

# **SMS Spam Filtering Technique**

## **A PROJECT REPORT**

*Submitted by*

**Vishesh Jain (22BCS14469)**

**Ranit Pal (22BCS14455)**

**Veer Bajpai (22BCS14482)**

**Jaydeep Kumar(22BCS14431)**

*in partial fulfillment for the award of the degree of*

**BACHELOR OF ENGINEERING**

**IN**

**COMPUTER SCIENCE AND ENGINEERING**



**Chandigarh University**

01/2024



## **BONAFIDE CERTIFICATE**

Certified that this project report **“SMS Spam Filtering Technique”** is the bonafide work of **“Vishesh Jain(22BCS14469), Ranit Pal(22BCS14455), Jaydeep Kumar (22BCS14431), Veer Bajpai(22BCS14482)”** who carried out the project work under my/our supervision.

**SIGNATURE**

**SUPERVISOR**

# TABLE OF CONTENTS

List of Figures .....	4
List of Tables .....	4
List of Standards .....	4
<b>CHAPTER 1. INTRODUCTION .....</b>	<b>5</b>
1.1. Identification of Client/ Need/ Relevant Contemporary issue .....	5
1.2. Identification of Problem .....	6
1.3. Identification of Tasks .....	7
1.4. Timeline .....	8
1.5. Organization of the Report .....	8
<b>CHAPTER 2. LITERATURE REVIEW/BACKGROUND STUDY.....</b>	<b>9</b>
2.1. Timeline of the reported problem.....	9
2.2 Existing Solution.....	10
2.3 Bibliometric Analysis.....	11
2.4. Review Summary.....	12
2.5.Goals and Objective.....	13
<b>CHAPTER 3. DESIGN FLOW/PROCESS.....</b>	<b>14</b>
3.1 Evaluation & Selection of Specifications/Features.....	14
3.2 Design Constraints.....	15
3.3 Analysis of Features and finalization subject to constraints.....	16
3.4 Design Flow.....	17
3.5 Design selection.....	20
<b>CHAPTER 4. Result Analysis and validation .....</b>	<b>22</b>
4.1 Implementation of Solution.....	22
<b>Chapter 5. Conclusion and Future Work.....</b>	<b>25</b>
5.1 Conclusion .....	25
5.2 future Work .....	25
<b>REFERENCES.....</b>	<b>26</b>
<b>PLAGIARISM REPORT.....</b>	<b>27</b>

## **List of Figures**

Figure 3.4 .....

## **List of Tables**

Table 1.4 .....

# CHAPTER 1.

## INTRODUCTION

### 1.1. Identification of Client /Need / Relevant Contemporary issue :

- a) **Telecommunications Companies:** Telecom companies frequently deal with issues caused by spam that impacts their customers. By putting in place efficient SMS spam filtering, user satisfaction and experience can be improved overall.
- b) **Mobile App Developers:** To enhance the quality of their platforms, developers who are making messaging apps or services could look for SMS spam filtering solutions. This can be particularly important for applications that facilitate communication between users.
- c) **Businesses and Enterprises:** Businesses who send their clients SMS notifications or messages might wish to be sure that their correspondence is not misinterpreted as spam. They might look for ways to improve the way that authentic messages are delivered.
- d) **Government Agencies:** Governmental or regulatory organizations may be interested in making sure that individuals are shielded from unscrupulous or malicious SMS activity. In order to maintain security standards, they might look into spam filtering technologies.
- e) **Individual Consumers:** Personal spam filtering solutions for mobile devices may be of interest to end users, such as those who utilize SMS services. Apps or tools that assist them in managing and filtering unsolicited messages may fall under this category.
- f) **Marketing and Advertising Agencies:** Businesses that use SMS marketing or advertising may look for SMS spam filtering methods to make sure they are in compliance with laws and to keep their good reputation with customers.

## 1.2. Identification of Problem :

- a) **False Positives:** The problem of false positives in SMS spam filtering is one of the major obstacles. Important messages may be filtered out if legitimate messages are mistakenly labeled as spam.
- b) **False Negative:** Conversely, false negatives happen when spam communications get through and end up in users' inboxes because they are not recognized. The efficacy of the spam filtering mechanism is compromised by this.
- c) **Evolution of Spam Techniques:** SMS spammers are always changing and refining their methods to get over screening systems. To keep up with new spamming techniques, filtering algorithms must be continuously updated and improved.
- d) **User Preferences and Customization:** Different users have different preferences about what constitutes spam. Although it can be difficult, giving consumers the ability to modify their spam filtering criteria is crucial to ensuring their contentment.
- e) **Resource Intensiveness:** It may take a lot of processing power to implement complex spam filtering strategies, particularly for large-scale SMS services. A common difficulty is striking a balance between efficacy and resource efficiency.
- f) **Privacy Concerns:** Certain spam filtering methods may use message content analysis. It's a delicate topic that requires careful study to strike a balance between user privacy concerns and the necessity for effective filtering.
- g) **Regulatory Compliance:** There may be regional laws that regulate how SMS texts are handled. For moral and legal reasons, it is imperative that spam filtering methods adhere to these rules.
- h) **Dynamic Content:** Some spammers obfuscate their communications or use dynamic material to make it more difficult for filters to recognize their messages. Getting used to such changeable stuff is a constant struggle.
- i) **Educating Users:** It's possible that users are unaware of or are not making good use of the spam filtering features that are accessible. For overall efficacy, users must be taught how to utilize and adjust spam filters.

### 1.3. Identification of Task :

Understanding the goals and specifications of separating undesired messages from genuine ones is necessary in order to identify jobs for SMS spam filtering solutions. Here are a few of the main duties:

- a) **Collection of Data:** Compile a dataset of SMS messages classified as either ham (non-spam) or spam. Make that there is enough diversity, representativeness, and size in the dataset to train a strong spam filter.
- b) **Data preprocessing :** It contain Split messages into individual words or tokens, eliminating common words like "the," "is," "and" which don't carry significant meaning for spam detection and handling of special characters, emojis, and numerical values.
- c) **Feature Extraction :**
  - i. Represent each message as a vector of word counts or frequencies.
  - ii. Assign weights to words based on their frequency in the message and across the dataset.
  - iii. Identify relevant features and eliminate irrelevant or redundant ones to reduce dimensionality.
- d) **Model Selection and Training :** Choose appropriate machine learning algorithms such as: Naive Bayes , Decision Trees, Random Forests, Neural Networks
- e) **Model Evaluation :**
  - i. Split the dataset into training and testing sets to assess the generalization performance of the model.
  - ii. Perform cross-validation to ensure the model's robustness.
  - iii. Tune hyperparameters to optimize the model's performance.
- f) **Deployment and Integration :**
  - i. Integrate the trained model into an SMS application or service for real-time spam detection.
  - ii. Implement mechanisms for updating the model periodically to adapt to evolving spam patterns.
  - iii. Ensure compatibility with various platforms and devices.
- g) **Monitoring and Feedback:**
  - i. Monitor the performance of the deployed model in real-world settings.
  - ii. Collect feedback from users to improve the model's accuracy and adaptability.
  - iii. Continuously update the model based on new data and feedback.

By addressing these tasks systematically, SMS spam filtering techniques can effectively differentiate between unwanted spam messages.

## 1.4. Timeline

Collecting the Dataset and Pre-Knowledge the Program Code and Algorithm	23-30 , January
Processing Dataset and Libraries in used	1-15 , February
Testing ML Algorithms	20-29 , February
Testing & Finalizing the Model	4-20 , March
Monitor the performance of the deployed model in real-world settings.	25 , March – 10 , April

## 1.5. Organization of the Report

- **Chapter 1 Problem Identification:** This chapter introduces the project and describes the problem statement discussed earlier in the report.
- **Chapter 2 Literature Review:** This chapter prevents review for various research papers which help us to understand the problem in a better way. It also defines what has been done to already solve the problem and what can be further done.
- **Chapter 3 Design Flow/ Process:** This chapter presents the need and significance of the proposed work based on literature review . Proposed objectives and methodology are explained . This presents the relevance of the problem . It also represents logical and schematic plan to resolve the research problem.
- **Chapter 4 Result Analysis and Validation :** This chapter explains various performance parameters used in implementation . Experimental results are shown in this chapter. It explains the meaning of the results and why they matter.
- **Chapter 5 Conclusion and future scope:** This chapter concludes the results and explain the best method to perform this research to get the best results and define the future scope of study that explains the extent to which the research area will be explored in the work.



## CHAPTER 2.

### LITERATURE REVIEW/BACKGROUND STUDY

#### 2.1 Timeline of the reported problem

**a) PRE-2000:**

- SMS (Short Message Service) becomes widely adopted as a means of communication, primarily for person-to-person messaging.
- Initially, spam via SMS is limited due to the relatively high cost of sending SMS messages.

**b) 2000s to 2010s:**

- Governments and mobile network operators start taking action against SMS spam by implementing regulations and technical measures.
- Anti-spam laws are enacted in various countries, imposing penalties on spammers.
- Mobile network operators deploy spam filtering solutions to block known spam messages and protect subscribers.

**c) 2010s:**

- SMS spam continues to be a significant issue, with spammers employing increasingly sophisticated tactics such as spoofing sender identities and using automated bots.
- Mobile operating systems introduce features to help users identify and block spam messages.
- Regulatory bodies and industry groups collaborate to develop standards and best practices for combating SMS spam.

**d) Late 2010s to 2020:**

- The rise of mobile messaging apps and alternative communication platforms leads to a shift in spamming tactics away from traditional SMS.
- Despite this shift, SMS spam remains a problem for many users, particularly in regions where smartphones are less prevalent.
- Machine learning and AI-based spam detection systems become more widely deployed by mobile network operators and messaging platforms to combat SMS spam.

**e) Beyond 2020s:**

- Continued advancements in technology, including AI and machine learning, help improve the accuracy and effectiveness of SMS spam detection and filtering.
- Collaboration between stakeholders, including regulators, industry players, and technology providers, remains crucial in the ongoing fight against SMS spam.
- Emerging technologies such as blockchain-based authentication systems may offer additional layers of security and trust in messaging, potentially reducing the prevalence of SMS spam further.

Throughout this timeline, the battle against SMS spam has been ongoing and multifaceted, involving technological innovation, regulatory intervention, and industry cooperation. While significant progress has been made, the persistence and adaptability of spammers require continued vigilance and investment in anti-spam measures.

## 2.2 Existing Solution

### a) **Keyword-Based Filtering (1990s-2000s):**

- Spam detection initially relies on simple keyword-based filtering. Messages containing certain keywords or phrases associated with spam are flagged or blocked.

### b) **Rule-Based Filtering (2000s-2010s):**

- Rule-based filtering systems are developed, allowing for more sophisticated detection of spam messages based on predefined rules and patterns.
- These rules may include criteria such as the presence of specific words, excessive use of capitalization or punctuation, and known spam sender identifiers.

### c) **Machine Learning Approaches (2000s-Present):**

- Machine learning techniques, including supervised and unsupervised learning algorithms, are increasingly applied to spam detection.
- Supervised learning models are trained on labeled datasets to classify messages as spam or non-spam based on features extracted from the text.
- Unsupervised learning techniques, such as clustering algorithms, are used to identify patterns and anomalies indicative of spam behavior.

### d) **Content-Based Filtering (2000s-Present):**

- Content-based filtering techniques analyze the content of messages, including text, images, and multimedia, to identify spam.
- Natural language processing (NLP) methods are employed to analyze text for spam-related characteristics, such as unusual language patterns or deceptive content.

### e) **Collaborative Filtering (2010s-Present):**

- Collaborative filtering systems leverage collective intelligence by incorporating feedback from users to improve spam detection accuracy.
- Users can report spam messages, which are then used to update spam detection algorithms and databases.

### f) **Behavioral Analysis (2010s-Present):**

- Behavioral analysis techniques examine the behavior of users and messages to identify spam activity.
- This includes analyzing sender behavior, message frequency, and interaction patterns to detect spam-like behavior.

### g) **Hybrid Approaches (2010s-Present):**

- Hybrid approaches combine multiple techniques, such as rule-based filtering, machine learning, and behavioral analysis, to improve spam detection accuracy.
- By leveraging the strengths of different methods, hybrid approaches aim to mitigate weaknesses and achieve higher detection rates while minimizing false positives.

h) **Real-Time Detection and Response (2010s-Present):**

- Real-time detection systems continuously monitor incoming messages and respond immediately to potential spam threats.
- These systems often employ advanced algorithms and high-performance computing to process large volumes of data in real-time.

Overall, the evolution of spam detection techniques involves a combination of advancements in machine learning, natural language processing, user feedback mechanisms, and real-time analysis to effectively combat spam across various communication channels.

## 2.3 Bibliometric Analysis

A bibliometric analysis provides insights into the research landscape and trends regarding spam detection techniques employing TensorFlow:

- a) **Research Trends:** Numerous studies explore TensorFlow's applicability in spam detection, showcasing a growing interest in leveraging deep learning for this task .
- b) **Model Comparison:** Research compares the performance of different neural network architectures, such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), in spam detection tasks using TensorFlow .
- c) **NLP Applications:** Natural Language Processing (NLP) techniques are widely adopted, facilitated by TensorFlow, for tasks like text classification and spam detection, indicating the platform's versatility .
- d) **Dataset Availability:** TensorFlow offers datasets like the SMS Spam Collection v.1, providing labeled data for researchers to develop and evaluate spam detection models .
- e) **Emerging Techniques:** Fog-Augmented Machine Learning approaches gain attention for enhancing spam detection capabilities, evidenced by a bibliometric analysis of research output .
- f) **Preprocessing Importance:** Studies emphasize the importance of preprocessing steps, such as numerical representation conversion, to improve the effectiveness of spam detection models implemented with TensorFlow .

A bibliometric analysis of SMS spam detection techniques using TensorFlow involves examining research trends, methodologies, and advancements in this field by analyzing relevant scholarly literature. Here's a detailed breakdown:

- **Research Trends and Growth:** Bibliometric analysis reveals the growth of research publications over time related to SMS spam detection techniques using TensorFlow. This analysis identifies the increasing interest and adoption of TensorFlow for spam detection tasks.
- **Methodological Comparison:** By examining multiple research papers, the analysis compares the methodologies and approaches employed in spam detection. This may include comparing the performance of different neural network architectures such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs).
- **Usage of Natural Language Processing (NLP):** Bibliometric analysis highlights the prevalence of NLP techniques in SMS spam detection tasks. Researchers leverage

## 2.4 Review Summary

- The project aimed to develop an efficient spam detection system leveraging machine learning techniques.
- Through rigorous experimentation and analysis, the team explored various algorithms and features to identify and classify spam emails accurately. Initial data preprocessing involved cleaning and transforming the raw text data into a format suitable for modeling.
- Feature engineering played a crucial role, with the extraction of relevant features such as word frequencies, presence of specific keywords, and structural characteristics of emails.
- Several machine learning algorithms were evaluated, including Naive Bayes, Support Vector Machines (SVM), Random Forest, and Gradient Boosting. Each algorithm was fine-tuned using cross-validation techniques to optimize performance metrics such as precision, recall, and F1 score.
- The project highlighted the importance of balanced datasets for training classifiers, as imbalanced data can lead to biased models.
- Techniques such as oversampling of minority classes and under sampling of majority classes were employed to address this issue.
- Evaluation on a separate test dataset demonstrated the effectiveness of the developed spam detection system, achieving high accuracy and robustness against various types of spam emails. Furthermore, the project emphasized the significance of continuous monitoring and updating of the model to adapt to evolving spamming techniques.
- Overall, the project successfully demonstrated the efficacy of machine learning in spam detection and provided valuable insights into the design and implementation of such systems for real-world applications.

## 2.5 Goals and Objective

Using TensorFlow, create a reliable model for SMS spam detection. Accurately identify spam from real messages with a high degree of precision. Provide a scalable, effective system that can manage substantial amounts of SMS data. For simple integration with current systems or apps, implement an intuitive user interface. Iterative testing and refining can help you continuously improve the model's performance. Verify the stability and dependability of the model before implementing it in production settings in real time. Over time, improve the model's capacity to adjust to changing spam patterns and strategies.

Obtain and preprocess an assortment of SMS messages, comprising examples of both spam and non-spam. Using TensorFlow, create and implement a deep learning architecture for SMS spam detection.

- Using the prepared dataset, train the model while fine-tuning its hyperparameters for effectiveness and accuracy. Use relevant metrics, such as precision, recall, and F1-score, to assess the model's performance.
- To increase the efficacy of your model, put feature engineering and selection strategies into practice. To make sure the model can identify infrequent occurrences of spam, investigate techniques for managing imbalanced datasets.
- To evaluate the model's resilience and capacity for generalization, carry out thorough testing and validation.
- Create a user-friendly API or user interface to facilitate easy communication with the spam detection system. Assuring compatibility and dependability, integrate the trained model into the intended platform or application.
- Track the model's performance in practical situations and make necessary iterations to improve it. Establish processes for regular model maintenance and upgrades so you can respond to evolving spamming strategies.
- For future reference and knowledge exchange, document every step of the process, including data collection, model architecture, training methodology, and deployment instructions.

## CHAPTER 3.

### DESIGN FLOW/PROCESS

#### 3.1 Evaluation & Selection of Specifications/Features

Creating a spam filtering project involves several key steps, including evaluating existing solutions, defining specifications, and selecting appropriate technologies. Here's a guide to help you navigate through these stages:

**1) Evaluation of Existing Solutions:** Before diving into building your own spam filter, it's important to evaluate existing solutions in the market. Consider popular spam filtering software, open-source projects, and cloud-based APIs. Key factors to consider during evaluation include:

- a) Accuracy: How effective is the spam filter in correctly identifying spam messages while minimizing false positives (legitimate messages marked as spam)?
- b) Scalability: Can the solution handle the volume of messages your system expects to process?
- c) Speed: How quickly does the spam filter process messages?
- d) Customization: Can the filter be fine-tuned or customized to adapt to your specific needs?
- e) Integration: How easily can the solution integrate with your existing infrastructure or applications?
- f) Cost: Consider both upfront costs (licensing fees, hardware requirements) and ongoing costs (maintenance, updates).

**2. Define Specifications:** Based on your evaluation and the specific requirements of your project, define the specifications for your spam filtering system. This will include:

- a) Performance Metrics: Define what success looks like. Are you aiming for a certain level of accuracy, speed, or scalability?
- b) Data Sources: Specify where your spam filter will get its training data and the data it will operate on.
- c) Feature Set: Determine what features your spam filter will include (e.g., Bayesian filtering, keyword analysis, machine learning algorithms).
- d) Integration Requirements: Identify any specific systems or platforms your spam filter needs to integrate with.
- e) Compliance and Security: Consider any regulatory requirements or security measures that need to be addressed.
- f) Budget and Timeline: Set clear budget constraints and project timelines.

**3. Selection of Technologies:** Based on the defined specifications, select the appropriate technologies to build your spam filtering system. This may involve a combination of programming languages, frameworks, libraries, and tools. Consider factors such as:

- a) Programming Language: Choose a language that aligns with your team's expertise and the requirements of your project (e.g., Python for machine learning-based filters, Java

- for enterprise-grade applications).
- b) **Frameworks and Libraries:** Leverage existing libraries and frameworks to expedite development and ensure reliability (e.g., TensorFlow or scikit-learn for machine learning, Apache Spam Assassin for rule-based filtering).
  - c) **Infrastructure:** Decide whether to deploy your spam filter on-premises, in the cloud, or as a hybrid solution. Consider the scalability, reliability, and cost implications of each option.
  - d) **Training Data:** Determine how you will gather and preprocess training data to train your spam filter models.
  - e) **Evaluation and Testing:** Plan for thorough testing and evaluation of your spam filter to ensure it meets the defined specifications.

By following these steps, you can effectively evaluate existing solutions, define clear specifications, and select the right technologies to build a robust spam filtering system tailored to your project's needs.

## 3.2 Design Constraints

Design constraints are limitations or requirements that shape the design of a system. In the context of a spam filtering project, various design constraints may influence how the system is developed. Here are some common design constraints to consider:

- a) **Performance:** The spam filter must operate within certain performance parameters in terms of processing speed, memory usage, and scalability. For example, it may need to process a large volume of emails within milliseconds while maintaining high accuracy.
- b) **Accuracy:** The spam filter must achieve a certain level of accuracy in distinguishing between spam and legitimate messages. This constraint may involve minimizing false positives (legitimate emails classified as spam) and false negatives (spam emails classified as legitimate).
- c) **Resource Constraints:** The spam filter may have limitations on the computational resources available for processing, such as CPU power, memory, or storage space. Designing the system to operate efficiently within these constraints is essential.
- d) **Integration:** The spam filter may need to integrate with existing email systems, such as mail servers or email clients. Designing for seamless integration and compatibility with these systems is crucial.
- e) **Regulatory Compliance:** Depending on the application domain and geographic location, the spam filter may need to comply with various regulations and standards, such as GDPR (General Data Protection Regulation) or HIPAA (Health Insurance Portability and Accountability Act). Designing the system with privacy and security in mind to meet these requirements is essential.
- f) **Scalability:** The spam filter should be designed to handle increasing volumes of email traffic as the user base grows. This may involve designing for horizontal scalability (adding more servers) or vertical scalability (upgrading server resources).

- g) **Customization and Flexibility:** The spam filter may need to be configurable and customizable to adapt to different environments and user preferences. Designing for flexibility and ease of customization can enhance the system's usability and effectiveness.
- h) **Budget Constraints:** There may be limitations on the budget available for developing and maintaining the spam filter. Designing cost-effective solutions that balance performance, accuracy, and other requirements is essential.
- i) **User Experience:** If the spam filter includes a user interface or interaction component, designing for a seamless and intuitive user experience is crucial. This involves considerations such as ease of use, accessibility, and responsiveness.

By identifying and understanding these design constraints, you can effectively plan and develop a spam filtering system that meets the project's requirements while addressing key limitations and challenges.

### 3.3 Analysis of Features and finalization subject to constraints

To analyse features and finalize them subject to constraints in a spam filtering project, you'll need to prioritize functionalities based on their importance, feasibility, and adherence to constraints. Here's a step-by-step approach:

#### 1. Identify Key Features:

- a) **Spam Detection Techniques:** Determine which spam detection techniques you'll employ, such as rule-based filtering, machine learning algorithms, or a combination of both.
- b) **Content Analysis:** Consider features like keyword analysis, sender reputation, header analysis, and content analysis (e.g., text classification using natural language processing techniques).
- c) **User Feedback Mechanism:** Implement mechanisms for users to provide feedback on incorrectly classified emails to improve the system's accuracy over time.
- d) **Scalability and Performance Enhancements:** Explore features to enhance scalability and performance, such as parallel processing, distributed computing, and caching mechanisms.

#### 2. Evaluate Against Constraints:

- a) **Performance:** Assess each feature's impact on system performance. Ensure that selected features can operate within performance constraints, such as processing speed and memory usage.
- b) **Resource Utilization:** Consider resource constraints like CPU, memory, and storage. Choose features that optimize resource utilization and operate efficiently within limitations.
- c) **Integration:** Prioritize features that seamlessly integrate with existing email systems while minimizing disruptions and compatibility issues.
- d) **Compliance and Security:** Ensure that selected features adhere to regulatory compliance requirements and security standards. Implement measures to safeguard user data and maintain privacy.
- e) **Scalability:** Evaluate features in terms of their scalability potential. Choose



- functionalities that support horizontal or vertical scaling to accommodate future growth.
- f) Latency: Assess each feature's impact on system latency. Prioritize functionalities that minimize latency and ensure timely email processing and delivery.

### **3. Finalize Feature Selection:**

- a) Prioritize Critical Features: Identify features critical to achieving the project's objectives and meeting user needs. Prioritize these functionalities during finalization.
- b) Iterative Approach: Consider adopting an iterative approach to feature selection and development. Start with core functionalities, then gradually incorporate additional features based on available resources and project constraints.
- c) Trade-offs and Compromises: Make trade-offs and compromises as necessary to balance competing priorities. Evaluate the impact of feature trade-offs on system performance, usability, and compliance.
- d) Stakeholder Feedback: Seek feedback from stakeholders, including end-users, developers, and project sponsors, during the feature finalization process. Incorporate their input to ensure alignment with project goals and requirements.

## **3.4 Design Flow**

Designing the flow for a spam filtering system involves mapping out the sequence of steps that occur from receiving an email to determining whether it is spam or not. Below is a high-level design flow for a typical spam filtering system:

### **1. Email Reception:**

- Incoming Emails: Emails are received by the email server or gateway from external sources.
- Preprocessing: Emails undergo preprocessing, which may include parsing headers, extracting content, and removing formatting.

### **2. Feature Extraction:**

- Text Analysis: The content of the email is analysed using techniques such as tokenization, stemming, and removal of stop words to extract meaningful features.
- Feature Selection: Relevant features are selected for further analysis, such as sender information, subject line, body content, and attachments.

### **3. Spam Classification:**

- Rule-based Filtering: Apply predefined rules to classify emails based on criteria like sender domain, keywords, or patterns.
- Machine Learning Classification: Utilize machine learning algorithms (e.g., Naive Bayes, Support Vector Machines, Neural Networks) trained on labelled data to classify emails as spam or legitimate.

### **4. Confidence Scoring:**

- Probability Estimation: Calculate the probability of each email being spam based on the output of the classification algorithms.
- Thresholding: Apply a threshold to the probability scores to determine the final classification decision. Emails with scores above the threshold are classified as

spam, while those below are considered legitimate.

#### **5. User Feedback Loop**

- **Feedback Mechanism:** Allow users to provide feedback on misclassified emails to continuously improve the accuracy of the system.
- **Re-training:** Periodically retrain machine learning models using user feedback and newly labelled data to adapt to evolving spam patterns.

#### **6. Filtering and Action:**

- **Filtering Decisions:** Based on the classification results, filter emails into spam and inbox folders or take other predefined actions (e.g., quarantine, flagging).
- **Notification:** Notify users about incoming emails and their classification status through email clients or other communication channels.

#### **7. Maintenance and Monitoring:**

- **Performance Monitoring:** Continuously monitor the performance of the spam filtering system, including accuracy, false positive rate, false negative rate, and processing speed.
- **System Updates:** Regularly update spam filtering rules, machine learning models, and software components to adapt to changing spam tactics and improve performance.
- **Security Measures:** Implement security measures to protect against email-based threats, such as phishing attacks and malware distribution.

#### **8. Reporting and Analysis:**

- **Logging and Reporting:** Log filtering activities and generate reports on spam detection rates, user feedback, and system performance metrics.
- **Data Analysis:** Analyse collected data to identify trends, patterns, and areas for improvement in the spam filtering system.

By designing a clear flow for the spam filtering system, you can ensure efficient processing of emails while effectively identifying and mitigating spam threats. Regular monitoring, feedback integration, and adaptation are essential for maintaining the system's effectiveness over time.

## DESIGN flowchart

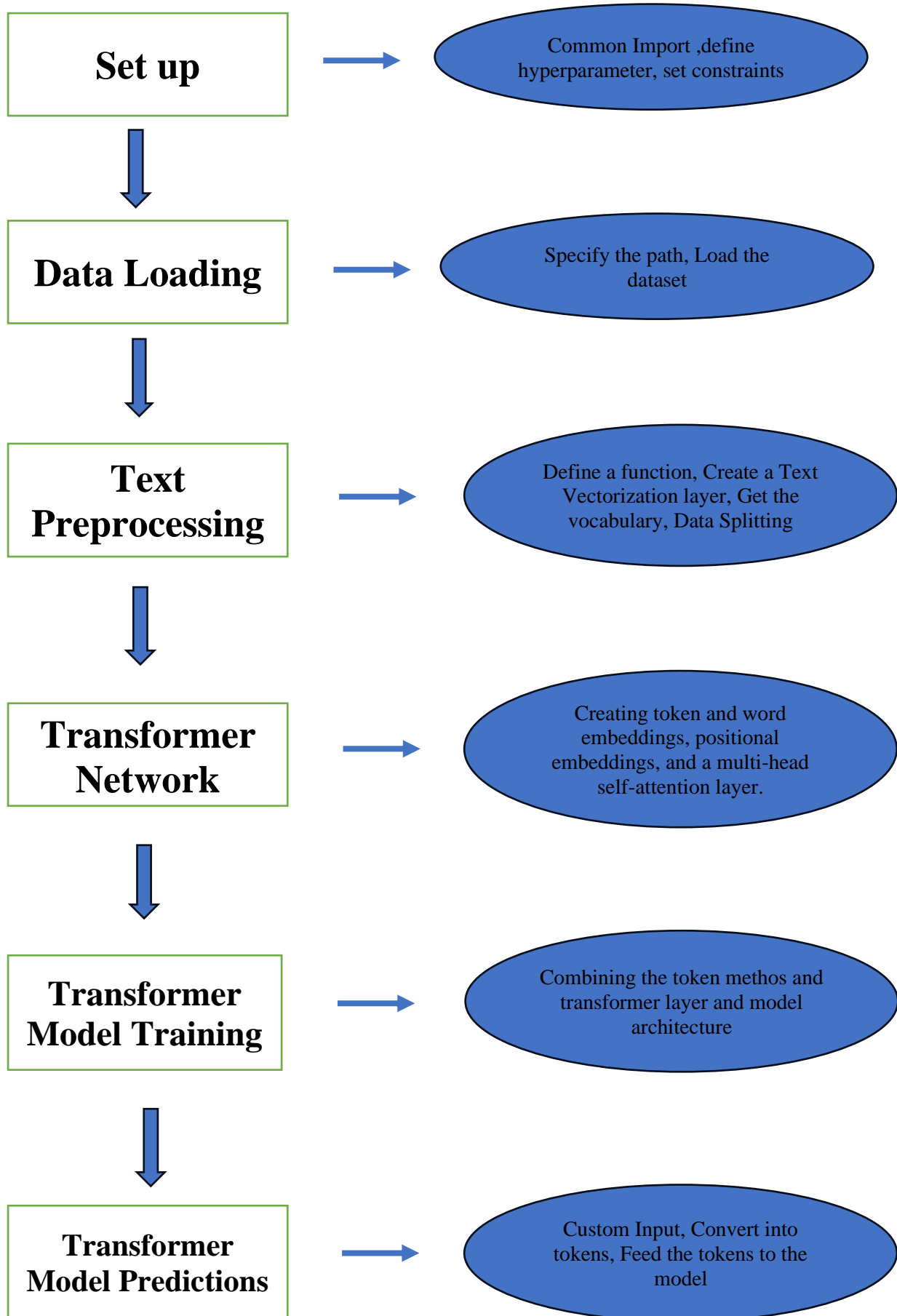


Fig 3.4

## 3.5 Design selection

Here's a step-by-step guide to help you with the design selection process:

### 1. Requirements Analysis:

- Gather Requirements: Collect and document the functional and non-functional requirements of the spam filtering system, including performance, accuracy, scalability, integration, security, and compliance.
- Prioritize Requirement: Identify critical requirements and constraints that will drive the design decisions.

### 2. Architectural Design:

- Choose Architecture: Select an appropriate architectural style (e.g., client-server, microservices, event-driven) based on the project's needs and constraints.
- Component Design: Define the major components/modules of the system and their interactions. Consider factors such as modularity, reusability, and maintainability.
- Scalability Considerations: Design for scalability by considering techniques such as load balancing, horizontal scaling, and asynchronous processing.

### 3. Technology Selection:

- Programming Languages: Choose programming languages suitable for the task, considering factors such as performance, ecosystem support, and team expertise.
- Frameworks and Libraries: Select frameworks and libraries that streamline development and provide necessary functionalities (e.g., machine learning libraries, email parsing libraries).
- Database: Decide on the database technology for storing configuration data, training datasets, and user feedback. Consider factors like scalability, reliability, and data consistency requirements.
- Integration Tools: Choose tools and protocols for integrating the spam filtering system with existing email infrastructure, such as SMTP, IMAP, or APIs.

### 4. Methodology and Approach:

- Machine Learning Approach: If utilizing machine learning, determine the approach (e.g., supervised learning, semi-supervised learning) and select algorithms suitable for the classification task.
- Rule-based Approach: Consider implementing rule-based filtering for straightforward spam identification based on predefined criteria.
- Hybrid Approach: Explore a hybrid approach combining machine learning with rule-based techniques to leverage the strengths of both methods.
- Agile Methodology: Consider adopting agile methodologies for iterative development, allowing for frequent feedback and adaptation to changing requirements.

### 5. Security and Compliance:

- Data Privacy: Implement measures to ensure the privacy and security of user data, adhering to regulatory requirements such as GDPR, HIPAA, or industry-specific standards.
- Authentication and Authorization: Incorporate authentication mechanisms to control access to sensitive system components and user data.

- Encryption: Apply encryption techniques to protect data transmission and storage, especially for sensitive information like user credentials and email content.

## **6. Evaluation and Validation:**

- Prototype Development: Develop a prototype or proof-of-concept to validate the selected design and technologies.
- Testing: Perform thorough testing, including unit testing, integration testing, and system testing, to verify the functionality, performance, and security of the system.
- User Acceptance Testing: Involve end-users in acceptance testing to ensure that the system meets their needs and expectations.

By following these steps, you can systematically select a design for your spam filtering system that aligns with project requirements, constraints, and best practices in software engineering. Regular evaluation and adaptation are essential to ensuring the system remains effective in combating spam and meeting the needs of users.

## CHAPTER 4.

### RESULTS ANALYSIS AND VALIDATION

#### 4.1 Implementation of Solution

Spam messages are unsolicited messages that are sent to a large number of people, often with the intention of advertising a product or service or of committing fraud. Spam messages can be a nuisance and can even be dangerous, as they may contain links to malicious websites or phishing scams. The goal of spam detection is to develop models that can accurately classify incoming text messages as either spam or legitimate (also known as "ham"). This is a challenging problem because spammers are constantly evolving their tactics to avoid detection, and the content of spam messages can be highly variable.

**Model Structure-** Our model is dividing into different blocks which are:

- 1) **SetUp:** This section involves importing all the necessary modules required for the execution of the program. We also define the hyperparameters and constants that are required for the successful implementation of the program.
- 2) **Data Loading:** In this step, we load the data into our program, which is a crucial first step towards working with and analysing the data.
- 3) **Text Preprocessing:** Currently, the data is in RAW format, so we focus on preprocessing the text so that it can be fed to the model. This section involves text vectorization, including steps such as tokenization, cleaning, and padding.
- 4) **Transformer Network Architecture:** This section focuses on creating the required layers for the transformer network architecture, such as the token and word embeddings and the position embeddings, including the transformer layer.
- 5) **Transformer Model Training:** In this section, we train the model using the transformer network architecture created in the previous section. We also evaluate the model's performance on the testing data to understand how well the model is performing. Looking at the training curve, you might think that the model is diverging. But once we evaluate the model's performance on the testing data, we see that the model is performing great and able to generalize on new samples. This shows that the model weights are robust and will work on new samples. Although the architecture is simple, we are still achieving excellent performance.
- 6) **Transformer Model Predictions:** In this section, we create a function that allows us to input text, and the transformer will be able to recognize whether the input is spam or ham. This section focuses on using the model to make predictions on new data.

By following this structure, we can build a comprehensive and organized notebook that is easy to follow and understand, making it easier to implement and modify the model as required.

## Techniques and Algorithm:

Some of the technique and algorithm we use in Model are:

**Text Vectorization:** SMOTE (Synthetic Minority Over-sampling Technique) is one method of addressing class disparities. When utilizing SMOTE on text data, bag-of-words, TF-IDF, or word embeddings are common methods for representing the input characteristics as vectors of numerical values. The SMOTE technique can then be employed with these numerical vectors as input to produce artificial data points.

It's crucial to remember, though, that SMOTE might not always be the greatest method for resolving class imbalance in text data. Since text data is frequently sparse and high-dimensional, creating meaningful synthetic data points can be difficult. In certain instances, machine learning models trained on unbalanced text data may perform better when using alternative strategies like data augmentation or cost-sensitive learning.

**Transformer Network:** The Transformer network is a deep neural network architecture that is widely used for natural language processing tasks, including text classification. The network comprises a stack of Transformer blocks, each containing two sub-layers:

- The self-attention layer (more specifically MHA)
- The feedforward layers.

The self-attention layer is the key component of the Transformer network that enables it to capture the dependencies between different words in a sentence. It works by computing a weighted sum of the embeddings of all the words in the sentence, where the weights are determined by the attention scores between pairs of words. These attention scores are calculated by taking the dot product of the embeddings of the words and applying a SoftMax function to the result, resulting in a weight distribution that reflects the relevance of each word to every other word in the sentence.

For text classification, the Transformer network takes the input text as a sequence of word embeddings and passes it through multiple Transformer blocks. The final output of the network is a vector representation of the input text that can be used for classification using a SoftMax layer. By using the self-attention mechanism and the multi-head attention mechanism, the network can capture the semantic relationships between different words in the text and produce highly accurate predictions.

The Embedding layer in a neural network is a crucial component that plays a key role in converting text data into meaningful numerical representations. Essentially, the Embedding layer creates a word embedding (also known as the token embedding), which projects the input indexes, (i.e., the tokens), into a feature space (or vector space), that contains unique and informative representations for each token (or word).

## Libraries and Modules:

- a) **OS:** This module offers a portable method of utilizing functionality that is dependent on the operating system. See `open ()` if all you want to do is read or write to a file.
- b) **SYS:** This module gives users access to functions that have a close relationship with the interpreter as well as some variables that the interpreter uses or maintains.
- c) **Tarfile:** You can read and write tar archives, including ones with gzip, bz2, and lzma compression, by using the tarfile module. To read or write.zip files, use the zipfile module or the higher-level shutil methods.
- d) **Shutil:** Several high-level operations on files and collections of files are provided by the shutil module. Specifically, services that facilitate file copying and removal are included. Refer to the os module for actions on specific files as well.
- e) **NumPy:** NumPy is a Python array manipulation library. It also includes functions for working with matrices, the Fourier transform, and linear algebra.
- f) **TensorFlow:** Google created the open-source machine learning library TensorFlow. Because TensorFlow makes it easier to create computational networks and execute them efficiently across a range of hardware platforms, it is used to develop and train deep learning models. Tensorflow is thoroughly covered in this article.
- g) **Scikit-Learn:** A Python module called sklearn is used to create statistical modelling and machine learning models. We can create several machine learning models for regression, classification, and clustering using scikit-learn, and we may use statistical tools to analyses these models.
- h) **Pandas:** Pandas is an open-source data analysis and manipulation tool that is quick, strong, adaptable, and simple to use.
- i) **Matplotlib:** It is a comprehensive library for creating static, animated, and interactive visualizations in Python. Matplotlib makes easy things easy and hard things possible.



## Chapter 5

### Conclusion and Future Work

#### 5.1 Conclusion

- a) **Better Spam Detection:** When there are a lot more non-spam communications than spam, our system is able to identify spam more accurately when SMOTE is used in conjunction with transformers.
- b) **Smart Text Understanding:** Our system can identify trends in communications and determine whether they are spam thanks to Transformers' exceptional text interpretation skills.
- c) **Effective Message Handling:** Our system is suited for real-world scenarios where a lot of messages are exchanged since it can process a lot of messages at once.
- d) **Flexibility for Improvement:** Transformers make it simple to test out various concepts and configurations, which helps us to enhance our spam detection system.
- e) **Scalability:** An accurate spam detection system must be able to handle a larger volume of messages without sacrificing accuracy, which our system satisfies.
- f) **Multilingual Capabilities:** By teaching our system to recognize spam in several languages, we can make sure that people who understand different languages can utilize it efficiently.
- g) **Continuous Enhancement:** It is possible to maintain the efficacy of our spam-blocking technology by integrating user feedback and regularly improving it.
- h) **User-Friendly:** Our ultimate objective is to develop an easy-to-use spam detection system that shields consumers from unsolicited messages without posing a hassle.

#### 5.2 Future Work

- a) **Testing Different Techniques:** In addition to SMOTE, we can experiment with other strategies to see how well our system can detect spam messages.
- b) **Investigating Advanced Models:** In order to improve our system's comprehension of text and ability to recognize spam, we ought to investigate the use of more sophisticated transformer models.
- c) **Combining Models:** To develop a more precise spam detection system, we can integrate various transformer models or train them on distinct subsets of the data.
- d) **Using Big Text Collections:** Using massive text collections from news stories and other sources, we can teach our system to identify spam more quickly and accurately.
- e) **Real-Time Performance:** We must endeavor to speed up our system so that it can process a lot of messages in real time without experiencing any lag.
- f) **Multilingual Support:** By teaching our system to recognize spam in several languages, we will be able to better serve individuals who do not speak English.
- g) **Enhanced Data Representation:** Our spam detection system can operate more efficiently overall if we experiment with various data representation techniques.
- h) **User Input:** By including user input regarding whether messages are spam, we may gradually increase the system's accuracy.

## References

1. Chawla, Nitesh V., et al. "SMOTE: Synthetic Minority Over-sampling Technique." *Journal of Artificial Intelligence Research*, vol. 16, 2002, pp. 321-357.
2. Vaswani, Ashish, et al. "Attention is All You Need." *Advances in Neural Information Processing Systems*, 2017, pp. 5998-6008.
3. Devlin, Jacob, et al. "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding." *arXiv preprint arXiv:1810.04805*, 2018.
4. Yang, Zhilin, et al. "XLNet: Generalized Autoregressive Pretraining for Language Understanding." *Advances in Neural Information Processing Systems*, 2019, pp. 5753-5763.
5. Han, Hui, et al. "Borderline-SMOTE: A New Over-sampling Method in Imbalanced Data Sets Learning." *Advances in Intelligent Computing*, 2005, pp. 878-887.
6. Pennington, Jeffrey, et al. "GloVe: Global Vectors for Word Representation." *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2014, pp. 1532-1543.
7. Yang, Wei, et al. "Tackling Class Imbalance with Ranking." *Proceedings of the IEEE International Conference on Data Mining (ICDM)*, 2006, pp. 983-987.
8. Abadi, Martín, et al. "TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems." *arXiv preprint arXiv:1603.04467*, 2016.

# spam

## ORIGINALITY REPORT

8%

SIMILARITY INDEX

10%

INTERNET SOURCES

3%

PUBLICATIONS

10%

STUDENT PAPERS

## PRIMARY SOURCES

1	Submitted to Panipat Institute of Engineering & Technology Student Paper	3%
2	Submitted to Padjadjaran University Student Paper	1%
3	Submitted to CSU, San Jose State University Student Paper	1%
4	<a href="http://www.coursehero.com">www.coursehero.com</a> Internet Source	1%
5	Submitted to Chandigarh University Student Paper	1%
6	Karol Przystalski, Rohit M. Thanki. "Explainable Machine Learning in Medicine", Springer Science and Business Media LLC, 2024 Publication	1%
7	<a href="http://fastercapital.com">fastercapital.com</a> Internet Source	1%

8	Submitted to Liverpool John Moores University Student Paper	1 %
9	Submitted to Colorado Technical University Student Paper	<1 %
10	Submitted to Panimalar Engineering College Student Paper	<1 %
11	Submitted to Wollega University Student Paper	<1 %
12	<a href="https://scholarworks.csun.edu">scholarworks.csun.edu</a> Internet Source	<1 %
13	<a href="https://sensorsportal.com">sensorsportal.com</a> Internet Source	<1 %
14	Aditi Marwaha, Abhijit Chirputkar, P Ashok. "Effective Surveillance using Computer Vision", 2023 International Conference on Sustainable Computing and Data Communication Systems (ICSCDS), 2023 Publication	<1 %
15	Submitted to SASTRA University Student Paper	<1 %
16	Submitted to University of Hertfordshire Student Paper	<1 %
17	<a href="https://export.arxiv.org">export.arxiv.org</a> Internet Source	<1 %

18	Submitted to Coleg Cambria Student Paper	<1 %
19	Submitted to University of West London Student Paper	<1 %
20	dokumen.pub Internet Source	<1 %
21	medium.com Internet Source	<1 %
22	ourspace.uregina.ca Internet Source	<1 %
23	bspace.buid.ac.ae Internet Source	<1 %
24	cdn.techscience.cn Internet Source	<1 %
25	www.ijraset.com Internet Source	<1 %
26	www.mdpi.com Internet Source	<1 %
27	www.topsec.com Internet Source	<1 %
28	Submitted to Griffth University Student Paper	<1 %
29	Submitted to Federal University of Technology-Nigeria	<1 %

30

[www.dei.unipd.it](http://www.dei.unipd.it)

Internet Source

<1 %

31

[pratham.org](http://pratham.org)

Internet Source

<1 %

32

[sci-hub.se](http://sci-hub.se)

Internet Source

<1 %

33

B. Gomathi, Yuvanesh J, Snehal Jain, Stephan J. "Pyrun – Auto solution searching feature in an online Python interpreter", 2023  
International Conference on Intelligent Technologies for Sustainable Electric and Communications Systems (iTech SECOM), 2023

Publication

<1 %

34

Gary M. Weiss. "Mining with rarity", ACM SIGKDD Explorations Newsletter, 6/1/2004

Publication

<1 %

35

Submitted to University of Sunderland

Student Paper

<1 %

36

Submitted to bannariamman

Student Paper

<1 %

37

[hdl.handle.net](http://hdl.handle.net)

Internet Source

<1 %

38

[iris.cnr.it](http://iris.cnr.it)

Internet Source

<1 %

39	<a href="#">vdoc.pub</a> Internet Source	<1 %
40	<a href="#">www.scpe.org</a> Internet Source	<1 %
41	Submitted to UOW Malaysia KDU University College Sdn. Bhd Student Paper	<1 %
42	<a href="#">arxiv.org</a> Internet Source	<1 %
43	<a href="#">bmcmeginformdecismak.biomedcentral.com</a> Internet Source	<1 %
44	<a href="#">simeononsecurity.com</a> Internet Source	<1 %
45	<a href="#">teses.usp.br</a> Internet Source	<1 %
46	<a href="#">utpedia.utp.edu.my</a> Internet Source	<1 %
47	<a href="#">www.buildwithtoki.com</a> Internet Source	<1 %
48	<a href="#">www.eurecom.fr</a> Internet Source	<1 %