

# Enhancing ECA (Event- Condition-Action) Rules: Fine-Tuning BERT for Security and Privacy Violation Detection

Hemavathi<sup>1</sup>, Kalpana R<sup>2</sup>, Kavitha M<sup>3</sup>, Swetha Rani L<sup>4</sup>

<sup>1</sup>Dept. of ECE, B.M.S. College of Engineering, Bengaluru-560019, Karnataka, India.

<sup>2</sup>Dept. of ECE, B.M.S. College of Engineering, Bengaluru-560019, Karnataka, India.

<sup>3</sup>Dept. of ECE, JSS Academy of Technical Education, Bengaluru-560060, Karnataka, India.

<sup>4</sup>Dept. of ECE, AMC Engineering College, Bengaluru-560083, Karnataka, India.

\*Corresponding Email: hemavathi.ece@bmsce.ac.in

## ARTICLE INFO

Received: 18 Dec 2024

Revised: 10 Feb 2025

Accepted: 28 Feb 2025

## ABSTRACT

In the world of Internet of Things, Event Condition Action rules are the secret sauce for smart device interaction. An Event triggers the rule. If a specific Condition is met, then an Action happens automatically. This article addresses trigger–action platforms, which empower users to define custom behaviors for IoT devices and web services through conditional rules. While these platforms enhance user creativity in automation, they also pose significant risks, such as unintentional disclosure of private information or exposure to cyber threats. The proposed solution leverages Natural Language Processing techniques to identify automation rules within these platforms that may compromise user security or privacy. Natural Language Processing based models are applied to analyze the semantic and contextual information of trigger–action rules, utilizing classification techniques on various rule features. Evaluation on the If-This-Then-That platform, using a dataset of 76,741 rules labeled through an ensemble of three semi-supervised learning techniques, demonstrates that the results from the Bidirectional Encoder Representations from Transformers based model training demonstrate promising outcomes, with an average validation accuracy of 89% over 2 epochs. The Test Accuracy of around 90.65% is achieved. Predicted outputs showcase the model's ability to categorize applets into different risk classes, including instances of cyber security threats and physical harm.

**Keywords:** NLP, IFTTT, BERT, IOT Platforms, Privacy and Security, Trigger–Action Rules.

## INTRODUCTION

The advent of Internet of Things technology has marked a significant paradigm shift, especially within residential settings. With an increasing number of users integrating IoT devices like connected cameras, smart locks, and smoke detectors into their homes, there's a palpable shift towards creating smarter living environments aimed at streamlining daily tasks and enhancing convenience. As these devices become more ubiquitous, they reshape the dynamics of household interactions, offering users unprecedented levels of control and efficiency.

This work represents a concerted effort to develop a sophisticated system poised at the nexus of NLP techniques and neural networks, with a singular focus on fortifying the security and privacy of smart home ecosystems. At its core, the work endeavors to autonomously pinpoint potential risks inherent in the intricate web of Event Condition Action rules governing the operation of IoT devices. Through the strategic deployment of advanced machine learning models, notably including BERT, the overarching aim is to furnish users with a meticulously safeguarded smart home

experience characterized by heightened levels of security and privacy assurance. By delving into the nuances of ECA rules and leveraging cutting-edge AI capabilities, the work seeks to provide a robust defense against emerging threats in the ever- evolving landscape of IoT security.

### Trigger-Action Platforms

Trigger-Action Platforms have garnered substantial popularity in the domain of IOT owing to their simplicity and adaptability. TAPs empower users to craft bespoke automation by delineating conditional rules that dictate the interactions between devices and services. These rules afford users the ability to delineate how and when their devices ought to react to specific triggers, thereby tailoring the automation experience to align with individual needs and preferences.

### ECA Rules

ECA rules serve as the cornerstone of TAPs, facilitating the automation of interactions among diverse IoT devices. ECA rules comprise three fundamental components: Event: The trigger instigating the rule's activation (e.g., motion detected by a security camera).

Condition: The explicit criteria mandating fulfillment for the action to be executed (e.g., time of day falls between 6 PM and 6 AM).

Action: The response effectuated when the event transpires and the condition is satisfied (e.g., illumination of outdoor lights).

This framework presents a straightforward yet potent mechanism for users to forge personalized smart environments without necessitating intricate technical expertise. In this example, the security camera continuously monitors the front entrance. When the camera detects motion (the event), the system checks the current time. If the motion occurs between 6 PM and 6 AM (the condition), the system activates the outdoor lights (the action) for 10 minutes. This setup ensures that the lights only turn on when it is likely to be dark outside, conserving energy while providing security and convenience. The homeowner does not need to manually control the lights or have advanced technical knowledge to implement this automation. The ECA rule seamlessly integrates the motion detection capability of the security camera with the lighting system, creating a smart environment that adapts to the homeowner's needs and enhances the safety and functionality of their home.

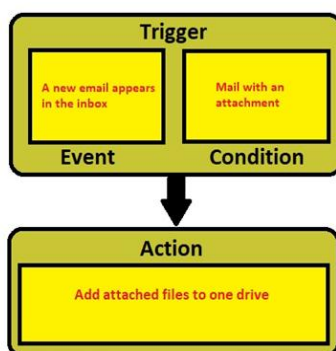


Fig. 1 Example of how an ECA rule works

### IFTTT

IF This Then That stands out as a prominent illustration of a TAP, providing a web- based service that empowers users to craft straightforward conditional statements termed as applets. These applets automate diverse tasks and interactions spanning across different web applications, devices, and online services, thereby augmenting the functionality and efficacy of users' digital and physical environments.

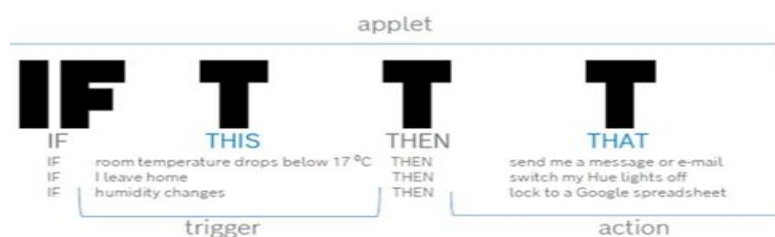


Fig. 2 IFTTT Breakdown

### How IFTTT Works

The operation of IFTTT revolves around users establishing applets through the configuration of conditional statements, wherein each applet comprises a trigger (the "If This" part) and an action (the "Then That" part). For instance, an applet might read as follows: "If the front door is unlocked, then turn on the hallway light." This user-friendly format renders it accessible to a wide audience, empowering even those devoid of technical prowess to capitalize on the potential of automation. Furthermore, IFTTT provides a vast library of pre-existing applets and supports the creation of custom applets tailored to individual preferences and requirements. Users can browse through a multitude of applets covering various functionalities, ranging from integrating smart home devices to automating social media tasks. The platform's intuitive interface facilitates the effortless creation and management of applets, enabling users to set up automation sequences with minimal effort.

The versatility of IFTTT extends beyond basic home automation, encompassing a wide array of use cases across different domains. Users can leverage IFTTT to streamline workflows, receive notifications, synchronize data between applications, and much more. For instance, an individual may utilize IFTTT to automatically save email attachments to cloud storage, receive weather updates via text message, or even control their smart thermostat based on their GPS location. Additionally, IFTTT fosters interoperability by offering seamless integration with numerous popular services and platforms, including social media networks, smart home ecosystems, productivity tools, and IoT devices. This interconnectedness enables users to orchestrate complex automation scenarios involving multiple services and devices, thereby enhancing their overall digital experience. In essence, IFTTT serves as a catalyst for simplifying and enhancing the integration of various digital services and devices, empowering users to customize and automate their interactions according to their preferences and needs. Whether it's streamlining routine tasks, optimizing productivity workflows, or orchestrating smart home environments, IFTTT provides a versatile and accessible platform for users to harness the power of automation.

### Role in IoT

In the domain of IoT, neural networks serve pivotal roles in processing the vast quantities of data generated by diverse sensors and devices. Through sophisticated algorithms and learning mechanisms, neural networks excel in discerning intricate patterns and correlations embedded within this data. This capability enables them to facilitate advanced automation within smart home systems, imbuing them with heightened responsiveness and intelligence. Technically, neural networks consist of layers of interconnected nodes, each node representing a computational unit akin to a neuron. These nodes are organized into input, hidden, and output layers, with each layer responsible for specific tasks. In the context of IoT, the input layer receives data from various sensors and devices, which is then processed through successive layers of hidden neurons. These hidden layers perform complex computations, extracting meaningful features and patterns from the raw data. Finally, the output layer generates predictions or actions based on the processed information, enabling the system to make decisions or initiate responses.

Moreover, neural networks employ sophisticated learning algorithms, such as backpropagation and gradient descent, to iteratively adjust the connections between neurons during training. This process, known as supervised learning, allows the network to refine its predictive capabilities by minimizing prediction errors. Additionally, advanced neural network architectures, such as convolutional neural networks (CNNs) and recurrent neural networks (RNNs), offer specialized capabilities for processing different types of data, such as images, time-series data. In the context of IoT applications, neural networks can perform a myriad of tasks, including anomaly detection, predictive maintenance,

adaptive control, and NLP. For example, in a smart home environment, neural networks can analyze sensor data to detect anomalies indicative of potential security breaches or equipment failures. They can also predict future system states based on historical data, enabling proactive maintenance to prevent costly breakdowns. Furthermore, neural networks can facilitate adaptive control algorithms that optimize energy consumption, comfort levels, and resource utilization within smart homes.

Overall, the integration of neural networks into IoT systems enhances their cognitive capabilities, enabling them to process vast amounts of data, extract valuable insights, and make informed decisions in real-time. This advancement paves the way for smarter, more efficient, and autonomous IoT applications, revolutionizing the way of interacting with the environments.

### **Application in IoT**

NLP is used to automatically identify potential security and privacy risks in individual ECA rules by deeply analyzing their semantics. This involves multiple steps: tokenizing the text to break down the rules, triggers, and actions into manageable units, parsing the grammatical structure to understand the relationships between different components, and applying machine learning algorithms to detect and classify potential risks. Tokenization splits the text into words and phrases, making it easier to handle, while parsing involves constructing a syntactic structure that represents the sentence grammatically. This structural understanding allows the system to comprehend context and nuances in the rules. The heart of the analysis lies in leveraging advanced NLP techniques, such as named entity recognition (NER) to identify specific entities mentioned in the rules, and dependency parsing to determine how different parts of the rule relate to each other. Machine learning models, particularly those based on neural networks, are trained on large datasets to recognize patterns associated with security and privacy risks. Techniques like vector embeddings are used to capture the semantic meaning of words and phrases, enabling the model to understand subtle variations in language.

Once potential risks are identified, the system categorizes them based on their severity and type. For instance, a rule that disables a security alarm when a door is unlocked might be flagged as a high-risk rule due to its potential security implications. The system then generates proactive guidance to help users mitigate these risks. This guidance could include suggestions for modifying the rule, adding additional conditions to enhance security, or providing alerts when certain high-risk actions are triggered. By integrating these NLP techniques, the work aims to significantly enhance the overall security and privacy of smart home environments, ensuring that automation not only adds convenience but also maintains user safety.

### **Bidirectional Encoder Representations from Transformers**

Bidirectional Encoder Representations from Transformers is a powerful NLP model developed by Google that has significantly advanced the state of the art in understanding the context of words in sentences. BERT's ability to consider the bidirectional context of words makes it particularly effective in comprehending and generating natural language.

### **Integration with IoT**

In this work, BERT is utilized to analyze the natural language descriptions of ECA rules. By leveraging BERT's deep understanding of language, the system can identify potential risks associated with specific automation rules and provide recommendations to mitigate these risks, thereby enhancing the overall security of smart home environments. The proliferation of IOT technology has transformed homes into smart, automated environments, offering numerous benefits to users. Smart devices, such as connected cameras, smart locks, and thermostats, are now commonplace, allowing users to create customized automation through ECA rules. These rules specify actions to be taken when specific events and conditions occur, simplifying daily tasks. For example, an ECA rule might state, "If the front door is unlocked, then turn on the hallway lights," providing both convenience and security by ensuring the home is well-lit upon entry. This ability to automate and customize how devices interact not only enhances the user experience but also allows for greater efficiency and energy management within the home.

However, this surge in IoT and ECA rule-based automation has given rise to a critical problem - the potential for security and privacy risks, as well as threats to cybersecurity, personal damage, and physical damage. Users, often lacking technical knowledge, may inadvertently create ECA rules that compromise their personal data, privacy, or the security of their smart environments. The intuitive nature of creating these rules can sometimes obscure the complex security implications they might entail. For instance, a rule like "If a motion is detected at the front door, then send a live feed to my smartphone" could potentially expose sensitive footage to unauthorized viewers if not properly secured. Some ECA rules may pose security risks, as exemplified by the rule "Open the shutters in the living room when the temperature of the house changes and goes above 25°C." While this rule seems convenient, it may inadvertently leave windows open during hot summer days when the house is empty, potentially providing an entry point for thieves. Such oversights can have serious security implications, turning what was meant to be a feature of comfort and efficiency into a vulnerability. Another example might be an ECA rule that disarms the security system when a Bluetooth device (such as a smartphone) is nearby, which could be exploited if the Bluetooth signal is spoofed.

Users often underestimate the risks associated with ECA rules, assuming that IoT device manufacturers are responsible for safeguarding their privacy. This misplaced trust can lead to a false sense of security, where users believe that their devices are inherently secure just because they come from reputable manufacturers. However, the responsibility for ensuring the security of these devices and the rules governing their operation often falls on the users themselves. This is particularly concerning given that many users may not have the necessary expertise to identify and mitigate these risks effectively. Consequently, the potential for cybersecurity threats, personal damage, and physical damage becomes a significant concern in the increasingly interconnected landscape of smart home environments. This work endeavors to create a robust solution for identifying security and privacy risks associated with IFTTT applets. This task is challenging due to several factors, including the diverse nature of applets, an imbalanced dataset, and the necessity to thoroughly grasp their contextual semantics. Applets can vary significantly in complexity and functionality, making it difficult to categorize them accurately. Additionally, the dataset may contain disproportionately more examples of certain risk categories, complicating the training process. Understanding the contextual nuances of applets is crucial for accurately assessing their potential risks.

The proposed solution takes a multi-faceted approach that leverages NLP and machine learning techniques to enhance the accuracy of IFTTT applet classification. The first step involves curating an extensive dataset comprising a diverse range of applets. These applets are then categorized into four risk classes: innocuous, personal harm, physical harm, and cybersecurity harm. This classification scheme allows us to capture the varying degrees of risk associated with different types of applets and tailor the analysis accordingly. Subsequently, data preprocessing is conducted to extract essential features from the applets, including trigger events, actions, channel information, and applet descriptions. This preprocessing step is crucial for preparing the data for analysis and ensuring that the model can effectively learn from it.

## **Deep Learning Model**

Deep learning models have become the cornerstone of many advanced AI applications due to their ability to learn complex patterns and representations from large datasets. For this work, state-of-the-art BERT is leveraged model to tackle the intricate task of classifying IFTTT applets based on their security and privacy risks. Deep learning, particularly models like BERT, excels in NLP tasks because of its capacity to understand and interpret the nuanced context of words within sentences. Unlike traditional machine learning models that rely on hand-crafted features, deep learning models automatically learn relevant features from raw data, making them highly effective for complex tasks.

BERT's architecture is particularly well-suited for this work because it processes input text bidirectionally, meaning it considers the entire context of a word by looking at both its preceding and succeeding words. This bidirectional approach contrasts with models like LSTM (Long Short-Term Memory), which process input sequentially and may not capture context as effectively. By leveraging BERT's deep understanding of language, the model can better interpret the semantics of IFTTT applet descriptions, triggers, and actions. This results in more accurate classifications and improved detection of potential security and privacy risks. Additionally, BERT's pre-trained



contextual embeddings allow it to perform well even with limited labeled data, as it leverages knowledge from vast amounts of text data it was trained on.

The BERT chosen over LSTM for several reasons. Firstly, BERT's bidirectional nature allows it to capture richer contextual information compared to LSTM, which processes sequences in one direction at a time. This results in more accurate and nuanced interpretations of the text, which is crucial for assessing the risk levels of applets. Secondly, BERT has been pre-trained on a large corpus of text, enabling it to generalize better and require less task-specific labeled data for fine-tuning. Lastly, BERT's transformer architecture allows for parallel processing of input data, making it more efficient and scalable compared to the sequential nature of LSTM models. This efficiency is particularly important given the large and diverse dataset of applets that are working with.

### **Raspberry Pi as Edge Solution**

The edge solution in the work leverages Raspberry Pi to provide localized, secure, and private predictions without relying on an internet connection. Raspberry Pi was chosen over Arduino due to its superior processing capabilities and the ability to run complex machine learning models like BERT. While Arduino excels in simple, real-time control tasks, it lacks the computational power and full-fledged operating system support required for intensive data processing and model inference. Raspberry Pi, on the other hand, can run a full Linux operating system and supports a wide range of software and programming languages, making it highly versatile for various applications.

By using Raspberry Pi, the sensitive data remains within the user's control, addressing concerns about data privacy and sovereignty. This is particularly important for users who operate in environments with limited connectivity or have strict data privacy requirements. The Raspberry Pi's compact size and low power consumption make it an ideal choice for edge computing applications, allowing users to deploy the solution in a variety of settings without significant infrastructure investments. Additionally, the extensive community support and documentation available for Raspberry Pi facilitate easy setup, troubleshooting, and customization, ensuring a smooth user experience.

### **Website as Cloud Solution**

In addition to the edge solution, a user-friendly website is created, that can be deployed as a cloud-based solution. This website will provide users with convenient access to the IFTTT applet risk prediction service from any device with an internet connection. Users can simply visit the website, input their applet details, and receive instant feedback on the potential security and privacy risks. The cloud-based approach offers several advantages, including scalability, ease of access, and centralized management of the prediction service. This ensures that users can benefit from the latest model updates and improvements without needing to manage the infrastructure themselves.

The website will feature a clean, intuitive interface designed to cater to both technical and non-technical users. It will provide detailed explanations of the predicted risks and offer actionable recommendations to mitigate potential threats. This approach not only enhances the security and privacy of IoT ecosystems powered by IFTTT but also empowers users to make informed decisions about their applet configurations. By offering both cloud-based and edge solutions, a wide range of user preferences and requirements is catered, ensuring that the work provides a comprehensive and flexible solution for enhancing the security and privacy of smart home environments.

## **BACKGROUND**

IoT ECA (Event-Condition-Action) rules are fundamental components in IoT systems that govern the behavior of devices based on specific events, conditions, and subsequent actions. When an event occurs, such as a sensor detecting a change in temperature, the system evaluates predefined conditions, like whether the temperature exceeds a certain threshold. If the conditions are met, corresponding actions are triggered, such as adjusting the thermostat to regulate the temperature. These rules enable automation and intelligent decision-making in IoT ecosystems, allowing devices to respond dynamically to changing environmental conditions or user-defined criteria. They form the backbone of many IoT applications, facilitating efficient resource management, enhanced user experiences, and the automation of various processes across diverse industries.

BERT (Bidirectional Encoder Representations from Transformers) is a cutting-edge natural language processing (NLP) model developed by Google. It revolutionized the field by introducing bidirectional context awareness, allowing it to capture deeper semantic meaning from text. Unlike previous models that processed text in a unidirectional manner, BERT considers both preceding and following words when encoding each word's representation. This bidirectionality enables BERT to understand the context of a word within a sentence more effectively. BERT achieved state-of-the-art performance across various NLP tasks, including question answering, sentiment analysis, and named entity recognition, by pre-training on vast amounts of text data and fine-tuning on specific tasks.

The Word embeddings are dense vector representations of words in a high-dimensional space, typically created using techniques like Word2Vec, GloVe, or FastText. In natural language processing (NLP), word embeddings serve as a fundamental tool for representing textual data in a format that is more easily interpretable by machine learning algorithms. Handling unbalanced datasets is crucial in machine learning to prevent biased models and ensure accurate predictions, especially in classification tasks where one class significantly outnumbers the others. Several methods exist to address this issue. One approach is resampling, which involves either oversampling the minority class or under sampling the majority class to achieve a more balanced distribution. Another technique is using different evaluation metrics such as precision, recall, F1-score, or area under the ROC curve (AUC) that are less sensitive to class imbalance.

The paper [1] by J. Cano et al. delves into the application of ECA rules within reactive systems, specifically in the context of the IOT. Acknowledging the complexity of rule interactions in such systems, particularly in highly distributed IoT applications, the authors advocate for runtime coordination and formal analysis to address potential side effects, crucial in critical applications. The paper presents a case study focusing on safe application development in IoT by extending the ECA semantic through discrete control. The proposed approach defines safety properties for interactions and introduces autonomous controllers to support the distribution of ECA rules, aiming to enhance the safety and reliability of IoT applications. Overall, this work provides valuable insights into designing secure and dependable IoT systems.

In the paper [2] authored by B. Breve, G. Cimino, and V. Deufemia, the authors address the potential risks associated with trigger-action platforms, which allow users to define custom behaviors for IoT devices and Web services through conditional rules. While these platforms foster user creativity in automation, they pose significant security and privacy concerns, such as unintentional data disclosure and vulnerability to cyber threats. The paper proposes leveraging NLP techniques to detect automation rules that may violate user security or privacy. The NLP models presented capture semantic and contextual information by applying classification techniques to various rule features. The evaluation, conducted on If-This- Then-That platform with a dataset of 76,741 labeled rules using semi-supervised learning techniques, shows that the BERT-based model achieves high precision and recall. The research highlights the efficacy of NLP in enhancing user safety and data protection in trigger-action IoT platforms.

The paper [3] authored by M. J. Jozani, É. Marchand, and A. Parsian, explores the use of balanced loss functions for estimating an unknown parameter  $\theta$ . The proposed loss functions, denoted as  $L_{\rho, \omega, \delta_0}(\theta, \delta)$ , take the form  $\omega \rho(\delta_0, \delta) + (1 - \omega) \rho(\theta, \delta)$ , where  $\rho(\theta, \delta)$  is an arbitrary loss function,  $\delta_0$  is a predefined "target" estimator of  $\theta$ , and  $\omega$  is a weight ranging from 0 to 1. The paper introduces weighted versions of these functions,  $q(\theta) L_{\rho, \omega, \delta_0}(\theta, \delta)$ , with  $q(\cdot)$  being a positive weight function. The authors discuss the general development of Bayesian estimators under  $L_{\rho, \omega, \delta_0}$ , establishing connections between such estimators and Bayesian solutions for the unbalanced case ( $\omega = 0$ ). The paper provides illustrative examples for various choices of  $\rho$ , including absolute value, entropy, linex, intrinsic (model-based), and a generalized form of squared error losses. This research contributes to the understanding and application of balanced loss functions in the context of parameter estimation.

The paper [4] authored by N. Reimers and I. Gurevych, addresses the computational challenges associated with BERT and RoBERTa in sentence-pair regression tasks, such as semantic textual similarity (STS). While BERT and RoBERTa exhibit state-of-the-art performance in STS, their requirement to process both sentences in a pair leads to significant computational overhead. The paper introduces Sentence-BERT (SBERT), a modification of the pretrained

BERT network. SBERT utilizes siamese and triplet network structures to generate semantically meaningful sentence embeddings, allowing for efficient cosine-similarity comparisons. This modification reduces the computational effort for finding the most similar pair from 65 hours with BERT/roBERTa to about 5 seconds with SBERT, while maintaining accuracy. The authors evaluate SBERT on STS tasks and transfer learning, demonstrating its superior performance compared to other state-of-the-art sentence embedding methods. The work addresses the efficiency challenges of BERT-based models in tasks requiring sentence similarity search and unsupervised clustering

The paper [5] authored by L. Wang, M. Han, X. Li, N. Zhang, and H. Cheng, provides an in-depth analysis of classification techniques tailored for unbalanced datasets. The paper begins by introducing unbalanced datasets and subsequently delves into a comprehensive examination of classification methods, considering various perspectives such as data sampling methods, algorithmic approaches, feature-level strategies, cost-sensitive functions, and the application of deep learning. The authors categorize data sampling methods into distinct technologies, including synthetic minority oversampling technology (SMOTE), support vector machine (SVM), k-nearest neighbor (KNN), among others, and compare the advantages and disadvantages of these approaches. The paper concludes by summarizing evaluation criteria for unbalanced dataset classifiers and offering insights into future research directions in this domain. This review contributes valuable insights for addressing the challenges posed by unbalanced datasets in the classification context.

In the paper [6] authored by S. Ghannay, B. Favre, Y. Esteve, and N. Camelin, addresses the application of word embeddings in various NLP and speech processing tasks. While word embeddings derived from neural networks have shown success, the paper identifies gaps in the existing literature regarding their evaluation. The study rigorously compares the performance of different word embeddings on diverse NLP and linguistic tasks, ensuring consistency in training data, vocabulary, dimensions, and other characteristics. The evaluation results align with previous literature, emphasizing that improvements observed in one task may not consistently translate to others. Consequently, the paper investigates and evaluates methods to combine word embeddings to leverage their complementarity, seeking effective embeddings that perform well across various tasks. In conclusion, the paper contributes new insights into the intrinsic qualities of popular word embedding families, offering perspectives that may differ from those previously published in scientific literature

In the paper [7] authored by B. A. Johnsson and B. Magnusson, the authors address the challenges of graphical user interface (GUI) development, which is generally complicated, time-consuming, and requires programming knowledge. In the context of the IOT (IoT), this work focuses on producing an efficient development approach that also supports non-experts. The authors introduce a novel “inverted” development approach that does not require program code to be written – a step towards supporting end-user development in the given context. The approach is realized as a language for describing GUIs, interpreters for rendering GUIs, and a graphical tool for creating and editing GUIs. The work is evaluated in a number of research projects in the domain of e-health, concluding that the GUI language is practically viable for building professional-grade GUIs. Furthermore, the presented editor is compared directly to a market-leading product in a controlled experiment. From this, the authors conclude that the editor is accessible to new users and can be more efficient to use than the commercial alternative.

In the paper [8] authored by J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, the authors introduce a new language representation model called BERT, which stands for Bidirectional Encoder Representations from Transformers. Unlike recent language representation models, BERT is designed to pre-train deep bidirectional representations from unlabeled text by jointly conditioning on both left and right context in all layers. As a result, the pre-trained BERT model can be fine-tuned with just one additional output layer to create state-of-the-art models for a wide range of tasks, such as question answering and language inference, without substantial task-specific architecture modifications. BERT is conceptually simple and empirically powerful. It obtains new state-of-the-art results on eleven NLP tasks, including pushing the GLUE score to 80.5% (7.7% point absolute improvement), MultiNLI accuracy to 86.7% (4.6% absolute improvement), SQuAD v1.1 question answering Test F1 to 93.2 (1.5 point absolute improvement) and SQuAD v2.0 Test F1 to 83.1 (5.1 point absolute improvement).

In the paper [9] authored by R. Wang, R. Ridley, W. Qu, X. Dai, and X. Su, the authors address the challenge of assigning a set of labels to a given document in multi-label text classification. Previous classifier-chain and sequence-



to-sequence models have demonstrated powerful abilities to capture label correlations; however, these models rely heavily on the label order, while labels in multi-label data are essentially an unordered set. The performance of these approaches is therefore highly variable depending on the order in which the labels are arranged. To avoid dependency on label order, the authors design a reasoning-based algorithm named Multi-Label Reasoner (ML-Reasoner) for multi-label classification. ML-Reasoner employs a binary classifier to predict all labels simultaneously and applies a novel iterative reasoning mechanism to effectively utilize the inter-label information, where each instance of reasoning takes the previously predicted likelihoods for all labels as additional input. This approach leverages the information between labels while avoiding the issue of label-order sensitivity. Extensive experiments demonstrate that ML-Reasoner outperforms state-of-the-art approaches on the challenging AAPD dataset. The authors also apply their reasoning module to a variety of strong neural-based base models and show that it is able to significantly boost performance in each case.

In the paper [10] authored by B. Breve, G. Desolda, V. Deufemia, F. Greco, and M. Matera, the authors address the need for laypeople to configure their IoT devices securely. Task Automation Systems (TASs) simplify the personalization of device behavior but often overlook security and privacy threats inherent to connected devices. The paper introduces a user-centered design approach that resulted in a visual paradigm, enabling end-users to understand and control security and privacy threats. This approach abstracts complex security configurations into manageable visual elements, allowing non-expert users to define and manage security rules effectively. The system alerts users to potential security issues and provides actionable insights for mitigation. Evaluations through user studies demonstrated that users could successfully manage security configurations using the proposed visual tools, significantly reducing the risk of security and privacy breaches in smart environments. This innovative approach makes security configuration accessible to a broader audience, enhancing the usability and security of IoT devices, and represents a significant contribution to IoT security and end-user development.

In the paper [11] authored by M. Saeidi, M. Calvert, A. W. Au, A. Sarma, and R. B. Bobba, the authors investigate the usage of trigger-action platforms like If-This-Then-That for creating applets to connect smart home devices and services. Despite their convenience, there are inherent risks associated with using such applets, including the potential leakage of sensitive information in certain contexts (such as the device's location or who can observe the resulting action), even in non-malicious scenarios. The study aims to understand end users' ability to assess these risks by exploring their concerns with using IFTTT applets and how these concerns vary based on contextual factors. Through a Mechanical Turk survey involving 386 participants and 49 smart-home IFTTT applets, the research reveals that prompting participants to consider different usage contexts prompts deeper reflection on associated risks and raises their concerns. Qualitative analysis demonstrates participants' nuanced understanding of contextual factors and their awareness of how these factors could lead to data leakage and unauthorized access to applets and data.

In the paper [12] authored by M. McCall et al., the authors address the intricacies of home automation rules established through trigger-action programs for connecting smart home devices. They highlight the challenge of ensuring these rules are free of undesirable behavior and advocate for tools to assist users in analyzing and verifying their rules. To understand users' needs for such tools, they conducted a user study where participants were provided with their custom analysis tool, SAFETAP, or not. The study revealed users' limited ability to identify issues in their TAP rules, despite perceiving the task as easy. Users expressed a desire to check their rules whenever they made changes. Consequently, the authors developed SAFETAP $\Delta$ , a novel incremental symbolic model checking (SMC) algorithm that extends the basic SMC algorithm of SAFETAP. SAFETAP $\Delta$  focuses solely on analyzing the addition or removal of rules and reports new violations. Performance evaluations demonstrate that incremental checking improves performance by 6X on average when adding new rules.

In the paper [13] authored by Wenshu Zha et al., the authors address the challenge of accurately predicting gas field production, crucial for reservoir engineers but hindered by numerous unknown reservoir parameters. They propose a CNN-LSTM model to provide a low-cost, intelligent, and robust prediction method for gas and water production in gas reservoirs. This model combines the feature extraction ability of convolutional neural networks (CNN) with the sequence dependence learning capability of long short-term memory networks (LSTM) to effectively capture the

changing trend of gas field production. Additionally, they introduce a new prediction strategy named partly unknown recursive prediction strategy (PURPS), estimating some input features using predicted gas and water production based on known equations. The study demonstrates that the CNN-LSTM model outperforms existing methods, with an average mean absolute percentage error (MAPE) for monthly gas production of 7.7%, surpassing models such as RNN, Random Forest, ARIMA, DNN, Support Vector Machine, CNN, and LSTM. This research contributes to advancing gas field production forecasting techniques, offering promising results for reservoir.

In the paper [14] authored by Mihai Masala, Stefan Ruseti, and Mihai Dascalu, the authors address the increasing ubiquity of deep pre-trained language models in NLP. These models, trained on vast amounts of unlabeled text data, learn contextualized representations and achieve state-of-the-art results across various NLP tasks through efficient transfer learning. However, for languages other than English, options for such models are limited, with most being trained solely on multilingual corpora. The authors introduce RoBERT, a Romanian-only pre-trained BERT model, and conduct a comparative analysis with different multilingual models on seven Romanian-specific NLP tasks categorized into sentiment analysis, dialect and cross-dialect topic identification, and diacritics restoration. Their findings demonstrate that RoBERT outperforms both multilingual models and another monolingual BERT implementation across all tasks. This highlights the effectiveness of language-specific pre-trained models like RoBERT in capturing the nuances and intricacies of Romanian language processing tasks, underscoring the importance of tailored solutions for diverse linguistic contexts in NLP.

In the paper [15] authored by García-Grao G and Carrera Á, the authors delve into the realm of workflow automation within software development systems. They highlight how automation, coupled with agile principles, has given rise to the DevOps paradigm, offering businesses enhanced efficiency, accelerated production, and adaptability to market changes. However, practitioners often encounter vendor lock-in, hindering their ability to switch tools or migrate to different platforms due to associated costs. To address this challenge, the authors propose standardizing service interfaces to facilitate integration, leveraging Linked Data for its versatility. The article focuses on extending the Open Services for Lifecycle Collaboration (OSLC) standard, specifically targeting event-based interoperable automation using the ECA model. This extension paves the way for semantic automation within OSLC services, enabling services to interact autonomously and with human users. The paper elucidates the fundamental concepts of the proposed model and provides real-life examples of its application in automation scenarios involving two services. Validation of the proposal is conducted using established ontology evaluation methods, including coverage and similarity metrics, as well as competency questions. This research contributes to advancing interoperability and automation standards, offering potential solutions to mitigate vendor lock-in concerns in software development systems.

In the paper [16] authored by Smagulova, K. and James, A.P., the authors delve into the realm of recurrent neural networks (RNN), focusing specifically on Long Short-Term Memory (LSTM) networks. RNNs are noted for their effectiveness in approximating dynamic systems dealing with time and order-dependent data such as video and audio. The paper explores the motivations behind the development of LSTM networks and provides a tutorial survey on existing LSTM methods. It elucidates the principles and functionalities of LSTM networks, highlighting their state memory and multilayer cell structure. Additionally, the authors discuss recent advancements in memristive LSTM architectures, emphasizing the use of memristor circuits for hardware acceleration. By examining the evolution of LSTM networks and discussing the latest developments in memristive LSTM architectures, the paper contributes to the understanding of neural network models and their hardware implementations. It serves as a valuable resource for researchers and practitioners interested in leveraging LSTM networks for various applications in dynamic systems approximation.

In the paper [17] authors delve into the realm of emotion recognition through facial expressions. The paper discusses how people use facial expressions as a form of communication to express their emotions, prompting researchers to delve into the field of emotion recognition. Understanding human emotions through facial expressions is crucial for interpreting individuals' current mood states. Therefore, the design of deep learning models plays a significant role in capturing facial gestures' patterns to interpret human emotions effectively. The authors propose a customized Convolutional Neural Network (CNN) tailored for emotion recognition, utilizing the Adaptive Moment Estimation (Adam) and Nesterov-accelerated Adaptive Moment Estimation (Nadam) optimizers. They vary parameters such as

the number of convolution layers, filters, filter sizes, and optimizers to optimize emotion recognition using the FER-2013 dataset. Emotions are classified using softmax activation in the output layer. Experimental results indicate that the proposed model achieves high accuracy, with 0.841 and 0.826 accuracy scores using Nadam and Adam optimizers, respectively. This research contributes to advancing emotion recognition technologies, offering promising results in accurately interpreting human emotions through facial expressions.

In the paper [18] authored by Yi, Dokkyun, Jaehyun Ahn, and Sangmin Ji, the authors tackle the challenge of optimizing machine learning models efficiently, particularly in the presence of non-convex cost functions. As the complexity of artificial intelligence structures and the volume of learning data increase, conventional gradient descent optimization methods encounter limitations due to the existence of local minima in the cost function landscape. To address this issue, the authors propose a novel optimization method aimed at enhancing the efficiency of machine learning. They introduce a modification to the parameter update rule of the ADAM optimization algorithm, incorporating additional terms in the cost function to prevent the model from converging to suboptimal solutions. The paper provides theoretical proofs of the convergence of the sequences generated by the proposed method and conducts numerical comparisons with traditional gradient descent methods such as GD, ADAM, and AdaMax. Through these comparisons, the authors demonstrate the superiority of their proposed optimization method in terms of efficiency and effectiveness for non-convex cost functions. This research contributes to advancing optimization techniques in machine learning, offering promising solutions for improving model performance and convergence.

In the paper [19] authored by Turki M. Alanazi tackles the challenge of text detection and recognition from natural scene images using an embedded system approach. The paper proposes a prototype system based on the Raspberry Pi 4 and a USB camera for text detection and recognition. The system integrates a text detection and recognition model developed in Python, leveraging the Efficient and Accurate Scene Text Detector (EAST) model for text localization and detection, and the Tesseract-OCR engine for text recognition. Controlled wirelessly via the Virtual Network Computing (VNC) tool, the prototype demonstrates promising results. Experimental findings indicate a high recognition rate of 99.75% for captured images, achieved with low computational complexity. Additionally, the prototype outperforms the Tesseract software in terms of recognition rate and exhibits similar recognition performance to EasyOCR software on the Raspberry Pi 4, with a significant reduction in execution time averaging 89%. This research contributes to advancing text detection and recognition technologies using embedded systems, offering practical implications for applications requiring efficient and accurate text processing from natural scene images.

In the paper [20] authored by Amin Biglari and Wei Tang, the authors provide an overview of embedded machine learning implementations across various hardware platforms, applications, and sensing schemes. The paper categorizes embedded machine learning implementations based on specific hardware devices such as NVIDIA Jetson and Raspberry Pi, along with less utilized embedded computers. It explores the applications of these devices in different fields and analyzes commonly implemented machine learning models and sensors used for input gathering. The authors conducted a comprehensive review of relevant literature. The selection criteria included the accuracy, power consumption, and inference time of the embedded machine learning systems. The review highlights the growing significance of embedded machine learning, driven by advancements in system performance, machine learning models, and increased affordability and accessibility. The authors note improvements in quality, power usage, and effectiveness of embedded machine learning systems, signaling an expansion in both scale and scope of applications.

### **MODEL IMPLEMENTATION**

The work is currently underway with the goal of developing a robust solution for identifying security and privacy risks associated with IFTTT applets. This endeavor presents a notable challenge, given the need to classify these applets into distinct risk categories. The complexity arises from their diverse nature, an imbalanced dataset, and the critical requirement to comprehend their contextual semantics.

The proposed solution is a multi-faceted approach that integrates natural language processing (NLP) and machine

learning techniques to advance the accuracy of IFTTT applet classification. The plan begins by curating an extensive dataset that encompasses a wide array of applets. These applets are categorized into four risk classes based on their potential for causing harm. Subsequently, data preprocessing is conducted, where features such as trigger, action, channel information, and applet descriptions are extracted. To address the challenge of imbalanced data, a weighted loss function is devised.

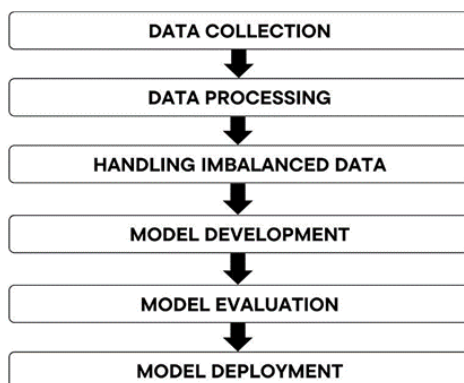


Fig. 3 Block Diagram of proposed work Implementation

In response to the imperative need for robust threat detection and classification systems, the team has successfully developed a cutting-edge BERT (Bidirectional Encoder Representations from Transformers) model. This model is designed to excel in the intricate task of multiclass classification across four distinct categories: innocuous, personal harm, physical harm, and cybersecurity harm. Leveraging the power of pre-trained contextual embeddings, BERT ensures a nuanced understanding of the semantic context within the textual data, enabling accurate and context-aware predictions. The training process primarily focuses on comprehending the crucial role of context in accurately classifying applets by employing established metrics such as accuracy, precision, recall, and F1 score. The proposed model is meticulously fine-tuned on a diverse and comprehensive dataset, encompassing a wide array of scenarios to ensure its adaptability and effectiveness in real-world applications.

### **Data Collection**

Data collection is crucial in deep learning, especially for harm classification using a BERT model, because the performance and accuracy of the model heavily depend on the quality, quantity, and diversity of the data. High-quality data ensures accurate labeling, minimizing noise and ambiguity, which directly impacts the model's ability to learn and generalize. A substantial and diverse dataset exposes the model to various contexts and types of harmful content, enhancing its robustness and adaptability to real-world scenarios. This diversity helps the model understand nuanced language and cultural variations, improving its overall effectiveness in identifying harmful content. The dataset used for training and testing the models for classifying harmful IFTTT applets was created by researchers from Indiana University Bloomington through a crawling process on the IFTTT.com website from November 2016 to May 2017. During this period, the authors took weekly “snapshots” of the applets published on the platform at that moment. These snapshots ensured that the dataset captured the dynamic nature of the applets over time, providing a comprehensive view of the types of applets available during the collection period.

The researchers maintained the same structure of the IFTTT paradigm, which involves the decomposition of the applet into trigger, action, and the corresponding channels. This structured approach allowed the researchers to systematically analyze the components of each applet, ensuring that the dataset accurately reflected the operational mechanics of IFTTT applets. By preserving this structure, the dataset offers detailed insights into how different triggers and actions are used across various channels, which is essential for training models to classify harmful content effectively.



The dataset for classifying harmful IFTTT applets is organized into several key headers, each representing crucial components of an applet. These headers provide structured information that is essential for understanding and analyzing the behavior and functionality of each applet. Here is a brief explanation of each header:

- **Title:** This is a string containing the title of the applet. It provides a concise summary or name of the applet, often giving an indication of its purpose or functionality.
- **Desc:** Short for description, this string details the applet's behavior. It explains what the applet does, outlining the interaction between the trigger and the action in a more descriptive manner.
- **TriggerTitle:** This string represents the name of the trigger that activates the applet. It specifies the condition or event that initiates the applet's action.
- **TriggerChannelTitle:** This string represents the name of the trigger channel chosen by the user. The trigger channel is the service or platform where the triggering event occurs, such as a social media platform, a weather service, or a device.
- **ActionTitle:** This string represents the name of the action that is executed in response to the trigger. It specifies what happens once the applet is activated.
- **ActionChannelTitle:** This string represents the name of the action channel chosen by the user. The action channel is the service or platform where the action takes place, such as sending a message, posting on social media, or controlling a smart device.
- **Target:** This column contains integer values (0, 1, 2, 3) representing the classes of harm. Each value corresponds to a specific level or type of harm, with 0 typically representing non-harmful applets and higher values indicating varying degrees or types of harmfulness.

These headers form the columns of the CSV file, structuring the dataset in a way that each row corresponds to a unique IFTTT applet with its associated trigger and action details. This structured approach facilitates the analysis and classification process, enabling the identification of potentially harmful applets based on their described behavior and interactions.

First few rows:

|   | triggerTitle          | triggerChannelTitle | actionChannelTitle      | \ |
|---|-----------------------|---------------------|-------------------------|---|
| 0 | New Popular photo     | 500px               | Update device wallpaper |   |
| 1 | New Popular photo     | 500px               | Update device wallpaper |   |
| 2 | New Popular photo     | 500px               | Update device wallpaper |   |
| 3 | New Popular photo     | 500px               | Update device wallpaper |   |
| 4 | New photo from search | 500px               | Update device wallpaper |   |

|   | actionTitle    | title   | \ |
|---|----------------|---|---|
| 0 | Android Device | If there is a new photo in 500px ,then but it ... |   |
| 1 | Android Device | Any 500px Popular photo on Wallpaper              |   |
| 2 | Android Device | Popular Images from 500px on Device               |   |
| 3 | Android Device | any popular photo on 500px on my mobile           |   |
| 4 | Android Device | New york  |   |

|   | desc  | target |
|---|---|--------|
| 0 | You can change the category                       | 3      |
| 1 | Turn notifications off                            | 3      |
| 2 | Just click and enable to have the 500 Px popul... | 3      |
| 3 | allows to have series of wonderful pictures as... | 3      |
| 4 | If new photo by anyone matching search new yor... | 3      |

Fig. 4 First few rows of the collected dataset

## Data Analysis and Preprocessing

Basic data analysis plays a pivotal role in NLP models, serving as a foundational step that significantly impacts the model's performance and outcomes. Understanding and scrutinizing the under lying data is crucial for identifying patterns, trends, and potential biases, which can profoundly influence the model's accuracy and generalization capabilities. By examining the distribution of data points across various categories, researchers can detect imbalances and address them appropriately, ensuring that the model does not favor any particular class. Additionally, data



analysis helps in assessing the quality of the dataset, identifying issues such as noise, missing values, or inconsistencies that need to be resolved before training the model. This step also provides insights into relevant features that should be included, guiding feature selection and engineering to enhance the model's effectiveness. Ultimately, thorough data analysis lays the groundwork for building robust and reliable NLP models that can perform well on diverse and real-world data. Various data analysis techniques for NLP have been employed in order to understand the obtained data set properly.

Desc: Describes the basic structure of the dataset considered. Around 79,000+ non null object in each column of the considered dataset.

Each class in the target column represents different forms of harm in the considered ECA rule:

- Innocuous (Class 0): Signifies actions or events that are harmless or benign, posing no threat or risk to individuals or systems.
- al Harm (Class 1): Indicates actions or events that may cause harm to individuals' privacy, reputation, or personal well-being, such as cyberbullying or harassment.
- Physical Harm (Class 2): Represents actions or events that have the potential to cause physical harm or damage to individuals or property, including threats of violence or unsafe activities.
- Cybersecurity Harm (Class 3): Denotes actions or events that pose risks to cybersecurity, such as data breaches, malware distribution, or exploitation of vulnerabilities in systems or networks.

In the context of input approach, which merges all columns except the target column, the presence of null values in the 'Desc' and 'Title' columns does not exert a substantial impact on the model's performance. While there are null values, accounting for just less than 100 instances in these columns, the merged input data still retains valuable information from other non-null columns, such as 'TriggerTitle', 'TriggerChannelTitle', 'ActionTitle', and 'ActionChannelTitle'. These columns provide substantial context for the model to analyze and classify harmful applets effectively. Given the relatively small number of null values and the comprehensive nature of the merged input data, the absence of information in 'Desc' and 'Title' does not significantly hinder the model's ability to learn patterns and make accurate predictions. Therefore, while null value handling is essential for dataset completeness, its impact on the model's performance within this specific input approach remains minimal.

Word length distribution can be utilized in NLP tasks to understand the complexity of text, identify language patterns, and inform feature engineering or model selection decisions.

## Word Cloud Representaion

Word cloud analysis holds significant importance in the realm of NLP models, providing a visually intuitive representation of the most frequent words in a corpus. This analysis serves as a powerful exploratory tool for gaining quick insights into the key themes, prevalent terms, and overall semantic structure of a text dataset. By visualizing word frequency distributions, word clouds enable researchers to identify prominent topics, discern patterns, and uncover underlying sentiment tendencies within the data. Additionally, they facilitate communication and interpretation of NLP results to stakeholders who may not have expertise in data analysis or machine learning, thereby fostering collaboration and informed decision-making.

## Data Preprocessing

Data processing is a critical step for the success of a BERT-based transfer learning model, as it directly influences the model's performance and accuracy. Effective data processing ensures that the input text is clean, consistent, and in a format suitable for analysis. This involves several steps, such as removing special characters to eliminate noise, converting text to lowercase to standardize it, and tokenizing the text to break it down into manageable units. Removing stopwords helps focus on the most meaningful words, while lemmatization ensures consistency by reducing words to their base forms. These preprocessing steps not only enhance the clarity and quality of the text

data but also help in reducing the complexity of the vocabulary, allowing the BERT model to better understand and learn from the data. Properly processed data enables the BERT model to accurately capture the underlying patterns and semantics, leading to more reliable and generalizable results in various NLP tasks. Consequently, thorough and precise data processing is indispensable for leveraging the full potential of BERT-based transfer learning models in NLP applications.

The following data processing methods are used to prepare input for a model utilizing transfer learning on a BERT model. Each step plays a crucial role in ensuring the text data is clean, standardized, and suitable for effective training and prediction:

- **Remove Special Characters:** This step uses regular expressions to replace non-alphanumeric characters with spaces. Removing special characters such as punctuation marks, symbols, and other non-standard characters helps clean the text by eliminating noise that can interfere with the model's understanding of the text. It ensures that only relevant textual data is retained.
- **Remove Single Characters:** Isolated single characters that appear after spaces are removed. Single characters often do not contribute meaningful information and can be artifacts of typing errors or shorthand notation. Removing them helps reduce noise and improves the clarity of the text data.
- **Remove Single Characters from the Start:** Single characters at the beginning of the text are removed to prevent any leading characters from affecting the analysis. Leading single characters might include bullet points or formatting markers that do not add value to the text content.
- **Substitute Multiple Spaces with Single Space:** Multiple consecutive spaces are replaced with a single space using regular expressions. This step ensures uniform spacing within the text, making it cleaner and easier to process. It helps in maintaining consistency and prevents the model from interpreting multiple spaces as different tokens.
- **Convert to Lowercase:** Converting the entire text to lowercase standardizes the text, reducing variability. For example, "Apple" and "apple" would be treated the same, which is crucial for accurate tokenization and subsequent analysis. This step helps in minimizing the number of unique tokens, making the model's vocabulary more efficient. Since a BERT-based uncased model is used, this step aligns perfectly with the model's design, as the uncased BERT model is already designed to ignore case differences. Therefore, converting text to lowercase ensures consistency with the model's expectations and contributes to its effective performance.
- **Tokenization:** The text is split into a list of words, known as tokens. Tokenization is crucial for breaking down the text into manageable units that the model can process. It allows the model to understand and analyze the text at the word level, which is essential for tasks like sentiment analysis and classification.
- **Remove Stopwords:** Common English stopwords, such as "and," "the," and "is," are removed from the text. Stopwords are frequently occurring words that typically do not carry significant meaning on their own. Removing them helps in focusing the analysis on more meaningful words that contribute to the overall context.
- **Lemmatization:** Lemmatization reduces words to their base or root form using a tool like the WordNet lemmatizer. For example, "running" is converted to "run." This step ensures consistency in word forms, which helps in reducing the complexity of the vocabulary and improving the model's ability to generalize from the data.
- **Join Tokens:** The processed tokens are joined back into a single string. This final step reconstructs the cleaned and processed text, ensuring it is in a suitable format for input into the BERT model. It helps in maintaining the integrity of the text while ensuring that it is standardized for the model's consumption.

Each of these preprocessing steps is essential for preparing high-quality input data, which directly impacts the performance and accuracy of the BERT model in various NLP tasks.

## Handling Imbalanced dataset

An imbalanced target column can significantly affect the final results of a BERT-based transfer learning model. When the target classes are not evenly distributed, the model tends to be biased towards the majority class, which can lead to poor performance in predicting minority classes. This imbalance means that the model might learn to prioritize

accuracy for the most frequent class at the expense of others, resulting in high overall accuracy but low precision, recall, and F1-score for the less frequent classes. Consequently, important patterns and signals from the minority classes may be overlooked or underrepresented during training, leading to a model that does not generalize well across all classes.

In practical terms, this can have serious implications, especially in applications where correctly identifying minority class instances is crucial, such as in detecting harmful content or cybersecurity threats. For example, if the model is used to classify types of harm in IFTTT applets and the minority class represents a critical harm type like cybersecurity threats, an imbalanced dataset could result in the model failing to detect these threats reliably. To mitigate the effects of class imbalance, techniques such as resampling (oversampling the minority class or under sampling the majority class), using class weights to give more importance to minority classes, or employing advanced methods like Synthetic Minority Over-sampling Technique (SMOTE) can be implemented. These approaches help in creating a more balanced training process, allowing the model to learn effectively from all classes and thereby improving its overall performance and robustness.

To handle the imbalanced distribution in the target column, Weighted Categorical Cross-Entropy and Stratified K-Fold Cross-Validation is employed. These methods are essential for ensuring that the model performs well across all classes, particularly when the data distribution is uneven.

### **Stratified K-Fold Cross-Validation**

Stratified k-fold cross-validation is a well-established technique tailored for scenarios where class imbalances pose significant challenges, especially in classification tasks. Unlike standard k-fold cross-validation, where the dataset is randomly partitioned into k folds, stratified k-fold ensures that each fold maintains the same class distribution as the original dataset. This preservation of class proportions across folds is crucial for obtaining reliable performance estimates, particularly when certain classes are underrepresented. By mitigating the risk of skewed class distributions in training and evaluation sets, stratified k-fold cross-validation provides a robust assessment of model performance across all classes.

### **Weighted Categorical Cross-Entropy**

In conjunction with stratified k-fold cross-validation, weighted categorical cross-entropy is employed as the loss function to address the imbalanced distribution in the target column. Traditional categorical cross-entropy treats each class with equal importance during model training, which might lead to suboptimal results when dealing with imbalanced datasets. Weighted categorical cross-entropy introduces a weighting factor, denoted as  $W(x)$ , that dynamically adjusts the contribution of each class to the overall loss based on its frequency in the dataset.

$W(x)$  Calculation:

The weighting factor,  $W(x)$ , is computed using the formula:

$$W(x) = \text{samples} / (\text{classes} \times \text{samples}(x))$$

Where:

- samples represent the total number of samples in the dataset.
- classes denote the total number of classes in the target variable.
- $\text{samples}(x)$  signifies the number of samples belonging to class  $x$ .

By incorporating this weighting factor into the categorical cross-entropy loss function, the contribution of minority classes is effectively prioritized, ensuring that the model is adequately trained to address the imbalanced nature of the dataset.

In summary, the combination of stratified k-fold cross-validation and weighted categorical cross-entropy offers a robust framework for evaluating Bert-based uncased transfer learning models on imbalanced datasets. By preserving class distributions across folds and dynamically adjusting loss contributions based on class frequencies, the approach provides a more accurate and representative assessment of model performance. This comprehensive strategy not only

enhances the reliability of performance metrics but also ensures the effectiveness of the trained models in real-world scenarios characterized by class imbalances.

## Model Development

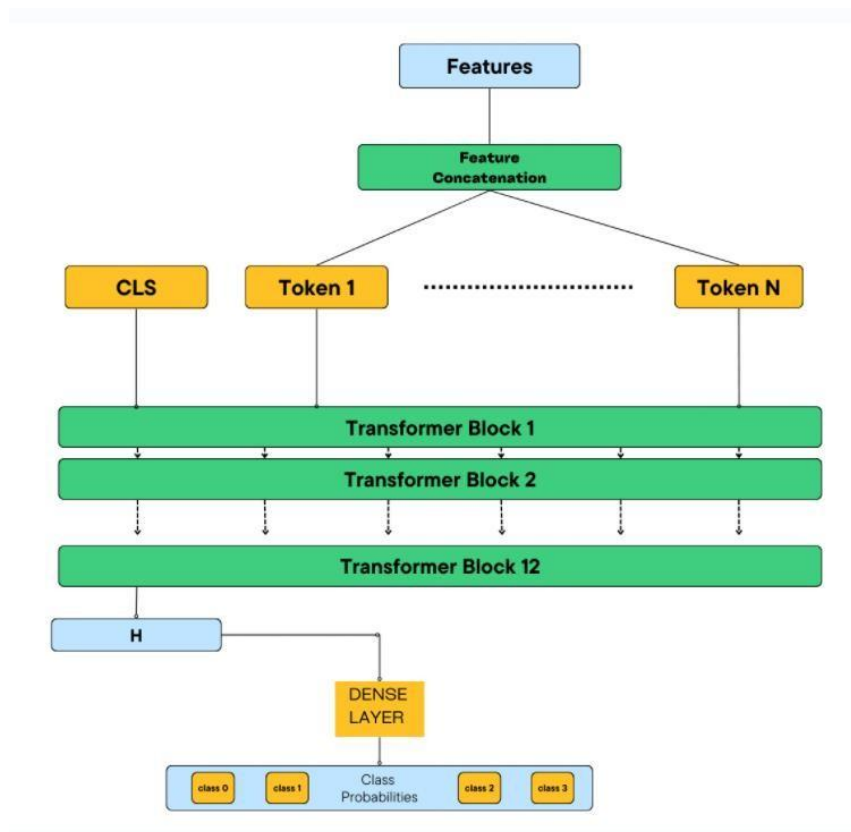


Fig. 5 Model development using BERT based transfer learning

In the landscape of NLP, the emergence of transformer-based models has revolutionized the field, offering unprecedented capabilities in understanding and processing language. Among these models, BERT stands as a beacon of innovation, transforming how to approach NLP tasks.

BERT owes its prowess to its foundation: the transformer architecture. Introduced in the seminal paper "Attention is All You Need" by Vaswani et al., transformers represent a departure from traditional recurrent neural networks (RNNs) by leveraging self-attention mechanisms for sequence processing.

Central to the transformer architecture is the self-attention mechanism, which allows models to capture contextual dependencies within input sequences. Unlike RNNs, which process sequences sequentially, BERT can consider all tokens simultaneously, discerning nuanced relationships between words irrespective of their positions. This bidirectional context modeling is pivotal in understanding the intricacies of natural language.

Transformers consist of encoder and decoder stacks, with BERT exclusively employing the encoder component. The encoder stack encodes input sequences into high-dimensional representations, which can be fine-tuned for various downstream NLP tasks. By leveraging self-attention mechanisms across multiple layers, BERT generates rich contextual embeddings that encapsulate semantic information crucial for tasks such as text classification, named entity recognition, and sentiment analysis. In essence, BERT's transformative power lies in its adeptness at capturing bidirectional context and leveraging it to generate contextually rich embeddings. This capability has propelled BERT to the forefront of NLP research and applications, enabling advancements in diverse domains such as question answering, language translation, and semantic understanding.

To effectively divide the imbalanced dataset of 79,214 ECA rules into training and test datasets while maintaining the imbalance ratio, there is a need to ensure that both subsets reflect the same class distribution. The test dataset should comprise 15,843 samples, and the remaining 63,371 samples will form the training dataset. stratified sampling is achieved by employing this, which ensures that each class is proportionately represented in both the training and test sets. This method is crucial for preserving the original distribution and preventing potential bias that could arise from an uneven class representation. By maintaining the class proportions, the model training and evaluation remain consistent and reflective of real-world scenarios is ensured, leading to more reliable and generalizable performance metrics. This stratified approach is particularly important in imbalanced datasets, where minority classes might otherwise be underrepresented in the training or test sets.

### Leveraging Transfer Learning with BERT for Classification Tasks

In the realm of NLP, the utilization of pre-trained language models has emerged as a cornerstone strategy for achieving state-of-the-art performance across various tasks. Among these models, BERT has garnered widespread acclaim for its exceptional ability to capture contextual information in text.

In the pursuit of harnessing the power of BERT for classification tasks, transfer learning approach is applied to capitalize on the wealth of knowledge encapsulated in pre-trained BERT models. The training process entails freezing the parameters of the pre-trained BERT layers while introducing an additional classification layer atop the architecture. This strategy ensures that during training, only the newly added classification layer is fine-tuned on the dataset, thereby maximizing efficiency and minimizing computational overhead.

Implementation-wise, TRANSFORMERS Python library is leveraged, specifically the Bert for Sequence Classification class, tailored for classification tasks. This class seamlessly integrates with the BERT model, incorporating a single linear layer designed for classification purposes.

The model configuration centers around the "bert-base-uncased" variant, chosen for its robustness and efficiency. This variant comprises 12 transformer blocks, 768 hidden units, and 12 self-attention heads, with a focus on processing lowercase letters for enhanced generalization.

In the input sequence, tokens are arranged, with a distinctive classification token ([CLS]) positioned at the outset. These token embeddings traverse through successive transformer blocks, each employing self-attention mechanisms and forwarding the output to a feedforward neural network. The representation associated with the [CLS] token serves as the input to the Dense Layer at the pinnacle of the architecture, orchestrating the final classification based on the processed token sequence.

During training, specific hyperparameters govern the optimization process. A batch size of 32, a learning rate of  $2e-5$ , and 2 epochs are employed, with an epsilon value of  $1e-8$  ensuring stability in the optimization process. The maximum sentence length is dynamically adjusted based on the features under consideration, with the initial feature combination denoted as BERT-1c, setting the maximum sentence length at 50. This configuration accommodates discrete features while allowing comprehensive coverage of textual content, crucial for accurate classification. Softmax as an Activation Function

$$s(x_i) = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$$

where:

- $s(x_i)$  is the softmax function applied to the  $i$ -th input value  $x_i$ .
- $e$  is the base of the natural logarithm (approximately equal to 2.71828).
- $x_i$  is the  $i$ -th element of the input vector  $x$ .
- $\sum_{j=1}^n e^{x_j}$  is the sum of the exponentials of all the elements in the input vector  $x$ , where  $n$  is



the number of elements (or classes).

In the realm of multiclass classification, the choice of activation function plays a pivotal role in shaping the output layer of the neural network model. Softmax activation function stands out as a fundamental component due to its ability to normalize the output into a probability distribution across multiple classes. Mathematically, softmax function transforms the raw output scores into probabilities by exponentiating each score and dividing by the sum of all exponentiated scores. This normalization ensures that the output values lie between 0 and 1, with the sum of probabilities across all classes equating to 1.

The significance of softmax in multiclass classification lies in its capacity to provide interpretable and actionable predictions. By generating probabilities for each class, softmax facilitates not only the identification of the most likely class but also quantifies the model's uncertainty in its prediction. Furthermore, softmax enables the model to output class probabilities, which are inherently more informative than raw scores, especially in scenarios where decision-making relies on confidence levels. Thus, by employing softmax activation function in the dense layers of the neural network model, by ensuring that the model's outputs are calibrated probabilities, conducive to informed decision-making in multiclass classification tasks.

## ADAM as an optimizer

In the realm of deep learning, optimizers serve as the driving force behind the training process, guiding the model parameters towards optimal solutions. Their role is particularly crucial in complex architectures like BERT, where the optimization landscape is intricate and high-dimensional. Optimizers are tasked with adjusting the model's parameters iteratively based on the gradients computed during backpropagation, aiming to minimize the loss function and improve model performance. The choice of optimizer significantly influences the efficiency, stability, and convergence speed of the training process, making it a critical decision in model development.

Mathematically:

$$w_{t+1} = w_t - \alpha m_t$$

Where,

$$m_t = \beta m_{t-1} + (1 - \beta) \left[ \frac{\partial L}{\partial w_t} \right]$$

- $m_t$  = Aggregate of gradients at time  $t$  [Current] (Initially,  $m_t = 0$ )
- $m_{t-1}$  = Aggregate of gradients at time  $t-1$  [Previous]
- $W_t$  = Weights at time  $t$
- $W_{t+1}$  = Weights at time  $t+1$
- $\alpha$  = Learning rate at time  $t$
- $\partial L$  = Derivative of Loss Function
- $\partial W_t$  = Derivative of weights at time  $t$
- $\beta$  = Moving average parameter (Constant, 0.9)

Among the plethora of optimization algorithms, Adam stands out as a popular choice, especially in transformer-based architectures like BERT. Adam's adaptive learning rate mechanism, which combines the benefits of both momentum and RMSprop techniques, proves to be highly advantageous in navigating the complex optimization landscape of BERT models. By dynamically adjusting the learning rates for individual parameters based on their past gradients and squared gradients, Adam effectively adapts to varying gradient magnitudes and optimizes the model's convergence trajectory. In the context of BERT, where attention mechanisms and transformer blocks introduce nonlinearities and gradient fluctuations, Adam's adaptive nature ensures stable and efficient optimization, ultimately leading to enhanced model performance.

In the development of a BERT-based transfer learning model for classifying multiclass harm in ECA rules, the choice of the Adam optimizer aligns strategically with the objectives. Given the complexity of BERT architectures and the

intricacies of language patterns inherent in ECA rules, employing Adam as the optimizer offers several key advantages. Its adaptive learning rate mechanism enables efficient exploration of the optimization landscape, mitigating issues like vanishing or exploding gradients commonly encountered in deep learning tasks. Moreover, Adam's ability to handle varying gradient magnitudes ensures stable training dynamics, facilitating the convergence of the model towards optimal solutions. Overall, the utilization of Adam as the optimizer underscores the commitment to optimizing the training process and maximizing the proficiency of the BERT-based transfer learning model in capturing nuanced language patterns relevant to multiclass harm classification in ECA rules.

## Exploring Model Magic: Building Brilliance with Google Colab

Google Colab offers a fantastic resource for college students embarking on machine learning projects, providing them with a powerful platform to build and train models effectively. One of the standout features of Google Colab is its availability of inbuilt GPUs, which can significantly accelerate the training process for deep learning models. This access to GPU resources allows students to experiment with complex architectures and large datasets without the need for expensive hardware investments. Furthermore, Google Colab provides a collaborative environment where students can easily share and collaborate on their projects, fostering knowledge exchange and collective learning. With its intuitive interface, seamless integration with popular libraries like TensorFlow and PyTorch, and access to GPU resources, Google Colab empowers college students to unleash their creativity and build sophisticated machine learning models, regardless of their hardware limitations.

## Library Luminaries: Fueling Model Mastery with Essential Tools

- NumPy: NumPy is pivotal for data manipulation tasks, aiding in the preprocessing of input data and facilitating array operations crucial for feeding data into the Bert- based model efficiently.
- Pandas: Pandas streamlines data handling, allowing seamless organization and manipulation of datasets. It aids in data exploration, assisting us in understanding the structure and characteristics of the dataset before model training.
- TensorFlow: TensorFlow serves as the backbone of the model, providing a robust framework for building and training neural networks. Its integration with Bert enables us to construct and fine-tune complex models for text classification tasks.
- Keras: Keras simplifies the implementation of neural networks, offering a user- friendly interface for constructing model architectures. By utilizing Keras's high-level API to design and configure the classification layers atop the Bert-based model.
- PyTorch: PyTorch complements the model development with its dynamic computation graph and intuitive interface. It facilitates efficient model training and experimentation, empowering us to iterate quickly and refine the Bert-based classifier.
- Transformers: Transformers, particularly the Hugging Face library, provides pre- trained Bert models and utilities for fine-tuning them on specific tasks. By leveraging this library to load pre-trained Bert embeddings and adapt them to the harm classification task through transfer learning.

By harnessing these libraries in conjunction with the Bert-based transfer learning approach, Streamlining the model development process empowers the classification system to effectively and efficiently identify harmful ECA rules with high accuracy.

## Model Evaluation

Model evaluation is the cornerstone of machine learning projects, as it assesses the performance and reliability of predictive models on unseen data. By analyzing various metrics, such as accuracy and precision, model evaluation provides actionable insights into a model's effectiveness, guiding decision-making and ensuring its suitability for real-world applications. This iterative process of scrutiny and refinement is essential for optimizing model performance, mitigating biases, and ultimately driving impactful outcomes in diverse domains.

In the iterative journey of refining machine learning models, the assessment of training loss and accuracy emerges as a cornerstone process. The training loss, serving as a pivotal metric, quantifies the disparity between predicted and actual values during the model's training phase. Its minimization signifies the model's proficiency in learning and adapting to the intricacies of the training data. Concurrently, monitoring training accuracy offers valuable insights into the model's ability to classify or predict instances accurately within the training set. Together, these metrics guide the optimization process, providing a clear view of the model's convergence and performance trajectory. Regular scrutiny of training loss and accuracy is indispensable for identifying potential overfitting or underfitting scenarios, enabling informed decisions on hyperparameter adjustments, and ensuring the model's adeptness in generalizing to new, unseen data. This ongoing evaluation process forms the bedrock for achieving robust and reliable machine learning outcomes.

In addition to tracking training loss and accuracy, the employment of stratified k-fold evaluation further enhances the robustness and reliability of the model. Stratified k-fold cross-validation plays a crucial role in mitigating biases and ensuring a representative distribution of classes across different folds. By preserving class proportions in each fold, this technique furnishes a more accurate estimate of the model's performance, which is particularly vital when confronted with imbalanced datasets. It facilitates an insightful examination of how well the model generalizes to diverse instances, offering a comprehensive assessment of its efficacy across various class distributions. Moreover, stratified k-fold evaluation aids in identifying potential weaknesses or biases inherent in the model's predictions, enabling targeted improvements and enhancing its overall robustness. Embracing this approach underscores the commitment to rigorous model evaluation and ensures the reliability of the machine learning outcomes.

As part of the model evaluation process, the confusion matrix—a quintessential tool is employed for assessing the performance of classification models. The confusion matrix provides a comprehensive summary of the model's predictions by tabulating true positive, true negative, false positive, and false negative instances across different classes. By visualizing the distribution of correct and incorrect predictions, the confusion matrix offers valuable insights into the model's strengths and weaknesses. Additionally, metrics derived from the confusion matrix, such as precision, recall, and F1-score, further illuminate the model's performance characteristics, enabling informed decision-making and iterative model refinement. Leveraging the confusion matrix alongside training metrics and stratified k-fold evaluation enriches the model assessment process, empowering us to build robust and reliable machine learning solutions tailored to real-world challenges.

### **Store model weights**

In the work pipeline for classifying harm in ECA rules, a pivotal step involves storing the learned parameters of the NLP model in a file named "pytorch\_model.bin". This file serves as a repository for the model's weights, encapsulating the knowledge acquired during the training phase. By saving these parameters, Ensuring the preservation and accessibility of the model's learned insights and optimizations facilitates seamless deployment and inference. The stored weights enable quick and efficient model reloading, eliminating the need for retraining from scratch. Additionally, separating the model's architecture from its learned parameters allows for greater flexibility in model management and version control. Utilizing the "pytorch\_model.bin" file not only streamlines the process of obtaining results but also underscores a commitment to efficient and scalable NLP practices.

### **Empowering Edge Intelligence: Deploying the model on Raspberry Pi 4**

Integrating edge solutions into IoT frameworks managing ECA rules heralds a paradigm shift, enriching system dynamics with unparalleled responsiveness and autonomy. By deploying models directly on edge devices like Raspberry Pi boards, the essence of real-time processing is epitomized, facilitating swift event detection and action execution. This local processing not only minimizes latency but also fortifies data privacy and security, a paramount concern in today's digital landscape. Moreover, the inherent resilience of edge solutions ensures seamless operation even in environments with intermittent internet connectivity, enhancing system reliability and adaptability. In essence, harnessing edge solutions in IoT applications imbued with ECA rules amplifies system efficiency, fostering innovation and unlocking new frontiers in diverse domains.

Opting for the Raspberry Pi over an expensive GPU for running the model stems from a combination of practicality, affordability, and versatility. While GPUs offer unparalleled processing power and efficiency for deep learning tasks, they come with a hefty price tag and often require additional infrastructure for setup and maintenance. In contrast, the Raspberry Pi provides a cost-effective and accessible solution, with its compact size, low power consumption, and relatively affordable price point making it an attractive choice for deployment in resource-constrained environments.

Moreover, the Raspberry Pi's versatility extends beyond just running the model—it can serve as a dedicated edge computing device, capable of executing a wide range of tasks beyond deep learning inference. Its ease of setup, portability, and connectivity options make it suitable for deployment in diverse scenarios, from IoT applications to educational projects and hobbyist endeavors. By leveraging the Raspberry Pi, one can not only reduce upfront costs but also gain flexibility and scalability in deploying the model across various use cases and environments.

Equipped with 4GB of RAM, the Raspberry Pi 4 delivers a compelling balance of performance and efficiency, making it an ideal choice for a wide range of projects and applications. This generous memory capacity enables smooth multitasking and seamless execution of resource-intensive tasks, ensuring responsive performance even in demanding scenarios. Whether used for running advanced software applications, hosting web services, or powering multimedia projects, the 4GB RAM variant of the Raspberry Pi 4 provides ample headroom for creativity and innovation. With its compact size and low power consumption, coupled with the versatility of the Raspberry Pi ecosystem, this iteration of the Raspberry Pi opens doors to endless possibilities for hobbyists, educators, and professionals alike.

To deploy the model efficiently on the Raspberry Pi, the 64-bit Raspberry Pi OS was chosen due to its compatibility with PyTorch, an essential component for running the model effectively. With a model size of approximately 420 MB, a robust operating system that supports PyTorch's functionalities was paramount. Additionally, a display was connected to the Raspberry Pi, enabling real-time visualization of the model's results. Using the terminal, a virtual environment was created to encapsulate the project dependencies, ensuring a clean and isolated environment for running the Python script. This meticulous setup aims to maximize the performance and reliability of the model on the Raspberry Pi, paving the way for seamless integration into the desired application environment.

Raspberry Pi OS, formerly known as Raspbian, is the official operating system for the Raspberry Pi single-board computers, including the Raspberry Pi 4. Designed specifically for the Raspberry Pi's ARM architecture, Raspberry Pi OS provides a lightweight and optimized Linux distribution tailored for embedded and IoT projects.

Deploying a deep learning model on a Raspberry Pi 4 using Raspberry Pi OS offers several advantages. Firstly, Raspberry Pi OS provides a familiar Linux environment, making it easy for developers to install and manage dependencies required for deep learning frameworks such as TensorFlow, PyTorch, or Keras. These frameworks are essential for training and deploying neural networks for various tasks, including image classification, object detection, and NLP.

Additionally, Raspberry Pi OS supports hardware acceleration through libraries like OpenCV (Open Source Computer Vision Library) and libraries optimized for the Raspberry Pi's GPU (Graphics Processing Unit). Leveraging hardware acceleration capabilities can significantly improve the performance of deep learning inference tasks on the Raspberry Pi 4, enabling real-time or near-real-time processing of data.

Moreover, Raspberry Pi OS provides access to the extensive Raspberry Pi ecosystem, including official and community-supported software repositories, documentation, and forums. This ecosystem fosters collaboration and innovation, allowing developers to share code, tutorials, and best practices for deploying deep learning models on the Raspberry Pi platform.

Furthermore, Raspberry Pi OS offers a range of system configurations and optimizations to maximize the Raspberry Pi 4's performance and stability. These include options for overclocking, memory allocation, power management, and software optimizations tailored for the Raspberry Pi's hardware architecture.

Overall, Raspberry Pi OS provides a reliable and versatile platform for deploying deep learning models on the Raspberry Pi 4. By harnessing the power of Raspberry Pi OS and the Raspberry Pi 4's hardware capabilities, developers can build cost-effective and energy-efficient solutions for edge AI (Artificial Intelligence) applications, IoT devices, robotics, and more.

**Fostering Cloud Collaboration: Elevating the Model to the Digital Stratosphere** Transitioning to a cloud solution holds immense potential for enhancing the efficacy of the classification system for identifying harm in ECA rules. By leveraging cloud computing resources, one can unlock a myriad of benefits that bolster the performance, scalability, and accessibility of the model. Cloud solutions offer unparalleled scalability, enabling us to seamlessly accommodate fluctuations in demand and handle large volumes of data with ease. Additionally, the inherent flexibility of cloud platforms allows for rapid deployment and updates, ensuring that the classification system remains agile and responsive to evolving needs. Furthermore, the accessibility of cloud-based solutions enables users to access the classification system from anywhere with an internet connection, fostering collaboration and expanding the reach of the solution to a global audience. Through the adoption of a cloud solution, one can maximize the impact of the classification system, driving innovation and empowering stakeholders to make informed decisions in navigating complex regulatory landscapes.

In the project's trajectory, a strategic leap has been taken by developing a robust model and building a user-friendly website to facilitate interaction with it. However, this website has not yet been deployed on a server. Looking ahead, utilizing this website as a cloud solution is recognized for its potential to extend the accessibility of the model to a broader audience. Deploying the website on a server aims to leverage cloud computing resources to enhance scalability, reliability, and performance. This transition to a cloud-based solution aligns with the vision of harnessing cutting-edge technologies to address real-world challenges effectively. Through this strategic evolution, new opportunities are anticipated, driving transformative impact in the domain of classifying harm in ECA rules.



Fig. 6 User interface of the Harmful ECA rule classification website

## Technology used to create the website

1) **Flask**: it is a lightweight and versatile Python web framework ideal for developing web applications. With its simplicity and flexibility, Flask allows developers to quickly create powerful web applications with minimal boilerplate code. It follows the WSGI (Web Server Gateway Interface) specification, making it compatible with various web servers and deployment options. Flask provides essential features for web development, including routing, templating, and handling HTTP requests and responses. Its modular design encourages the use of extensions, enabling developers to add functionalities like authentication, database integration, and RESTful APIs seamlessly. Whether building a small-scale project or a complex web application, Flask empowers developers to create elegant and efficient solutions with ease.

2) **HTML and CSS**: HTML (Hypertext Markup Language) and CSS (Cascading Style Sheets) are foundational technologies essential for developing modern web applications. HTML provides the structure and content of web pages, while CSS adds style and presentation to make them visually appealing and user-friendly.

HTML serves as the backbone of web development by defining the structure of a web page through various elements like headings, paragraphs, lists, links, forms, and more. These elements allow developers to organize and present



information in a logical and hierarchical manner, creating a seamless user experience. HTML also facilitates accessibility by providing semantic markup that helps screen readers and search engines understand the content of web pages, enhancing usability for all users.

### Model Evaluation

Specifically, the assessment of training loss and accuracy assumes critical significance in the iterative process of refining machine learning models. The training loss serves as a pivotal metric, quantifying the disparity between predicted and actual values during the model's training phase. Its minimization is essential, indicating the model's capacity to learn and adapt effectively to the training data. Concurrently, the tracking of training accuracy provides valuable insights into the model's ability to accurately classify or predict instances within the training set. These metrics collectively guide the optimization process, offering a lens into the model's convergence and performance. Regular scrutiny of training loss and accuracy is indispensable for identifying potential overfitting or underfitting scenarios, facilitating informed decisions on hyperparameter adjustments, and ensuring the model's adeptness in generalizing to new, unseen data. This ongoing evaluation process is fundamental for achieving robust and reliable machine learning outcomes. In addition to tracking training loss and accuracy, employing stratified k-fold evaluation enhances the robustness and reliability of machine learning models. Stratified k-fold cross-validation is crucial for mitigating biases and ensuring a representative distribution of classes across different folds. By maintaining class proportions in each fold, this technique provides a more accurate estimate of the model's performance, especially vital when dealing with imbalanced datasets. It helps identify how well the model generalizes to diverse instances, offering a more comprehensive assessment of its efficacy.

Furthermore, evaluating the loss on a separate validation set is pivotal for understanding how well the model performs on unseen data. The evaluation loss serves as a key indicator of the model's ability to generalize beyond the training set. A meticulous examination of stratified k-fold cross-validation and evaluation loss collectively ensures a thorough understanding of a model's performance, contributing to its reliability in real-world applications.

### RESULTS

Training loss:



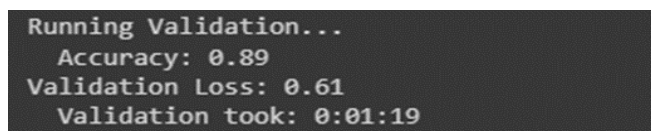
Fig. 7 Training loss in various epochs

In the context of the model, training loss refers to a metric used to quantify the disparity between the model's predictions and the actual labels during the training phase. Specifically, it measures how well the model is performing on the training data by calculating the error or deviation between the predicted outputs and the ground truth labels. A lower training loss indicates that the model's predictions are closer to the true values, signifying better performance and improved learning.

For instance, in epoch one of the training process, the average training loss of 0.53 suggests that, on average, the model's predictions deviated by approximately 0.53 units from the actual labels in the training dataset. As training progresses, this loss is minimized, reflecting the model's gradual improvement in capturing the underlying patterns and relationships in the data. In epoch two, the average training loss decreases to 0.49, indicating a further refinement in the model's performance and a reduction in prediction errors.

#### ● Validation Accuracy:

Validation accuracy is a critical metric used to evaluate the performance of a machine learning model on a separate validation dataset. Specifically, it measures the proportion of correctly predicted instances out of the total number of



instances in the validation dataset. In the context, achieving a validation accuracy of 89% indicates that the model correctly classified approximately 89% of the instances in the validation dataset.

Fig. 8 Validation Accuracy and Validation loss

This high validation accuracy signifies that the model demonstrates strong performance and generalization capabilities, as it is able to accurately classify instances it hasn't seen during the training phase. A validation accuracy of 89% suggests that the model effectively captures the underlying patterns and relationships in the data, enabling it to make accurate predictions on unseen data.

Monitoring validation accuracy is crucial in assessing the robustness and reliability of the model, as it provides insights into how well the model is likely to perform in real-world scenarios. By achieving a validation accuracy of 89%, the model demonstrates promising performance and instills confidence in its ability to classify harm in ECA rules accurately.

### Test Accuracy and Confusion Matrix

In the multiclass classification of harmful ECA rules, a confusion matrix serves as a fundamental tool for evaluating the performance of the model. Essentially, a confusion matrix is a table that allows us to visualize the performance of a classification algorithm by tabulating the actual and predicted classes for each instance in the dataset. In the context of the classification task, the confusion matrix would have rows corresponding to the actual classes of harmful ECA rules and columns corresponding to the predicted classes generated by the model. Each cell in the matrix represents the number of instances that belong to a specific actual class and were predicted to belong to a specific predicted class.

The value in each cell of the confusion matrix provides insights into the model's performance across different classes. For instance, cells along the diagonal of the matrix represent instances that were correctly classified, where the actual class matches the predicted class. On the other hand, off-diagonal cells indicate misclassifications, where the actual class differs from the predicted class. By analyzing the distribution of correct and incorrect predictions across the matrix, one can identify patterns of misclassification and assess the model's strengths and weaknesses in classifying different types of harmful ECA rules. Additionally, summary statistics derived from the confusion matrix, such as accuracy, precision, recall, and F1-score, offer comprehensive measures of the model's performance across all classes. The confusion matrix stands out as the best way to evaluate the model due to its ability to provide detailed and granular insights into classification performance. Unlike single-value metrics like accuracy, which provide an overall assessment of model performance, the confusion matrix allows us to delve deeper into the specific types of errors made by the model. This granularity is particularly crucial in the domain, where accurately identifying harmful ECA rules is of utmost importance. Leveraging the information provided by the confusion matrix allows for fine-tuning the model, prioritizing classes that are prone to misclassification, and ultimately enhancing the effectiveness of the classification system in identifying and mitigating harm in ECA rules.

In addition to its tabular format, the confusion matrix can also be visualized through a heatmap representation, which offers a more intuitive and visually appealing way to interpret the classification performance of the model. In this representation, each cell in the confusion matrix is color-coded based on its value, with brighter colors indicating higher counts and darker colors representing lower counts. By visualizing the confusion matrix as a heatmap, patterns of correct and incorrect classifications become readily apparent, allowing us to identify clusters of misclassifications and areas of high accuracy more intuitively. The heatmap representation not only facilitates easier interpretation of the confusion matrix but also enhances the ability to identify trends and insights at a glance, making it a valuable tool for assessing and improving the performance of the model in classifying harmful ECA rules.

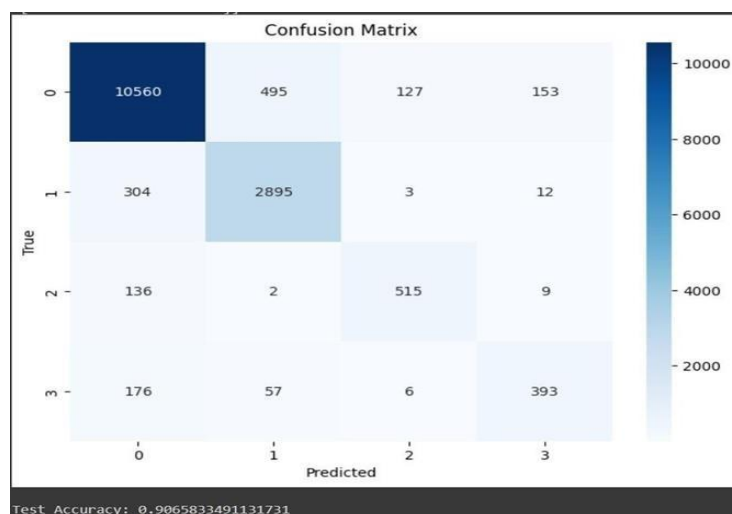


Fig. 9 Heatmap Representation

Class 0 (Row 1):

- True Positives (TP): 10560
- False Negatives (FN):  $495 + 127 + 153 = 775$
- False Positives (FP):  $304 + 136 + 176 = 616$
- True Negatives (TN): Sum of all other elements = 13408

Class 1 (Row 2):

- TP: 2895
- FN:  $304 + 3 + 12 = 319$
- FP:  $495 + 2 + 57 = 554$
- TN: 12203

Class 2 (Row 3):

- TP: 515
- FN:  $136 + 2 + 9 = 147$
- FP:  $127 + 3 + 6 = 136$
- TN: 14045

Class 3 (Row 4):

- TP: 393
- FN:  $176 + 57 + 6 = 239$
- FP:  $153 + 12 + 9 = 174$
- TN: 15037

## Classification Report

The classification report includes precision, recall, and F1-score for each class, along with the overall accuracy, macro average, and weighted average metrics.

### Precision

Precision is the ratio of correctly predicted positive observations to the total predicted positives.

$$\text{Precision} = \text{TP} / (\text{TP} + \text{FP})$$

### Recall (Sensitivity)

Recall is the ratio of correctly predicted positive observations to the all observations in the actual class.

$$\text{Recall} = \text{TP} / (\text{TP} + \text{FN})$$

## F1-Score

The F1-score is the harmonic mean of precision and recall, providing a single metric that balances both.

$$\text{F1-Score} = 2 * (\text{Precision} * \text{Recall}) / (\text{Precision} + \text{Recall})$$

|                                      |              |           |        |          |         |
|--------------------------------------|--------------|-----------|--------|----------|---------|
| Confusion Matrix (Testing Set):      |              |           |        |          |         |
| [[                                   | 10560        | 495       | 127    | 153]     |         |
| [                                    | 304          | 2895      | 3      | 12]      |         |
| [                                    | 136          | 2         | 515    | 9]       |         |
| [                                    | 176          | 57        | 6      | 393]]    |         |
| Classification Report (Testing Set): |              |           |        |          |         |
|                                      |              | precision | recall | f1-score | support |
|                                      | 0            | 0.94      | 0.93   | 0.94     | 11335   |
|                                      | 1            | 0.84      | 0.90   | 0.87     | 3214    |
|                                      | 2            | 0.79      | 0.78   | 0.78     | 662     |
|                                      | 3            | 0.69      | 0.62   | 0.66     | 632     |
|                                      | accuracy     |           |        | 0.91     | 15843   |
|                                      | macro avg    | 0.82      | 0.81   | 0.81     | 15843   |
|                                      | weighted avg | 0.91      | 0.91   | 0.91     | 15843   |

Fig. 10 Confusion Matrix and various other evaluation parameters

## Support

Support is the number of actual occurrences of the class in the dataset.

## Overall Metrics

Accuracy: The ratio of correctly predicted instances to the total instances.

Accuracy = Number of correct predictions/Total number of predictions = 0.91

Macro Average: The arithmetic mean of precision, recall, and F1-score across all classes. It treats all classes equally regardless of their support.

Macro Precision =  $(0.94 + 0.84 + 0.79 + 0.69) / 4 = 0.82$  Macro

Recall =  $(0.93 + 0.90 + 0.78 + 0.62) / 4 = 0.81$  Macro F1-Score =

$(0.94 + 0.87 + 0.78 + 0.66) / 4 = 0.81$

Weighted Average: The weighted mean of precision, recall, and F1-score, taking into account the support of each class.

Discussions:

### Confusion Matrix Analysis:

Class 0:

- High TP count indicates the model accurately identifies the majority of class 0 instances.
- The FN count (775) suggests some instances of class 0 are being misclassified as other classes, but relatively few compared to the large support of 11335.
- FP count (616) indicates some non-class 0 instances are incorrectly classified as class 0, but this is relatively small compared to the TN count.
- Overall, the model performs very well for class 0, as seen in the high precision (0.94) and recall (0.93).

Class 1:

- The model correctly classifies a large number of class 1 instances.
- The FN count (319) indicates some class 1 instances are being misclassified, but this is small relative to the support of 3214.
- FP count (554) shows some misclassification into class 1, which is manageable.

- Precision (0.84) and recall (0.90) are high, indicating good performance for class 1.

Class 2:

- The model identifies most class 2 instances correctly.
- The FN count (147) shows a noticeable amount of misclassification, significant given the smaller support of 662.
- FP count (136) indicates some confusion with other classes.
- Precision (0.79) and recall (0.78) are slightly lower, reflecting moderate performance.

Class 3:

- The model struggles more with class 3, having the lowest TP count.
- The FN count (239) indicates significant misclassification, which is notable given the support of 632.
- FP count (174) shows moderate confusion with other classes.
- Precision (0.69) and recall (0.62) are the lowest, indicating the model has difficulty accurately identifying class 3 instances.

### **Classification Report Analysis**

Precision, Recall, and F1-Score

Class 0:

- Precision (0.94) and recall (0.93) are very high, indicating the model is excellent at identifying class 0 instances with minimal misclassification.
- High F1-score (0.94) confirms the balanced high precision and recall.

Class 1:

- Precision (0.84) and recall (0.90) are also high, suggesting strong performance with slight room for improvement in reducing FP.
- F1-score (0.87) reflects this good performance balance.

Class 2:

- Precision (0.79) and recall (0.78) are moderately high, indicating the model is reasonably good at identifying class 2, but with more misclassifications than classes 0 and 1.
- F1-score (0.78) is slightly lower, pointing to a need for improvement in balancing precision and recall.

Class 3:

- Precision (0.69) and recall (0.62) are lower, highlighting difficulties the model has with accurately identifying class 3 instances.
- F1-score (0.66) suggests significant room for improvement in model performance for this class.

Overall Metrics:

Accuracy (0.91): Indicates that 91% of the total instances are correctly classified. This high accuracy suggests the model performs well overall.

Macro Average:

Precision (0.82), recall (0.81), and F1-score (0.81) give an unweighted average, treating all classes equally. The lower scores compared to the weighted averages suggest that the smaller classes (like class 3) negatively impact these averages.

Weighted Average:

Precision (0.91), recall (0.91), and F1-score (0.91) take into account the number of instances for each class, indicating overall excellent performance but with the larger classes having more influence.



## Implications

- **High Accuracy:** The model is generally reliable across the dataset.
- **Strong Performance for Major Classes (0 and 1):** The model excels in identifying the more frequent classes, with high precision, recall, and F1-scores.
- **Moderate Performance for Minor Classes (2 and 3):** Performance decreases for less frequent classes, especially class 3, indicating a need for further model tuning or perhaps additional training data to better learn these classes.
- **Balanced Metrics:** High weighted averages show the model performs well overall, but lower macro averages suggest that improvements are needed to ensure better performance across all classes.

CLASS 0 : NO HARM



Fig. 11 Predicting an ECA rule belonging to class 0

The ECA rule "If I meet my daily step target, update a file with the statistics on my mobile" poses no personal, physical, or cybersecurity harm. It simply tracks the user's steps and records the data locally on the mobile device, ensuring privacy as no sensitive personal information is shared externally. The rule encourages physical activity, promoting health without compelling unsafe or excessive exercise. Additionally, as the data update occurs within the secure environment of the mobile device, it minimizes cybersecurity risks such as data breaches or unauthorized access. Overall, the rule supports health monitoring while maintaining user privacy and security.

CLASS 1 : PERSONAL DAMAGE HARM



Fig. 12 Predicting an ECA rule belonging to class 1

The ECA rule "Retweet everything that person X tweets" poses personal harm. Automatically retweeting all content can expose the user to unwanted or inappropriate material, potentially damaging their online reputation and personal relationships. It reduces control over their social media presence, leading to a cluttered timeline and misrepresentation of their views. This rule also risks privacy by revealing patterns of behavior. Overall, it compromises the user's control and privacy, leading to significant personal harm.

## CLASS 2: PHYSICAL HARM



Fig. 13 Predicting an ECA rule belonging to class 2

The ECA rule "Open the windows if the temperature goes up by 25 degrees Celsius" poses physical harm. Automatically opening windows can be dangerous during high winds, storms, or heavy rainfall, allowing water and debris to enter and cause damage or injury. It also creates safety hazards for young children and pets, increasing the risk of falls. Additionally, open windows compromise home security, making unauthorized entry easier. These factors highlight the potential physical risks associated with this rule.

## CLASS 3: CYBER SECURITY HARM

The ECA rule "Share any attached file in email to drive" poses a significant cybersecurity threat. Automatically transferring all email attachments to a drive can expose sensitive or confidential information to unauthorized access, especially if the drive is not adequately secured. This rule can also inadvertently upload malware or phishing attachments, spreading harmful software and compromising the security of the drive and other connected systems. Furthermore, without proper filtering, this action may violate privacy policies or data protection regulations, leading to potential legal and financial repercussions. Overall, this rule significantly increases the risk of data breaches and cybersecurity threats.



Fig. 14 – Predicting an ECA rule belonging to class 3

## Output in Raspberry Pi board:

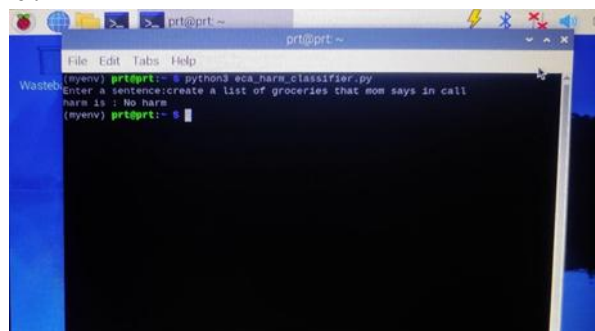


Fig. 15 – Predicting an ECA rule belonging to class 0 in Raspberry pi

In the exploration of ECA rules and their potential impacts, the rule: "If mom says groceries in a call, then create a

list of groceries," is evaluated. This automation is beneficial and causes no harm due to its focus on non-sensitive data. By transcribing grocery lists from calls, it enhances efficiency and convenience, ensuring no items are forgotten and reducing the need for multiple store trips. Operating within a secure environment, it minimizes cybersecurity risks as it does not interact with potentially malicious external inputs. The automated list is accurate and clear, eliminating the risk of miscommunication. Furthermore, the process is transparent and traceable, allowing for easy verification and correction of any errors. This example underscores the importance of thoughtfully designed ECA rules that improve user experience without compromising privacy or security.

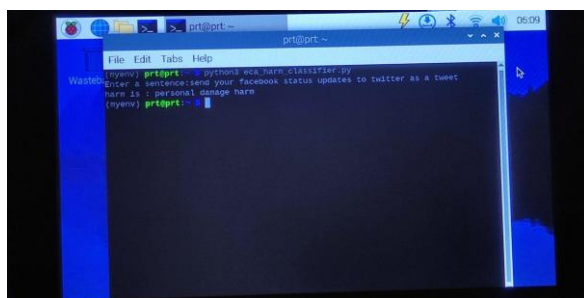


Fig. 16 Predicting an ECA rule belonging to class 1 in Raspberry pi

Automatically sharing Facebook status updates on Twitter through the ECA rule raises significant concerns regarding personal harm. By broadcasting private thoughts and activities to a wider audience without discretion, individuals risk exposing sensitive information to potential stalkers or malicious actors. This indiscriminate sharing of personal content on a public platform like Twitter removes control over who can access and interpret the shared information, leading to privacy breaches, harassment, or even threats to personal safety. Additionally, without context or consideration for the diverse audience on Twitter, individuals may inadvertently disclose details that could be exploited by individuals with harmful intentions, highlighting the inherent risks associated with this ECA rule.

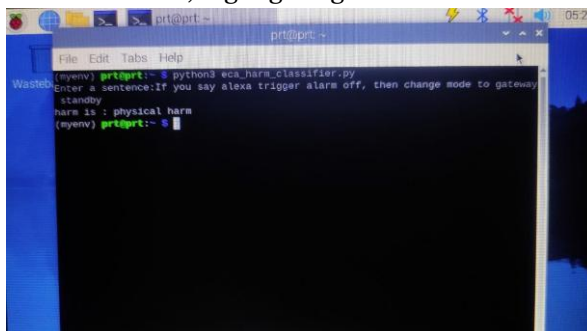


Fig. 17 Predicting an ECA rule belonging to class 2 in Raspberry pi

The command "Alexa, trigger alarm off then change the mode to gateway standby" can cause physical harm by compromising the security of a space. Disabling an alarm and switching to a standby mode can leave a property unprotected, making it vulnerable to unauthorized access or intrusion. This action can lead to potential physical harm to the inhabitants if intruders exploit the security lapse. For example, in a scenario where the alarm is a critical component of a home security system, turning it off can expose residents to the risk of burglary or other physical threats. The failure to maintain a secure environment directly endangers the physical safety of individuals, highlighting the severe implications of such an ECA rule.

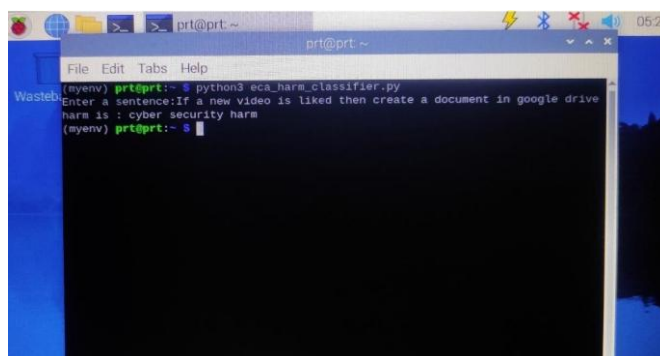


Fig. 18 Predicting an ECA rule belonging to class 3 in Raspberry pi

In the exploration of ECA rules and their potential impacts, the cybersecurity implications of the rule: "If a new YouTube video is liked, then create a document in Google Drive." This automation, while seemingly harmless and aimed at organizing content, poses significant cybersecurity risks. Automatically creating documents based on external interactions, like liking a YouTube video, could expose the Google Drive to unnecessary vulnerabilities. For instance, if an account is compromised, an attacker could manipulate this rule to create numerous documents, leading to potential misuse of storage and clutter. Moreover, this could serve as a gateway for phishing attacks, where malicious links or content could be embedded in the automatically created documents, thereby compromising the security of other files in Google Drive. This highlights the critical need for evaluating the security ramifications of integrating different services through ECA rules, ensuring that such automations do not inadvertently open.

### CONCLUSION

This work addresses the security and privacy risks associated with IoT automation using IFTTT applets by proposing a comprehensive solution. The study identifies the proliferation of IoT technology in creating smart homes but highlights the potential security vulnerabilities and privacy concerns posed by ECA rules. Leveraging NLP and machine learning techniques, a multi-faceted approach is employed to classify applets based on their security and privacy risks. A deep learning model, particularly BERT, is utilized for accurate classification, while Raspberry Pi serves as the edge solution for localized predictions, ensuring data privacy and sovereignty. Additionally, a user-friendly Ibsite deployed as a cloud solution offers convenient access to the risk prediction service, providing users with actionable insights to mitigate potential threats. By offering both edge and cloud solutions, the work aims to enhance the security and privacy of IoT ecosystems, empowering users to make informed decisions about their automation configurations. In the study, the dataset was divided into training and testing sets with an allocation of 80% and 20% respectively. This deliberate distribution ensured that the model learned from a diverse and representative sample of data, contributing to its ability to generalize effectively to unseen instances. During the training phase, the model achieved an impressive validation accuracy of approximately 89% and exhibited a notably low training loss of 0.53. These results indicated the model's proficiency in minimizing errors and capturing underlying patterns in the data.

Subsequently, the trained model underwent evaluation on the separate testing set, serving as a crucial benchmark for assessing its generalization capability. The model demonstrated remarkable performance with a testing accuracy of 90.65%, reaffirming its reliability and robustness. Furthermore, the validation loss on unseen validation data remained low at 0.61, indicating consistent generalization even on previously unseen instances. Additionally, the model's performance was assessed using metrics such as F1 scores and precision scores, providing insights into its accuracy across multiple classes. The F1 scores for each class (0, 1, 2, and 3) were 94%, 87%, 78%, and 66% respectively, while the precision scores were 94%, 84%, 79%, and 69% respectively. These metrics highlighted the model's proficiency in accurately distinguishing between different classes.

In summary, the study showcased the efficacy of the training approach and the robustness of the model in real-world applications. The commendable validation accuracy and low training loss underscored the model's ability to learn

effectively. Moreover, the impressive testing accuracy and additional performance metrics validated the model's proficiency and provided a promising foundation for practical deployment in real-world scenarios.

## REFERENCES

- [1] J. Cano, E. Rutten, G. Delaval, Y. Benazzouz, and L. Gurgun, "ECA Rules for IoT Environment: A Case Study in Safe Design," in 2014 IEEE Eighth International Conference on Self-Adaptive and Self-Organizing Systems Workshops, London, UK, 2014, pp. 116–121, doi: 10.1109/SASOW.2014.32.
- [2] B. Breve, G. Cimino, and V. Deufemia, "Identifying Security and Privacy Violation Rules in Trigger-Action IoT Platforms With NLP Models," in IEEE IOT Journal, vol. 10, no. 6, pp. 5607–5622, March 15, 2023, doi: 10.1109/JIOT.2022.3222615.
- [3] M. J. Jozani, É. Marchand, and A. Parsian, "On estimation with weighted balanced-type loss function," Stat. Probability Lett., vol. 76, no. 8, pp. 773–780, 2006.
- [4] N. Reimers and I. Gurevych, "Sentence-BERT: Sentence embeddings using siamese BERT-networks," in Proc. 9th Int. Joint Conf. Nat. Lang. Process., 2019, pp. 3982–3992.
- [5] L. Wang, M. Han, X. Li, N. Zhang, and H. Cheng, "Review of Classification Methods on Unbalanced Data Sets," in IEEE Access, vol. 9, pp. 64606–64628, 2021, doi:10.1109/ACCESS.2021.3074243.
- [6] S. Ghannay, B. Favre, Y. Esteve, and N. Camelin, "Word embedding evaluation and combination," in Proc. 10th Int. Conf. Lang. Resour. Eval., 2016, pp. 300–305.
- [7] B. A. Johnsson and B. Magnusson, "Towards end-user development of graphical user interfaces for IOT," Future Gener. Comput. Syst., vol. 107, pp. 670–680, Jun. 2020.
- [8] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in Proc. Conf. North Amer. Chapter ACL Human Lang. Technol., vol. 1, 2019, pp. 4171–4186.
- [9] R. Wang, R. Ridley, W. Qu, X. Dai, and X. Su, "A novel reasoning mechanism for multi-label text classification," Inf. Process. Manage., vol. 58, no. 2, 2021, Art. no. 102441.
- [10] B. Breve, G. Desolda, V. Deufemia, F. Greco, and M. Matera, "An enduser development approach to secure smart environments," in Proc. 8th Int. Symp. End-User Develop., 2021, pp. 3652.
- [11] M. Saeidi, M. Calvert, A. W. Au, A. Sarma, and R. B. Bobba, "If this context then that concern: Exploring users' concerns with IFTTT applets," Proc. Privacy Enhanc. Technol., vol. 2022, no. 1, pp. 166–186, 2022.
- [12] M. McCall et al., "SAFETAP: An efficient incremental analyzer for trigger– action programs," Carnegie Mellon University, Pittsburgh, PA, USA, Rep. 14792271, 2021.
- [13] Wenshu Zha a, Yuping Liu a, Yujin Wan b, Ruilan Luo b, Daolun Li a, Shan Yang b, Yanmei Xu b, "Forecasting monthly gas field production based on the CNN-LSTM model" in July 2022.
- [14] Mihai Masala, Stefan Ruseti, and Mihai Dascalu. 2020. RoBERT – A Romanian BERT Model. In Proceedings of the 28th International Conference on Computational Linguistics, pages 6626–6637, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- [15] García-Grao G, Carrera Á. Extending the OSLC Standard for ECA-Based Automation. Electronics. 2023; 12(14):3043.
- [16] Smagulova, K., James, A.P. A survey on LSTM memristive neural network architectures and applications. Eur. Phys. J. Spec. Top. 228, 2313–2324 (2019).
- [17] K. P, R. S. P, H. P, D. M and C. Iwendi, "Customized CNN with Adam and Nadam Optimizers for Emotion Recognition using Facial Expressions," 2023 International Conference on Wireless Communications Signal Processing and Networking (WiSPNET), Chennai, India, 2023.
- [18] Yi, Dokkyun, Jaehyun Ahn, and Sangmin Ji. 2020. "An Effective Optimization Method for Machine Learning Based on ADAM" Applied Sciences 10, no. 3: 1073. <https://doi.org/10.3390/app10031073>
- [19] Alanazi, Turki M. "Embedded System Based Raspberry Pi 4 for Text Detection and Recognition." Intelligent Automation & Soft Computing 36, no. 3 (2023).
- [20] Biglari, Amin, and Wei Tang. "A review of embedded machine learning based on hardware, application, and sensing scheme." Sensors 23, no. 4 (2023): 2131.