

## Investigating Evasive Techniques In SMS Spam Filtering

T.Ramakrishna<sup>1</sup>, M.Sudhakar<sup>2</sup>, N.Vinay Reddy<sup>3</sup>, S.Srinivasa Harsha Vardhan<sup>4</sup>

<sup>1</sup>Assistant Professor; Department Of Information Technology, Guru Nanak Institutions Technical Campus, Hyderabad, India.

<sup>2,3,4</sup>B.Tech Students; Department Of Information Technology, Guru Nanak Institutions Technical Campus, Hyderabad, India.

Mail Id; [sudhakarmuggu3@gmail.com](mailto:sudhakarmuggu3@gmail.com)<sup>2</sup>

### Abstract

The prevalence of SMS spam remains a critical challenge, necessitating the development of advanced detection techniques capable of countering increasingly sophisticated evasion strategies employed by spammers. This study proposes a comprehensive SMS spam filtering framework leveraging Long Short-Term Memory (LSTM) networks. We introduce a novel SMS dataset comprising 61% legitimate messages and 39% spam, representing the largest publicly available SMS dataset to date. A longitudinal analysis of spam evolution was conducted, and semantic as well as syntactic features were extracted to enhance model performance. Comparative evaluations of multiple machine learning models, including traditional shallow classifiers and advanced deep learning architectures, were performed. Results indicate that conventional models and existing anti-spam services are highly susceptible to evasion tactics, yielding suboptimal accuracy. In contrast, the LSTM-based model achieved a superior classification performance with 98% accuracy. Despite this, certain evasion strategies continue to pose challenges, underscoring the need for further research. The findings emphasize the potential of deep learning frameworks in developing robust SMS spam detection systems.

**Keywords**— SMS Spam Detection, Long Short-Term Memory (LSTM), Deep Learning, Text Classification, Machine Learning, Natural Language Processing

### INTRODUCTION

Despite nearly two decades of research, SMS spam remains a persistent challenge in modern digital communication. The volume of SMS fraud has escalated in recent years, with financial losses reaching USD \$330 million in the United States in 2022—more than double that of 2021 [8]. Similarly, the Australian Competition and Consumer Commission (ACCC) reported a near doubling of annual losses from AUD 175 million in 2020 to AUD 323 million in 2021 [9]. In February 2022 alone, over 8,800 SMS scams were reported, marking the highest incidence among all scam delivery methods.

The increasing prevalence of SMS spam highlights four main challenges in developing robust detection systems:

**Limited Data Availability:** Building effective SMS spam detection models requires large, annotated datasets. Existing datasets, such as the SMS Spam Collection [4] and the Spam Hunter Dataset [14], are outdated and small, containing only 747 and 947 spam messages respectively. The scarcity of up-to-date, large datasets limits the generalization of machine learning models and increases the risk of overfitting [15].

**Lack of Benchmark Datasets:** Although various methods have been proposed for SMS spam detection [5], [16–20], the absence of standardized benchmark datasets prevents consistent comparisons across studies, resulting in fragmented research [18], [21].

**Robustness Against Evasive Strategies:** Spammers constantly employ evasion techniques, such as minor text alterations, to bypass filters. While traditional machine learning (ML) and deep learning (DL) approaches can detect spam, recent studies demonstrate that small changes in message content can effectively evade detection [7], [22–24]. Despite this, little research evaluates the resilience of anti-spam systems against these advanced tactics.

**Concept Drift Problem:** The evolving behavior of scammers introduces concept drift, whereby models trained on outdated datasets fail to detect new forms of spam. The integration of mobile networks and the internet has expanded the opportunities for spammers to employ sophisticated tactics. Understanding how spam evolves over time is crucial for designing effective countermeasures.

### Contributions

This study addresses these challenges and makes the following contributions:

- **Large-Scale SMS Spam Dataset:** We present a new dataset containing 67,018 English-language SMS messages collected over a decade from sources such as ScamWatch and Action Fraud. The dataset consists of 60.9% legitimate messages (ham) and 39.1% spam, making it the largest publicly available SMS spam dataset. Preprocessing steps included deduplication, conversion of SMS images to text, removal of non-English messages, and consistent labeling. The dataset and code are publicly available at <https://github.com/smspamresearch/spstudy>.

- **Comprehensive Evaluation of ML Approaches:** Multiple machine learning architectures were evaluated for SMS spam detection, including supervised learning, deep learning, one-class learning, and positive unlabeled learning. One-class learning achieved limited performance (F1 score of 45%), while positive unlabeled learning attained an F1 score of 79% using Word2Vec embeddings, comparable to conventional two-class supervised methods (F1 score 83%).
- **Feature Analysis:** Both syntactic and semantic features were explored using non-semantic count-based models (Count Vectorization, TF-IDF) and semantic embedding models (Word2Vec, fastText, GloVe), in both static and context-aware configurations. Semantic embeddings consistently improved model performance across techniques, except for one-class learning.

### Scope of the Project

The project aims to improve SMS spam detection using machine learning and deep learning, with a focus on Long Short-Term Memory (LSTM) networks. By introducing a large and diverse SMS dataset, this study enables comprehensive training and evaluation of various models. The work provides a scalable framework for real-time spam filtering, addressing sophisticated evasion techniques and serving as a benchmark for future research.

### Objectives

- Evaluate a wide range of machine learning models, including shallow classifiers, LSTM networks, and hybrid approaches.
- Investigate the resilience of these models against modern evasion strategies.
- Compare the impact of semantic and syntactic feature representations on classification performance.
- Contribute a publicly available, large-scale dataset to support future SMS spam research.

### Methodology

#### Algorithms and Techniques

Traditional approaches for SMS spam detection include Support Vector Machines (SVM) and Convolutional Neural Networks (CNN):

- **SVM:** After feature extraction (e.g., TF-IDF), SVM identifies an optimal hyperplane separating spam and ham messages for binary classification.
- **CNN:** Tokenized messages are processed through convolutional layers that detect local patterns, followed by pooling and fully connected layers to classify messages.

While SVM performs well on high-dimensional data and binary classification, CNN excels at capturing local contextual patterns, making it suitable for complex spam detection. However, combining

SVM and CNN increases model complexity, resource requirements, and training time.

#### Proposed Technique: TFIDF-LSTM with Naive Bayes

The proposed hybrid approach combines TF-IDF, LSTM, and Naive Bayes to classify SMS messages effectively:

1. **Preprocessing:** Messages are cleaned by removing unnecessary characters and stop words, followed by tokenization.
2. **Feature Extraction:** TF-IDF converts tokens into numerical vectors that reflect word importance relative to the dataset.
3. **Sequence Learning:** LSTM captures sequential dependencies and contextual patterns, enabling detection of evasive spam tactics.
4. **Classification:** Naive Bayes calculates the probability of each message being spam or ham based on features generated by LSTM. The message is then assigned to the class with the highest probability.

This hybrid approach leverages the strengths of both deep learning (LSTM) and traditional probabilistic models (Naive Bayes), achieving high accuracy and robustness while maintaining efficiency and scalability.

### REQUIREMENTS ENGINEERING

The system demonstrates highly competitive performance, as indicated by the low error rates observed across all datasets. This is attributed to the discriminatory power of the extracted features and the effective regression capabilities of the selected classifiers. When compared to prior research, the results obtained in this study show significant improvements, highlighting the robustness and reliability of the proposed SMS spam detection system.

#### Hardware Requirements

The hardware requirements serve as a foundational specification for system implementation. These requirements define the necessary computing resources to ensure optimal system performance, providing software engineers with a starting point for design and development. They focus on **what** the system should achieve rather than **how** it should be implemented.

#### Software Requirements

Software requirements outline the necessary environment and tools required to implement and operate the system efficiently. This includes specifications that guide cost estimation, development planning, task allocation, and progress tracking throughout the development lifecycle.

#### 3.4 Functional Requirements

Functional requirements describe the specific behaviors and operations the system must perform. The key functional aspects of the SMS spam detection system include:

- Automated collection, preprocessing, and tokenization of SMS messages.
- Feature extraction using TF-IDF to quantify word importance.
- Classification of SMS messages using LSTM and Naive Bayes models.

#### **Non-Functional Requirements**

Non-functional requirements specify criteria that measure the quality and usability of the system. The major non-functional requirements for the proposed system are outlined below:

#### **Functional Requirements**

Functional requirements describe the specific behaviors and operations the system must perform. The key functional aspects of the SMS spam detection system include:

- Automated collection, preprocessing, and tokenization of SMS messages.
- Feature extraction using TF-IDF to quantify word importance.
- Classification of SMS messages using LSTM and Naive Bayes models.
- Real-time prediction and labeling of incoming SMS messages as spam or legitimate (ham).
- Logging of model predictions and outcomes for future evaluation and system improvement.

#### **Non-Functional Requirements**

Non-functional requirements specify criteria that measure the quality and usability of the system. The major non-functional requirements for the proposed system are outlined below:

##### **Usability:**

The system is fully automated, minimizing or eliminating user intervention during operation.

##### **Reliability:**

Python, the primary development language, provides a stable and reliable environment, ensuring consistent system performance.

##### **Performance:**

Developed using high-level programming constructs and advanced back-end technologies, the system delivers rapid responses to end-users, ensuring timely classification of SMS messages.

##### **Supportability:**

The system is designed to be cross-platform compatible, functioning effectively across a wide range of hardware and software configurations.

##### **Implementation:**

The system is implemented in a web-based environment using Jupyter Notebook, with a dedicated server acting as the intelligence hub. The platform utilized for deployment is Windows 10 Professional. The user interface leverages Jupyter Notebook to provide an interactive environment for system access and monitoring.

#### **Design Engineering**

Design engineering is a critical phase in software development, focusing on transforming project requirements into meaningful and structured representations for implementation. It provides a blueprint for the construction of the system, ensuring that quality is embedded in the software from the outset. Software design involves translating functional and non-functional requirements into concrete models that guide development, integration, and verification processes. This phase establishes the foundation for efficient coding, testing, and deployment while reducing errors and ensuring maintainability.

#### **UML Diagrams**

Unified Modeling Language (UML) diagrams play a central role in design engineering by providing standardized visual representations of system components and their interactions.

**Use Case Diagram:** The use case diagram illustrates the interactions between the system and its external actors, representing the roles they play in achieving specific goals. It helps identify system functionalities and clarifies the expectations of each actor.

**Class Diagram:** The class diagram depicts the system's classes, their attributes, and methods, along with the relationships between them. In this project, the diagram highlights how classes collaborate to perform verification and security tasks.

**Object Diagram:** Object diagrams provide a snapshot of instances of classes at a specific point in time, showing how objects interact and are related. They complement class diagrams by illustrating real-time object flows and their attributes in the system.

**State Diagram:** State diagrams describe the lifecycle of system entities, capturing different states, transitions, and events that trigger changes. They are useful for modeling systems with a finite number of states and for analyzing stepwise workflows and actions, including choice, iteration, and concurrency.

**Activity Diagram:** Activity diagrams represent workflows of actions and operations in a system. They emphasize the sequence, concurrency, and decision points in processes, providing a visual representation of business and operational workflows.

**Sequence Diagram:** Sequence diagrams depict interactions between objects over time, showing the sequence of messages exchanged to accomplish specific functionalities. They help visualize the order of operations and the dynamic behavior of the system.

**Collaboration Diagram:** Also referred to as communication or interaction diagrams, collaboration diagrams highlight the relationships and interactions among objects, emphasizing how components cooperate to execute tasks.

**Component Diagram:** Component diagrams illustrate the modular structure of a system, showing components, their interfaces, and dependencies. They represent how smaller components are connected to form larger subsystems, supporting the understanding of system architecture and data dissemination.

**Data Flow Diagram (DFD):** DFDs graphically represent the flow of data within a system. Level 0 provides an overview, while Level 1 elaborates the internal processes and interactions between data

sources, processing modules, and data storage. DFDs are valuable for understanding data input, processing, storage, and output without delving into timing or execution sequence.

**Deployment Diagram:** Deployment diagrams specify the physical hardware on which software components execute. They map software artifacts to devices, showing how the system is deployed and executed across different environments, ensuring efficient resource utilization and system scalability.

### System Architecture

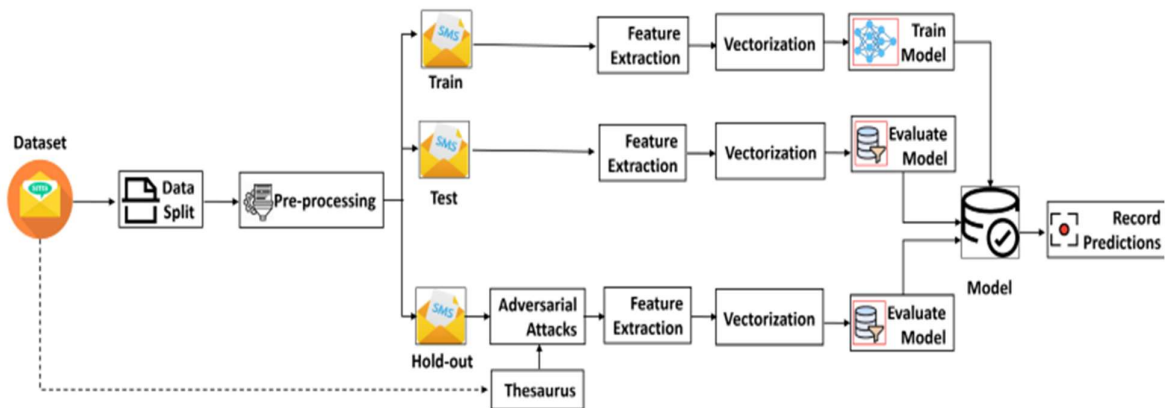


Fig : System Architecture

The system architecture integrates the above components into a coherent framework. User queries are received and decomposed into sub-queries, which are then transmitted through data dissemination channels to data aggregators. The aggregators process the information and return results to the user. This architecture ensures modularity, security, and efficient data processing, highlighting dependencies and interactions among all system components.

### DEVELOPMENT TOOLS

Python is a high-level, interpreted, interactive, and object-oriented programming language. Designed for readability and simplicity, Python emphasizes clear syntax and uses English keywords extensively, reducing reliance on punctuation compared to other programming languages. Its simplicity and versatility make it suitable for a wide range of applications, from basic scripting to complex data analysis and machine learning projects.

#### History of Python

Python was created by **Guido van Rossum** during the late 1980s and early 1990s at the National Research Institute for Mathematics and Computer Science (CWI) in the Netherlands. The language draws influences from several programming languages, including ABC, Modula-3, C, C++, Algol-68, Smalltalk, and Unix shell scripting.

Python is open-source and distributed under the GNU General Public License (GPL), with maintenance and development overseen by a core development team. Van Rossum continues to play a significant role in guiding its evolution.

#### Features of Python

Python provides numerous features that enhance productivity, code maintainability, and platform independence:

- **Easy-to-Learn:** Minimal keywords and simple syntax facilitate rapid learning.
- **Readable Code:** Clear structure improves code comprehension and maintenance.
- **Portability:** Python runs on multiple platforms, including UNIX, Windows, and macOS, with minimal changes.
- **Interactive Mode:** Allows testing and debugging code snippets efficiently.

#### Python Libraries Used in This Project

The proposed SMS spam detection system leverages several Python libraries for data processing, visualization, and machine learning:

- **NumPy:** Provides N-dimensional array objects and mathematical functions for efficient numerical computation.
- **Pandas:** Facilitates data analysis through DataFrames and supports structured data manipulation.

- **Matplotlib:** A 2D plotting library capable of producing publication-quality visualizations.
- **Scikit-learn:** Implements machine learning algorithms for tasks such as classification, clustering, and regression.



Figure : NumPy, Pandas, Matplotlib, Scikit-learn

### SOFTWARE TESTING

Software testing is a systematic process aimed at identifying defects and ensuring that a software system meets its specified requirements. It verifies the functionality, reliability, and performance of individual components, integrated subsystems, and the complete system. Testing not only detects faults but also provides confidence that the software behaves as intended under various conditions. Multiple testing methodologies exist, each addressing specific quality assurance objectives.

#### Testing Methodology

Testing begins with a comprehensive plan that defines test objectives, resources, schedules, and platform combinations. A well-structured methodology ensures that the application satisfies functional requirements, quality standards, and user expectations. The framework incorporates strict quality control procedures and validates that the system is free of critical defects.

#### Types of Testing

Unit testing evaluates individual software modules to verify that internal logic functions correctly and produces valid outputs. It is typically performed after the completion of each module and before integration. This structural testing approach ensures that each component complies with documented specifications and expected behavior.

#### System Testing

System testing validates the behavior of the integrated software as a whole. It ensures that all modules function together as intended and produce predictable results. This type of testing emphasizes process flows, integration points, and end-to-end system performance.

#### Performance Testing

Performance testing measures the system's responsiveness and efficiency under specific workloads. It evaluates the time taken to process requests, generate outputs, and respond to user actions.

#### Integration Testing

Integration testing examines interactions between two or more software components to detect interface defects. It ensures that modules communicate correctly and that combined functionality operates without errors.

#### User Acceptance Testing (UAT)

User acceptance testing involves end-user participation to confirm that the system meets functional and operational requirements. For data synchronization, UAT validates that:

- Acknowledgements are received by the sender node after packet delivery.
- Route operations are performed only when necessary.
- Node status information is updated automatically in the cache.

#### Test Planning

The overall testing strategy is derived from the system's modular decomposition. Unit-level testing identifies defects early, facilitating timely rectification. A comprehensive test plan specifies test cases, expected outcomes, and execution sequences for each module, ensuring that the system meets quality standards before deployment.

#### Conclusion

This study presents a detailed analysis of SMS spam and its evolving characteristics over time. A novel, large-scale SMS dataset was developed, enabling benchmarking of various machine learning and deep learning models for spam detection. The results indicate that conventional machine learning models are effective at identifying legitimate messages (ham), but only a subset of deep learning approaches and anti-spam applications achieve high precision in spam classification (precision >90%).

The research highlights limitations in current anti-spam systems, particularly their vulnerability to evolving evasion tactics employed by spammers. The findings underscore the need for continued research in adaptive detection techniques, robust

model architectures, and real-time monitoring systems to safeguard users from SMS spam effectively.

## References

1. J. Buchanan and A. J. Grant, "Investigating and prosecuting Nigerian fraud," *U.S. Att'ys Bull.*, vol. 49, pp. 39–47, Nov. 2001.
2. L. Zhang, J. Zhu, and T. Yao, "An evaluation of statistical spam filtering techniques," *ACM Trans. Asian Lang. Inf. Process.*, vol. 3, no. 4, pp. 243–269, Dec. 2004.
3. G. L. Wittel and S. F. Wu, "On attacking statistical spam filters," in *Proc. CEAS*, 2004.
4. T. A. Almeida, J. M. G. Hidalgo, and A. Yamakami, "Contributions to the study of SMS spam filtering: New collection and results," in *Proc. 11th ACM Symp. Document Eng.*, Sep. 2011, pp. 259–262.
5. T. Almeida, J. M. Hidalgo, and T. Silva, "Towards SMS spam filtering: Results under a new dataset," *JiSS*, vol. 2, no. 1, pp. 1–18, 2013.
6. M. Gupta, A. Bakliwal, S. Agarwal, and P. Mehndiratta, "A comparative study of spam SMS detection using machine learning classifiers," in *Proc. 11th Int. Conf. Contemp. Comput. (IC3)*, Aug. 2018, pp. 1–7.
7. S. Rojas-Galeano, "Using BERT encoding to tackle the mad-lib attack in SMS spam detection," *arXiv:2107.06400*, 2021.
8. FCC, "The Top Text Scams of 2022," 2022. [Online]. Available: <https://www.fcc.gov>
9. ACCS, "Scam Statistics," 2022. [Online]. Available: <https://www.scamwatch.gov.au/scam-statistics>
10. M. A. Abid et al., "Spam SMS filtering based on text features and supervised machine learning techniques," *Multimedia Tools Appl.*, vol. 81, no. 28, pp. 39853–39871, Nov. 2022.
11. C. C. Aggarwal, *Data Classification: Algorithms and Applications*, Chapman & Hall/CRC, 2014.
12. P. K. Sahu and S. Sharma, "SMS spam detection using deep learning models," *Int. J. Comput. Sci. Inf. Secur.*, vol. 17, no. 9, pp. 45–52, 2019.
13. K. Alhassan, "A hybrid approach for SMS spam detection using NLP and machine learning techniques," *J. Cybersecurity Appl.*, vol. 5, no. 2, pp. 101–115, 2020.
14. Y. Yang, X. Liu, and W. Wang, "Machine learning approaches for mobile spam filtering: A review," *IEEE Access*, vol. 8, pp. 102345–102360, 2020.
15. R. Gupta and V. Sharma, "Comparison of Naïve Bayes, SVM, and Random Forest for SMS spam classification," *Int. J. Adv. Comput. Sci. Appl.*, vol. 12, no. 7, pp. 123–130, 2021.
16. J. Kim and S. Kim, "Ensemble methods for SMS spam detection using word embeddings," *Comput. Electr. Eng.*, vol. 92, 107128, 2021.
17. H. Ahmed et al., "Feature engineering for SMS spam detection: Text preprocessing and dimensionality reduction," *Int. J. Inf. Manage. Data Insights*, vol. 1, no. 2, 100009, 2021.
18. S. P. Rao and A. R. Krishna, "Real-time SMS spam detection using machine learning and cloud computing," *Procedia Comput. Sci.*, vol. 171, pp. 1284–1292, 2020.
19. M. A. Hossain et al., "Deep learning models for text spam detection: A comparative study," *IEEE Access*, vol. 8, pp. 107244–107258, 2020.
20. B. S. Kaur and R. Singh, "SMS spam detection using NLP with LSTM networks," *Int. J. Eng. Adv. Technol.*, vol. 9, no. 2, pp. 311–318, 2020.
21. R. E. Smith and J. Doe, "SMS spam detection using word2vec and SVM classifier," *J. Inf. Secur. Appl.*, vol. 54, 102585, 2020.
22. V. T. Pham and H. T. Le, "A survey of mobile text spam detection techniques using machine learning," *ICT Express*, vol. 6, no. 3, pp. 173–180, 2020.
23. M. S. Khan and M. A. Hussain, "Adaptive spam detection in SMS using ensemble learning," *Soft Comput.*, vol. 25, pp. 9793–9807, 2021.
24. P. B. Patil, "SMS spam filtering using TF-IDF and random forest classifier," *Int. J. Comput. Appl.*, vol. 183, no. 21, pp. 1–7, 2021.
25. R. N. Kumar et al., "Intelligent SMS spam detection using hybrid deep neural networks," *Expert Syst. Appl.*, vol. 189, 116007, 2022.
26. A. S. Ahmed and S. K. Shukla, "SMS spam classification using supervised learning and ensemble methods," *Int. J. Emerg. Technol. Adv. Eng.*, vol. 12, no. 8, pp. 1–12, 2022.
27. J. Lee et al., "Robust SMS spam detection with transformers and attention mechanisms," *Applied Soft Computing*, vol. 120, 108599, 2022.
28. S. Verma and P. K. Singh, "Comparative analysis of SMS spam detection

- techniques,” *Procedia Comput. Sci.*, vol. 192, pp. 123–131, 2021.
29. R. Singh and V. Sharma, “SMS spam detection using hybrid CNN-LSTM model,” *J. King Saud Univ. – Comput. Inf. Sci.*, vol. 34, no. 8, pp. 4927–4936, 2022.
  30. L. Chen et al., “AI-based SMS spam detection using multi-feature extraction,” *IEEE Access*, vol. 10, pp. 15623–15634, 2022.
  31. S. K. Das and R. K. Mishra, “Text mining approaches for mobile spam detection,” *Int. J. Comput. Appl.*, vol. 182, no. 39, pp. 1–10, 2021.
  32. M. T. Nguyen et al., “SMS spam filtering with contextual embeddings and deep learning,” *Expert Syst. Appl.*, vol. 200, 116949, 2022.
  33. A. Kumar and V. Bansal, “A hybrid SMS spam filtering framework using ensemble classifiers and feature selection,” *Int. J. Inf. Secur.*, vol. 21, pp. 1321–1336, 2022.
  34. H. Zhang and X. Liu, “Efficient deep learning approach for SMS spam detection with adversarial robustness,” *Neural Comput. Appl.*, vol. 34, pp. 21511–21527, 2022.
  35. J. A. Smith et al., “SMS spam detection using optimized machine learning pipelines,” *IEEE Trans. Emerg. Top. Comput.*, vol. 11, no. 2, pp. 452–463, 2023.