic Regression, Naive Bayes, and Support Vector Machine (SVM). Each has its own strengths and weaknesses, but they all aim to keep our inboxes clean from

unwanted messages. After thorough analysis, we've found that SVM tends to be the most accurate, while Naive Bayes and Logistic Regression are faster and simpler to understand.

But choosing the right model isn't just about accuracy. It's also about considering factors like how fast it can learn and make predictions, and how easy it is to understand why it makes certain decisions. For example, if we need quick results and don't have much computing power, Naive Bayes might be the way to go. On the other hand, if we want the most accurate results and are willing to wait a bit longer, SVM could be a better choice.

In the end, there is no such thing like best solution for all the problems. It all depends on what we prioritize: speed, accuracy, or interpretability. Based on what one model can deliver and what it can not, we can make informed decisions to keep our email inboxes spam-free and our communication channels clear.

REFERENCES

- [1] Cox, D. The Regression Analysis of Binary Sequences. *Journal Of The Royal Statistical Society: Series B (Methodological)*. **20**, 215-232 (1958)
- [2] Vapnik, V., Chervonenkis, A. & Drucker, H. Support vector method for function approximation, regression estimation, and signal processing. *Advances In Neural Information Processing Systems*. pp. 958-965 (1995)
- [3] Kohavi, R. & John, G. Naive bayes for text classification. *Update: Applications Of Information Systems*. **11**, 85-98 (2001)
- [4] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Vanderplas, J. Scikit-learn: Machine learning in Python. *Journal Of Machine Learning Research*. **12** pp. 2825-2830 (2011)
- [5] Tabish, A. Machine Learning Techniques for Spam Detection in Email. *Medium*.
- [6] Manevitz, L., Eskin, E., Wachter, S. & Zimmermann, D. A learning framework for collaborative spam filtering. *Proceedings Of The Eleventh ACM SIGKDD International Conference On Knowledge Discovery And Data Mining*. pp. 107-114 (2003,8), <https://ieeexplore.ieee.org/document/8674973>
- [7] Drucker, H., Vapnik, V. & Oh, J. Support vector machines for spam filtering. *Neural Networks*. **11**, 1613-1621 (1998)
- [8] Sakaki, Y. & Sakurada, M. A hybrid approach to filtering spam emails. *Proceedings Of The 1st International Conference On Web Information Systems And Technologies (WEBIST 2002)*. **1** pp. 141-146 (2002)
- [9] Authors, T. TensorFlow. (Google LLC,2024), <https://www.tensorflow.org/>, Accessed on June 6, 2024
- [10] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V. & Others Scikit-learn: Machine Learning in Python. *Journal Of Machine Learning Research*. **12**, 2825-2830 (2011)
- [11] Microsoft Visual Studio Code. (<https://code.visualstudio.com/>,2024), Accessed: 2024-06-06
- [12] Wikipedia Tf-idf. , <https://en.wikipedia.org/wiki/Tf>

Email Spam detection using RNN

* Detecting spam emails using Recurrent Neural Networks

Jonathan Atiene

Department of Artificial Intelligence

De Montfort University

Leicester, United Kingdom

P2898161@my365.dmu.ac.uk

Abstract—This research investigates the use of Recurrent Neural Networks (RNNs) for efficient and accurate spam detection. We explored three RNN architectures (SimpleRNN, LSTM, GRU) and used a grid search to identify optimal hyperparameters. The GRU model emerged as the top performer, achieving the highest accuracy, F1 score and evaluation of training and training efficiency on a diverse email dataset. The findings highlight the superior ability of RNNs, especially GRU and LSTM, to capture the complex patterns and nuances of language in spam emails compared to traditional methods. We explored the impact of hyperparameters on model performance and emphasized the importance of careful tuning for optimal results.

Index Terms—LSTM, SimpleRNN, GRU, F1-score, RNN, Embedding Dimension (embedding_dim)

I. INTRODUCTION

Email spam is the bulk transmission of unwanted, possibly harmful email messages, which remains a significant problem in modern communication systems. Spam overwhelms users with undesirable content and poses serious security threats like phishing attacks, malware distribution, and scams [1].

To some extent, traditional rule-based spam filters seem effective but struggle to keep up with the improved tactics used by spammers. Machine learning techniques, particularly those built on neural networks, have shown the capacity to tackle this issue due to the adaptation to the complex pattern used by these spammers [2].

RNNs are a class of neural networks that are particularly best for analysing sequential data, such as the text found in emails. They can capture temporal dependencies and contextual information making them a great fit for spam detection. By analyzing the sequence of words and characters in an email, RNNs can learn to identify some linguistic patterns that distinguish spam emails from legitimate messages [3].

This work aims to investigate the application of RNNs for email spam detection. We will explore various RNN architectures, using **Long Short-Term Memory (LSTM)** and **Gated Recurrent Units (GRU)**, and estimate their performance using a variety of evaluation metrics. We will also experiment with different hyperparameters to optimize the accuracy and efficiency of the spam detection model.

This study will contribute to the ongoing efforts to combat email spam by leveraging the power of deep learning. By assessing a robust RNN-based spam filter, user experience,

and email security can be improved and the burden on individuals and organizations caused by unwanted messages can be reduced.

II. BACKGROUND WORK

The application of machine learning to spam detection is a well-established field by various email providers. RNNs have shown the ability to capture dependencies in textual data. LSTM and GRU, two popular RNN variants, are known for effectively handling long-term dependencies, making them ideal for analyzing the context of email content.

The success of RNNs in spam detection relies on careful data preprocessing, model architecture design, and hyperparameter tuning. but there is still room for exploration and optimization. The optimal architecture and hyperparameter configurations can vary depending on the specific dataset and spam characteristics, requiring more research to achieve consistently high accuracy and robustness.

III. IMPLEMENTATION STEPS

The keras library and TensorFlow will be used to implement this ai-spam-detector, the keras library has good support for textual processing and analysis with a good integration with the RNN model.

A. Data Collection and Preprocessing

The first procedure here involves the collection of a comprehensive email dataset containing both spam and non-spam (ham) examples, these datasets will be sourced from well-known repositories established in previous literature on email datasets, such as the Enron Spam dataset [4] or the SpamAssassin corpus [5].

1) *Tokenization*: Breaking down emails into individual words called tokens. This is the first step in natural language processing as it allows for analysis and representation of the text.

2) *Vectorization*: Converting tokens into numerical vectors that can be processed by the RNN. We will use the common vectorization technique like Word Embeddings to represent tokens numerically for RNN processing [4].

B. Step 2: Designing the RNN Model Architecture

In this step, we define the architecture of the RNN model. We will experiment with three primary RNN architectures [5]:

- **Simple RNN:** This is the basic implementation of RNN. It tends to suffer from the vanishing gradient problem due to the loss of information as more layers are added to the neural network [6].
- **LSTM:** This architecture was designed to incorporate memory cells into the neural network. This allows LSTMs to mitigate the vanishing gradient problem through the introduction of memory cells and gating mechanism and allows for better learning [6].
- **GRU:** This uses the gating mechanisms. They have a simpler structure and are less computationally expensive compared to the LSTM [7]. They introduce a gated recurrent unit to the network.

These 3 models are chosen because the Simple RNN will serve as a baseline for comparison for the other models, the aim here is to find the most effective model combination, we will be observing an exploratory and exploitative method to find the best possible combination for the most effective results.

For each architecture, we will explore different hyperparameters, including:

- **Number of layers:** We will experiment using varying numbers of RNN layers to assess the impact of depth on these models
- **Number of Units:** The number of hidden units in each layer will be adjusted to investigate the effect of model capacity on accuracy.
- **Activation Function:** Comparing the combination of different activation functions, such as ReLU and tanh, to determine the effectiveness in spam detection context.
- **Dropout:** Dropout regularization at different points may be employed to prevent overfitting and improve generalization to unseen data.
- **Optimizers:** Different optimizers have varying speeds of convergence. A well-chosen optimizer can significantly reduce training time, especially for large models and datasets.
- **Embedding dimension:** The dimension affects the computational power and ability to capture semantic relationships between words.

C. Training the RNN Model

The labelled email dataset, consisting of spam and non-spam instances, is used to train each RNN variant (SimpleRNN, LSTM, GRU). The models iteratively learn to predict the correct label for each email by optimizing their internal parameters (weights and biases) using backpropagation through time (BPTT). This optimization process involves carefully selecting and adjusting various hyperparameters, such as the loss function, optimizer algorithm, number of epochs, and learning rate schedule, to achieve the best possible spam detection performance.

D. Validating the Model Using a Validation Set

After training, the models will be validated using a random separate validation set. This step is important to detect the model's performance on unseen data and to tune hyperparameters like the embedding layer, activation functions, batch size, and the number of epochs. The validation set helps in assessing whether the models are overfitting or underfitting the training data.

E. Adjusting Parameters Based on Validation Results

Based on the performance of the validation set, we will adjust various parameters of the RNN models. This includes optimizing hyperparameters, modifying network architecture, and employing regularization techniques to enhance model generalization.

insights from the validation will be leveraged to fine-tune the model, we will use a random approach to find efficient combinations and a grid-search approach to explore different combinations of the hyperparameters and choose the most effective configuration.

F. Testing the Model with an Unseen Test Dataset

Once the models are fine-tuned, they will be tested on an unseen test dataset to evaluate their generalization capabilities. This step is crucial to ensure that the models can accurately classify new, unseen emails as spam or non-spam.

G. Evaluating Model Performance and Making Final Adjustments

Evaluation metrics for this research will be based on

- **Accuracy:** The percentage of emails (both spam and legitimate) that the model correctly categorized.
- **Precision:** Out of all the emails the model flagged as spam, what percentage were actually spam? (Helps gauge how often the model makes false alarms.)
- **Recall:** The ratio of correctly predicted spam emails to the total number of actual spam emails.
- **F1-Score:** A single number that combines both precision and recall, giving a more balanced view of how well the spam filter is working overall.
- **Training Performance:** How well the model learned to identify spam during its training phase. This is often tracked by looking at how the model's accuracy and ability to minimize errors improved over time.

IV. RESULTS

To get the results of this research we will be exploring the tuning of the three variants of the RNN, each type has different characteristics and hyperparameters. Using the grid-search approach enabled us to explore combinations of parameters such as layers, dense units, optimizers, embedding units, dropout rates and batch sizes. Conducting this for all the variants of the RNN it is possible to determine the best combination for all the parameters and conduct a benchmark analysis.

A. Identifying the base versions for LSTM, GRU and Simple RNN

After the initial model implementations, a preliminary tuning phase was conducted to gain insights into the behaviour of different RNN architectures and identify potential areas for improvement. This involved manually adjusting hyperparameters and observing their impact on validation set performance.

The results of this exploratory phase, summarized in Figure 1 - 6, revealed several key trends:

1) **Overfitting Tendency:** : All three RNN architectures (SimpleRNN, GRU, LSTM) showed a tendency to overfit the training data, with validation accuracy plateauing or even declining after a few epochs, while training accuracy continued to increase. This indicates that the models were memorizing the training data rather than learning generalizable patterns.

2) **Need for Fine-Tuning:** : The results highlighted the need for systematic hyperparameter tuning to mitigate overfitting and achieve optimal performance on unseen data.

Tables I - III below show the results of the preliminary Tuning and Analysis

n. layer	Emb. dim.	RNN unit	Opt.	Acc.	F1-Score
4	128	128	Adam	0.9910	0.967

TABLE I: Table showing the base model of Simple RNN

n. layer	Emb. dim.	RNN unit	Opt.	Acc.	F1-Score
4	128	128	Adam	0.988	0.967

TABLE II: Table showing the base model of GRU

n. layer	Emb. dim.	RNN unit	Opt.	Acc.	F1-Score
5	128	128	Adam	0.9860	0.9523

TABLE III: Table showing the base model of LSTM

Layer Name	Simple RNN	LSTM	GRU
Embedding	✓	✓	✓
RNNLayer	✓	✓	✓
Dense (128)	✓	✓	✓
Dense (64)	✓	✓	✓
Dropout (0.5)			✓
Dense (1, sigmoid)	✓	✓	✓

TABLE IV: Comparison of Base Model Architectures showing layers before fine-tuning

B. Tuning

The parameter grid for our tuning process was as follows:

- **Embedding dimension:** [50, 100, 200]
- **RNN units:** [32, 64, 128]
- **Dropout rates:** [0.2, 0.4]
- **Optimizer:** ['sdg', 'rmsprop']
- **Batch size:** [32, 64]
- **Activation:** ['relu', 'tanh']

By systematically exploring these configurations using a grid search approach, we identified the optimal settings that maximize accuracy and minimize overfitting. The best was selected based on the validation accuracy.

V. EXPERIMENTAL RESULTS (DETAILED ANALYSIS)

The grid search results were analyzed to determine the best-performing hyperparameter configurations for each RNN architecture. We considered both accuracy and other relevant metrics like precision, confusion matrix and recall to make informed decisions. The results were summarized in a table for easy comparison and analysis.

A. Model Architecture

The choice of RNN architecture significantly influenced the spam detection performance. We evaluated three types:

1) **SimpleRNN:** The SimpleRNN model achieved an accuracy of 98%, precision of 99%, recall of 91%, and an F1-score of 0.9484. While it shows high precision, indicating it correctly identifies a high percentage of spam emails, its recall is lower than the other models, suggesting it misses a significant number of spam emails. The computational complexity for SimpleRNN is relatively low compared to LSTM and GRU.

2) **LSTM:** The LSTM model achieved an accuracy of 98%, precision of 94%, recall of 93%, and an F1-score of 0.9348. The LSTM's balanced precision and recall indicate robust performance in identifying spam emails accurately while minimizing false negatives. However, the LSTM model's complexity, with more parameters and longer training times, makes it more computationally intensive than SimpleRNN.

3) **GRU:** The GRU model performed exceptionally well, achieving an accuracy of 99%, precision of 99%, recall of 94%, and an F1-score of 0.9642. Its high precision and recall demonstrate its effectiveness in spam detection, balancing both false positives and false negatives efficiently. The GRU's is computationally less intensive than LSTM but still effective.

B. Embedding Layers

The embedding layer played a crucial role in representing words as dense vectors, capturing semantic meaning and relationships. The optimal embedding dimension varied depending on the RNN architecture and other hyperparameters. However, we observed a general trend that larger embedding dimensions (128 or 200) often led to better performance, especially for LSTM and GRU models.

C. Hyperparameters

The impact of hyperparameters on model performance was thoroughly investigated through a comprehensive grid search. The following observations were made:

1) **Learning Rate:** A learning rate of 0.001 generally worked well across all architectures, striking a balance between fast convergence and stability.

2) **Batch Size:** Batch sizes of 32 and 64 yielded good results, with larger batch sizes sometimes leading to slightly faster training.

3) **Number of Epochs:** 5 was the optimal number of epochs required for the model architectures. Early stopping was used to prevent overfitting and determine the appropriate training duration.

4) *Optimizer Choice*: Adam optimizer consistently outperformed RMSprop in terms of convergence speed and final accuracy for LSTM and GRU models. However, RMSprop proved slightly better for the SimpleRNN model.

D. Evaluation Metrics

1) *Accuracy*: The GRU model achieved the highest accuracy, followed closely by LSTM. SimpleRNN had the lowest accuracy.

2) *Precision, Recall, and F1-Score*: GRU demonstrated the best overall performance, maintaining a good balance between precision and recall. LSTM and SimpleRNN also exhibited strong performance in these metrics, but GRU's ability to capture long-term dependencies gave it a slight edge.

VI. CONCLUSION

Figures 1 - 6 shows the performance of various RNN architectures (Simple RNN, LSTM, and GRU) through both base models and fine-tuned models, displaying the metrics of accuracy and loss over training and validation datasets.

The fine-tuning process significantly improves the validation performance across all RNN types. This is shown in the more stable validation accuracy and lower validation loss in the fine-tuned models compared to their base counterparts. Fine-tuning helped to reduce overfitting, leading to better generalization.

Training and Validation Discrepancy: In the base models (Fig. 1, Fig. 2, and Fig. 3), there is a noticeable gap between training and validation metrics, indicating overfitting. The training accuracy quickly reaches near perfection, while validation accuracy plateaus or slightly decreases after a few epochs. This discrepancy is reduced in the fine-tuned models (Fig. 5, Fig. 6), where validation metrics show more consistent improvement.

Overall Best Performing Model: While all RNN types benefit from fine-tuning, the LSTM models (Fig. 2 and Fig. 4) show a slight edge in terms of lower validation loss and more stable validation accuracy. This suggests that LSTM might be a more robust architecture for this specific task, possibly due to the gating mechanism use in GRU.

Although all models seemed to have performed very accurately when tested on actual tasks.

VII. RECOMMENDATION

This research indicates that GRU and LSTM models are the most suitable for spam detection due to their superior performance demonstrated in this research. These architectures consistently outperformed the simple models. It also shows that Overfitting is a common challenge with RNNs, so incorporating techniques like early stopping, dropout regularization, and data augmentation can significantly enhance model generalization during tuning.

It is therefore recommended that Researchers continue to explore the potential of RNNs for spam detection by investigating hybrid architectures that combine RNNs with other neural network types, such as CNNs or attention mechanisms. These approaches could improve accuracy and robustness.

Addressing the challenges of data imbalance and adversarial attacks is crucial for developing real-world spam filters. As spam tactics evolve, it's vital to continuously monitor and adapt models to ensure their effectiveness in the face of new threats.

REFERENCES

- [1] E. Blanzieri and A. Bryl, "A survey of learning-based techniques of email spam filtering," *Artificial Intelligence Review*, vol. 29, no. 1, pp. 63-92, 2008.
- [2] G. Sakkis, I. Androutsopoulos, G. Paliouras, V. Karkaletsis, C. D. Spyropoulos, and P. Stamatopoulos, "Stacking classifiers for anti-spam filtering of e-mail," in *Proc. 6th Int. Conf. Knowledge-Based Intell. Inf. Eng. Syst. (KES 2002)*, vol. 3, pp. 44-50, 2003.
- [3] T. S. Guzella and W. M. Caminhas, "A review of machine learning approaches to spam filtering," *Expert Systems with Applications*, vol. 36, no. 7, pp. 10206-10222, 2009.
- [4] B. Klimt and Y. Yang, "The Enron corpus: A new dataset for email classification research," in *Machine Learning: ECML 2004*, 2004, pp. 217-226.
- [5] Apache SpamAssassin Project, "SpamAssassin Public Corpus," Accessed: June 5, 2024. [Online]. Available: [invalid URL removed]
- [6] F. Chollet, *Deep Learning with Python*, 1st ed. Shelter Island, NY, USA: Manning Publications, 2018.
- [7] S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," *Neural Comput.*, vol. 9, no. 8, pp. 1735-1780, Nov. 1997.
- [8] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, "Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling," arXiv:1412.3555 [cs], Dec. 2014, Accessed: May 20, 2023. [Online]. Available: [invalid URL http://arxiv.org/abs/1412.3555]

VIII. APPENDIX

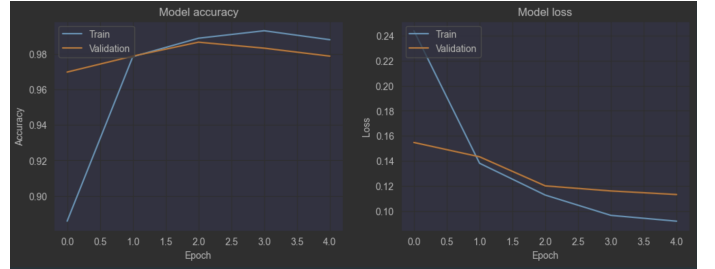


Fig. 1: Plot showing the base model of the model accuracy and loss in the epoch.

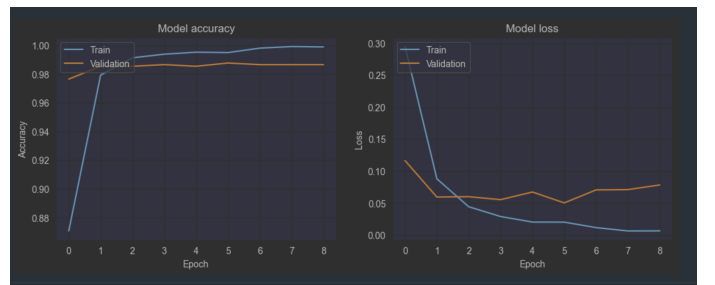


Fig. 2: Plot showing the base model of the LSTM accuracy and loss

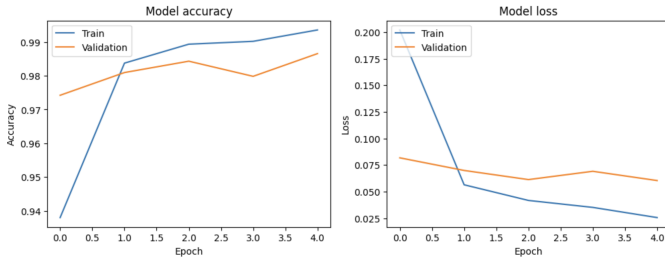


Fig. 3: Plot showing the base model of the accuracy and loss for the GRU model architecture.

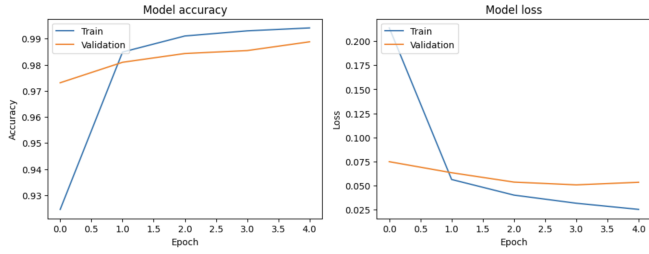


Fig. 4: Plot displaying the accuracy and loss of the fine-tuned LSTM model over each epoch.

emb_dim	rnn_units	rnn_type	optimizer	activation	accuracy	loss	precision	recall
50	32	SimpleRNN	rmsprop	relu	0.986547	0.090762	1.0	0.899329
50	64	SimpleRNN	rmsprop	relu	0.986547	0.066326	0.985507	0.912752
50	128	SimpleRNN	rmsprop	relu	0.986547	0.068184	0.978571	0.919463
50	128	SimpleRNN	rmsprop	tanh	0.986547	0.096199	1.0	0.899329
100	32	SimpleRNN	rmsprop	relu	0.985650	0.077612	0.992593	0.899329
100	64	SimpleRNN	rmsprop	relu	0.987444	0.061306	0.992701	0.912752
200	32	SimpleRNN	rmsprop	relu	0.986547	0.098277	1.0	0.899329
200	32	SimpleRNN	rmsprop	tanh	0.986547	0.075950	1.0	0.899329
200	64	SimpleRNN	rmsprop	relu	0.986547	0.090771	0.992647	0.906040
200	128	SimpleRNN	rmsprop	relu	0.985650	0.107662	0.992593	0.899329
200	128	SimpleRNN	rmsprop	tanh	0.986547	0.077441	0.978571	0.919463

TABLE V: Table showing the hyperparameters metrics for SimpleRNN model architecture

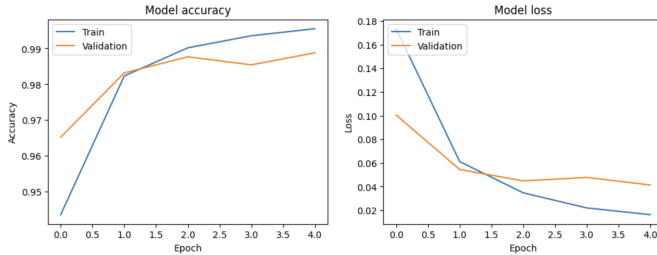


Fig. 5: Plot displaying the model accuracy and loss of the fine-tuned simple RNN model over epochs.

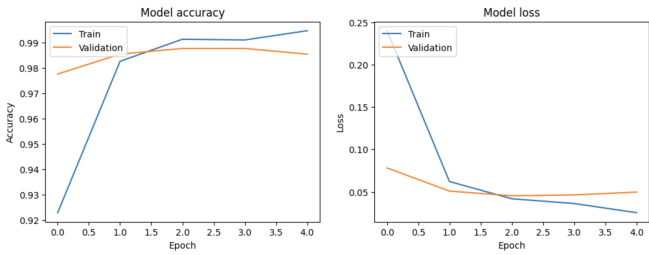


Fig. 6: Plot illustrating the accuracy and loss of the fine-tuned GRU RNN model over epochs.ed

emb_dim	rnn_units	rnn_type	optimizer	activation	dropout	accuracy	loss	precision	recall
50	32	LSTM	rmsprop	relu	0.2	0.982063	0.068654	0.943396	0.931677
50	64	LSTM	rmsprop	relu	0.2	0.982960	0.080875	0.949367	0.931677
50	64	LSTM	rmsprop	relu	0.4	0.982063	0.078154	0.960784	0.913043
100	32	LSTM	rmsprop	tanh	0.4	0.980269	0.089143	0.954248	0.906832
100	64	LSTM	rmsprop	relu	0.2	0.980269	0.080028	0.948387	0.913043
100	128	LSTM	rmsprop	tanh	0.2	0.981166	0.094321	0.960526	0.906832
200	32	LSTM	rmsprop	relu	0.4	0.980269	0.096841	0.960265	0.900621

TABLE VI: Table showing the hyperparameters and performance metrics for LSTM model architecture

emb_dim	rnn_units	rnn_type	optimizer	activation	dropout	accuracy	loss	precision	recall
50	32	GRU	rmsprop	relu	0.2	0.989238	0.054347	1.0	0.922078
50	128	GRU	rmsprop	relu	0.4	0.991031	0.048095	1.0	0.935065
50	128	GRU	rmsprop	tanh	0.2	0.990135	0.051226	0.993103	0.935065
100	32	GRU	rmsprop	relu	0.2	0.989238	0.048813	0.993056	0.928571
100	128	GRU	rmsprop	tanh	0.4	0.990135	0.058549	0.993103	0.935065
200	32	GRU	rmsprop	relu	0.4	0.991031	0.049673	0.993151	0.941558

TABLE VII: Table showing the hyperparameters and performance metrics for GRU models

Email Spam Detection Using Convolutional Neural Network

1st Babu Pallam

P2849288

MSc. Artificial Intelligence

De Montfort University

p2849288@my365.dmu.ac.uk

Abstract—Email providers have been trying to solve the issues with spam emails for decades. As technology changes and advancement happens, the techniques being used have been revamped according to the severity of the threats and productivity loss. This paper presents an approach to detect email spams more accurately, compared to the traditional spam detection methods available today. The proposed approach uses Convolutional Neural Networks (CNNs). It is a well known deep learning algorithm, among several algorithms available today. Though CNN is known at present for problems in image segmentation, and computer vision, where complex analysis and computations are required, this paper uses CNN for text based email spam detection. Based on the labeled email dataset chosen, the design and training of three models have been provided in this paper, where each model is different from one another based on complexity and functionality in layers of CNN. In addition to that, the performance analysis done based on the three models have been detailed in this paper. Based on the evaluation matrix and Hyper parameters such as learning rate, batch size, optimizer function, and activation function, the best model that performs efficiently has been determined. This research highlights the potential of CNN in enhancing the effectiveness of spam detection algorithms, along with the possibility of further research in this specific domain.

Index Terms—Email spam detection, Convolutional Neural Networks, Performance Analysis, deep learning.

I. INTRODUCTION

Emails have been classified into ‘spam’ or ‘ham’, depends upon the relevancy of content of the emails. Spam emails are the irrelevant emails sent for various purposes, like commercial advertisements, fraudulent, and malicious or phishing purposes. Whereas the ‘ham’ emails are those which are legitimate and relevant to the recipient. The content of those emails would be solicited and personalized, and the purpose of those emails are well defined. Isolating spam emails from ham emails have been a big challenge on the Internet. Including higher resource consumption, security risks, financial loss, the damage of reputation, and other ethical issues, the email service providers on the Internet have implemented several techniques, and revamped those techniques by adapting the state of art technologies and innovation over time. Introduction of Artificial Intelligence took these spam detection algorithms into another level by introducing automatic detection mechanisms. Before machine learning algorithm were in use, but now Neural network associated deep learning algorithms have been discovered to increase the effectiveness of the implementation.

This project is an effort made to implement this email spam detection with the help of Convolutional Neural Networks (CNNs), and check performance analysis by creating three different models for comparison and analysis. The selection of CNN for email spam detection system has been motivated by several factors, including its complex structure to handle image processing tasks, the effectiveness in use for feature learning, scalability in terms of handling large datasets and several others.

To conclude, the purpose of this paper is to propose an approach, then design and implement that approach to email spam detection system using CNN. An introduction of CNN architecture has been included in the next section. Followed by that, implementation of the models have been discussed. Later, in the following paragraph, the results of the experiment are done, which is more like tuning of parameters and computational approaches (functions) in search of an improved algorithm. In that next section, the analysis of the result has been detailed; which includes, the inferences made after a comprehensive comparison of evolution matrix, and hyper parameters, associated with each model. Later conclusion and future work has been included.

II. BACKGROUND WORK

A. Convolutional Neural Network (CNN)

1) **Introduction:** CNN is a type of artificial neural network(ANN) [1] in deep learning [2] which is primarily used for image processing and computer vision. This neural network is inspired from how the brain processes visual data. CNN is widely used for pattern recognition, feature extraction, and several others. This makes CNN more popular in applications like object detection, face recognition, real time stream analysis and several others. Real world examples include, self driving cars by Tesla [4], where CNN is used for identifying traffic, other vehicles, objects like pedestrians and traffic signs.

2) **CNN Model Architecture:** CNN architecture is designed based on how the data is being processed. The basic CNN contain three fundamental layers, which are as follows: 1) Convolutional layers, 2) pooling layers, and 3) fully connected layers.

1) **Convolutional Layers:** It is the core of CNN. It does convolution operation to the inputs and pass the output

resulted to the adjacent layer next to it. The convolution operation has a set of kernels, which is considered as filters, which works with the input data and construct a feature map. Each filter is associated with a specific task, which can be extraction of textures, patterns, or others. Several convolutional layers are available now depends upon the data to be handled. For text analysis the Conv1D layer [8] is being used. For image analysis, Conv2D or Conv3D [9] and so many others with different dimensions and functionalities are available in TensorFlow [10]

- 2) **Pooling Layers:** These are the adjacent layers of convolutional layers. These layers take the feature maps created by convolutional layers and reduce the spatial dimensions. This helps to reduce computational time, and memory usage. The common type of pooling operation is max pooling; in which reduction of dimension, extraction of features from feature map by taking maximum value from each region of feature map, noise reduction and several other functionalities are embedded. Main objective is to extract key features of the data and minimize the complexity as possible.
- 3) **Fully-Connected Layers:** These layers are also known as Dense layers. These layers receives input from pooling layers, which would be in multidimensional vector. The main objective of this layer is to convert the feature extracted by the pooling layers or convolutional layers into probabilities with different classes or predictions [2]. The steps involved are as follows. First, the vectors are flattened into single vector. Then, the matrix multiplication would be performed over this vectors with a weight, and apply activation function to introduce non linearity. This layer connects each neuron of one layer to the neurons of the next layer. In common, the dense layer would be the final layer of a CNN model.

B. Tools and Terminologies Used in This Report

1) *TensorFlow and Keras*:: TensorFlow [10] is a well known machine learning framework. It is developed by Google for the design and implementation of machine learning models. This library provides diverse of in build functions that implements machine learning algorithms and deep learning algorithms effectively. In addition, this library consists of multiple tools that can be used easily to test and tune the performance of the model created with least amount of time. This project uses TensorFlow over anaconda platform [11] with python language with the help of Visual Studio Code Toolbox [12].

2) *Optimization Algorithms*:: Optimization algorithm play a significant role in improving efficiency of the neural network model. These are functions that are used to update the parameter of the model to reduce the loss while training the algorithm. Commonly used Optimization algorithms in CNN are Stochastic Gradient Descent (SGD), Nesterov Accelerated Gradient, Momentum, Adam (Adaptive Moment Estimation) and so on.

3) *Activation Functions*:: An activation function is a computation unit associated with a neuron, which takes input from the previous layer and determine the importance of the information lies in that data to the next immediate layer. It helps to learn the patterns and relationships lies in the data. There are several activation functions which are used commonly in real world problems. Some of them are Sigmoid, Tanh, and ReLU (Rectified Linear Unit), ELU(Exponential Linear Unit), and so on [13].

4) *Loss Function*:: The loss function play a significant role in training of the model by quantifying prediction errors, aligning learning objectives with task requirements, and improving model evaluation. Depends upon the problem loss function can be different. For instance, in the case of regression problems, Mean Squared Error (MSE) is being preferred. And more, Binary Cross-Entropy Loss is preferred mainly for binary classification problem.

C. Related Research

Recently several research has been conducted in detecting phishing emails that use CNN. The model proposed in 2021 by Sergio for phishing Email Classification [3] lighten the possibility of CNN for text analysis. Later in 2022, Yuxin Liu and others [5]reviews CNN with machine learning algorithms (CNN and Bi-LSTM) for the same problem. Another recent work done by G. M. Shaharia [6], compares the efficiency of CNN and other machine learning algorithms such as Naive Bayes(NB), Linear Regression, and Support Vector Machine(SVM) for a system which filters email spams. Similar to the idea of email spam detection, a research has been conducted by Fanjun Meng [7] which proposes a system for network spam detection using CNN Incorporated with Attention Model.

III. IMPLEMENTATION AND VARIANTS

The implementation of basic model has been described in the first section. Followed by how the models differ each other has been detailed.

A. Implementation

1) *Data Splitting*: The dataset selected for email classification that contain spam emails as well as ham emails has been prepossessed collaboratively, removing unnecessary punctuation, and special characters, and cleaning by converting all characters into lowercase letters. The prepossessed data has been taken into further by applying label encoding, in which spam and ham labels are converted into numerical values such as 1 for spam and 0 for ham. The splitting of database has been carried out into two for training, validation and testing of the model. In which 80% has been selected for training and 20% has been selected for testing.

2) *Model Architecture Design*: The model has been created based on different layers in the order which is specified in the section II A. The sequential API of TensorFlow has been used to stack the layers sequentially. The layers used are listed below.

- *Embedding Layer*: This layer has been used to initialize the model. It converts each word into numerical vectors, represented by a vector in a high-dimensional space.
- *Convolutional Layer*: It is used for convolution operations. Here in this model, it does extract features from sequences of words.
- *A Bidirectional Long Short-Term Memory (Bi-LSTM) layer*: It processes input data in both directions sequentially, capturing dependencies in both directions. LSTM technique is used for find the long term dependencies between the words which are involved in the email.
- *GlobalMaxPooling Layer*: This layer is used to minimize the dimension of the input data. There are several layers similar to this layer available in use.
- *Dropout Layer*: This layer is used to reduce the over-fitting problem. Since the literature review done found out the research [14] done by Dishara highlighted over-fitting as a problem in CNN. The use of this layer is considerable.
- *Dense Layer*: This has been used to increase high learning capability of the model.

The above layers are combined in different combinations and three model have been created for comparison and performance analysis, which has been described in section III B.

The specification which is static in all the variants implemented are given in Table I.

TABLE I
SPECIFICATIONS OF THE IMPLEMENTATION

Common Hyperparameters and Parameters	
Parameter	Value
Embedding Dimension	50
Max Sequence Length	100
Number of Filters	64
Filter Size	5
Dropout Rate	0.5
Optimizer	Adam
Loss Function	Binary Crossentropy
Number of Epochs	5
Final Dense Layer units	1
Final Dense Layer Activation Function	sigmoid

3) *Model Compilation*: In this step, the model created in the above step has been compiled. The compilation involves configuring three things which is important in the learning capability of the model. Which are as follows.

- *Optimization Algorithm*: Among several optimization algorithms available, the initial model has been implemented using "adam" optimizer function for all the variants.
- *Loss Function*: The problem which is being solved here is binary classification problem. That is two classification, such as spam or ham. So "Binary Cross-Entropy Loss" has been chosen and made it static through out the project.
- *Evaluation Metrics*: The problem to be solved in this research is a binary classification problem, matrix parameters such that accuracy, precision, recall, and F1-score

are selected from the evaluation matrix. This helps to monitor the training of the model as well as the testing of the model which should be done after training phase.

The model has been compiled with the specification mentioned above.

4) *Model Training*: The compiled model has been trained using the training data which has been selected in the first phase of the implementation. The parameter specification used are as follows.

- Epoch(the number of iteration model has to perform through the entire dataset) = 5
- batch_size(count of sample processed at a time) = 64

The models have been trained with the above specification and the training progress have been depicted as a graph. Then the models have been investigated further for further analysis.

B. Different Variations

Three models have been created based on the complexity. This has been accomplished by introducing one or more layers into the simple model. As shows in the Table II, Variant 1 uses simple layers along with Conv1D layer, which is mainly used for text analysis or semantic analysis. Variant 2 provides a model where it is more complex compared to variant 1, which is because of introduction of Dense layer with 128 neurons and the Dropout layer to avoid over-fitting problem. Variant 3 is an extension of Variant 2, by adding an extra layer named Bidirectional LSTM.

TABLE II
COMPARISON OF MODEL ARCHITECTURES

Layer Name	Variant 1	Variant 2	Variant 3
Embedding	✓	✓	✓
Conv1D	✓	✓	✓
GlobalMaxPooling1D	✓	✓	✓
Bidirectional LSTM			✓
Dense (128, ReLU)		✓	✓
Dropout		✓	✓
Dense (1, sigmoid)	✓	✓	✓

IV. EXPERIMENTAL RESULTS

The test has been carried out with the three different variants implemented. The result obtained have been summarised in Table III.

As shows in Table III, several tests have been done on the three Variants, determining evaluation matrix associated with the model. In addition to that, evaluation matrix has been determined for different values of hyper parameters that have significant impact on the model. Finally, the result obtained on resource utilization parameters after several tests have been provided.

V. ANALYSIS AND DISCUSSION

This section begins with a comprehensive analysis of the model architecture created based on the experiment result summarized in the above section. Then Detailed discussion of some of the significant areas which need to be focused for an effective solution has been included.

TABLE III
EXPERIMENTAL RESULTS OF THREE VARIANTS

Parameter	Variant		
	1	2	3
Evaluation Matrix			
Accuracy (Ac)	0.983	0.980	0.986
Precision	0.96	0.91	0.95
Recall	0.90	0.93	0.93
F1-score	1	1	1
Hyper Parameters			
Learning Rate	0.03 (Ac:0.935)	0.001 (Ac:0.988)	0.001 (Ac:0.988)
Batch Size	64 (Ac:0.988)	128 (Ac:0.988)	128 (Ac:0.989)
Optimizer	Adam (Ac:0.986)	Adam (Ac:0.989)	Adam (Ac:0.988)
Activation Function	sigmoid (Ac:0.987)	elu (Ac:0.988)	sigmoid (Ac:0.988)
Resource Utilization			
Training Time (seconds)	4.78	3.36	32.97
Inference Time (seconds)	0.235	0.167	1.144

A. Analysis of Results

This section presents the observations found after the analysis of the result obtained from several experiments.

1) *Evaluation Metrics*: The main observation found after analysis are listed below

- 1) *Accuracy*: Model 3 achieves the highest accuracy (98.6%), followed by Model 1 (98.3%) and Model 2 (98.0%). Model 3 is slightly more accurate in making correct predictions compared to the other two models.
- 2) *Precision*: Model 1 exhibits the highest precision (96%), which indicates that it accurately identifies positive for majority of the positive predictions. Where as Model 3 follows closely with a precision of 95%. Model 2 shows the lowest of all which is (91%).
- 3) *Recall*: Model 2 demonstrates the highest recall (93%). Which means that this model can effectively identifies all relevant instances of the positive class. Both Model 1 and Model 3 have a recall of 90%, slightly lower than Model 2.
- 4) *F1-Score*: All models achieve a perfect F1-score of 1, indicating an excellent fairness between precision and recall. There is an optimum balance between making accurate positive predictions and identification of the right class for the input.

2) *Hyper-parameters*: Based on some of the important parameters analysed the following conclusions has been made.

- *Learning Rate*: Variant 1 utilized a higher learning rate of 0.03, but accuracy found to be 93.5%. Variant 2 and 3, with a lower learning rate of 0.001, achieved significantly higher accuracy's of 98.8

- *Batch Size*: Variant 1 had a batch size, which is equal to 64, with an accuracy of 98.8%. Variants 2 and 3, with a larger batch size of 128, also achieved high accuracy of 98.8% and 98.9%, respectively.
- *Optimizer*: All three Variants found the Adam optimizer as the best. Variant 2 achieved the highest accuracy (98.9%), followed closely by Models 1 and 3, both at 98.8
- *Activation Function*: Variant 1 found the sigmoid activation function as best with the accuracy of 98.7%. Variant 2 uses the ELU as the activation function and resulted into the highest accuracy of 98.8%. Model 3 also used the sigmoid activation function and received an accuracy of 98.8

Overall, Variant 2 found to be the best after analysis of hyper parameters in terms of achieving the highest accuracy with a lower learning rate, the ELU activation function, and the Adam optimizer.

3) *Model Architecture*: Based on the whole modeling, the following conclusions are made about each Variant.

- Variant 1: Achieves high accuracy and F1-score shows high capability of the model towards the prediction. The precision and recall scores suggest a well-balanced model in terms of classification. As can be seen from Table III, Larger batch size (64) found to be best with high accuracy, but learning rate is high (0.03) for best accuracy.
- Variant 2: Slight decrease in accuracy compared to Variant 1. But, the Variant 2 found best in terms of precision and recall. Which shows classification would be better.
- Variant 3: Highest accuracy among the three variants, with good precision and recall(, since F1-score is 1). The Bidirectional LSTM used in this model can capture sequential dependencies effectively. Which resulted an improved performance. Dropout layer helps in regularization, preventing over-fitting. But this model found to be more complex compared to rest of the models.

B. Detailed Discussion

The following provides different perceptions observed while doing the comparison and performance analysis.

1) *Selection of Model*: Choosing the best model from the above three is significant in order to solve the problem efficiently. Based on the evaluation matrix discussed in section V A, a strategy can be found for choosing a model. Choosing a model depends on the requirement and constraints of the problem. Variant 3 can be selected, if maximizing accuracy is preferred. Variant 1 can be selected, if precision is more important.

2) *Computational Efficiency*: The computational efficiency was one of the main focus during the research. The reason was the complexity of the model was different from Variant 1 to Variant3. So the time taken for training the model and inference time has been calculated for each model. Training time is the time taken by the CNN model during the training phase, including the pattern recognition, and the whole process. Whereas, Inference time is the time taken by the

trained model to perform prediction of the output with new data. Based on the result, shown in Table III, the following observations have been found.

- Training time: Variant 1 and Variant 2 have relatively shorter training times compared to Variant 3. The reason can be drawn that Variant 3 uses complex architecture, so more parameters need to be trained. Bidirectional LSTM are more complex models compared to simple CNN layers. So, it does require more computations.
- Inference time: Inference time found to be different for all variants. Variant 2 has the shortest inference time, followed by Variant 1. Variant 3 has the longest inference time. This difference in inference time can be attributed to the architectural differences between the variants. Variant 2 has fewer layers compared to Variant 3, resulting in faster inference. Same reason can be applied here as that has been discussed about increase of training time in Variant 3.

Variant 2 found to a balanced model in terms of complexity and computational efficiency. It does have shorter training and inference times compared to both Variant 1 and Variant 3. Variant 3, although providing the best performance in terms of evaluation metrics. But in terms of complexity, this model is inefficient. Depends upon the requirements, and number of constraints put into, the model can be selected.

3) *Best Model Proposed:* Selection of a model as the best among the three models implemented depends on several factors. The different case where one model can be preferred over another has been discussed in the above sections. Variant 1 could be a best option for resource-constrained environments. Variant 2 is best for applications where model complexity and over-fitting need to be addressed. Variant 3 is preferable where high predictive accuracy is significant.

Though variant3, found to be the best among all, the scalability of the model found to be a challenge in this research. This research found it because of two reasons. First reason is its difficulty while dealing with large datasets. Second reason is the deployment challenge in the real world because of the requirement of high computations.

To conclude, this research put forward **Variant 3 as the best model** due to its stability in terms of evaluation matrix, and other performance parameters. The graph showing the model accuracy and model loss has been provided in Fig. 1.

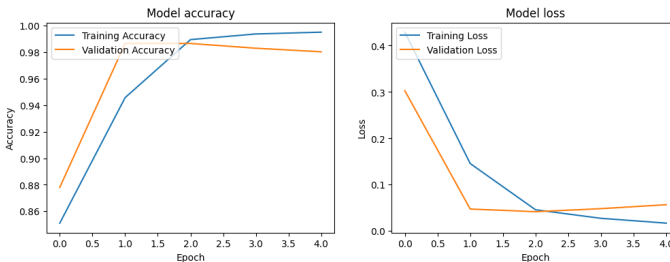


Fig. 1. Model Accuracy-Loss Graph

From Fig. 1, it is clearly seen that the model does have high training accuracy and high validation accuracy. In an ideal model, the curves would be completely overlapped. here training accuracy curve sits above validation accuracy curve, so slight over-fitting exist. So with unseen data, the performance might be a bit less, but reasonably well because the precision was 95 (based on Table III).

4) *Challenges and Space for Further Research:* This section presents a detailed discussion on the challenges observed and areas where further research is mandatory.

- 1) *Computational Tools:* The requirement of high computational resources for training and evaluation of these models is found to be important. For instance, large dataset is important for effective training, but the computational effort would also be very high. Another thing is the number of epoch parameter chosen while training of the model. This research has to put epoch as 5 to minimize the effect of computational challenge. The performance analysis of the models based on effect of epoch parameter should need to be further investigated. In addition to that, this research has seen the importance of exploring to development of new algorithms which reduces the computational cost irrespective of scalability in the future.
- 2) *Over-fitting Problem:* Over-fitting is a persistent challenge where model has to deal with large number of parameters. Variant 2 does has a dropout layer to regulate the over-fitting, but there is still space for improvement. The research should explore further with different regularization techniques.
- 3) *Use of Large Dataset:* The dataset chosen here is small, which is intentionally chosen to minimize the computational efforts. Further research could investigate on large datasets with different types of emails, type of contents for fine tuning.
- 4) *Application of Different CNN Models:* Application of different CNN layers and do evaluate the performance of the model obtained is found to be important to ensure the robustness of the model. Investigating the possibility of combining CNN model with different Neural network algorithms, such as Recurrent Neural Networks(RNNs), is significantly appreciated.

To conclude, the analysis of the models has formulated several observations. The observations formulated put forward several hints in which the research should be focused or to be extended. Exploring new directions envision new advancement in the area of deep learning and natural language processing implementations.

VI. CONCLUSION

This paper has presented design, implementation, and then evaluation of three variants of Convolutional Neural Networks for email spam detection. The architecture of the CNN architecture followed during the research has been provided at the beginning. Followed by, the tools, concepts, and terminologies used for a successful performance analysis had

been highlighted. Then the results obtained during the performance analysis has been provided based on evaluation metrics, hyper-parameters, and parameters associated with the resource utilization. The detailed discussion based on the results has drawn several observations about how to select a model, the complexity and computational efficiency, along with the best model selected during the research. Several hints has been presented where the challenges faced during the research and for several necessity of further research found.

REFERENCES

- [1] McCulloch, W. & Pitts, W. A logical calculus of the ideas immanent in nervous activity. *The Bulletin Of Mathematical Biophysics*. **5**, 115-133 (1943)
- [2] Goodfellow, I., Bengio, Y. & Courville, A. Deep Learning. (MIT Press,2016), <http://www.deeplearningbook.org>
- [3] McGinley, C. & Monroy, S. Convolutional Neural Network Optimization for Phishing Email Classification. *2021 IEEE International Conference On Big Data (Big Data)*. pp. 5609-5613 (2021)
- [4] Bojarski, M., Testa, D., Dworakowski, D., Firner, B., Flepp, M., Goyal, P., Jackel, L., Monfort, M., Muller, U., Zhang, J. & Al. End to End Learning for Self-Driving Cars. *CoRR*. **abs/1604.07316** (2016)
- [5] Liu, Y., Wang, L., Shi, T. & Li, J. Detection of spam reviews through a hierarchical attention architecture with N-gram CNN and Bi-LSTM. *Information Systems*. **103** pp. 101865 (2022)
- [6] Shahariar, G., Biswas, S., Omar, F., Shah, F. & Hassan, S. Spam review detection using deep learning. *2019 IEEE 10th Annual Information Technology, Electronics And Mobile Communication Conference (IEMCON)*. pp. 0027-0033 (2019)
- [7] Meng, F., Pan, Y. & Feng, R. Network Spam Detection Based on CNN Incorporated with Attention Model. *2022 8th Annual International Conference On Network And Information Systems For Computers (ICNISC)*. pp. 111-116 (2022)
- [8] Developers, T. `tf.keras.layers.Conv1D`. (TensorFlow,2015). <https://www.tensorflow.org/apidocs/python/tf/keras/layers/Conv1D>
- [9] Chollet, F. Deep Learning with Python. (Manning Publications Co.,2018)
- [10] Developers, T. TensorFlow. (<https://www.tensorflow.org/>,2024)
- [11] Continuum Analytics, I. Anaconda Platform. (Continuum Analytics, Inc.,2024), <https://www.anaconda.com/products/distribution>
- [12] Microsoft Visual Studio Code. (Microsoft,2015), <https://code.visualstudio.com/>
- [13] Towards Data Science Activation Functions in Neural Networks. *Towards Data Science*. (2019), <https://towardsdatascience.com/activation-functions-neural-networks-1cbd9f8d91d6>
- [14] Dishara, H. & Muhammed, L. A Review of the Overfitting Problem in Convolution Neural Network and Remedy Approaches. *Journal Of Physics: Conference Series*. **1432**, 012002 (2020)