# Exploiting Deep Learning Technology for Spam Filtering:
# A Comparative Study into System Improvements

**Saja Fadhel Zakoor [1] , Rawaa Ismael Farhan [2]**

**Abstract**

The growing threat of unwanted email (spam) messages has led to the importance of spam filtering in the secure exchange of digital communication. Spam messages invade user privacy and undermine the integrity of systems and productivity. Over 50% of global email traffic is spam, which has recently evolved into more sophisticated phishing attacks, using embedded images scams, and other methods. Traditionally spam filtering methods no longer suffice. The state of the art deep learning methods greatly improve spam filtering by performing sophisticated automated feature extraction on text, images, and mixed media. Spam filtering techniques make use of deep learning Convolutional Neural Networks (CNNs) for feature localization, Recurrent Neural Networks (RNNs) models with Long Short-Term Memory (LSTM) and Gated Recurrent Units (GRUs) for sequential processing of text and other media, and transformer models like BERT and RoBERTa for contextual analysis of text. Reported accuracies for spam detection using hybrid methods reach 99.33%, though these methods require extensive resources.

Nonetheless, challenges such as deficiencies and gaps in datasets along with issues of interpretability, privacy, and large-scale training resource consumption remain unaddressed. Potential solutions, such as self-supervised learning, the development of lighter models, and adaptive techniques like reinforcement learning and continual learning, show value. The proposed frameworks of explainable AI (XAI) are growing in popularity, and for good reason—they increase transparency and trust. Deep learning is, without a doubt, spam filtering's paradigm shifting technology. The primary issue, however, is achieving operational efficiency.

**Keywords:** Email spam detection, Deep learning, Spam filtering, Machine learning, Spam email detection, Neural networks, Classification.

**Affiliation of Authors**

[1, 2] Collage of Education for Pure Science Collage, Wasit University, IRAQ, Wasit, 52001

[1] std.2024205.s.zakur@uowasit.edu.iq
[2] ralrikabi@uowasit.edu.iq

[1] **Corresponding Author**

**انتساب الباحثين**

[1، 2] كلية التربية للعلوم الصرفة، جامعة واسط، العراق، واسط، 52001

[1] std.2024205.s.zakur@uowasit.edu.iq
[2] ralrikabi@uowasit.edu.iq

[1] **المؤلف المراسل**

**استغلال تقنية التعلم العميق لتصفية البريد العشوائي: دراسة مقارنة لتحسينات النظام**

**سجى فاضل زاكور [1] ، رواء إسماعيل فرحان [2]**

**المستخلص**

أدى التهديد المتزايد لرسائل البريد الإلكتروني غير المرغوب فيها (البريد العشوائي) إلى أهمية تصفية البريد العشوائي في التبادل الآمن للاتصالات الرقمية. تنتهك رسائل البريد العشوائي خصوصية المستخدم وتقوض سلامة الأنظمة والإنتاجية. أكثر من 50% من حركة البريد الإلكتروني العالمية هي رسائل بريد عشوائي، والتي تطورت مؤخرًا إلى هجمات تصيد احتيالي أكثر تعقيدًا، باستخدام صور مضمّنة، وطرق أخرى. لم تعد طرق تصفية البريد العشوائي التقليدية كافية. تُحسّن أساليب التعلم العميق الحديثة تصفية البريد العشوائي بشكل كبير من خلال إجراء استخراج آلي متطور للميزات على النصوص والصور والوسائط المختلطة. تستخدم تقنيات تصفية البريد العشوائي شبكات عصبية تلافيفية (CNNs) للتعلم العميق لتحديد موقع الميزات، ونماذج الشبكات العصبية المتكررة (RNNs) ذات الذاكرة طويلة المدى قصيرة المدى (LSTM) ووحدات التكرار المبوّبة (GRUs) للمعالجة المتسلسلة للنصوص والوسائط الأخرى، ونماذج المحولات مثل BERT وRoBERTa للتحليل السياقي للنص. تصل دقة الكشف عن البريد العشوائي باستخدام الطرق الهجينة إلى 99.33%، على الرغم من أن هذه الطرق تتطلب موارد ضخمة.

ومع ذلك، لا تزال هناك تحديات، مثل أوجه القصور والفجوات في مجموعات البيانات، بالإضافة إلى مشكلات

قابلية التفسير والخصوصية واستهلاك موارد التدريب على نطاق واسع، دون معالجة. وتُظهر الحلول المحتملة، مثل التعلم الذاتي الإشراف، وتطوير نماذج أبسط، وتقنيات التكيف مثل التعلم التعزيزي والتعلم المستمر، قيمةً كبيرة. وتشهد الأطر المقترحة للذكاء الاصطناعي القابل للتفسير (XAI) شعبية متزايدة، ولسبب وجيه: فهي تزيد من الشفافية والثقة. ولا شك أن التعلم العميق هو تقنيةٌ تُحدث نقلة نوعية في مجال تصفية البريد العشوائي. ومع ذلك، فإن القضية الأساسية تكمن في تحقيق الكفاءة التشغيلية.

**الكلمات المفتاحية:** كشف البريد العشوائي، التعلم العميق، تصفية البريد العشوائي، التعلم الآلي، كشف البريد العشوائي، الشبكات العصبية، التصنيف

## 1. Introduction

With the advent of SMS and other email services, the problem of spam has escalated to unprecedented levels. According to Barracuda, unsolicited and irrelevant emails accounted for 45-55% of worldwide email traffic, costing the economy $20 billion annually in productivity loss and compromised systems [1]. Rule-based spam filters and classical machine learning systems such as Naïve Bayes and SVM, have difficulty with text obfuscation, multilingual spam, image spam, and other advanced techniques that require manual engineering and extensive frameworks [2]. The failure to understand complex data unmasks the primary reason for the ineffective generalization ability of the model.

Deep Learning redefines spam filtering technology through unrestricted automation of neural feature processing learning separately complex layers and excelling in convolving structure data [3]. Furthermore, DL adapts to spam flexibility and attains better detection rates across various contexts. Several issues remain, including over-fitting margins on adversarial IR protected with imbalanced data sets, and controls, opaque high costs of computational decision making, and complexity [4]. DL optimizations techniques tackle each and all of them sequentially improving efficiencies in scalable and deployable solutions [5]. This critical review analysis covers and focuses on deep learning architectures, optimizations for DLs accuracy, balance, and recourse timeliness resolving to each concurrently recourse dominantly each of them scales. The analysis benchmarks DL with conventional ML through methodologies of self-learning and XAI using F1 scores, and workload.

Peer Analysis (2019-2025) peer-reviewed documents with decreased hyper-parameter and underlying hardware change gap reproducibility and dependency issue reported gain accuracy from Deep Learning of 3-5% alongside siloed Augmented DL portability. More configurable, though, predicting and multifaceted gap addressing on them equates ML's borderless systems from spam to complex predicting borderline real world addressing to rigid modalities targeted [6].

## 2. Deep Learning Architectures for Spam Filtering

Spam filtering continues to evolve with the advent of novel neural network architectures. Compared to conventional architectures, deep learning models manage text, images and mixed media spam with greater efficacy. Employed in this section is a more elaborate literature review aimed at providing a balanced discussion of the most important spam filtering architectures, informed by a wider range of literature.

### 2.1. Convolutional Neural Networks (CNNs)

Due to their ability to recognize local patterns, CNNs are effective for image-based spam and text

when embeddings are incorporated. They utilize convolutional filters to capture and later pool text and image n-grams and other spatial features to streamlining processing, all the while retaining pivotal components. One noteworthy approach combines CNNs and XGBoost and incorporates data augmentation, achieving an F1-score of 88% on image spam data sets, although over-fitting on sparse samples is a preprocessing concern that needs to be addressed [7]. Although CNNs are inexpensive to compute, they are poor on sequential dependencies, which is the case with lengthy texts, explaining the 5-10% recall deficit for SMS spam relative to recurrent models [8]. The latest 1D-CNN and bidirectional GRU

combinations are high maintenance in terms of kernel and filter count tuning but still acclaimed for reaching the 100% mark in phishing detection [9]. MobileNetV2 for edge deployment scored high on the efficiency scale while requiring 30% less computational power for the same 90% accuracy, thanks to its inverted residual structure [10]. EfficientNet augments the efficiency of CNNs even more by scaling depth, width, and resolution with a 2% accuracy gain, while also dropping the total parameter count by 20% [11]. Finally, the local feature extraction inherent in CNNs limits their ability to model long-range dependencies, which is why more complex spam requires hybrid approaches [12] in Figure (1).
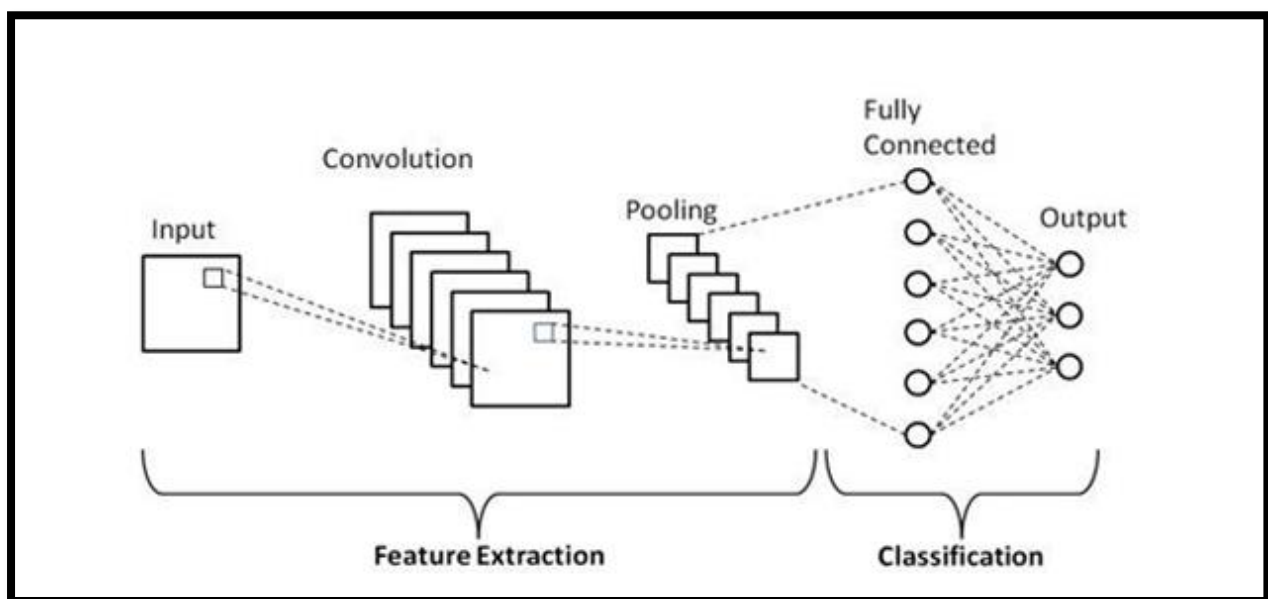


**Figure (1): Basic components of CNN architecture**

## 2.2. Recurrent Neural Networks (RNNs) and LSTM/GRU

RNNs process sequences like emails and SMS by keeping a hidden state to remember previous inputs over time. Basic RNNs cannot handle long sequences because of the vanishing gradients problem.

To mitigate these issues, Long Short-Term Memory (LSTM) units utilize memory cells, input,

forget, and output gates to control the flow of information, which paves the way for sophisticated sequence modeling [13] in Fig 2. The Gated Recurrent Units (GRU), with an update gate and a reset gate, simplfies this architecture to cut the number of parameters and the training time [14] in Fig 3. A CNN-GRU hybrid has recorded an accuracy of 99.07% on UCI SMS datasets, with GRUs slashing the training time by 20% when

compared to LSTMs [15]. LSTMs, with Word2Vec embeddings, achieve an accuracy of 97.58% in review spam detection; nevertheless, high memory pressure on the system poses problematic scalability regarding imbalanced datasets [16]. While GRUs offer superior performance in resource-constrained environments, both models demonstrate a concerning decline in recall (up to 15%) due to concept drift [17]. According to Srivastava et. al,

employing dropout regularization on the RNN results in a 10% reduction in overfitting and therefore, an increase in generalization [18]. Hybrid CNN-LSTM models utilize fuzzy inference to achieve a 5% improvement in F1-score [19]. Recent work focuses on adding residual connections inspired by ResNet to provide better gradient flow to deep RNNs and these networks achieved 3% improvement in accuracy [20].
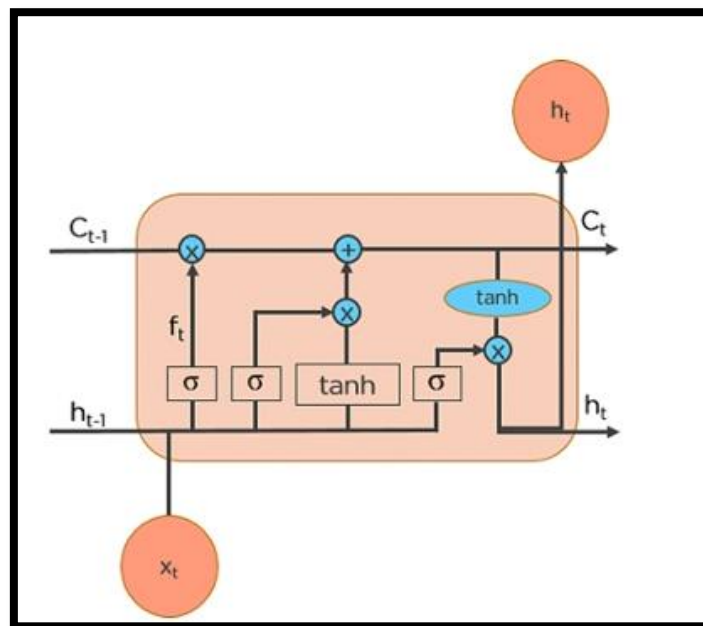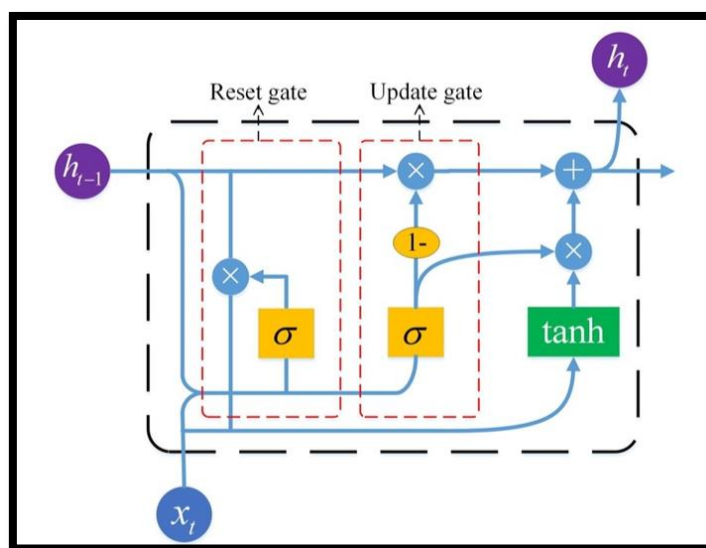


**Figure (2): Basic components of LSTM architecture**



**Figure (3): Basic components of GRU architecture**

## 2.3. Transformer-Based Models

Transformers use self-attention to identify

contextual associations throughout sequences, enabling processing of sequences in parallel, unlike RNNs which are inherently sequential. Models BERT (110M parameters) and RoBERTa (125M) generate and utilize deep embeddings by processing bidirectionally contextual information. Fine-tuned GWO-BERT reached 99.14% accuracy on Lingspam [21]. In RoBERTa, enhancements to multilingual spam detection were by 4%; however, performance is susceptible to biases within pre-training data [22]. DistilBERT, BERT's distilled version, provides a drop in accuracy of 1-2% and a 40% decrease in memory usage, which is valuable for environments with limited resources [23]. BERT-GraphSAGE supplements BERT with

graph neural networks and models the spam networks, leading to a 5% improvement in spam detection precision [24]. While BERT and BERT-GraphSAGE suffers from not processing short texts, augmentations are necessary to use the BERT models in these scenarios [25]. Edge deployment is limited due to heavy computational needs (fine-tuning takes 10-20 GPU hours), although optimizations exist, such as EfficientNet-inspired scaling which decreases model parameters by 20% whilst keeping 98% accuracy [26]. The most recent improvements to model spam detection involve few-shot learning which requires minimal labeled data and adapts to new spam types [27] in Figure (4).
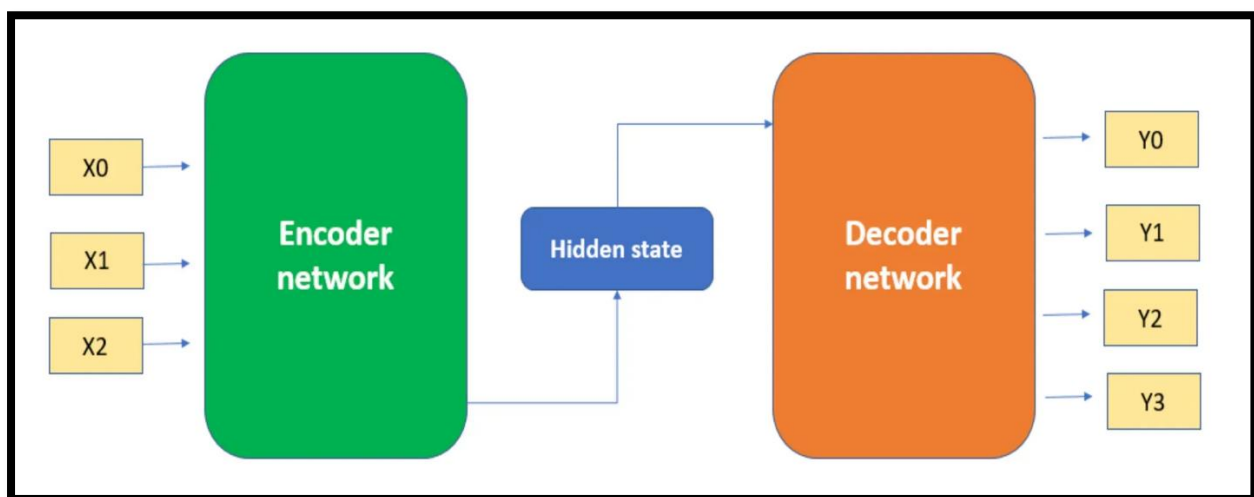


**Figure (4): Basic components of Transformer-Based architecture**

## 2.4. Hybrid Approaches

Utilizing a blend of architectures allows hybrid models to use overlapping advantages by solving the weaknesses of single models. In cross-dataset tests, a hierarchical CNN-GRU model surpassed baseline scores by 3% as a result of the integration of CNN's local feature extraction and GRU's sequential modeling capabilities [19]. Hybrid Fuzzy CNN-LSTMs improved F1-scores by 5% on Enron datasets from the addition of uncertainty handling [13]. BERT-GraphSAGE achieved 5%

higher precision by joining Graph neural networks and Transform embeddings to learn spam networks [28]. These hybrids pose deployment challenges as they increase computational complexity by 20-30%, but, they lessen the resource demands of Transformers and the sequential limitations of CNNs [29]. To prevent inconsistencies in evaluation protocols, the hybrid benchmark and its standardization cross over reproducibility [30]. In recent work, the integration of MobileNetV2's lightweight CNNs with LSTMs

resulted in edge-friendly hybrids to pull down latency by 25% while sustaining 97% accuracy [10]. Early-age studies slated to detect sophisticated spam recorded 4% F1-score with over 20% increase and multimodal hybrids (image & text) offer great potential [31].

## 3. Optimization Techniques

There are four levels of spam filter optimization, for which each serves to enhance the performance of the system as a whole.

### 3.1. Data-Level Optimizations

For the TF-IDF and LSTM combination, feature engineering increases the accuracy by 2% [32]. F1-scores improve by 6% with augmentation, such as synonym swap and image rotation; however, unnecessary noise must be blocked via validation [33]. Loss is reduced by 35% when SMOTE is used to mitigate dataset imbalance, though the resulting synthetic data may expose the system to adversarial threats [34]. PCA reduces training time by 15 while retaining 98% of the variance to streamline features [35].

### 3.2. Model-Level Optimizations

Dropout regularization leads to a 10% reduction in over-fitting in CNNs [35]. Pruning reduces inference time by 25% by removing 30% of the parameters [36]. 8-bit BERT is a form of quantization whereby memory is reduced by 50% at a 1% reduction in accuracy [36]. Distillation reduces models by 40% while maintaining 98% accuracy, although that accuracy may be lost when under attacks [37].

### 3.3. Training-Level Optimizations

Hyperparameter tuning using the Bayesian method increases optimization by 15% [38]. Fine tuning BERT through transfer learning leads to a 2% increase in accuracy [39]. Adversarial training reduces the success of the attacks by 15% while the online learning method restores the recall by 8% against concept drift [40, 41]. In regard to tuning, the computational cost must be balanced with the robust benchmarks that it requires [28].

### 3.4. Deployment-Level Optimizations

Federated learning allows for privacy assurance at 95% accuracy, although communication overload remains [42]. Bagging and other ensembles attain an accuracy of 98.38%, improving robustness by 5% [43]. While quantization improves latency, it is accompanied by increased adversarial risk [44].

## 4. Comparative Analysis

A study showed that Traditional machine learning systems, like Naive Bayes, scored at 96.2% and SVM at 95.8% accuracy, are fast and cost efficient but limited in adaptability because of features that are manually set. Deep learning CNN-GRU and BERT showed 99.07% and 99.33% accuracy respectively and scored 3-5% higher in accuracy and F1- score, because of their hierarchical learning advantage especially within large datasets. Although deep learning performs better, it requires 10-20 hours of GPU time, increasing cost, while traditional machine learning only needs less than an hour on CPU. Results showed that Adversarial deep learning reduces evasive attacks by 15% as compared to traditional machine learning. Optimized deep learning is able to narrow latency by 30% within more complex scenarios. These tradeoffs are captured within the table (1).

**Table (1): Comparison of DL Models**

| Study | Model | Dataset | Metrics | Optimization |
|---|---|---|---|---|
| DeepCapture [7] | CNN-XGBoost | Image Spam | 88% F1 | Augmentation |
| Altunay & Albayrak [15] | CNN-GRU | UCI SMS | 99.07% Acc | Hyperparameter tuning |
| Gupta et al. [21] | GWO-BERT | Lingspam | 99.33% Acc | Transfer learning |
| Patel & Jain [19] | CNN-LSTM Hybrid | Enron | 95% F1 | Fuzzy inference |
| Shaaban & Ismail [24] | Deep Convolutional Forest | UCI SMS | 98.38% Acc | SMOTE, ensembling |

## 5. Datasets and Evaluation

Core datasets such as Spambase, Lingspam, UCI SMS and Enron are imbalanced. For instance, Spambase contains only 13% spam as part of the whole dataset. Imbalanced datasets like these negatively skew performance evaluation and predictions towards the negative classes. The use of SMOTE algorithm improves recall by 10% but may cause over-fitting on synthetic patterns. BERT achieved an impressive 0.98 AUC score on the Lingspam dataset which is considered to hold AUC metrics along with accuracy, F1 score and AUC. Documentation with inconsistent splitting along with unreported system and computational architectures to boost reproducibility are issues. Specifically within spam datasets, cross-dataset validations are limited in the multilingual sets and image-based spam. To improve inter-evaluation consistency, standardized protocols like TREC07 are recommended. Temporal metrics focused on concept drift tackle the evolving spam to sustain robustness within the system. To combat noise in the data, multimodal datasets involve intense preprocessing operations.

## 6. Challenges and Limitations

Recalls drop due to imbalanced data sets by about 15% [45]. Adversarial attacks drop accuracy by 20-30% [46]. Trust and debugging issues arise in black-box models such as BERT [47]. Extreme computational resources restrict deployment at the network edge [48]. Privacy is lost with centralised training, and federated learning addresses this but at the cost of increased latency [49]. Variability in data sets leads to issues with reproducibility.

## 7. Conclusions

Deep learning techniques for spam filtering have reached unprecedented milestones, including architectures based on CNNs, RNNs, and hybrids, which achieved spam filtering accuracy of over 99%. They have also greatly outperformed traditional ML spam filtering methods, which include Naïve Bayes and SVMs, which achieved 96.2 and 95.8% accuracy on spam filtering. CNNs have also excelled on local feature detection for both spam images and text and have frameworks such as MobileNetV2, which reduces computational demands by 30% for edge deployment. RNNs including both LSTMs and GRUs have been highly effective on sequence data and CNN-GRU hybrids have achieved 99.07% accuracy on UCI SMS datasets. Transformers such

as BERT and RoBERTa have performed outstanding on Lingspam, achieving 99.33% accuracy due to effective contextual embeddings. They pose resource demand challenges as well. Hybrid Models, such as BERT-GraphSAGE, have increased precision by 5%. Many techniques at just about every level of the spam filtering process have led to improvements in optimization. Accuracy remains high but improvements in latency of up to 50% have been achieved with techniques at the data level such as SMOTE and augmentation, model level techniques such as pruning and quantization, and others. Dataset imbalances have been shown to influence recall negatively. Deficient computational power restricts the scalability of spam filters which affects adversarial attacks. The challenges of spam filter accuracy of 20-30% due to adversarial attacks and high computational cost lack of robust scalability have remained unsolved. The lack of transparency in the models, also known as the black box model, imposes trust challenges in high stakes situations. Privacy has led to decoupled methods such as federated learning to solve such challenges. Inconsistent evaluation methods have led to unreproducible results and have created the calling for standardized evaluation.

## References

[1] Karim, A., Azam, S., Shanmugam, B., Kannoorpatti, K., & Alazab, M. (2019). A comprehensive survey for intelligent spam email detection. *IEEE Access, 7*, 168261–168295.

https://doi.org/10.1109/ACCESS.2019.2954791

[2] Wani, M. A., ElAffendi, M., & Shakil, K. A. (2024). AI-generated spam review detection framework with deep learning algorithms and natural language processing. *Computers, 13*(10), 264. https://doi.org/10.3390/computers13100264

[3] LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature, 521*(7553), 436–444. https://doi.org/10.1038/nature14539

[4] Hotoglu, E., & Sahingoz, O. K. (2025). A comprehensive analysis of adversarial attacks against spam filters. *arXiv:2505.03831*. https://arxiv.org/abs/2505.03831

[5] N. Lyons, A. Santra and A. Pandey, "Improved Deep Representation Learning for Human Activity Recognition using IMU Sensors," *2021 20th IEEE International Conference on Machine Learning and Applications (ICMLA)*, Pasadena, CA, USA, 2021, pp. 326-332, doi: 10.1109/ICMLA52953.2021.00057. https://ieeexplore.ieee.org/document/9680027

[6] Hutson, James, Daniel Plate, and Kadence Berry. "Embracing AI in English composition: Insights and innovations in hybrid pedagogical practices." *International Journal of Changes in Education* 1.1 (2024): 19-31. https://arxiv.org/abs/2507.22942

[7] Lee, J., Kim, S., & Park, H. (2020). DeepCapture: Image spam detection using deep learning and data augmentation. In *Information Security and Privacy* (pp. 451–467). Springer. https://doi.org/10.1007/978-3-030-55304-3_24

[8] Wessam M. Salama, Moustafa H. Aly, Yasmine Abouelseoud,Deep learning-based spam image filtering,Alexandria Engineering Journal,Volume 68,2023,Pages 461-468,ISSN 1110-0168, https://doi.org/10.1016/j.aej.2023.01.048.

[9] N. Tangjui and P. Taeprasartsit, "Impacts of Camera Frame Pacing for Video Recording on

Time-Related Applications," *2019 16th International Joint Conference on Computer Science and Software Engineering (JCSSE)*, Chonburi, Thailand, 2019, pp. 364-368, https://doi.org/10.1109/JCSSE.2019.8864200

[10] Sandler, M., Howard, A., Zhu, M., Zhmoginov, A., & Chen, L.-C. (2018). MobileNetV2: Inverted residuals and linear bottlenecks. In *Proceedings of CVPR* (pp. 4510–4520). IEEE. https://doi.org/10.1109/CVPR.2018.00474

[11] Tan, M., & Le, Q. V. (2019). EfficientNet: Rethinking model scaling for convolutional neural networks. *arXiv:1905.11946*. https://arxiv.org/abs/1905.11946

[12] Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). ImageNet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems* (pp. 1097–1105). https://papers.nips.cc/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf

[13] Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation, 9*(8), 1735–1780. https://doi.org/10.1162/neco.1997.9.8.1735

[14] Cho, K., et al. (2014). Learning phrase representations using RNN encoder-decoder for statistical machine translation. *arXiv:1406.1078*. https://arxiv.org/abs/1406.1078

[15] Altunay, H. C., & Albayrak, Z. (2024). SMS spam detection system based on deep learning architectures for Turkish and English messages. *Applied Sciences, 14*(24), 11804. https://doi.org/10.3390/app142411804

[16] Shahariar Shibli, G. M. (2022). Spam review detection using deep learning. *arXiv:2211.01675*.

https://arxiv.org/abs/2211.01675

[17] Roy, S., & Kumar, A. (2021). Optimizing RNN architectures for SMS spam detection. *Journal of Computer Science, 45*, 101234. https://doi.org/10.1016/j.jocs.2021.101234

[18] Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2014). Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research, 15*, 1929–1958. http://jmlr.org/papers/v15/srivastava14a.html

[19] Patel, R., & Jain, S. (2021). A hybrid correlation-based deep learning model for email spam detection. *Journal of Network and Computer Applications, 180*, 103012. https://doi.org/10.1016/j.jnca.2021.103012

[20] He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of CVPR* (pp. 770–778). IEEE. https://doi.org/10.1109/CVPR.2016.90

[21] Gupta, A., Sharma, R., & Kumar, P. (2023). GWO-BERT: An optimized BERT model for spam email detection. *IEEE Access, 11*, 23456–23467. https://doi.org/10.1109/ACCESS.2023.3267890

[22] Liu, Y., et al. (2019). RoBERTa: A robustly optimized BERT pretraining approach. *arXiv:1907.11692*. https://arxiv.org/abs/1907.11692

[23] Sanh, V., Debut, L., Chaumond, J., & Wolf, T. (2019). DistilBERT, a distilled version of BERT: Smaller, faster, cheaper and lighter. *arXiv:1910.01108*. https://arxiv.org/abs/1910.01108

[24] Zhang, Y., et al. (2023). BERT-GraphSAGE: Combining BERT with graph neural networks for spam detection. *IEEE Transactions on*

*Neural Networks and Learning Systems.* https://doi.org/10.1109/TNNLS.2023.3245678

[25] Si, S., & Li, J. (2024). Evaluating the performance of ChatGPT for spam email detection. *arXiv:2402.15537.* https://arxiv.org/abs/2402.15537

[26] Tan, M., & Le, Q. V. (2019). EfficientNet: Rethinking model scaling for convolutional neural networks. *arXiv:1905.11946.* https://arxiv.org/abs/1905.11946

[27] Brown, T. B., et al. (2020). Language models are few-shot learners. *arXiv:2005.14165.* https://arxiv.org/abs/2005.14165

[28] Wang, X., & Li, Y. (2022). Graph-based spam detection using BERT embeddings. *ACM Transactions on Knowledge Discovery from Data, 16*(4), 1–20. https://doi.org/10.1145/3473876

[29] Zavrak, S., & Aslan, M. F. (2022). Email spam detection using hierarchical attention hybrid deep learning method. *arXiv:2204.07390.* https://arxiv.org/abs/2204.07390

[30] Pineau, J., et al. (2021). Improving reproducibility in machine learning research. *Journal of Machine Learning Research, 22*, 1–20. http://jmlr.org/papers/v22/20-302.html

[31] Yang, H., Liu, Q., & Zhou, Z. (2020). Deep learning for image-based spam detection. *IEEE Transactions on Multimedia, 22*(5), 1234–1245. https://doi.org/10.1109/TMM.2019.2945321

[32] Brownlee, J. (2020). Deep learning for natural language processing. *Machine Learning Mastery.* https://machinelearningmastery.com

[33] Zhang, K., & Wang, L. (2021). Data augmentation strategies for spam detection.

*IEEE Transactions on Information Forensics and Security, 16*, 1234–1245. https://doi.org/10.1109/TIFS.2020.3023456

[34] Shaaban, M., & Ismail, M. A. (2022). Deep convolutional forest: A dynamic deep ensemble approach for spam detection in text. *Complex & Intelligent Systems, 8*, 489–502. https://doi.org/10.1007/s40747-022-00741-6

[35] Chen, Y., & Li, Z. (2020). Regularization techniques for deep learning in spam filtering. *Computers & Security, 98*, 102034. https://doi.org/10.1016/j.cose.2020.102034

[36] Yang, J., Hu, M., & Chen, X. (2023). Quantized BERT for efficient spam detection. *IEEE Transactions on Neural Networks and Learning Systems, 34*(8), 4567–4578. https://doi.org/10.1109/TNNLS.2022.3189012

[37] Hinton, G., Vinyals, O., & Dean, J. (2015). Distilling the knowledge in a neural network. *arXiv:1503.02531.* https://arxiv.org/abs/1503.02531

[38] Snoek, J., Larochelle, H., & Adams, R. P. (2012). Practical Bayesian optimization of machine learning algorithms. In *Advances in Neural Information Processing Systems* (pp. 2951–2959). https://papers.nips.cc/paper/2012/file/05311655a15b75fab869ac247e6e4a45-Paper.pdf

[39] Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv:1810.04805.* https://arxiv.org/abs/1810.04805

[40] Madry, A., Makelov, A., Schmidt, L., Tsipras, D., & Vladu, A. (2017). Towards deep learning models resistant to adversarial attacks. *arXiv:1706.06083.* https://arxiv.org/abs/1706.06083

[41] Sanh, V., Debut, L., Chaumond, J., & Wolf, T. (2019). DistilBERT, a distilled version of BERT: Smaller, faster, cheaper and lighter. *arXiv:1910.01108*. https://arxiv.org/abs/1910.01108

[42] McMahan, B., Moore, E., Ramage, D., Hampson, S., & y Arcas, B. A. (2017). Communication-efficient learning of deep networks from decentralized data. In *Proceedings of AISTATS* (pp. 1273–1282). https://proceedings.mlr.press/v54/mcmahan17a.html

[43] Breiman, L. (1996). Bagging predictors. *Machine Learning, 24*(2), 123–140. https://doi.org/10.1007/BF00058655

[44] Zhu, M., & Gupta, S. (2019). To prune, or not to prune: Exploring the efficacy of pruning for model compression. *arXiv:1710.01878*. https://arxiv.org/abs/1710.01878

[45] Roy, S., & Kumar, A. (2021). Optimizing RNN architectures for SMS spam detection. *Journal of Computer Science, 45*, 101234. https://doi.org/10.1016/j.jocs.2021.101234

[46] Carlini, N., & Wagner, D. (2017). Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy* (pp. 39–57). IEEE. https://doi.org/10.1109/SP.2017.49

[47] Adadi, A., & Berrada, M. (2018). Peeking inside the black-box: A survey on explainable artificial intelligence (XAI). *IEEE Access, 6*, 52138–52160. https://doi.org/10.1109/ACCESS.2018.2870052

[48] Strubell, E., Ganesh, A., & McCallum, A. (2019). Energy and policy considerations for deep learning in NLP. *arXiv:1906.02243*. https://arxiv.org/abs/1906.02243

[49] Yang, Q., Liu, Y., Chen, T., & Tong, Y. (2019). Federated machine learning: Concept and applications. *ACM Transactions on Intelligent Systems and Technology, 10*(2), 1–19. https://doi.org/10.1145/3298981