# Interpretability and Transparency-Driven Detection and Transformation of Textual Adversarial Examples (IT-DT)

Bushra Sabir[1,2], M. Ali Babar[1], and Sharif Abuadbba[2]

[1] University of Adelaide
[2] CSIRO's Data61

**Abstract.** Transformer-based text classifiers, such as BERT, Roberta, T5, and GPT-3, have achieved impressive performance in Natural Language Processing (NLP). However, their vulnerability to adversarial examples poses a significant security risk. Existing defense methods often lack interpretability and transparency, making it difficult to understand the reasoning behind adversarial classifications and identify vulnerabilities in the models. To address these limitations, we propose an approach that focuses on interpretability and transparency in the detection and transformation of textual adversarial examples. Our framework, titled Interpretability and Transparency-Driven Detection and Transformation (IT-DT), aims to provide insights into the decision-making process of the models and enable effective mitigation of adversarial attacks. The IT-DT framework leverages techniques such as attention maps, integrated gradients, and model feedback to achieve interpretability in the detection phase. By visualizing the attention and gradient information, we can identify the salient features and perturbed words that contribute to adversarial classifications. This interpretability enhances our understanding of the model's vulnerabilities and helps us identify potential adversarial inputs. In the transformation phase, the IT-DT framework utilizes pre-trained embeddings and model feedback to generate optimal replacements for the perturbed words. By finding suitable substitutions, we aim to convert the adversarial examples into non-adversarial counterparts that align with the model's intended behavior. This transformation process ensures that the modified inputs no longer deceive the model while preserving the overall meaning and context of the text. Moreover, the IT-DT framework emphasizes transparency by involving human experts in the loop. Human intervention plays a crucial role in reviewing and providing feedback on the detection and transformation results. This collaborative approach enhances the system's decision-making process, particularly in complex scenarios where automated methods may face challenges. Furthermore, the IT-DT framework generates valuable insights and threat intelligence that empower security analysts to identify vulnerabilities and improve model robustness. The framework aims to bridge the gap between the technical aspects of adversarial attacks and the human understanding required for effective mitigation. Through comprehensive experiments, we demonstrate the effectiveness of the IT-DT framework in detecting and transforming textual adversarial examples.

Our approach enhances interpretability, provides transparency into the decision-making process, and enables accurate identification and successful transformation of adversarial inputs. By combining technical analysis and human expertise, the IT-DT framework significantly improves the resilience and trustworthiness of transformer-based text classifiers against adversarial attacks.

**Keywords:** Human-centric, Interpretability, Transperancy, Proactive Defence against Adversarial Examples, Textual Adversarial Examples

## 1   Introduction

Transformer-based Large Language Models (TLLMs) and their variants have significantly transformed the field of Natural Language Processing (NLP) [13, 26, 48]. These models have excelled in user review analysis, hate speech detection, content moderation, cyber-attack detection, and news classification on online platforms [13, 26, 48]. Prominent industry leaders, including Google, Facebook, and Amazon, have been actively leveraging Language Models (LLMs) for proactive text moderation purposes. These companies have implemented LLM-based models, such as Google's Perspective API [3], Facebook's CS2T model for hate speech detection [4], and Amazon Comprehend [5], to effectively address text moderation challenges.

However, these models, like their other DNN counterparts, are not foolproof and are vulnerable to carefully designed Adversarial Examples (AEs) [19, 24]. AEs are generated by perturbing the original example in a way that retains its original semantics but induces test-time misclassification, resulting in evading the target text classifier. For example, transforming "The Fish N Chips are excellent" to "The Fish N Chips are first-class" changes the target model output (BERT for restaurant review classification) from positive to a negative restaurant review. Adversaries utilize several obfuscation methods to generate AEs ranging from substituting to inserting critical words (e.g., HotFlip [14], TextFooler [19], characters (e.g., DeepWordBug [16], TextBugger [23] phrases (e.g., [21, 57]) in the text to fool the target model. Among these techniques, Word Substitution Attacks (WSA) are a popular approach to creating adversarial examples against transformer-based models [19]. By replacing a critical word or multiple words with mis-spelt versions, synonyms, or contextually similar words, these attacks can bypass the security measures of transformer-based models while maintaining the original text's meaning [23]. The ease with which WSA can be generated using pre-existing tools and libraries further enhances their significance [32].

Recent studies have shown that transformer-based models can misclassify over 90% of adversarial examples generated using WSA, highlighting the potential impact on users who rely on the model's output. Therefore, to ensure

---

[3] https://www.perspectiveapi.com/

[4] https://ai.facebook.com/blog/cs2t-a-novel-model-for-hate-speech-detection/

[5] https://aws.amazon.com/comprehend/

the reliability and integrity of these models, it is crucial to focus on developing defences against WSA. Failure to secure these models against WSA can lead to incorrect information being presented to users, affecting their satisfaction and trust in the model's output [20]. Therefore, techniques must be developed to protect transformer-based models from WSA, ensuring their reliability and safety in real-world applications [20].

The landscape of adversarial attacks against transformer-based models has been evolving rapidly, necessitating the development of sophisticated defenses to mitigate the associated risks [5, 34, 54, 62]. While existing defense mechanisms, including Adversarial Training, Synonym Encoding, Frequency-Guided Word Substitution, and Discriminate Perturbations, have shown promising results in countering Word Substitution Attacks (WSA), they suffer from a lack of transparency and interpretability, thereby limiting their efficacy in real-world applications. These defenses often operate as black boxes, impeding their transparency and interpretability. In practical scenarios, it is crucial to have actionable intelligence, such as raising alerts upon detecting adversarial examples, attempting autocorrection, and involving security analysts equipped with enriched quality logs that aid in identifying the attack type, its source, and the vulnerabilities exploited by the attacker [5].

Human-centric AI defenses play a crucial role in bridging the gap in existing defenses by prioritizing transparency and understandability for humans [36, 56]. These defenses aim to provide security analysts with actionable intelligence, enabling them to interpret attack types and sources and determine the need for human intervention [37]. Integrating comprehensive threat intelligence and fostering collaboration enhances preemptive capabilities. Strengthening defenses through alarms, transparency, and collaboration is vital for detecting and mitigating new attacks. Logging relevant information enables forensic analysis and the development of proactive countermeasures, aiding in attack investigations [44]. Human-centric AI defenses leverage human abilities to improve reliability, reduce false positives and negatives, and establish trust in AI systems. Developing such defenses is crucial for ensuring the integrity and reliability of transformer-based models in real-world applications. While previous studies have explored interpretability to demonstrate the significance of their methodologies [18], to the best of our knowledge, none have achieved the detection and transformation of adversarial examples in a transparent and human-understandable manner. Additionally, existing defenses often operate as black boxes, offering limited insights into the nature of the attack or whether human intervention is necessary to address it.

In this study, we propose a novel framework called **I**nterpretibility and **T**ransparency driven **D**etection and **T**ransformation (IT-DT) that provides interpretable detection, identification, and correction of adversarial examples for transformer-based text classification models. The IT-DT framework leverages explainability features, such as attention maps and integrated gradients, to automate the detection and identification of adversarial examples and pinpoint the perturbed words in Word Substitution Attacks at multiple levels of granularity, including char-

acters, words, and multiple word substitutions. By incorporating explainability, the framework offers insights into the distinguishing features and patterns that differentiate adversarial examples from legitimate ones. Furthermore, it identifies specific words that require transformation, enabling the conversion of adversarial examples into non-adversarial ones. The transformation process in the IT-DT framework utilizes pre-trained embeddings, attention weights, word frequency analysis, and model feedback to determine the words that need to be modified for effective transformation. Additionally, the framework logs intermediate information and identifies scenarios that require human intervention, providing valuable threat intelligence and interpretability to security analysts for identifying vulnerabilities in the target model and enhancing resilience against Word Substitution Attacks.

We comprehensively evaluate our proposed framework and its components on two transformer-based architectures, namely BERT and ROBERTA, trained on four state-of-the-art text classification datasets: IMDB, YELP, AGNEWS, and SST2. Additionally, we assess our framework against seven state-of-the-art Word Substitution Attacks at three levels of granularity: character-level, word-level, and multi-level. Our experimental results demonstrate that our approach significantly enhances the model's resistance to adversarial word substitution attacks, improving its reliability and effectiveness in practical applications. We achieve a median detection performance, measured by Matthew Correlation Coefficient (MCC), of 0.846 and an F-score of 92% across the four datasets. Our method dramatically improves the median accuracy of the considered models on adversarial examples from zero to 92.98% against the state-of-the-art adversarial word substitution attacks.

*Contributions* The main contributions of this work can be summarized as follows:

1. We develop a novel Adversarial Detector ($D_{adv}$) by extracting features from attention maps, integrated gradients, and frequency distribution to effectively differentiate between adversarial and clean examples.
2. We devise a hybrid approach that combines attention, frequency, and model feedback to identify perturbed words accurately.
3. We propose a novel transformation mechanism to find optimal replacements for perturbed words, converting adversarial examples into non-adversarial forms. Additionally, we identify specific scenarios that require human intervention.
4. We comprehensively evaluate our proposed detector, our identification and transformation module over four datasets, two transformation-based architectures and seven SOTA adversarial attacks.
5. Our framework provides valuable threat intelligence and actionable insights to security analysts facilitating human-AI collaboration. These insights can enhance the robustness of the defence mechanism and identify vulnerabilities in the target model.

*Significance and Implications* The proposed method, the Interpretability and Transparency driven Detection and Transformation of Adversarial Examples

(IT-DT) framework, has several implications for both research and practitioners in the field of cyber-security, natural language processing (NLP) and machine learning. These implications are as follows:

– The proposed method, the Interpretability and Transparency driven Detection and Transformation of Adversarial Examples (IT-DT) framework, has several implications for both researchers and practitioners in the field of cyber-security, natural language processing (NLP) and machine learning. These implications are as follows:

– (i) Advancement in Adversarial Defense Research: The IT-DT framework contributes to advancing research in adversarial defence techniques targeted explicitly at transformer-based models. By addressing the vulnerability of transformers to adversarial word substitution attacks, it offers a comprehensive defence mechanism that combines explainability, detection, and transformation methods. In addition, this research can inspire further investigations into enhancing the security and reliability of deep learning models against adversarial attacks.

– (ii) Enhanced Transparency and Interpretability: The IT-DT framework emphasizes the importance of transparency and human involvement in understanding and defending against adversarial attacks. By incorporating explainability-based features and involving human intervention during the transformation process, the framework provides practitioners with insights into the reasons behind adversarial examples and allows them to intervene when necessary. In addition, it promotes a more transparent and interpretable approach to defence, enabling practitioners to gain deeper insights into the vulnerabilities and robustness of their models.

– (iii) Practical Defense Mechanism: The IT-DT framework offers a practical defence mechanism that can be implemented by practitioners working with transformer-based models. It provides a step-by-step process, including training an adversarial detector (TAD) and performing test-time detection and transformation (TDT). This practical guidance allows practitioners to effectively safeguard their models against adversarial attacks, enhancing the reliability and security of their NLP systems in real-world applications.

– (iv) Generalization and Real-World Evaluation: The IT-DT framework has been evaluated extensively on various state-of-the-art word substitution attacks, popular text datasets and transformer architectures. This evaluation demonstrates the generalization ability of the defence mechanism across different attack scenarios and datasets, highlighting its effectiveness in countering adversarial manipulations at various granularity levels. Furthermore, the real-world evaluation conducted in the study provides practitioners with empirical evidence of the framework's performance, instilling confidence in its applicability in real-world scenarios.

– (v) Practical Application in Industry: The IT-DT framework has implications for practitioners heavily relying on NLP applications, such as sentiment analysis, content filtering, and language understanding. By implementing the proposed defence mechanism, organizations can enhance the security and reliability of their transformer-based models, mitigating the risks associated

**Input: Input sequence** $X = [x_1, x_2, ..., x_N]$**, where** $x_i \in R^d$ **Output:**
**Probability distribution over the possible class labels** $P(y|x)$ **for**
$l = 1$ *to* $L$ **do**

    *Compute self-attention*
    $Z^{self} \leftarrow MultiHeadSelfAttention(X)$
    *Apply residual connection and layer normalization*
    $Z \leftarrow LayerNorm(X + Z^{self})$
    *Compute feed-forward sub-layer*
    $Z^{ffn} \leftarrow FeedForward(Z)$
    *Apply residual connection and layer normalization*
    $O \leftarrow LayerNorm(Z + Z^{ffn})$
    *Update input for next layer*
    $X \leftarrow O$
**end**
*Compute class probabilities*
$P(y|x) \leftarrow softmax(W^T O_N + b)$

**Algorithm 1:** Transformer-based text classifier

    with adversarial attacks and ensuring more reliable and robust NLP systems
    for their customers and users.

Overall, the IT-DT framework contributes to the research landscape by addressing the vulnerability of transformer-based models to adversarial attacks and providing a practical defense mechanism. It also offers transparency, interpretability, and generalization, making it a valuable resource for practitioners seeking to secure their NLP systems in real-world applications.

## 2    Background and Related Work

The following section aims to provide a comprehensive overview of the relevant background and related research in the field.

### 2.1    Preliminaries

**Transformer-based Text Classification Models** The transformer architecture was first introduced in the paper "Attention is All You Need" [48], and has since become a popular choice for Natural Language Processing (NLP) tasks. The Transformer-based text classifier $T_m$ is a machine learning algorithm that takes as input a sequence $X$ of word embeddings and outputs a probability distribution over the possible class labels $y$. The algorithm for $T_m$ is given by Algorithm 1. $T_m$ consists of multiple Transformer layers, each of which has a self-attention sub-layer and a feed-forward sub-layer. The $T_m$ takes an input sequence $X$ as input, where each element $x_i$ is a word embedding of dimension $d$. The then proceeds to iteratively process the input through the Transformer layers. Each layer first computes a self-attention sub-layer, which involves computing a set

of $H$ self-attention heads. The outputs from all $H$ heads are concatenated and passed through a linear projection to produce the final output for the sub-layer. The result is then passed through a residual connection and layer normalization operation before being processed by the feed-forward sub-layer. After each layer, the output is updated to be the input for the next layer. Once all the layers have been processed, the final output representation $O_N$ is passed through a linear projection and a softmax function to produce a probability distribution over the possible class labels. During training, the model is optimized to minimize the cross-entropy loss between the predicted class probabilities and the true class labels.

**Adversarial Examples (AEs)** Given a real example $x$ with label $y$, an adversarial example $AE$ against $x$ can be defined as an input $x' \leftarrow x + \imath$, where $\imath$ is a minor perturbation. This perturbation is constrained to satisfy a set of perturbation constraints $P_{\text{const}}$. Furthermore, the classifier $F \in T_m$ predicts an incorrect label for the adversarial example, i.e., $y' \leftarrow F(x')$ such that $y' \neq y$.

**Word Substitution Attacks (WSA)** Given an original text $E$ with a set of important words $W$, WSA results in a new text $E'$ where some of the words in $W$ have been replaced with similar or mis-spelt words $W'$. The goal is to produce an altered text $AE$ that modifies the output of an NLP model such that $F(E) \neq F(AE)$, where $AE = E_{\text{subst}}(W, W')$ and $E_{\text{subst}}$ represents the word substitution function.

**Prediction Confidence Score (PCS)** $\text{PCS}(x,y)$ of model $F$ depicts the likelihood of an input $x \in X$ having a label $y \in Y$. The smaller $\text{PCS}(x,y)$ suggests $F$ has low confidence that $x$ has a label $y$.

**Explainability** Explainability refers to the capability of comprehending and interpreting how a machine learning model reaches its decisions or forecasts. The importance of explainability lies in transparency, accountability, and trust in the decision-making process of machine learning models [7]. More formally, let $F$ be a machine learning model that takes an input $X$ and produces an output $y$. We can represent this relationship as $y = F(x)$, where $F$ is a complex function that maps the input space to the output space. To assess the explainability of $F$, we need to examine how it makes its predictions. One way to do this is to decompose $F$ into simpler functions that can be more easily understood. For example, we can use feature importance analysis to identify which input features have the strongest influence on the output. Another approach is to use model-agnostic methods such as LIME (Local Interpretable Model-Agnostic Explanations) [40] or SHAP (SHapley Additive exPlanations) [27] to generate local explanations for individual predictions. These methods identify which input features were most important for a particular prediction, and how changes to those features would affect the output.

**Pre-Trained Embedding** Pre-trained embeddings refer to the use of pre-existing word vectors that have been trained on large datasets using unsupervised machine learning techniques such as word2vec, GloVe, or fastText. These embeddings are trained on large corpora of text data, and the resulting word vectors capture semantic and syntactic relationships between words. By using pre-trained embeddings, it is possible to leverage the knowledge captured in these embeddings to improve the performance of natural language processing (NLP) tasks such as text classification, sentiment analysis, and machine translation, without the need to train a new embedding from scratch.

## 2.2    Related Work

In NLP, the research community has dedicated significant efforts to developing models that are robust against WSA. These efforts can be broadly categorized into two mainstreams, which we summarize in this section.

(a) **Robustness Enhancement defenses (RED)**: the first stream of research focus on developing resilient models with a goal to learn a robust model that can achieve high performance on both clean and adversarial test inputs [25]. These defenses include Adversarial Training [50], Synonym Encoding [53], Certifiable Robustness [51] However, these defences have four main limitations [20]: (i) they often require retraining the DNN with adversarial data augmented via known attack strategies, which can cause substantial retraining overhead and is computationally expensive. Implementing such defenses often demands a considerable amount of time and resources. (ii) defenses often have limited applicability. They may only be effective against specific types of attacks, making them less useful against more sophisticated and varied attacks. For example, certifiable defenses may be effective against simple word substitution attacks, but may not be effective against attacks that involve multiple word substitutions or more complex obfuscation methods. (iii) External information, such as knowledge graphs, is required to enhance the model's robustness. (iv) These defenses often rely on unrealistic assumptions, such as assuming that the attacker's embedding method is known, which is not always practical or feasible.

(b) **Detection, Blocking and Recovery Strategy (DBR)**: While robustness enhancement is a widely studied defense mechanism against word substitution attacks, there is another defense mechanism that involves identifying and recovering adversarial examples to prevent such attacks. However, this area of research is not as well-established compared to robustness enhancement defenses. One notable study in this category is DIScriminate Perturbation (DISP) [62], which utilizes a discriminator to detect suspicious tokens and a contextual embedding estimator to restore the original word. Frequency-Guided Word Substitution [34] is another approach that leverages the frequency properties of replaced words to recover the original example. More recently, Word-level Differential Reaction (WDR) by Mosca et al. [33] investigated the variation of logits in the target classifier when perturbing the input text and trained an adversarial detector. In another study, Yoo et al. [59] extracted features from the last out-

put layer of the model and adapted Robust Density Estimation (RDE) to detect adversarial samples.

### 2.3   Comparison

We present a comparison between our study and other relevant studies in Table 1. Each row represents a specific topic or feature addressed in the studies, while each column corresponds to a particular study. A checkmark ($\checkmark$) indicates that the corresponding study covers the topic, while a cross ($\times$) indicates the topic is not addressed in that study. Our defence framework, IT-DT, falls under the category of detection and correction mechanisms. It incorporates explainability features inspired by TextShield [43], but with notable distinctions. Unlike TextShield, which relies on computationally intensive saliency factors [9] and lacks the utilization of attention cues inherent in transformer architecture, IT-DT is specifically designed to tackle the challenges posed by large and complex transformer-based models. In terms of transformation, IT-DT focuses on finding optimal replacements rather than relying solely on rule-based substitutions, thereby enhancing its effectiveness in generating non-adversarial counterparts. Additionally, we comprehensively evaluate our identification and transformation models, setting our study apart from previous research. This thorough assessment enables us to determine the efficacy and performance of our approach. Transparency and interpretability are fundamental aspects of our framework. We go beyond by incorporating a human-in-the-loop component by integrating rich threat intelligence logs. This aspect facilitates the active engagement of security analysts, empowering them to make informed decisions based on the insights provided by our framework. Furthermore, compared to DISP [62], our approach eliminates the need to train a discriminator and generator model directly from the embedding. Instead, we utilize explainability distributions to generate interpretable justifications for detection, identification of perturbed words, and transformation of adversarial examples. Additionally, we extensively evaluate our approach to seven SOTA WSA. We have compared our detection mechanism with the SOTA WDR. However, we could not find the replication package for TextShield for a direct comparison. Additionally, we tried to compare our results with DISP, but we faced challenges in replicating their experiments due to its high memory requirement for transformer-based methods.

## 3   IT-DT Framework

### 3.1   Motivating Example

To illustrate the significance of the proposed IT-DT framework, consider a scenario where an organization employs a machine learning-based text classification system to filter incoming emails for potential cybersecurity threats. Attackers employ sophisticated techniques such as carefully crafted language patterns and evasion strategies to bypass the email filters. These adversarial emails, if undetected, can result in security breaches, data leaks, or system compromise.

**Table 1.** Comparison of our Study with other studies

| Topic Covered | DISP [62] | FGWS [34] | WDR [33] | RDE [59] | TextShield [43] | SHAP [27] | IT-DT |
|---|---|---|---|---|---|---|---|
| Detection | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Learning-based | ✓ | ✗ | ✓ | ✗ | ✓ | ✓ | ✓ |
| Utilize Explainability | ✗ | ✗ | ✗ | ✗ | ✓ | ✓ | ✓ |
| Provides Transformation | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✓ |
| Provides Interpretability | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ |
| Comprehensive Evaluation | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✓ |
| Provides Threat Intelligence | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ |



**Fig. 1.** Overview of IT-DT Framework

By integrating the IT-DT framework into the organization's text classification system, the capabilities of the email filters are significantly enhanced. The explainability component of IT-DT allows analysts to gain insights into the linguistic cues and syntactic structures that differentiate adversarial emails from legitimate ones. This understanding enables security professionals to adapt the classification model and update the filters to counter emerging attack patterns

effectively. Moreover, the identification stage of the IT-DT framework provides valuable information about the specific attack vectors employed, such as Word Substitution Attacks or Sentence Modification Attacks. Armed with this knowledge, security analysts can develop targeted defense mechanisms and refine their detection strategies to stay one step ahead of the attackers. The transformation capabilities of the IT-DT framework play a vital role in neutralizing adversarial emails. However, in cases where the transformation process proves challenging, the framework triggers human intervention. Security analysts, equipped with their expertise, assess and classify the difficult examples manually, ensuring accurate classification and preventing potential security risks associated with misclassifications. In summary, the IT-DT framework addresses the challenges posed by adversarial examples in text classification for cybersecurity applications. By incorporating explainability, precise identification, effective transformation, and human intervention when needed, IT-DT provides a comprehensive solution to enhance the resilience of text classification systems against adversarial attacks. This framework empowers analysts, facilitates vulnerability identification, strengthens the overall security posture, and fosters collaboration between automated systems and human experts in the cybersecurity landscape.

## 3.2   Overview of IT-DT Framework

Figure 1 shows the overview of our proposed framework. Our framework consists of two main phases: Training Adversarial Detector (TAD) and Test-Time Detection and Adversarial Example Transformation (TDT).

*TAD: Training Adversarial Detector* The TAD phase aims to construct a machine learning (ML) model that discriminates between clean and adversarial examples. To achieve this objective, TAD takes the original and adversarial datasets as inputs and transfers them to the Explainability and Frequency-based Features Extraction (EFFE) stage. In the EFFE stage, the original and adversarial datasets are utilized as inputs to generate explainability maps, specifically attention maps ($A$map) and integrated gradient ($I$grad) distributions, for each example. This process employs a surrogate of the transformer model ($ST_m$), which acts as a shadow copy of the original transformer model ($OT_m$). The surrogate enables the detection of adversarial examples without impeding the regular functioning of the original model. Statistical features are extracted from the attention maps and integrated gradient distributions, and these features are combined with frequency-based statistical features to construct a feature set for training the ML models. Subsequently, we employ several supervised traditional ML classifiers to obtain the Adversarial Detector ($D_{adv}$). These classifiers undergo fine-tuning using K-fold cross-validation, and the most effective models are selected for each ML model type. Finally, an unknown test dataset is employed to evaluate the selected models, and the ML model exhibiting the highest validation performance is designated as the Trained Adversarial Detector ($D_{adv}$).

**Fig. 2.** Example of Attention Map for $E$ from Yelp-polarity Dataset [Green and Red color depict importance of words with respect to correct and incorrect classification label $y_{pred}$ of $E$ respectively.]

*TDT: Test-Time Detection and AE Transformation* The objective of the TDT phase encompasses three key aspects: (i) Differentiating clean examples ($CE$) from adversarial examples ($AE$) using the test-time Adversarial Detector ($D_{\mathrm{adv}}$) when the model is deployed. (ii) Transforming $AE$ into an optimal non-adversarial variant ($TF_E$) that retains semantic similarity to avoid evasion attacks. (iii) Generating threat intelligence logs and identifying the need for human intervention in cases where the transformation of $AE$ to a non-adversarial form fails. To accomplish this, the EFFE step first converts test-time data into a feature set and employs the $D_{\mathrm{adv}}$ model to detect $AE$ and $CE$ examples. Next, examples identified as $CE$ are processed through the spelling-based transformation module and then passed to the $T_m$ model for label prediction. On the other hand, the $AE$ examples are directed to the identification step, where Perturbation Identification (PI) module determines the candidate perturbed words, denoted as $P_{\mathrm{cand}}$, within the $AE$. Then, these candidates are inputted into the Embedding-based Transformation selection step, where a set of $N$ replacement options ($S_{\mathrm{cand}}$) is generated using pre-trained embeddings. Subsequently, the $AE$ is reverse-engineered by selecting an optimal substitution for each $P_{\mathrm{cand}}$, utilizing the outputs from both the $T_m$ and $D_{\mathrm{adv}}$ models. The transformed $AE$, denoted as $TF_E$, is further processed by the spelling-based transformation module to correct typos. Finally, $TF_E$ is sent to the $T_m$ model for label detection. In addition to generating the transformed $AE$ ($TF_E$), our framework produces analytical logs ($A_{\mathrm{logs}}$) that provide information about the success or failure of the conversion process from $AE$ to $TF_E$. These logs offer valuable insights for security analysts to enhance the model's robustness and counter word-substitution attacks.

We discuss the details of each stage in the subsequent subsections.

### 3.3   Explainability and Frequency based Feature Extraction (EFFE)

EFFE is responsible for exacting features for building ML models. In this study, we have explored explainability to distinguish adversarial examples from clean. This stage has two sub-phase as described below:

**Input:** Pre-trained model $T_m$, training dataset $T_{data}$ with $k$ examples;
**Output:** $A_{\text{map}}$, $I_{\text{grad}}$;
Initialize $A_{\text{map}}$ and $I_{\text{grad}}$ to $I^{\text{nxn}}$;
**forall** $E$ *in* $T_d ata$ **do**
    $y_{\text{pred}} = T_m(E)$
    **forall** $l$ *in* $T_m.layers$ **do**
        Compute attention weights using the inputs query $q$ and key $k$
        vectors with $d$ embedding dimensions;
        $A_{\text{w}}^l \leftarrow \frac{softmax(q.k^T)}{\sqrt{d}}$
        Compute gradient of attention with respect to $y_{\text{pred}}$;
        $\nabla A_{\text{w}}^l = \frac{\partial y_{\text{pred}}}{\partial A_{\text{w}}^l}$
        Compute weighted averaged cross all the heads in $l$ range
        [0,max], clamping negative contributions;
        $A_{\text{map}} = A_{\text{map}} + \overline{(A_{\text{w}}^l \cdot \nabla A_{\text{w}}^l)}$
        $I_{\text{grad}} = I_{\text{grad}} + \overline{(\nabla A_{\text{w}}^l)}$
    **end**
**end**
**return** $A_{\text{map}}^{dxk}$, $I_{\text{grad}}^{dxk}$
**Algorithm 2:** Explainability and Frequency based Feature Extractor

***Generate Explainability Maps*** In this step, we compute the $A_{\text{maps}}$ and $I_{\text{grad}}$ distributions for all words $n$ in an input sequence $E$. The $A_{\text{maps}}$ distribution measures the importance of each word based on how much attention it receives from the other words in the input sequence $E$. It represents the distribution of attention weights across the input sequence for a given prediction $y_{pred}$. On the other hand, the $I_{\text{grad}}$ distribution estimates the significance of each input word based on how much it affects the output of the model. Its utility is to identify which words are most sensitive to changes in the input and which words are most important for making predictions. Overall, these distributions provide useful information for understanding how the Transformer model processes input sequences and how it generates its predictions.

Our method for computing $A_{\text{maps}}$ and $I_{\text{grad}}$ is inspired by the latest seminal work in computer vision [8]. In that work, the authors captured the contribution of each token (word/ sub-word) to the predicted label $y_{pred}$ to explain the prediction of the $T_m$ models, as shown in Figure 2. Figure 2 highlights the relevance of each words with to the predicted label $y_{pred}$ obtained by this method. We have chosen this explainability method because it is computationally efficient and does not require additional training or modification of the original transformer model. However, unlike this work, we have explored the $A_{\text{maps}}$ and $I_{\text{grad}}$ distributions across all words to differentiate original versus adversarial examples instead of explainability. Moreover, we have only considered self-attention-based transformer models due to their popularity in text classification tasks [47]. Nevertheless, our work can be extended to other transformer architectures. Algorithm

2 summarizes the process of generating explainability maps: Firstly, we initialize $A_{\mathrm{maps}}$ and $I_{\mathrm{grad}}$ with an identity matrix $Inxn$, where $n$ denotes the number of words in $E$. It signifies that initially, each word is important to itself (the basis of the self-attention model). After that, we obtain attention weights $A\mathrm{weights}^l$ across each layer of the $T_m$ model, which signify the relative importance of each word or sub-word in $E$ for each other word or sub-word in $E$. These attention weights $A_{\mathrm{w}}^l$ are computed by taking the dot product of a query vector $q$ with a set of key vectors $k$ that represent each word or sub-word in $E$, and then applying a softmax function to normalize the resulting scores. Here, the $q$ vectors represent the current location the model is attending to, whereas the $k$ vectors represent different positions or contexts in the input sequence that may be relevant for the current task. After that, the attention gradients $\nabla A_{\mathrm{w}}^l$ are quantified by taking the derivative of the output $y_{\mathrm{pred}}$ with respect to the attention weights $A_{\mathrm{w}}^l$. Finally, $A_{\mathrm{maps}}$ is determined using the Hadamard product of attention weights and attention gradients, which are then averaged across all the heads in $l$ after removing negative contributions. Similarly, $I_{\mathrm{grad}}$ is calculated by averaging $\nabla A_{\mathrm{w}}^l$ across all the heads in $l$, and negative values are clamped. The outputs of this module are $dxk$ dimensional, $A_{\mathrm{maps}}$ and $I_{\mathrm{grad}}$ for all examples $k$.

***Extract Statistical*** We utilize statistical measures [30] such as minimum, maximum, mean, median, mode, variance, skewness, kurtosis, entropy, mean of gradient, and the sum of peaks to extract information regarding the central tendency, variability, and shape for four distinct types of distribution. Let $E$ denote a text sample.

The first distribution is the Attention Map ($A_{map}$), representing the overall attention received by all words in the text sample $E$. The second distribution is the Integrated Gradient ($I_{grad}$) Distribution, which enables us to compare the relative importance of individual words across different predicted labels $y_{pred}$ and to capture the difference between expected and adversarial samples. The third set of features is the Overall Frequency distribution, which characterizes the frequency distribution of words in $E$ with respect to their frequency in the training dataset. The inclusion of these features is motivated by the results reported in [34]. Finally, we utilize the Outlier Frequency (OF) distribution, representing the frequency of outliers in $A_{map}$. To compute the OF distribution, we first identify potential outliers in $A_{map}$ using the inter-quartile range (IQR) method. Let $q_1$ and $q_3$ be the first and third quartiles of $A_{map}$, respectively, and let $iqr$ be the interquartile range, which is defined as $iqr = q_3 - q_1$. A token $t_i$ having attentional $A_{map}$ is considered a potential outlier if it satisfies the following condition:

$$w_i > q_3 + 1.5 \times iqr and w_i \notin StopWords \tag{1}$$

where $w_i$ is the value of $t_i$ in $A_{map}$. This condition checks whether a value $w_i$ is more than 1.5 times the IQR away from the third quartile. If a value satisfies this condition, it is considered a potential outlier. Once the potential outliers are identified, we obtain the frequency of each outlier in the set by counting the

number of occurrences of the value in $A_{map}$. This frequency distribution is the OF distribution, which represents the frequency of outliers in $A_{map}$.

### 3.4   Buliding ML Models

To effectively detect adversarial examples, we developed machine learning (ML) models using the following steps:

*Clean Feature Set* After extracting the statistical features, the Feature set ($F_{set}$) was cleaned to eliminate irrelevant or redundant features that may have negatively impacted the models' performance. Any NaN (Not a Number) or infinite values in the statistical features were removed to ensure consistency and completeness. Moreover, we removed any duplicate values that might have been present to improve feature quality. It ensures that each feature contributes unique information to the ML models, eliminating redundancy.

*Train Model* We trained multiple ML models on the cleaned $F_{set}$ to identify the best-performing model for detecting adversarial examples. To do so, we fine-tuned the models by adjusting their hyperparameters.

*Tune Model* In order to optimize the performance of our adversarial attack detection system, we optimized the hyperparameters using various commonly used performance measures for imbalanced dataset problems, such as the Matthews Correlation Coefficient (MCC). MCC is a popular performance metric for binary classification tasks, particularly for imbalanced datasets [42]. Additionally, to ensure the models' generalization ability and avoid overfitting, we used 10-fold cross-validation with a stratified split. This involved splitting the dataset into ten equal parts, training the models on nine parts, and testing their performance on the remaining part. This process was repeated ten times, with a different part used for testing each time. We identified the best-performing model for our dataset by fine-tuning the ML models using hyperparameter optimization and cross-validation, thereby improving the accuracy and effectiveness of our adversarial attack detection system. We selected the model with the best cross-validation performance as the optimal ML model $D_{adv}$ for detecting AEs.

*Test Model* Finally, we evaluated the effectiveness of the best optimal model in distinguishing AE from CE for a unseen dataset.

### 3.5   Explainability-based AE Detection (EAD)

After identifying the optimal adversarial detector $D_{adv}$ during the training phase, the next step is to test the detector during the test phase. During this phase, the EFFE module extracts features from a given example $X$, which are then input into $D_{adv}$ for classification as either an adversarial $AE$ or a clean $CE$ example. The $AE$ examples are then processed through the identification phase, which is responsible for detecting potential perturbed words, while the $CE$ examples

are directed to undergo the Spelling-based transformation step, which applies spelling-based transformations to the input to correct any spelling errors that may exist.

### 3.6   Identification

The purpose of this module is to identify potential perturbed words within the adversarial example that have influenced the model's output. This is accomplished through a combination of explainability-based and frequency-based identification techniques, as outlined below:

*Explainability-based Perturbation Identification (EPI)* : Let $E = w_1, w_2, \ldots, w_n$ be an input sequence, where $w_i$ denotes the $i$-th word in the sequence. Let $AE$ be an adversarial example detected by $D_{adv}$. $AE$ is fed into the EPI module, which computes explainability scores $A_{map}$ for each word $w_i$ in $AE$. The words $w_i$ with attention scores greater than a threshold *thres* are identified as the potential perturbed candidates ($P_{cand}$) that have the most influence on the incorrect decision of the model i.e. $y_{pred}$ in the case of adversarial examples. Formally, the set of influential perturbed words is defined as:

$P_{cand} = w_i \mid A_{map}(w_i) > \text{thres} and w_i \notin \text{Stopwords}$

Here, $A_{map}(w_i)$ denotes the explainability score of word $w_i$ (section 3.3). Additionally, the words in $P_{cand}$ are sorted in descending order according to attention score such that the word with the highest attention becomes the first candidate.

*Model-based Identification (MPI)* : The SPI approach uses the model score to determine the importance of words in the decision-making process. This method follows the BERT-Attack [24] strategy and replaces each word in the adversarial example ($AE$) with a '[MASK]' token. If the Probability Confidence Score (PCS) of the surrogate model ($ST_m$) on the predicted label $y_{pred}$ decreases more than a certain threshold, the word is considered important. More specifically, the set of important words $P_{cand}$ is defined as:

$$P_{cand} \cup (w_i \mid PCS(AE)[y_{pred}] - PCS(AE \setminus w_i)[y_{pred}] > sc_{thres}) \qquad (2)$$

where $P_{cand}$ is the set of potential perturbed candidates, $w_i$ is a word in the adversarial example $AE$, $y_{pred}$ is the predicted label, $PCS(AE)[y_{pred}]$ is the Per-Class Score of the adversarial example for the predicted label $y_{pred}$, $PCS(AE \setminus w_i)[y_{pred}]$ is the Per-Class Score of the adversarial example with the $i$-th word replaced by a '[MASK]' token for the predicted label $y_{pred}$, and $sc_{thres}$ is the threshold that determines the level of importance.

*Frequency-based Identification (FPI)* : The FPI module is designed to detect potential perturbed words in adversarial examples at the character and multi-level (word and character) levels by considering frequency as a critical element. To achieve this, we build upon the work of Mozes et al. [34], which demonstrated

that the frequency of perturbed words in an adversarial example is typically lower than the frequency of the original words in its non-perturbed counterpart. Specifically, our FPI module selects words whose frequency is less than a defined threshold $fq_{thres}$ and are not pronouns or nouns, and concatenates them with the set of potential perturbed candidates $P_{cand}$.

$$P_{cand} \cup (w_i \mid Freq(w_i) < fq_{thres}, w_i \notin pronouns) \qquad (3)$$

This represents the set of potential perturbed candidates ($P_{cand}$) selected by the Frequency-based Identification (FPI) module, which includes words whose frequency is less than a defined threshold $fq_{thres}$ and are not pronouns. The $\cup$ symbol denotes set union, and the curly braces  are used to define a set. By incorporating frequency-based identification into our overall approach, we can improve the accuracy and reliability of our adversarial example detection methods.

### 3.7   Transformation

*Embedding-based Transformation (ET) Generation* The objective of this module is to generate a set of potential candidate substitutions, referred to as $S_{cand}$, for each candidate word *cand* that exists within a given set of words $P_{cand}$. This is achieved through the utilization of three distinctive embedding techniques. Firstly, synonym embedding is employed, which utilizes the vast lexical database, WordNet, to identify synonyms for a given word $w_i$. WordNet's synsets group words into sets of cognitive synonyms, which are indicative of distinct concepts. Secondly, contextual embedding is utilized, which involves masking the target word and utilizing its surrounding context to predict a potential replacement. This technique takes into account contextual information from surrounding words, thereby capturing nuanced variations in meaning that are dependent on context. Lastly semantic pre-trained word embedding models such as Word2Vec or GloVe are used to identify semantically related words to the target word. The output of this process, $S_{cand}$, is generated as a key-value pair, where the key is the substitution $S_{w_i}$ for $w_i$, and the value is a similarity score $simi_{PCS}$, indicative of the degree of absolute value of similarity between $S_{w_i}$ and $w_i$.

*Model-based Transformation (MT) Selection* :
   Algorithm 3 presents the Model-based Transformation (MT) Selection module, which aims to transform an original adversarial example ($AE$) into an optimal transformed adversarial example ($TF_E$) and determine if human intervention is required. To achieve this goal, the algorithm takes $AE$ as input, along with other variables, and produces $TF_E$ and a boolean value $Human$ indicating if human verification is necessary.
   The algorithm first generates various transformed examples ($TF_{cand}$) by replacing each word in the original phrase $P_{cand}$ with all possible substitutions in $S_{cand}$. Then, for each transformed example, the module queries the $T_m$ model to

**Input:** $AE$, $P_{cand}$, $S_{cand}$, $org_{pcs}$, $org_{label}$, $adv_{flag}$, $adv_{org}$

**Output:** $TF_E$, $Human$, $HumanMsg$

**for** $try$ $in$ $Tries$ **do**

    **for** $cand$ $in$ $P_{cand}$ **do**

        $TF_E \leftarrow AE$

        Generate possible transformation of AE

        $TF_{cand}$ = replace $cand$ in $TF_E$ with all its substitutions $S_{cand}$

        Obtain replacement score using target model

        $r_{PCS}, r_{label} = T_m(TF_{cand})$

        $r_{score} = abs(r_{PCS}[orglabel] - orgpcs[orglabel])$

        Inquire $D_{adv}$ score

        $freq_{score} \leftarrow frequency(S_{cand})$

        $opt_{score} = r_{score} + simi_{score} + freq_{score}$

        $sel_{id} == argmax(opt_{score})$

        **if** $r_{score}[sel_{id}] > threshold$ $or$ $r_{label}[sel_{id}] \neq orglabel$ **then**

            select the optimal transformation of AE

            $TF_E \leftarrow TF_{sel_{id}}$

            $cand_{opt} = S_{cand}[sel_{id}]$

            $adv_{label}, adv_{PCS} = D_{adv}(TF_E)$

            **if** $adv_{label} = 0$ $and$ $r_{label}[sel_{id}] \neq orglabel$ **then**

                $adv_{flag} \leftarrow False$

                $Human \leftarrow False$

                $HumanMsg \leftarrow$ Converted to non-Adv

                **return** $TF_E$, Human,HumanMsg

            **end**

        **end**

        **if** $adv_{score} > adv_{org}$ **then**

            $Human \leftarrow True$

            $HumanMsg \leftarrow$ Got more Adv

            **return** $TF_E$, Human,HumanMsg

        **end**

        **if** $no$ $replacement$ $done$ **then**

            $Human \leftarrow True$

            $HumanMsg \leftarrow$ No optimal replacement was found

            **return** $TF_E$, Human,HumanMsg

        **end**

    **end**

**end**

**if** $adv_{\text{flag}} = True$ **then**

    $Human \leftarrow True$

    $HumanMsg \leftarrow$ Couldn't convert to non-Adv

    **return** AE, Human,HumanMsg

**end**

**Algorithm 3:** Model-based Transformation (MT) Selection

**Table 2.** Dataset details

| Dataset | Task | Classes | Median Length | No of Training Samples | No of Test Samples |
|---|---|---|---|---|---|
| IMDB [29] | Movie Review | 2 | 161 | 25k | 25k |
| AG-News [60] | Headline Topic | 4 | 44 | 120k | 7.6k |
| SST-2 [45] | Movie Review | 2 | 16 | 2.7k | 1821 |
| YELP [60] | Restaurant Review | 2 | 152 | 100k | 38k |

obtain its PCS and label $r_{PCS}$ and $r_{label}$, and computes the replacement score ($r_{score}$), which is the difference in PCS of the $T_m$ model before and after the replacement. The module also computes the frequency score of each substitution based on its frequency in the $T_m$ model's training dataset. Next, the algorithm computes the optimization score ($opt_{score}$) by adding the $r_{score}$, the similarity score between the original and transformed examples ($simi_{score}$), and the frequency score of the substitution in $S_{cand}$, and selects the candidate $cand_{opt}$ that has the maximum $opt_{score}$. Then, the module selects the transformed example with the highest optimization score ($sel_{id}$). If the replacement score of the selected example exceeds a certain threshold or if the label of the selected example is not the original label, the algorithm replaces $AE$ with the selected example.

Subsequently, the module uses the adversarial detector $D_{adv}$ to determine if $AE$ is detected as a clean example, i.e., if its label is zero and the label of the selected example is not the original label. If so, the module sets $adv_{flag}$, $Human$ to false and $humanmsg$ to converted to non-Adversarial (Adv), indicating that the MT module has converted the example into a non-adversarial version and, therefore, human verification is unnecessary. If the adversarial PCS of the selected example is more significant than $adv_{org}$, the module sets $Human$ to true, indicating that human verification is necessary as the substitution has made the example more adversarial. Finally, the algorithm repeats these steps for a fixed number of attempts, and if $adv_{flag}$ is still True after all attempts, the module sets $Human$ to true, and $HumanMsg$ indicates that human verification is necessary.

*Spelling-based Transformation (ST)* The Spelling-based Transformation (ST) phase is the final step in the Transformation module, where spelling correction is performed on every input example $E$ for words that are not pronouns or present in the training dataset of $T_m$. This process involves identifying typos involving homoglyphs, such as replacing 'l' with '1' in "Hel10", followed by utilizing the SymSpell library [17] to search for a closely resembling word within an IT-DT distance of $ed_d s$. If a suitable replacement is found, the typo is replaced with the correct word, resulting in the transformed example $TF_E$.

### 3.8  Final Label Prediction (FLP)

The FLP stage uses the transformed example $TF_E$ as input to $T_m$ to predict the final label $y'_{pred}$.

**Table 3.** Research Questions and their Motivation

| ID | Research Questions | Motivation |
|---|---|---|
| **RQ1** | To what extent does explainability discriminate between genuine and adversarial examples? | This research question aims to assess the effectiveness of explainability in differentiating between genuine and adversarial examples, which is critical to the success of the proposed framework. By answering this question, we can better understand the capability of explainability in mitigating the impact of adversarial examples. |
| **RQ2** | How effective are different machine learning classifiers in detecting adversarial examples? **RQ1.1)** How successful is selected adversarial detector $D_{adv}$ in detecting particular type of Word Substitution attack? | The first part of this RQ explores the overall performance of multiple classifiers in detecting adversarial examples generated through various attack methods on transformer models. It also determines the most effective classifier to detect adversarial examples. On the other hand, RQ 1.1) strives to analyze the performance of our chosen optimal adversarial detector on each attack. This research's results intend to contribute to developing more robust machine learning models and improving their security against adversarial attacks. |
| **RQ3** | How proficiently does the proposed TDT module operate as a whole? **RQ3.1)** How successful our identification module is to detect the perturbed words? **RQ3.2)** How effectively does the Transformation module convert adversarial instances into non-adversarial ones? **RQ3.3)** To what extent can the TDT module identify the necessity for human intervention? | This RQ assesses the TDT module's performance in accurately detecting and transforming adversarial examples to their original non-adversarial labels at test time. We further subdivide this research question into three parts: RQ3.1) aims to identify the overlap between identified Potential perturbed candidates versus actual perturbed words, RQ3.2) aims to evaluate the performance of the proposed Transformation module in converting detected adversarial examples into non-adversarial ones, whereas RQ3.3) assesses the ability of the proposed method to identify situations where human intervention is required to restore the original non-adversarial label of adversarial examples, which is critical to reducing false alarms for the security analyst while involving humans in the loop when necessary. |

## 4  Experimental Design and Setup

This section elaborates on the experimental setup we adopted for each phase of our proposed IT-DT framework.

Moreover, we undertook a comprehensive assessment of the performance of the IT-DT framework by addressing three principal research questions, as outlined in Table 3. The subsequent sections elaborate the experimental and evaluation setup used for each of the research question.

### 4.1  RQ1: Explainability Effectiveness

In order to evaluate the effectiveness of explainability in discriminating between legitimate and adversarial examples, we conducted a statistical analysis. Initially, we calculated the distributions of explainability measures $A_{map}$ and $I_{grad}$ for all instances in the training (Tr), testing (Te), and adversarial (Ad) datasets, as described in section 3.3. The testing dataset was utilized to avoid any potential bias resulting from the fact that the explainability measures may differ for unseen examples. Subsequently, we computed the logarithm of the standard deviation of these measures for each distribution, and utilized Bayesian hypothesis testing [41] to evaluate the distributions. Specifically, we calculated the Bayes factor $BF_{10}$, which measures the degree to which the data support the alternative hypothesis (H1) that the distributions of the genuine datasets (Tr, Te) for $A_{map}$ and $I_{grad}$ differ from those of the adversarial dataset (Ad) and Tr and Te have no such difference, over the null hypothesis (H0) the converse. A higher $BF_{10}$ value indicates stronger evidence in favor of H1. Additionally, Cohen's d effect sizes were calculated for all mean explainability measures $A_{map}$ and $I_{grad}$.

### 4.2  RQ2: Building and Evaluating Adversarial Detector

To train the adversarial detector, we implemented the following experimental setup.

*Baseline Transformer Models and Datasets* For our experimental setup, we chose two preeminent classifiers, namely BERT [13] and ROBERTA [26], which have been trained on four benchmark datasets: IMDB, AGNEWS, YELP, and SST2. These pre-trained models were procured from the HuggingFace platform [15]. BERT and ROBERTA have been widely used in various machine learning applications and have demonstrated exceptional accuracy in text classification [55].

Furthermore, these models are currently being adapted in the cybersecurity domain to detect a range of attacks, including phishing and spam [1,3,49]. Recent studies have shown that these models are effective in detecting phishing emails [22]. Cybersecurity researchers often use these models in the context of adversarial ML to evaluate their proposed adversarial attacks or defenses, given that these models exhibit high accuracy and are vulnerable to attacks, as evidenced by their high attack success rates [2,10,61]. Hence, the selection of these models is appropriate for evaluating the effectiveness of the proposed framework in the context of text classification and cybersecurity.

Lastly, the datasets we used in our experiments covered diverse topics and had varying lengths, thereby ensuring the generalizability of our approach. Table 2 provides more information about the datasets used in our experiments.

*Adversarial Datasets* This study explores Word Substitution attacks across three granularity levels: char-level, word-level, and multi-level. We employ the TextAttack Library [32] and utilize an adversarial dataset from a recent study [59] to generate adversarial examples for each model. The following attacks are employed in this study:

1. **Char-level attack:** For this attack, we employed the state-of-the-art Deep-WordBug attack [52]. *DeepWordBug (DWB)* is a character-level adversarial attack method that utilizes a generative model based on a deep recurrent neural network to generate perturbations at the character level of the input text. [32].

2. **Word-level attack:** We generate the adversarial examples at this level of granularity through four distinct word-level substitution attacks: (i) *TextFooler (TF)* [19]: a word-level adversarial attack that generates semantically similar but different sentences to deceive a text classification model. It uses semantic and syntactic techniques to generate adversarial examples. (ii) *Bert_Attack (BAE)* [24]: a word substitution attack that utilizes the BERT language model to generate adversarial samples. The attack aims to replace words in the input sentence with similar words that would not change the meaning of the sentence. BAE employs a combination of gradient-based word importance scoring and masking to identify words that are most likely to be changed without changing the sentence's meaning. (iii) *Probability Word Saliency Substitution (PWWS)* [39]: a word substitution attack method that selects the most salient words in the input text based on their importance for the model's prediction. PWWS then substitutes these words with semantically similar words to create an adversarial example that can fool the model while maintaining semantic meaning. (iv) *TextFooler Adjusted (TF_adj)* [59]: a variation of the TF algorithm that incorporates a more effective synonym replacement strategy. TF_adj replaces only the top-K most essential words in the input text by considering the probability of word saliency. It helps to generate more efficient and potent adversarial examples with fewer word substitutions, leading to better attack success rates.

3. **Multi-level attack:** For the multi-level attack, we employed the *TextBugger attack (TB)* [23]. TextBugger is a multi-level attack approach that can generate adversarial examples by modifying text at the character, word, and phrase levels. It utilizes strategies to generate adversarial examples, including synonym replacement, word deletion, and insertion. It is effective against a wide range of NLP models, including those based on RNN, CNN, and BERT architectures. Additionally, it has been shown to have a high attack success rate while maintaining a low distortion rate, making it a powerful tool for evaluating the robustness of NLP models. .

*Explainability Computation* LIME and SHAP are the most widely used methods that provide explainability. However, these methods can be computationally and memory-intensive for large and complex Transformer models such as BERT and RoBERTa. This is due to the fact that these models have a large number of parameters, necessitating a significant amount of memory and processing power to evaluate input data. Additionally, the large number of synthetic samples generated by LIME and SHAP to approximate the model's behavior and compute the feature importance scores can become computationally infeasible for BERT and RoBERTa models due to the size and complexity of the models and input data.

**Table 4.** Adversarial Dataset Details with number of adversarial samples

| Dataset | Model | DWB | TF | BAE | PWWS | TF-ADJ | TB | Total |
|---------|-------|-----|----|----|------|--------|----|-------|
| IMDB | **BERT** | 714 | 9164 | 5959 | 9082 | 1415 | 854 | 27188 |
| | **ROBERTA** | 514 | 10212 | 7289 | 10338 | 870 | 863 | 30086 |
| SST2 | **BERT** | 728 | 1636 | 1039 | 1491 | 80 | 471 | 5445 |
| | **ROBERTA** | 765 | 1650 | 1062 | 1501 | 61 | 409 | 5448 |
| AGNEWS | **BERT** | 576 | 5895 | 1034 | 4140 | 343 | 483 | 12471 |
| | **ROBERTA** | 541 | 6179 | 1190 | 4720 | 367 | 497 | 13494 |
| YELP | **BERT** | 697 | 4588 | 2613 | 4572 | 489 | 822 | 13781 |
| | **ROBERTA** | 1509 | 1734 | 2301 | 1655 | 106 | 1200 | 8505 |

This can lead to memory crashes or out-of-memory errors when processing large datasets. To address these challenges, researchers have proposed various optimizations and approximations for explainability methods that are more suited to large and complex models like BERT and RoBERTa. For instance, attention-based methods like Attention-based Integrated Gradients and Attention Rollout can be used to compute feature importance scores for BERT and RoBERTa models by focusing on the attention weights of the model instead of individual parameters. These methods can be faster and more memory-efficient than LIME and SHAP for BERT and RoBERTa models. Therefore, our calculation of explainability scores for an input example $x$ is motivated by Chefer.et.al's [8] explainability method for text classification, which is more appropriate for large and complex models like BERT and RoBERTa.

*Building Machine Learning (ML) models* To develop the Adversarial Detector, we utilized five popular and effective classifiers, namely XGB, LGBOOST, RandomForestTree (RF), Decision Tree, and Logistic Regression, as reported in [6]. Given the class imbalance present in our classification problem, we employed the Matthew Correction Coefficient (MCC) as the primary performance criterion for training and fine-tuning the classifiers. To ensure a balanced class distribution in each fold, we employed attack-based stratified sampling with 10-fold cross-validation, and early stopping criteria to identify the optimal parameters. MCC was used as the model selection criteria, as it considers all classes explicitly, and is resilient against imbalanced datasets [4, 11, 28]. MCC values range from -1 to 1, with values near -1, 0, and 1 indicating a poor model (misclassifying both classes), a random model (classifying both classes randomly), and a good model (classifying both classes correctly), respectively. We used 10-fold cross-validation and Bayesian Optimization to fine-tune the models. We evaluated the performance of the classifiers using various metrics, including Accuracy (ACC), Balanced Accuracy (Bal_ACC), F1 Score, Precision, Recall, Area Under the Curve (AUC), False Positive Rate (FPR), and False Negative Rate (FNR), as reported in [4, 11, 28]. These metrics provided a comprehensive analysis of the classifiers' performance. Furthermore, we presented the accuracy of the trained

**Table 5.** Accuracy of the Transformer Models on Clean Dataset

| Dataset | Model | $ACC_{clean}$ |
|---------|---------|-----------|
| AGNEWS | BERT | 0.965 |
| | ROBERTA | 0.960 |
| IMDB | BERT | 0.953 |
| | ROBERTA | 0.953 |
| SST2 | BERT | 0.873 |
| | ROBERTA | 0.893 |
| YELP | BERT | 0.990 |
| | ROBERTA | 0.989 |

classifiers in detecting specific attacks to provide insights into their strengths and limitations in various attack scenarios.

### 4.3   RQ3: TDT Module

To assess the effectiveness of the TDT module, a random sampling technique was employed to select 200 adversarial examples for each of the six attacks described in Section 4.2, along with 200 clean (no attack) examples. In addition to the known attacks, we also included unseen word-level attack, *Attack to Training (A2TY)* [58]. A2T is a recent attack technique that utilizes gradient-based search to identify essential words in the text and replaces them with similar words using word embeddings. This process ensures that the Part of Speech (POS) and semantic similarity are preserved by employing Distil-Bert similarity scores. It is important to note that the adversarial examples generated by this attacks were not utilized during the training and validation of the $D_{Adv}$ model. Instead, it represent zero-day attack that enable the assessment of the generalizability of our framework during testing. These adversarial examples were misclassified by the models, indicating zero accuracy for the $ST_m$ model whereas the performance of the model on clean examples in shown in Table 5. To avoid over-fitting bias, these examples were unseen by both the $ST_m$ and $D_{adv}$ models. Next, these examples were fed into the TDT stage for test-time evaluation. The TDT module's evaluation metric considers both the $D_{Adv}$ detection error and the accuracy achieved in correctly detecting the label, as well as human intervention. Thus, it evaluates the TDT module as a whole. The overall accuracy of the TDT was measured using the following equation 4 Here, $n$ represents the total number of examples, $GT_i$ represents the ground-truth label of example $i$, $TF_i$ represents the final predicted label of example $i$, and $H_i$ is a Boolean variable indicating whether human intervention is required or not. This metric measures the end-to-end performance of the TDT phase by considering both correctly classified examples and examples flagged for human intervention. In addition, the accuracy of the $D_{adv}$ modules on the randomly selected examples, denoted as $ACC_Det$, and the accuracy achieved by the TDT module in correctly identifying the label, denoted as $ACC$, are also provided to obtain more clarity of the results. Furthermore, to provide a more accurate picture of the strengths and weaknesses of

the framework, identify areas for improvement in each component of the system and assess the impact of each component on the overall performance of the system. We evaluate identification and transformation steps using the experimental setup below [6]

$$[tb!]Acc_{all} = \frac{\sum_{i=1}^{n}[(GT_i \neq TF_i \wedge H_i) \vee (GT_i = TF_i)]}{n} \tag{4}$$

*RQ3.1: Identification* To identify perturbations based on explainability, we set a threshold (*thres*) above the third quartile (q3) of the attention distribution. Any attention score above this threshold was considered an influential word for the model prediction. On the other hand, we set the $sc_{thres}$ to 0.3 and $fq\_thres$ to 0.001. A pilot study was conducted to determine these optimal thresholds, where we varied the q3 threshold from q3 to q3 plus the inter-quartile range (IQR), which is a measure of the spread of the attention scores. The IQR is the difference between the third quartile (q3) and the first quartile (q1) of the attention distribution. For $sc_{thres}$, we explored the range of 0.1 to 0.5, and for $fq_{thres}$, we set the threshold from 0.001 to 0.01. Based on our observations, we discovered that the q3, 0.3, and 0.001 were the optimal thresholds for all the datasets for EPI, SPI, and FPI modules, respectively. As a result, we utilized these thresholds consistently across all datasets and models.

**Evaluation Metrics** The performance of the identification module is evaluated by analyzing key statistical measures including Recall, Precision, F1-Score, False Positive Rate (FPR), and False Negative Rate (FNR). These measures provide a comprehensive evaluation of the algorithm's ability to accurately detect perturbed words while minimizing false positives and false negatives. Recall indicates the algorithm's capacity to capture perturbed words, while Precision calculates the proportion of correctly identified perturbed words out of all the words identified as perturbed. The F1 score combines Precision and Recall into a single metric, providing a balanced measure of the algorithm's performance in detecting perturbed words, particularly in imbalanced data scenarios. FPR measures the proportion of non-perturbed words that are incorrectly classified as perturbed, assessing the rate of false alarms or false positives. On the other hand, FNR measures the proportion of actual perturbed words that are incorrectly classified as non-perturbed, indicating the rate of missed detections or false negatives. These measures collectively provide an accurate evaluation of the algorithm's performance in identifying perturbed words. Also, note that here we have not considered MCC because instead of overall performance we wanted to more targeted assessment of the method's ability to identify perturbed words [46] and therefore the selected measures are appropriate choice.

*RQ3.2: Transformation* To implement the Embedding-based Transformation technique, we utilized multiple methods for word replacement. Specifically, for

---

[6] It is essential to note that the TDT module only processes examples that are detected as adversarial, except for the spelling-based correction sub-module. Therefore, we report the performance based on detected adversarial examples only.

synonym replacement, we leveraged the WordNet synset [31]. To avoid bias to-
wards recent evasion attacks, we refrained from using the synonym embedding
provided by counter-fitted [35] [25]. For contextual replacement, we used the
BERT Masking model [13]. BERT is a state-of-the-art neural network architec-
ture for natural language processing tasks, and it generates high-quality contex-
tualized word embeddings.

To generate contextually similar words, we used the BERT Masking tech-
nique, which involves masking a word in a sentence and predicting the masked
word based on the surrounding context. For semantic replacement, we utilized
the GloVe word2vec, a pre-trained word embedding model developed by Stanford
University. Specifically, we used the "glove-wiki-gigaword-50" model, trained on
a combination of Wikipedia 2014 and Gigaword 5 data, which has a vocabulary
size of 400,000 and consists of 6 billion tokens. Each word is represented as a
50-dimensional vector in the embedding space [38]. To ensure a wide range of
options for replacement, we selected 15 similar candidates from each embedding.
However, for model-based transformation, we set a threshold of 0.1. This means
that if the optimal candidate, represented as $cand_{opt}$, has a minimum change
of 0.1 in PCS or a change in the label, we replace the candidate. Otherwise,
we ignore the replacement. The BERT model's ability to generate contextually
similar words makes it a powerful tool for tasks that involve text transforma-
tion and data augmentation, and its effectiveness in natural language processing
tasks is widely recognized.

**Evaluation Metrics** To evaluate the performance of the transformation
module, we utilize, the accuracy of transformation $Acc_{tf}$ (equation 5) which
measures how effectively the transformation module converts adversarial exam-
ples to non-adversarial examples,

$$Acc_{tf} = \frac{N_{ct}}{N_{det}} \tag{5}$$

Here, the notation $N_{ct}$ refers to the number of examples that were correctly
transformed by the TDT module and $N_{det}$ is the number of examples detected
as adversarial by the detector module.

*RQ3.3: Human Intervention* We evaluate $Acc_{human}$ (equation 7) to measure
the transformation module's ability to flag examples for human intervention
when it fails to convert adversarial examples to non-adversarial labels (i.e.,
$Transform_{Error}$, equation 6). These metrics provide feedback to security ana-
lysts and improve the system's overall robustness.

$$Transform_{Error} = \frac{(N_{det} \wedge N_{in})}{N_{det}} \tag{6}$$

$$Acc_{human} = \frac{(N_{det} \wedge N_{in} \wedge HI)}{(N_{det} \wedge N_{in})} \tag{7}$$

Here, $HI$ is a Boolean variable that indicates whether the example requires
human intervention, $N_{in}$ represents the number of examples incorrectly classified

by the Final Label Prediction module (FLP), and $N_{det}$ is the number of examples detected as adversarial by the detector module.

***Generating Logs*** During the TDT phase, a critical step is to continuously gather insights from the Identification and Transformation phases. The purpose of this step is to collect various relevant information for security analysts. Specifically, we log the example ID, the score of the adversarial perturbation $D_{adv}$ for the given example, $P_{cand}$ and its corresponding word replacements, the transformed text, the ground truth label, the original prediction $y_{pred}$ of the model before transformation, the final prediction $y'_{pred}$ of the model after transformation, the confidence on the final prediction, the human intervention flag, the human intervention message and the detected adversarial flag. We use following human intervention messages . (i) "DETECTED AS ADVERSARIAL EXAMPLE, But NO REPLACEMENT DONE, Requires Human Intervention" where no replacement was found against $P_{cand}$ that transforms it into non-adversarial example. (ii) "Converted to non-ADVERSARIAL EXAMPLE", if our algorithm was successful in transforming detected adversarial example to non-adversarial. (iii) "GOT MORE ADVERSARIAL after substitutions Requires Human Intervention", if our algorithm transformed the example but adversarial detector have more confidence on it being adversarial. (iv) "Detected as non Adversarial", if the example is detected as non-adversarial example by the $D_{adv}$/ (v) "STILL ADVERSARIAL EXAMPLE After n tries, Requires Human Intervention", if the algorithm failed to convert a detected adversarial example to non-adversarial after n tries. We considered n=3 in our experiment.

We believe, this information is essential for analysts to understand the attack scenarios, assess the severity of the attack, and develop appropriate countermeasures. By gathering these insights throughout the TDT phase, analysts can gain a comprehensive understanding of the adversarial attack and improve the overall security of the model.

## 5   Evaluation Results and Analysis

The results of our investigation are presented in the subsequent subsections.

### 5.1   RQ1: Explainability Analysis

Table 6 presents the results of the comparison between the AGNEWS, IMDB, SST2, and YELP datasets and the BERT and ROBERTA models used in the study, with the effect size computed based on the Cohen's d value. A small effect size is between 0.2 and 0.5, a medium effect size is between 0.5 and 0.8, and a large effect size is greater than 0.8 [12]. Overall both Cohen's d, effect size and $BF_{10}$ values indicate that Tr and Te (legitimate) examples are different from Ad (Adversarial) example based on both $A_{map}$ and $I_{grad}$ distribution across all the datasets and models considered. This suggests that there is a strong evidence that the distributions of explainability scores are different for genuine and adversarial examples. Specifically, the Cohen's d effect size is high and $BF_{10}$ values

**Table 6.** Comparison of different models and datasets based on Explainability ($\log \sigma A_{map}$) and ($\log \sigma I_{grad}$) using Cohen's d, Effect size, Baye's Factor ($BF_{10}$)[Here Tr refers to training dataset distribution, Te refers to testing dataset distribution and Ad refers to Adversarial dataset distribution ]

| Dataset | Model | Distributions | Attention Map ($\log \sigma A_{map}$) | | | Integrated Gradients ($\log \sigma I_{grad}$) | | |
|---|---|---|---|---|---|---|---|---|
| | | | **d** | **Effect Size** | $BF_{10}$ | **d** | **Effect Size** | $BF_{10}$ |
| AGNEWS | BERT | Tr-Te | 0.038 | low | 1.338 | 0.129 | low | 1.23E+15 |
| | | Tr-Ad | 1.962 | high | inf | 3.571 | high | inf |
| | | Te-Ad | 2.329 | high | inf | 2.179 | high | inf |
| | ROBERTA | Tr-Te | 0.088 | low | 2.242E+06 | 0.116 | low | 6.957E+10 |
| | | Tr-Ad | 1.985 | high | inf | 2.319 | high | inf |
| | | Te-Ad | 2.112 | high | inf | 1.852 | high | inf |
| IMDB | BERT | Tr-Te | 0.078 | low | 3.113E+13 | 0.147 | low | 2.08E+53 |
| | | Tr-Ad | 1.610 | high | inf | 3.835 | high | inf |
| | | Te-Ad | 1.523 | high | inf | 3.269 | high | inf |
| | ROBERTA | Tr-Te | 0.040 | low | 169.115 | 0.071 | low | 1.107E+11 |
| | | Tr-Ad | 1.618 | high | inf | 2.579 | high | inf |
| | | Te-Ad | 1.564 | high | inf | 2.426 | high | inf |
| SST2 | BERT | Tr-Te | 0.105 | low | 122.059 | 0.224 | low | 4.201E+11 |
| | | Tr-Ad | 1.399 | high | inf | 1.910 | high | inf |
| | | Te-Ad | 1.385 | high | inf | 1.497 | high | inf |
| | ROBERTA | Tr-Te | 0.061 | low | 0.564 | 0.148 | low | 1.013E+05 |
| | | Tr-Ad | 1.282 | high | inf | 1.884 | high | inf |
| | | Te-Ad | 1.327 | high | inf | 1.581 | high | inf |
| YELP | BERT | Tr-Te | 0.033 | low | 2871.167 | 0.071 | low | 1.10E+22 |
| | | Tr-Ad | 2.293 | high | inf | 3.692 | high | inf |
| | | Te-Ad | 2.341 | high | inf | 3.276 | high | inf |
| | ROBERTA | Tr-Te | 0.0514 | low | 21120000000 | 0.0297 | low | 100.298 |
| | | Tr-Ad | 2.053 | high | inf | 4.445 | high | inf |
| | | Te-Ad | 2.165 | high | inf | 4.142 | high | inf |

are infinte for all the datasets and models. A $BF_{10}$ value of 1 indicates anecdotal evidence in favor of the alternative hypothesis, while values greater than 3 indicate substantial evidence, values greater than 10 indicate strong evidence, and values greater than 30 indicate very strong evidence. Thus, the $BF_{10}$ values suggest that there is strong evidence in favor of the alternative hypothesis that the distributions of explainability scores are different for genuine and adversarial examples. In addition, when comparing the distributions of explainability scores for genuine examples between training and testing datasets, no significant difference was found. The Cohen's d effect sizes for this comparison ranged from 0.033 to 0.105, indicating a low effect size for all models and datasets. However, the $BF_{10}$ values, while still high, were lower than the values obtained for the previous comparison between adversarial and genuine examples. In fact, for AG-NEWS BERT and SST2 ROBERTA models, the $BF_{10}$ values were less than 3, specifically 1.338 and 0.564, respectively, which indicates no significant difference between training and testing distributions.

We believe that the lack of difference between the training and testing datasets is due to the fact that both datasets contain genuine examples. Additionally, we observed that the Cohen's d effect sizes for $A_{map}$ and $I_{grad}$ differed for each distribution pair, indicating that these measures contribute differently to distinguishing between adversarial and legitimate examples.

**Summary RQ1**

The analysis suggests that the distributions of both $A_{map}$ and $I_{grad}$ can be utilized for distinguishing between adversarial and legitimate examples.

**Table 7.** Comparison of various classifiers in detecting WSA

| Dataset | T_m | Classifier | MCC | ACC | BalACC | Precision | Recall | F1Score | AUC | FPR | FNR |
|---------|-----|-----------|-----|-----|--------|-----------|--------|---------|-----|-----|-----|
| AGNEWS | BERT | **XGB** | **0.844** | **0.975** | **0.905** | **0.940** | **0.905** | **0.921** | **0.991** | **0.009** | **0.095** |
| | | **LGBM** | **0.844** | **0.975** | **0.902** | **0.942** | **0.902** | **0.921** | **0.991** | **0.009** | **0.098** |
| | | RFT | 0.813 | 0.971 | 0.879 | 0.937 | 0.879 | 0.905 | 0.987 | 0.009 | 0.121 |
| | | DT | 0.743 | 0.958 | 0.875 | 0.868 | 0.875 | 0.871 | 0.875 | 0.024 | 0.125 |
| | | LR | 0.706 | 0.938 | 0.919 | 0.797 | 0.919 | 0.844 | 0.977 | 0.057 | 0.081 |
| | ROBERTA | XGB | 0.776 | 0.954 | 0.862 | 0.917 | 0.862 | 0.886 | 0.978 | 0.016 | 0.138 |
| | | **LGBM** | **0.783** | **0.955** | **0.864** | **0.921** | **0.864** | **0.889** | **0.980** | **0.015** | **0.136** |
| | | RFT | 0.752 | 0.950 | 0.843 | 0.913 | 0.843 | 0.873 | 0.973 | 0.015 | 0.157 |
| | | DT | 0.651 | 0.924 | 0.828 | 0.823 | 0.828 | 0.826 | 0.828 | 0.045 | 0.172 |
| | | LR | 0.646 | 0.901 | 0.883 | 0.772 | 0.883 | 0.812 | 0.957 | 0.093 | 0.117 |
| IMDB | BERT | XGB | 0.896 | 0.951 | 0.952 | 0.945 | 0.952 | 0.948 | 0.984 | 0.050 | 0.048 |
| | | **LGBM** | **0.898** | **0.952** | **0.953** | **0.945** | **0.953** | **0.949** | **0.984** | **0.050** | **0.047** |
| | | RFT | 0.897 | 0.952 | 0.953 | 0.945 | 0.953 | 0.948 | 0.984 | 0.051 | 0.047 |
| | | DT | 0.826 | 0.919 | 0.913 | 0.913 | 0.913 | 0.913 | 0.913 | 0.063 | 0.087 |
| | | LR | 0.883 | 0.945 | 0.947 | 0.936 | 0.947 | 0.941 | 0.980 | 0.062 | 0.053 |
| | ROBERTA | XGB | 0.853 | 0.929 | 0.931 | 0.922 | 0.931 | 0.926 | 0.977 | 0.078 | 0.069 |
| | | **LGBM** | **0.858** | **0.931** | **0.934** | **0.925** | **0.934** | **0.928** | **0.978** | **0.076** | **0.066** |
| | | RFT | 0.848 | 0.926 | 0.929 | 0.919 | 0.929 | 0.923 | 0.975 | 0.083 | 0.071 |
| | | DT | 0.761 | 0.887 | 0.880 | 0.881 | 0.880 | 0.881 | 0.880 | 0.091 | 0.120 |
| | | LR | 0.820 | 0.912 | 0.916 | 0.904 | 0.916 | 0.909 | 0.965 | 0.101 | 0.084 |
| SST2 | BERT | XGB | 0.890 | 0.949 | 0.945 | 0.945 | 0.945 | 0.945 | 0.987 | 0.040 | 0.055 |
| | | **LGBM** | **0.894** | **0.951** | **0.947** | **0.947** | **0.947** | **0.947** | **0.988** | **0.039** | **0.053** |
| | | RFT | 0.875 | 0.942 | 0.938 | 0.937 | 0.938 | 0.937 | 0.984 | 0.047 | 0.062 |
| | | DT | 0.831 | 0.921 | 0.919 | 0.912 | 0.919 | 0.915 | 0.962 | 0.074 | 0.081 |
| | | LR | 0.835 | 0.922 | 0.924 | 0.911 | 0.924 | 0.916 | 0.975 | 0.082 | 0.076 |
| | ROBERTA | XGB | 0.880 | 0.943 | 0.940 | 0.940 | 0.940 | 0.940 | 0.985 | 0.046 | 0.060 |
| | | **LGBM** | **0.883** | **0.945** | **0.941** | **0.942** | **0.941** | **0.941** | **0.985** | **0.043** | **0.059** |
| | | RFT | 0.864 | 0.936 | 0.932 | 0.932 | 0.932 | 0.932 | 0.981 | 0.052 | 0.068 |
| | | DT | 0.814 | 0.912 | 0.908 | 0.906 | 0.908 | 0.907 | 0.908 | 0.075 | 0.092 |
| | | LR | 0.806 | 0.907 | 0.907 | 0.899 | 0.907 | 0.903 | 0.970 | 0.093 | 0.093 |
| YELP | BERT | XGB | 0.866 | 0.972 | 0.934 | 0.932 | 0.934 | 0.933 | 0.991 | 0.016 | 0.066 |
| | | **LGBM** | **0.870** | **0.973** | **0.936** | **0.934** | **0.936** | **0.935** | **0.991** | **0.016** | **0.064** |
| | | RFT | 0.857 | 0.970 | 0.932 | 0.924 | 0.932 | 0.928 | 0.990 | 0.019 | 0.068 |
| | | DT | 0.810 | 0.959 | 0.918 | 0.893 | 0.918 | 0.905 | 0.977 | 0.028 | 0.082 |
| | | LR | 0.772 | 0.940 | 0.947 | 0.834 | 0.947 | 0.878 | 0.982 | 0.062 | 0.053 |
| | ROBERTA | XGB | 0.928 | 0.988 | 0.964 | 0.964 | 0.964 | 0.964 | 0.997 | 0.006 | 0.036 |
| | | **LGBM** | **0.929** | **0.989** | **0.964** | **0.964** | **0.964** | **0.964** | **0.997** | **0.006** | **0.036** |
| | | RFT | 0.925 | 0.988 | 0.962 | 0.963 | 0.962 | 0.963 | 0.996 | 0.007 | 0.038 |
| | | DT | 0.884 | 0.981 | 0.954 | 0.930 | 0.954 | 0.942 | 0.954 | 0.014 | 0.046 |
| | | LR | 0.778 | 0.955 | 0.952 | 0.835 | 0.952 | 0.882 | 0.989 | 0.044 | 0.048 |

## 5.2   RQ2: Performance Evaluation of Adversarial Detector

Table 7 presents the overall performance of multiple classifiers in detecting adversarial examples, addressing our research question. The results indicate that
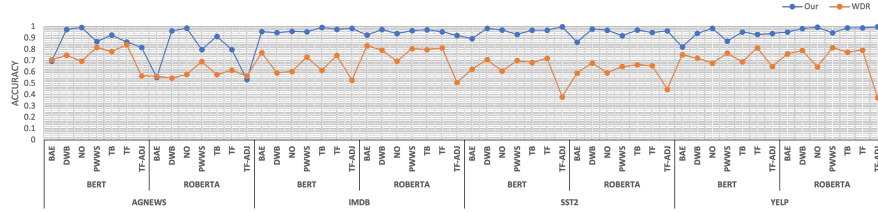
**Fig. 3.** Performance comparison of Our detector $D_{adv}$ versus WDR [33] Accuracy on Various Word Substitution Attacks: NO (No Attack), TB (TextBugger), TF (TextFooler), DWB (DeepWordBug), PWWS (Probability Word Saliency Substitution), and TF-ADJ (TF-Adjusted)

the LGBM and XGB models consistently outperform other classifiers, achieving an average MCC of 0.87, F1 score greater than 93.5%, and a balanced accuracy of 93.0%. Notably, the LGBM model achieved impressive results, with the highest MCC value of 0.929 for ROBERTA models trained on the YELP dataset. Furthermore, for the LGBM classifier, the False Negative Rate (FNR) ranged from 3.6% to 13.6%, while the False Positive Rate (FPR) ranged from 0.6% to 5%. On the other hand, the Logistic Regression (LR) and Decision Tree (DT) classifiers demonstrated weak performance in most cases, with an average MCC of 0.78 and 0.789, respectively. The Random Forest Tree (RFT) showed good performance (average MCC of 0.854) on some datasets, but it was outperformed by XGB and LGBM classifiers in most cases. Therefore, the ensemble classifiers (XGB, LGBM, RFT) are better for detecting adversarial examples. Additionally, there is no clear winner regarding the transformer models used, as BERT and ROBERTA perform similarly across the datasets. However, the classifiers performed comparatively better for the BERT model than ROBERTA, attaining an average MCC of 0.843 and 0.822, respectively. Moreover, it is crucial to highlight that accuracy (ACC) often overestimates performance, especially for imbalanced datasets like ours. For example, the LGBM model achieved an ACC of 95.5% for the AGNEWS ROBERTA model, while the balanced accuracy was only 86.4%. Thus, ACC alone is not a reliable measure to demonstrate the overall performance of the models in the detection task. Based on these findings, we have selected the LGBM classifier as our adversarial detector $D_{Adv}$.

*Analysis of $D_{Adv}$ with attack type (RQ1.1)* Fig 3 visually compares our proposed method and the WDR technique on the test dataset. Our method outperforms WDR in detecting adversarial attacks in most cases (53 out of 56 combinations of dataset, model, and attack). However, there are a few instances where our method shows slightly lower performance, specifically in BAE attacks against AGNEWS (BERT and ROBERTA) and TF-ADJ attacks against AGNEWS (ROBERTA), with an average accuracy reduction of 0.017 compared to WDR. Our method achieves high accuracy (¿90%) in detecting all attacks except BAE. Additionally, it maintains a fidelity of 97.6% in correctly identifying
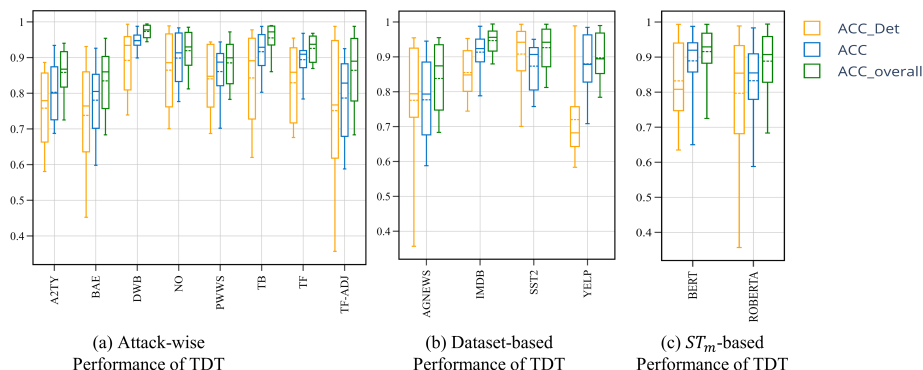
(a) Attack-wise
Performance of TDT

(b) Dataset-based
Performance of TDT

(c) $ST_m$-based
Performance of TDT

**Fig. 4.** Analysis of TDT module based on various Word Substitution Attacks ( refer to section 4.3 for specifics; the dotted line represents the mean and the solid line represents the median of the distribution.), Datasets and Transformer models

clean examples. In comparison, WDR performs best with an average accuracy of 75% for TF and PWWS attacks, but our method surpasses it with an average accuracy of 91%. When considering different datasets, our method demonstrates highly accurate detection for models trained on IMDB, YELP, and SST2 datasets, with an average accuracy of 95%. However, for AGNEWS, our method achieves a lower average accuracy of 83.5%. This discrepancy may be due to AGNEWS being a multi-class problem, where the distribution of adversarial examples can overlap with some class data. Regarding target models, our method performs equally well, with an average accuracy of 93.2% for BERT and 91.77% for ROBERTA across all datasets and attacks. Interestingly, the same trend is observed for WDR, with 68.82% for BERT and 66.55% for ROBERTA.

> **Summary RQ1**
>
> The LGBM classifier performed superior in detecting WSA, achieving an impressive average accuracy of 97%. Consequently, the LGBM model was selected as the adversarial detector ($D_{Adv}$) for the TDT phase. Moreover, our detector consistently outperformed WDR across all attack types.

### 5.3   RQ3: Overall Performance of TDT

The analysis of Fig 4(a) reveals that the TDT module achieved a median accuracy of over 86% for all attacks, indicating its efficacy in mitigating attacks and retaining the accuracy of clean examples. Furthermore, the TDT module acquired the highest median accuracy on DWB and TB attacks, with 97.7% and 97.2% accuracy, respectively, followed by TF, No attack, and PWWS, with accuracy levels of 93.7%, 92.9%, and 90.0%, respectively. These results indicate that

the TDT module effectively mitigates various attacks while maintaining high accuracy levels for clean examples. Furthermore, the TDT module demonstrated moderate accuracy in defending against 88%, 87% and 86% of TF-ADJ, A2TY and BAE attacks, respectively. Particularly noteworthy is the effectiveness of the TDT module against the A2TY attack, which represents a zero-day attack in this evaluation. The module achieved an 87% accuracy, indicating its capability to not only prevent known attack types but also generalize well to unknown attacks.. Moreover, we also observation that the overall accuracy of the TDT module is directly proportional to the adversarial detection accuracy, as the $D_{adv}$ was able to detect more than 80% of the examples correctly. For instance, the method successfully identified a median of 93.5% of DWB examples and accurately identified 88.5% of clean examples. In terms of accurately identifying the final prediction label of the example, the TDT module achieved a median accuracy of more than 80% for all attacks, including the no-attack scenario. In addition, it demonstrated the highest accuracy levels of 94%, 92.9%, 91.3%, and 90.9% on DWB, TB, NO, and TF attacks, respectively. These findings suggest that the TDT module is a reliable and effective method for defending against a wide range of attacks while ensuring high accuracy levels for clean examples.

Furthermore, an analysis of Fig 4(b) and (c) presents a performance evaluation of the TDT module on the considered datasets and transformer models. The TDT module attained a high overall median accuracy of 89.26% on all the datasets. Notably, the YELP dataset demonstrated a performance of 87.83% in correctly identifying the accurate label for the examples, despite only 68.20% of examples being correctly detected by the $D_{adv}$. Additionally, the TDT module exhibited the highest overall accuracy of 95.60% on IMDB, with a $ACC_{D}et$ of 84.82% and correct labelling of 92.37% examples by the FLP module. However, the AGNEWS dataset's overall accuracy was relatively poor, achieving only 87.41% $ACC_{overall}$, which may be attributed to its multi-class label. Regarding the surrogate transformer models, BERT and Roberta models both attained an approximately equal overall median accuracy of 92.91% and 90.73%, respectively. Interestingly, the median detection accuracy of attacks on the Roberta model was 6% higher than the BERT model. In comparison, the accuracy of obtaining a correct label by the FLP module was 6.47% higher for BERT than for Roberta. These results suggest that while attacks on the Roberta model are more easily detectable by the $D_{adv}$ model than the BERT model (4.62%), transforming the example to a non-adversarial example is more difficult for the Roberta model than the BERT model.

**RQ3.1: Performance of Identification Module** Fig 5 (a-c) displays a boxplot of the performance evaluation results of our identification module across different attack methods, datasets and $ST_m$ model (section 4.3). Figure 5 (a) presents an analysis of our algorithm's performance in identifying perturbed words under different attack methods. Our algorithm exhibits varying levels of effectiveness across these attacks.
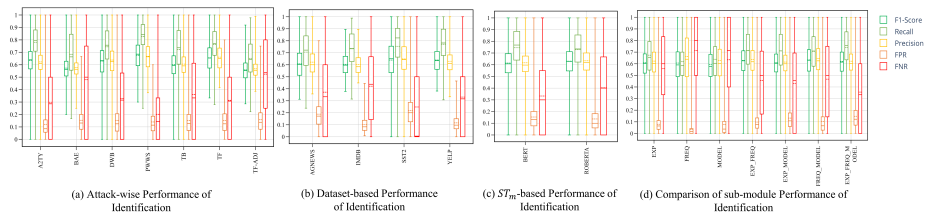
(a) Attack-wise Performance of Identification
(b) Dataset-based Performance of Identification
(c) $ST_m$-based Performance of Identification
(d) Comparison of sub-module Performance of Identification

**Fig. 5.** Performance of Identification module on various Word Substitution Attacks ( refer to section 4.3 for specifics; the dotted line represents the mean and the solid line represents the median of the distribution), datasets and transformer models. It also provide a comparison of the performance of different sub-module (Exp, Freq, Model refer to Explainability, Frequency and Model-based Perturbation Identification respectively

Notably, our algorithm demonstrates remarkable proficiency in detecting perturbed words in the PWWS, TF, A2TY, DWB, and TB attacks, as indicated by their higher median F1-scores (69%, 65%, 63.98%, 62%, 59.44% respectively) and median recalls (84%, 77%, 79.67%, 73.66%, and 76% respectively). The higher recall rates suggest that the method successfully identifies the perturbed words in over 70% of cases for these attacks. Furthermore, the algorithm exhibits lower median False Positive Rates (FPR) and False Negative Rates (FNR) below 14% and 31% respectively, indicating its ability to minimize both types of errors in these attack scenarios. Conversely, our algorithm did not perform as effectively in the BAE and TF-ADJ attacks, as it exhibited higher FNR rates exceeding 50%. These results highlight opportunities for improvement in these specific attack scenarios. Overall, our findings underscore the effectiveness of our algorithm in identifying perturbed words, particularly in the PWWS, TF, and DWB attacks. However, they also suggest the need for further enhancements in the BAE and TF-ADJ scenarios.

In our assessment of the performance of our algorithm in identifying perturbed words across multiple datasets, several essential deductions can be drawn from the results presented in Fig 5 (b). Firstly, it is apparent that the algorithm displays its most robust performance in the SST2, Yelp, and AGNEWS datasets, where it achieves a median F1-score and FNR of greater than 60% and less than 33%, respectively, indicating its effectiveness in accurately identifying perturbed words in these contexts. However, there is room for improvement in the IMDB dataset, where the algorithm exhibits higher rates of false negatives (a median of 43.75%), suggesting some missed detections. Notably, the algorithm maintains a relatively low FPR in all datasets, particularly for IMDB and Yelp datasets, where it is 8.76% and 9.90% , respectively, indicating its ability to minimize incorrect identifications. These findings underscore the algorithm's potential and offer insights for further refinement to improve its detection capabilities by reducing the FNR to ensure the precise identification of perturbed words across diverse datasets.

Fig 5 (c) presents the performance of our algorithm in detecting perturbed words generated against the BERT and RoBERTa models. The algorithm demonstrates commendable recall values for both models, with BERT achieving a median recall of 76.81% and RoBERTa achieving a median recall of 73.50%. It indicates the algorithm's capability to capture a significant portion of the perturbed words. Furthermore, the precision values are reasonable, with BERT achieving a median precision of 59.75% and ROBERTA achieving 61.84%. These results suggest that a substantial proportion of the words identified as perturbed by the algorithm are indeed correct detections. A key area for improvement lies in reducing the FNR for ROBERTA, which stands at a median value of 40.00%. Enhancements targeting the reduction of FNR would enhance the algorithm's ability to detect all relevant perturbed words for ROBERTA.

Further, to conduct a thorough analysis, we deconstructed the identification module into its sub-components (features) and evaluated their performance, as depicted in Fig 5 (d). presents the performance results of various feature combinations in detecting perturbed words within adversarial examples. Among the evaluated feature combinations, the "EXP_FREQ_MODEL" combination stood out as the most effective in terms of recall, achieving a value of 75.59%. This indicates that the "EXP_FREQ_MODEL" combination has a strong ability to correctly identify a significant portion of the perturbed words present in the adversarial examples.

To further analyze the performance of our individual features in comparison to the "EXP_FREQ_MODEL" combination, we can observe the following: EXP: The "EXP" feature exhibited an F1-score of 60.16% and a recall of 65.56%. While these values are relatively close to those of the "EXP_FREQ_MODEL" combination, it is evident that the combination of additional features leads to improved performance. On the other hand, the frequency "FREQ" feature demonstrated an F1-score of 59.01%, a recall of 58.33%. Comparing these values to the "EXP_FREQ_MODEL" combination, it becomes apparent that incorporating other features, such as "MODEL," contributes to better performance. The "MODEL" feature combination achieved an F1-score of 57.76% and a recall of 60.33%. These values are lower than those of the "EXP_FREQ_MODEL" combination, indicating that the inclusion of frequency-based and experimental features enhances the overall performance. The analysis of individual features suggests that the combination of experimental, frequency-based, and model-based features in the "EXP_FREQ_MODEL" combination leads to improved detection of perturbed words. The inclusion of these additional features enables a more comprehensive analysis, which ultimately results in a higher recall rate.

However, it is important to note that the "EXP_FREQ_MODEL" combination still exhibits room for improvement. The FPR of 12.22% and FNR of 33.33% indicate the presence of false positives and false negatives, respectively. This highlights the need for further research and refinement to minimize these errors and enhance the overall performance of the feature combination.

In conclusion, the "EXP_FREQ_MODEL" combination outperformed individual features in terms of recall, indicating its effectiveness in identifying per-
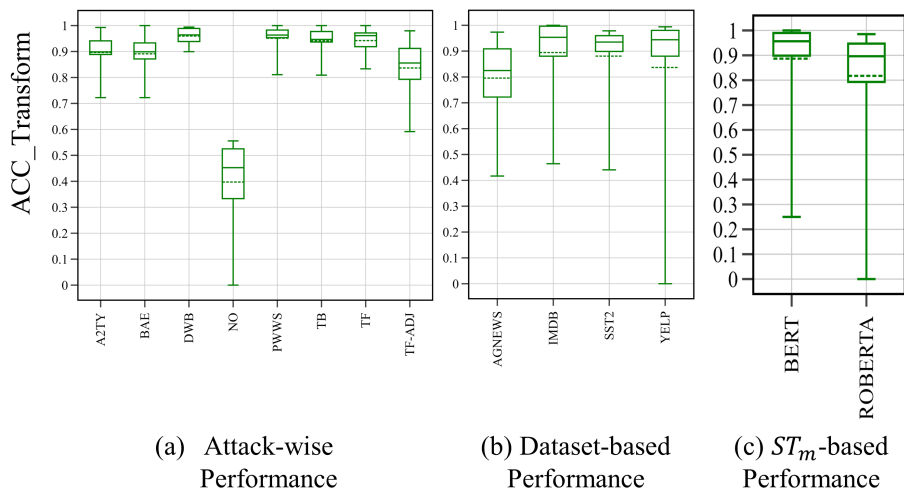
(a) Attack-wise Performance

(b) Dataset-based Performance

(c) $ST_m$-based Performance

**Fig. 6.** Performance of transformation module on various Word Substitution Attacks ( refer to section 4.3 for specifics; the dotted line represents the mean and the solid line represents the median of the distribution.), datasets and models

turbed words within adversarial examples. However, continued research efforts should focus on optimizing the feature combination to reduce false positives and false negatives, thereby further improving its performance in practical applications.

*Our identification module demonstrated superior effectiveness in detecting perturbed words in PWWS, TF, and DWB attacks, with higher median F1-scores and recalls. The algorithm showcased robust performance in the SST2, Yelp, and AGNEWS datasets, accurately identifying perturbed words with median F1-scores above 60% and FNR below 33%. Furthermore, the algorithm achieved commendable recall values for BERT and RoBERTa models, indicating its ability to capture a significant portion of perturbed words while maintaining reasonable precision.*

**RQ3.2: Performance of Transformation Module** Fig 6 (a) presents the box-plot representation of the performance evaluation of our transformation module using the ACC_Transform metric, which offers profound insights into the efficacy of our transformation module in converting adversarial examples to non-adversarial labels across different attack types. Among the assessed attacks, the DWB attack exhibited the highest median ACC_Transform score of 96.42%. This remarkable achievement underscores the exceptional ability of our transformation module to successfully convert a substantial proportion of DWB adversarial examples into their benign counterparts, thereby effectively neutralizing the adversarial perturbations. In contrast, the NO attack, characterized by examples falsely identified as adversarial (false positives), displayed a compara-

tively lower median ACC_Transform of 45.25%. This observation indicates that our transformation module adeptly preserved the inherent non-adversarial characteristics of these instances, aligning with its intended objective. The attainment of a lower median ACC_Transform for the NO attack is highly desirable as it signifies the transformation module's astute identification of these examples as non-adversarial and its prudent avoidance of unnecessary label alterations. Overall, our transformation module demonstrated commendable performance across the evaluated attacks, with median ACC_Transform scores ranging from 85% to 96.33% for all attacks. These outcomes substantiate the remarkable efficacy of our transformation module in mitigating the adversarial nature of the attacks on average, thereby reinforcing its practical value as a potent defense mechanism.

Additionally, the performance analysis of the transformation module was conducted based on the datasets and transformer models, as illustrated in Fig 6 (b-c). Regarding the datasets on which the transformer models were trained, distinct accuracy were obtained. The AGNEWS dataset demonstrated a median accuracy of 82.49%, indicating the substantial effectiveness of the transformation module in mitigating adversarial influences on models trained on the AGNEWS dataset, albeit with a slightly lower $ACC_{Transform}$ compared to other datasets. For the IMDB, SST2, and YELP datasets, notably high median ACC_Transform scores of 95.35%, 93.51%, and 94.41% were achieved, respectively. These exceptional results underscore the robust performance of the transformation module in successfully restoring the non-adversarial nature of instances from these datasets, thereby affirming its efficacy in countering adversarial attacks and preserving the integrity of the original labels. Furthermore, in terms of the adversarial examples targeted at the models, our method achieved an impressive ACC_Transform score of 95.66% for the BERT module. This highlights the remarkable capability of the transformation module in effectively neutralizing the adversarial perturbations introduced to BERT-based models, thereby accurately restoring the integrity of the original non-adversarial labels. Similarly, for the ROBERTA model, a slightly lower but still commendable ACC_Transform score of 89.62% was observed. This underscores the efficacy of the transformation module in mitigating the adversarial impact on ROBERTA-based models, facilitating the accurate recovery of the underlying non-adversarial semantics with notable success.

*The evaluation results demonstrate the transformation module's impressive performance (median 94.65%) in converting adversarial to non-adversarial, fortifying the robustness of models.*

**RQ3.3: Performance of Human Intervention** Lastly, we analyse the performance of our TDT module in detecting situations where human intervention is required. We present the results with transformation error and $ACC_{human}$ (see section 4.3) as shown by boxplot illustration in Fig 7. Lastly, we analyse the performance of our TDT module in detecting situations where human intervention is required. We present the results with transformation error and $ACC_{human}$ (see section 4.3) as shown by boxplot illustration in Fig 7. Regarding the attacks (Fig 7(a)) , we observed that the TDT module achieved a median $ACC_{Human}$
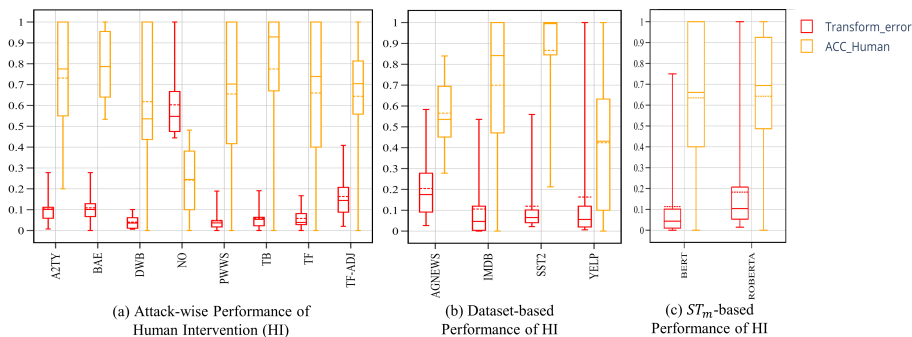
**Fig. 7.** Performance of Human Intervention on various Word Substitution Attacks ( refer to section 4.3 for specifics; the dotted line represents the mean and the solid line represents the median of the distribution.), datasets and models

of 78.61% for the BAE attack, indicating its effectiveness in flagging cases that need human intervention. The DWB attack showed a slightly lower $ACC_{Human}$ of 53.57%, suggesting a moderate level of accuracy in identifying examples requiring human intervention. Notably, the NO attack exhibited a significantly lower $ACC_{Human}$ of 24.49%, indicating that the TDT module correctly identified a large portion of non-adversarial examples without the need for human intervention. This outcome aligns with the objective of minimizing false positives for adversarial examples. For the PWWS, TB, TF, and TF-ADJ attacks, the TDT module demonstrated reasonable $ACC_{Human}$ values of 70.24%, 92.86%, 73.88%, and 70.43%, respectively. In terms of $Transform_{Error}$, higher values indicate a lower rate of false positives for adversarial examples. The results revealed that the TDT module exhibited a relatively low $Transform_{Error}$ of 3.58% for the DWB attack, indicating its effectiveness in correctly converting adversarial examples to non-adversarial ones. Similarly, the PWWS, TB, and TF attacks showed $Transform_{Error}$ values of 3.67%, 5.35%, and 3.88%, respectively, indicating the module's ability to successfully transform adversarial examples. The NO attack had a significantly higher $Transform_{Error}$ of 54.75%, suggesting that the TDT module correctly identified a majority of non-adversarial examples without unnecessary human intervention. Analyzing the results across datasets (Fig 7 (b)), we found that the TDT module achieved diverse $ACC_{Human}$ values. For the AGNEWS dataset, the module attained an $ACC_{Human}$ of 53.57%, suggesting its ability to identify cases requiring human intervention in this context. The IMDB dataset exhibited a higher $ACC_{Human}$ of 84.17%, indicating the module's effectiveness in flagging examples that need manual inspection and transformation. Remarkably, the SST2 dataset achieved a perfect $ACC_{Human}$ of 100.00%, indicating that the TDT module accurately identified all cases that required human intervention. The YELP dataset showed a relatively lower $ACC_{Human}$ of 43.08%, suggesting the need for improvement in identifying cases requiring human intervention. Across datasets, the AGNEWS dataset exhibited the highest

Transform_Error of 17.51%, indicating a relatively higher rate of false positives for adversarial examples. The IMDB, SST2, and YELP datasets demonstrated $Transform_{Error}$ values of 4.65%, 5.48%, and 5.59%, respectively, indicating the module's ability to successfully transform adversarial examples while minimizing false positives. Considering the models (Fig %reffig:TDT-perf-human (c)) , both BERT and ROBERTA achieved similar $ACC_{Human}$ values of 66.11% and 67.33%, respectively. These results suggest that the TDT module performed comparably in flagging examples for human intervention for both models. For the BERT and ROBERTA models, the TDT module achieved $Transform_{Error}$ values of 4.34% and 10.03%, respectively. These results suggest that the module performed slightly better in converting adversarial examples for the BERT model compared to the ROBERTA model. Our TDT module generate logs and threat intelligence report to help the analyst understand the behaviour of our framework and also examine the logs to identify the vulnerabilities in the underlying target model. A sample of threat intelligence report is shown in **??**.

*In summary, the TDT module exhibited satisfactory performance in identifying cases requiring human intervention and effectively converting adversarial examples to non-adversarial ones. The results varied across attacks, datasets, and models, indicating the module's adaptability and effectiveness in different scenarios. These findings highlight the TDT module's value as a robust defense mechanism against adversarial attacks and its potential for deployment in real-world applications.*

---

**Overall Summary RQ3**

The TDT module exhibits formidable defensive capabilities in countering adversarial attacks, effectively preserving high levels of accuracy while thwarting malicious attempts. In addition, it demonstrates exceptional proficiency in mitigating DWB and TB attacks, showcasing its effectiveness in neutralizing sophisticated adversarial perturbations. Moreover, its commendable performance transcends multiple datasets, with noteworthy efficacy observed, particularly in the challenging IMDB dataset. These salient features position the TDT module as an indispensable tool for researchers and industry practitioners, providing a reliable and robust solution to safeguard against adversarial threats.

---

## 6    Discussion and Threats to Validity

In this section, we discuss the overall insights gathered by conducting this study.

*Resiliency to Adaptive Attacks.* The proactive and transformative nature of the IT-DT framework poses a significant challenge for adaptive attacks, particularly score-based attacks that heavily rely on the model's feedback during the attack process. By actively transforming adversarial examples and disrupting the attacker's feedback loop, the framework hinders the attacker's ability to
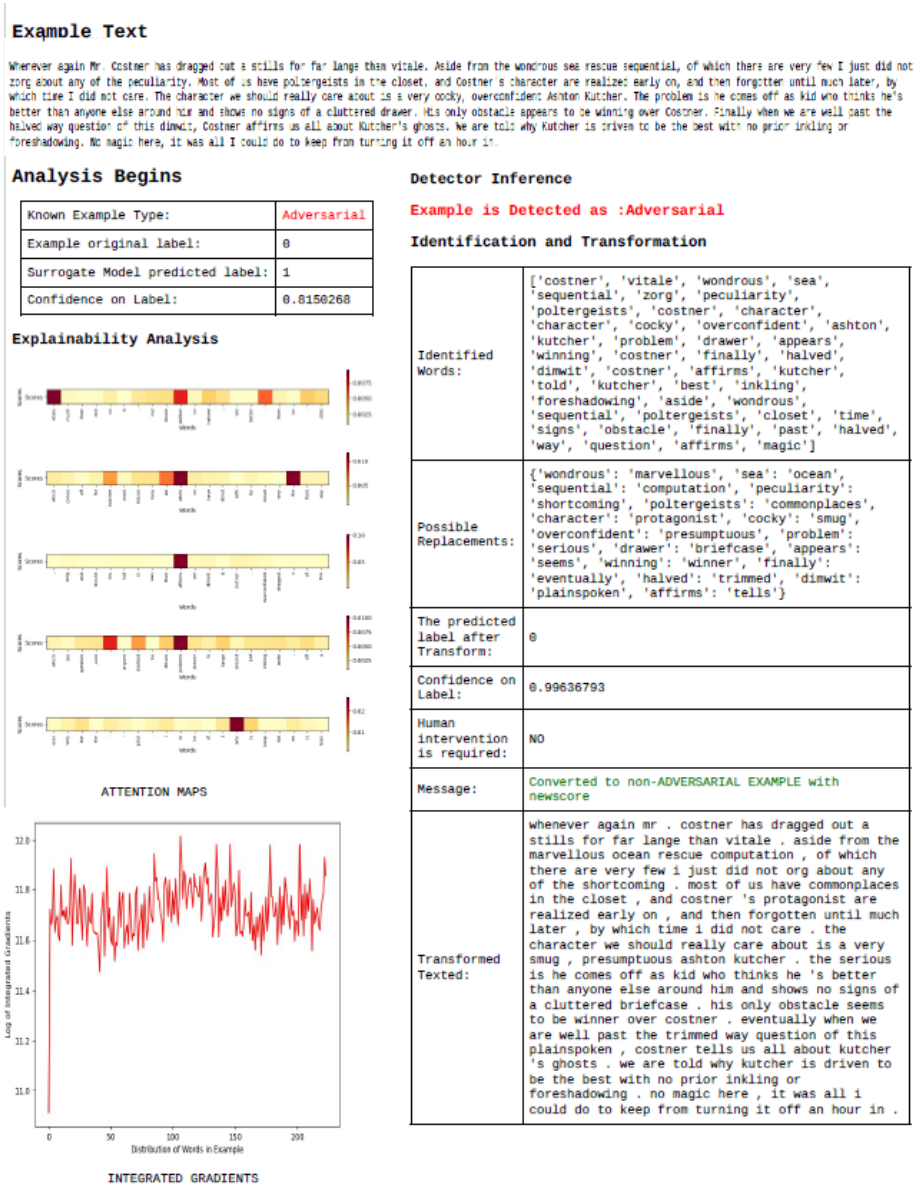
## Example Text

Whenever again Mr. Costner has dragged out a stills for far lange than vitale. Aside from the wondrous sea rescue sequential, of which there are very few I just did not zorg about any of the peculiarity. Most of us have poltergeists in the closet, and Costner's character are realized early on, and then forgotten until much later, by which time I did not care. The character we should really care about is a very cocky, overconfident Ashton Kutcher. The problem is he comes off as kid who thinks he's better than anyone else around him and shows no signs of a cluttered drawer. His only obstacle appears to be winning over Costner. Finally when we are well past the halved way question of this dimwit, Costner affirms us all about Kutcher's ghosts. We are told why Kutcher is driven to be the best with no prior inkling or foreshadowing. No magic here, it was all I could do to keep from turning it off an hour in.

## Analysis Begins

| Known Example Type: | Adversarial |
|---|---|
| Example original label: | 0 |
| Surrogate Model predicted label: | 1 |
| Confidence on Label: | 0.8150268 |

### Explainability Analysis



ATTENTION MAPS



INTEGRATED GRADIENTS

## Detector Inference

**Example is Detected as :Adversarial**

**Identification and Transformation**

| | |
|---|---|
| Identified Words: | ['costner', 'vitale', 'wondrous', 'sea', 'sequential', 'zorg', 'peculiarity', 'poltergeists', 'costner', 'character', 'character', 'cocky', 'overconfident', 'ashton', 'kutcher', 'problem', 'drawer', 'appears', 'winning', 'costner', 'finally', 'halved', 'dimwit', 'costner', 'affirms', 'kutcher', 'told', 'kutcher', 'best', 'inkling', 'foreshadowing', 'aside', 'wondrous', 'sequential', 'poltergeists', 'closet', 'time', 'signs', 'obstacle', 'finally', 'past', 'halved', 'way', 'question', 'affirms', 'magic'] |
| Possible Replacements: | {'wondrous': 'marvellous', 'sea': 'ocean', 'sequential': 'computation', 'peculiarity': 'shortcoming', 'poltergeists': 'commonplaces', 'character': 'protagonist', 'cocky': 'smug', 'overconfident': 'presumptuous', 'problem': 'serious', 'drawer': 'briefcase', 'appears': 'seems', 'winning': 'winner', 'finally': 'eventually', 'halved': 'trimmed', 'dimwit': 'plainspoken', 'affirms': 'tells'} |
| The predicted label after Transform: | 0 |
| Confidence on Label: | 0.99636793 |
| Human intervention is required: | NO |
| Message: | Converted to non-ADVERSARIAL EXAMPLE with newscore |
| Transformed Texted: | whenever again mr . costner has dragged out a stills for far lange than vitale . aside from the marvellous ocean rescue computation , of which there are very few i just did not org about any of the shortcoming . most of us have commonplaces in the closet , and costner 's protagonist are realized early on , and then forgotten until much later , by which time i did not care . the character we should really care about is a very smug , presumptuous ashton kutcher . the serious is he comes off as kid who thinks he 's better than anyone else around him and shows no signs of a cluttered briefcase . his only obstacle seems to be winner over costner . eventually when we are well past the trimmed way question of this plainspoken , costner tells us all about kutcher 's ghosts . we are told why kutcher is driven to be the best with no prior inkling or foreshadowing . no magic here , it was all i could do to keep from turning it off an hour in . |

**Fig. 8.** Threat Intelligence Report

accurately assess the model's vulnerabilities or refine their attack strategy. The transformed examples yield different feedback or prediction outcomes, making it difficult for adversaries to adapt and circumvent the system's defenses. This proactive approach increases the complexity and cost of launching successful attacks, enhancing the overall resiliency and security of the defended system.

*Human-in-the-Loop.* The IT-DT framework takes a step towards a human-in-the-loop approach by generating alerts for human intervention and leveraging security analysts' expertise to enhance robustness. Through providing logs and interpretability, analysts can perform case analysis, identify attack patterns, and propose countermeasures, enabling the framework to adapt to new attack types. This collaborative process fosters trust, improves interpretability, and facilitates continuous monitoring and evaluation. Involving security analysts in the feedback loop is crucial to enhance the framework further. Establishing a seamless workflow for analysts to review alerted examples, provide feedback, and refine detection and transformation mechanisms will bolster the framework's robustness.

*Scalability and Efficiency.* The IT-DT framework features a fast and efficient detection phase (<1 sec for an example), enabling quick identification of adversarial examples. However, the transformation process, which involves replacing perturbed words, can be time-consuming. The complexity of the transformation mechanism is O(nxd), where n represents the number of words to replace and d denotes the possible replacements for each word. While parallelization techniques can accelerate the transformation using multiple GPUs, the availability of computational resources can limit the extent of parallelization. Despite the time constraints, the transformation step is crucial for converting adversarial examples into non-adversarial ones and improving the model's robustness. Future research can focus on optimizing the transformation process and exploring efficient parallelization methods to reduce the time required.

### 6.1   Strengths

Our framework demonstrates the capability to detect and transform adversarial examples (AEs) into non-adversarial variants, while also incorporating human intervention when necessary, as demonstrated in the previous section 4. Unlike a simplistic approach of discarding AEs upon detection [59], our framework recognizes the importance of maintaining a positive user experience. It actively transforms AEs and provides meaningful responses to users, only resorting to blocking examples that require human intervention. This user-centric approach ensures that users receive appropriate feedback, minimizing frustration and confusion. Furthermore, our framework generates comprehensive logs and analytical reports, offering valuable insights into the characteristics of AEs, the effectiveness of the transformations, and potential vulnerabilities. These logs enable security analysts to gain a deeper understanding of adversarial attacks, refine defense strategies, and enhance the overall security posture of the system. Moreover, by actively transforming AEs and providing appropriate responses, our framework acts as a deterrent against adversaries. This proactive approach makes it more challenging for adversaries to adapt and circumvent the system's defenses. For instance, in countering score-based attacks that heavily rely on the model's feedback during the attack process to craft more effective AEs [32], our framework may poses challenges by actively transforming AEs and disrupting the attacker's feedback loop, it hinders the successful execution of score-based attacks. The

transformed examples yield different feedback or prediction outcomes, making it difficult for attackers to accurately assess the model's vulnerabilities or refine their attack strategy based on the model's responses. Lastly, the IT-DT framework is designed to address adversarial examples in real-world scenarios across various domains. Its effectiveness extends beyond theoretical settings, making it suitable for deployment in practical applications where robust defenses against adversarial attacks are crucial.

## 6.2   Threat to Validity

The limitations and threats to the validity of our approach are as follows:

Firstly, our approach focuses solely on defending transformer-based models against adversarial attacks. In the future, we aim to extend this approach to other model architectures such as convolutional neural networks (CNNs) or long short-term memory (LSTM) models. By incorporating explainability methods that do not rely solely on attention mechanisms, we aim to generalize our framework and enhance its applicability across different model types. Additionally, our results indicate that certain attacks, such as TF-ADJ and BAE, pose challenges for the detection of adversarial examples. Moreover, certain datasets, such as AGNEWS, may require improvements in order to achieve better detection performance. To address these limitations, we intend to explore additional features and techniques that can enhance the detection capabilities of our framework and improve its robustness against a wider range of adversarial attacks. Furthermore, our experiments were conducted solely on offline adversarial examples. In the future, we plan to deploy our framework in a client-server setup and evaluate its performance in real-time scenarios, particularly for query-based attacks. By simulating real-time attack scenarios, we can assess the effectiveness and efficiency of our framework in detecting and transforming adversarial examples in practical, dynamic environments. Lastly, although our framework generates security logs, their utility and effectiveness in assisting security analysts have not been thoroughly evaluated. Future work will involve conducting utility evaluations and gathering feedback from security analysts to assess the practical value of the generated logs and converted adversarial examples to further enhance the framework's usability and effectiveness in real-world settings.

Addressing these limitations and exploring future directions will contribute to the continued development and improvement of our framework, enabling it to provide robust defense against adversarial attacks across various model architectures, datasets, real-time scenarios, and practical security analysis.

## 7   Conclusion

In conclusion, this study introduces the IT-DT framework, which prevents transformer-based models from evading adversarial examples. Our approach involves training an adversarial detector, utilizing explainability and frequency-based features to distinguish between clean and adversarial examples. We then propose the

Test-Time Detection and AE Transformation (TDT) method, which deploys the adversarial detector at test-time to identify and transform detected adversarial examples into non-adversarial variants. Our experimental results demonstrate the effectiveness of our approach, achieving a median detection MCC of 0.846 and an overall median accuracy of 92.98% for TDT. Furthermore, our framework successfully identifies perturbations with a median recall of 74%, transforms adversarial examples with a median accuracy of 94.65%, and accurately involves human intervention with a median accuracy of 70.43%. Moving forward, we aim to extend our framework to non-transformer-based models, enhance the utility of generated logs by involving security analysts, and evaluate the framework's performance in real-time scenarios. These future endeavours will further strengthen the capabilities of our framework and broaden its applicability in defending against adversarial attacks.

# References

1. Akbar, K.A., Halim, S.M., Hu, Y., Singhal, A., Khan, L., Thuraisingham, B.: Knowledge mining in cybersecurity: From attack to defense. In: Data and Applications Security and Privacy XXXVI: 36th Annual IFIP WG 11.3 Conference, DBSec 2022, Newark, NJ, USA, July 18–20, 2022, Proceedings. pp. 110–122. Springer (2022)
2. Azizi, A., Tahmid, I.A., Waheed, A., Mangaokar, N., Pu, J., Javed, M., Reddy, C.K., Viswanath, B.: T-miner: A generative approach to defend against trojan attacks on dnn-based text classification. arXiv preprint arXiv:2103.04264 (2021)
3. Bayer, M., Kuehn, P., Shanehsaz, R., Reuter, C.: Cysecbert: A domain-adapted language model for the cybersecurity domain. arXiv preprint arXiv:2212.02974 (2022)
4. Bergstra, J., Yamins, D., Cox, D.: Making a science of model search: Hyperparameter optimization in hundreds of dimensions for vision architectures. In: International conference on machine learning. pp. 115–123. PMLR (2013)
5. Biggio, B., Roli, F.: Wild patterns: Ten years after the rise of adversarial machine learning. In: Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security. p. 2154–2156. CCS '18, Association for Computing Machinery, New York, NY, USA (2018). https://doi.org/10.1145/3243734.3264418, https://doi.org/10.1145/3243734.3264418
6. Bonaccorso, G.: Machine Learning Algorithms: Popular algorithms for data science and machine learning. Packt Publishing Ltd (2018)
7. Burkart, N., Huber, M.F.: A survey on the explainability of supervised machine learning. Journal of Artificial Intelligence Research **70**, 245–317 (2021)
8. Chefer, H., Gur, S., Wolf, L.: Generic attention-model explainability for interpreting bi-modal and encoder-decoder transformers. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 397–406 (2021)
9. Chefer, H., Gur, S., Wolf, L.: Transformer interpretability beyond attention visualization. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 782–791 (2021)
10. Chen, X., Salem, A., Chen, D., Backes, M., Ma, S., Shen, Q., Wu, Z., Zhang, Y.: Badnl: Backdoor attacks against nlp models with semantic-preserving improvements. In: Annual Computer Security Applications Conference. pp. 554–569 (2021)

11. Chicco, D., Jurman, G.: The advantages of the matthews correlation coefficient (mcc) over f1 score and accuracy in binary classification evaluation. BMC genomics **21**(1), 6 (2020)

12. Cohen, J.: Statistical power analysis for the behavioral sciences. Academic press (2013)

13. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805 (2018)

14. Ebrahimi, J., Rao, A., Lowd, D., Dou, D.: HotFlip: White-box adversarial examples for text classification. In: Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers). pp. 31–36. Association for Computational Linguistics, Melbourne, Australia (Jul 2018). https://doi.org/10.18653/v1/P18-2006, `https://www.aclweb.org/anthology/P18-2006`

15. Face, H.: Transformers: State-of-the-art natural language processing. `https://huggingface.co` (accessed 2023)

16. Gao, J., Lanchantin, J., Soffa, M.L., Qi, Y.: Black-box generation of adversarial text sequences to evade deep learning classifiers. In: 2018 IEEE Security and Privacy Workshops (SPW). pp. 50–56. IEEE (2018)

17. Garbe, W.: wolfgarbe/symspell: Symspell: 1 million times faster through symmetric delete spelling correction algorithm, `https://github.com/wolfgarbe/SymSpell`

18. Huber, L., Kühn, M.A., Mosca, E., Groh, G.: Detecting word-level adversarial text attacks via shapley additive explanations. In: Proceedings of the 7th Workshop on Representation Learning for NLP. pp. 156–166 (2022)

19. Jin, D., Jin, Z., Zhou, J.T., Szolovits, P.: Is bert really robust? a strong baseline for natural language attack on text classification and entailment. In: Proceedings of the AAAI conference on artificial intelligence. vol. 34, pp. 8018–8025 (2020)

20. Le, T., Park, N., Lee, D.: Shield: Defending textual neural networks against multiple black-box adversarial attacks with stochastic multi-expert patcher. In: Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). pp. 6661–6674 (2022)

21. Le, T., Wang, S., Lee, D.: Malcom: Generating malicious comments to attack neural fake news detection models. In: 2020 IEEE International Conference on Data Mining (ICDM). pp. 282–291. IEEE (2020)

22. Lee, J., Tang, F., Ye, P., Abbasi, F., Hay, P., Divakaran, D.M.: D-fence: A flexible, efficient, and comprehensive phishing email detection system. In: 2021 IEEE European Symposium on Security and Privacy (EuroS&P). pp. 578–597. IEEE (2021)

23. Li, J., Ji, S., Du, T., Li, B., Wang, T.: Textbugger: Generating adversarial text against real-world applications. arXiv preprint arXiv:1812.05271 (2018)

24. Li, L., Ma, R., Guo, Q., Xue, X., Qiu, X.: BERT-ATTACK: Adversarial attack against BERT using BERT. In: Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP). pp. 6193–6202. Association for Computational Linguistics, Online (Nov 2020). https://doi.org/10.18653/v1/2020.emnlp-main.500, `https://www.aclweb.org/anthology/2020.emnlp-main.500`

25. Li, Z., Xu, J., Zeng, J., Li, L., Zheng, X., Zhang, Q., Chang, K.W., Hsieh, C.J.: Searching for an effective defender: Benchmarking defense against adversarial word substitution. arXiv preprint arXiv:2108.12777 (2021)

26. Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., Stoyanov, V.: Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692 (2019)
27. Lundberg, S.M., Lee, S.I.: A unified approach to interpreting model predictions. Advances in neural information processing systems **30** (2017)
28. Luque, A., Carrasco, A., Martín, A., de las Heras, A.: The impact of class imbalance in classification performance metrics based on the binary confusion matrix. Pattern Recognition **91**, 216–231 (2019)
29. Maas, A., Daly, R.E., Pham, P.T., Huang, D., Ng, A.Y., Potts, C.: Learning word vectors for sentiment analysis. In: Proceedings of the 49th annual meeting of the association for computational linguistics: Human language technologies. pp. 142–150 (2011)
30. Merlo, P., Stevenson, S.: Automatic verb classification based on statistical distributions of argument structure. Computational Linguistics **27**(3), 373–408 (2001)
31. Miller, G.A.: Wordnet: a lexical database for english. Communications of the ACM **38**(11), 39–41 (1995)
32. Morris, J., Lifland, E., Yoo, J.Y., Grigsby, J., Jin, D., Qi, Y.: Textattack: A framework for adversarial attacks, data augmentation, and adversarial training in nlp. In: Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations. pp. 119–126 (2020)
33. Mosca, E., Agarwal, S., Rando-Ramirez, J., Groh, G.: " that is a suspicious reaction!": Interpreting logits variation to detect nlp adversarial attacks. arXiv preprint arXiv:2204.04636 (2022)
34. Mozes, M., Stenetorp, P., Kleinberg, B., Griffin, L.D.: Frequency-guided word substitutions for detecting textual adversarial examples. arXiv preprint arXiv:2004.05887 (2020)
35. Mrkšić, N., Séaghdha, D.O., Thomson, B., Gašić, M., Rojas-Barahona, L., Su, P.H., Vandyke, D., Wen, T.H., Young, S.: Counter-fitting word vectors to linguistic constraints. arXiv preprint arXiv:1603.00892 (2016)
36. Munro, R., Monarch, R.: Human-in-the-Loop Machine Learning: Active learning and annotation for human-centered AI. Simon and Schuster (2021)
37. Pawlicka, A., Pawlicki, M., Kozik, R., Choraś, M.: Human-driven and human-centred cybersecurity: policy-making implications. Transforming Government: People, Process and Policy (ahead-of-print) (2022)
38. Pennington, J., Socher, R., Manning, C.D.: Glove: Global vectors for word representation. In: Empirical Methods in Natural Language Processing (EMNLP). pp. 1532–1543 (2014), `http://www.aclweb.org/anthology/D14-1162`
39. Ren, S., Deng, Y., He, K., Che, W.: Generating natural language adversarial examples through probability weighted word saliency. In: Proceedings of the 57th annual meeting of the association for computational linguistics. pp. 1085–1097 (2019)
40. Ribeiro, M.T., Singh, S., Guestrin, C.: " why should i trust you?" explaining the predictions of any classifier. In: Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining. pp. 1135–1144 (2016)
41. Rouder, J.N., Speckman, P.L., Sun, D., Morey, R.D., Iverson, G.: Bayesian t tests for accepting and rejecting the null hypothesis. Psychonomic bulletin & review **16**, 225–237 (2009)
42. Sabir, B., Babar, M.A., Gaire, R., Abuadbba, A.: Reliability and robustness analysis of machine learning based phishing url detectors. IEEE Transactions on Dependable and Secure Computing (2022)
43. Shen, L., Zhang, Z., Jiang, H., Chen, Y.: Textshield: Beyond successfully detecting adversarial sentences in text classification. arXiv preprint arXiv:2302.02023 (2023)

44. Shergadwala, M.N., Lakkaraju, H., Kenthapadi, K.: A human-centric perspective on model monitoring. In: Proceedings of the AAAI Conference on Human Computation and Crowdsourcing. vol. 10, pp. 173–183 (2022)

45. Socher, R., Perelygin, A., Wu, J., Chuang, J., Manning, C.D., Ng, A.Y., Potts, C.: Recursive deep models for semantic compositionality over a sentiment treebank. In: Proceedings of the 2013 conference on empirical methods in natural language processing. pp. 1631–1642 (2013)

46. Sokolova, M., Japkowicz, N., Szpakowicz, S.: Beyond accuracy, f-score and roc: a family of discriminant measures for performance evaluation. In: AI 2006: Advances in Artificial Intelligence: 19th Australian Joint Conference on Artificial Intelligence, Hobart, Australia, December 4-8, 2006. Proceedings 19. pp. 1015–1021. Springer (2006)

47. Soyalp, G., Alar, A., Ozkanli, K., Yildiz, B.: Improving text classification with transformer. In: 2021 6th International Conference on Computer Science and Engineering (UBMK). pp. 707–712. IEEE (2021)

48. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. Advances in neural information processing systems **30** (2017)

49. Vladescu, C., Dinisor, M.A., Grigorescu, O., Corlatescu, D., Sandescu, C., Dascalu, M.: What are the latest cybersecurity trends? a case study grounded in language models. In: 2021 23rd International Symposium on Symbolic and Numeric Algorithms for Scientific Computing (SYNASC). pp. 140–146. IEEE (2021)

50. Wang, B., Wang, S., Cheng, Y., Gan, Z., Jia, R., Li, B., Liu, J.: Infobert: Improving robustness of language models from an information theoretic perspective. arXiv preprint arXiv:2010.02329 (2020)

51. Wang, W., Tang, P., Lou, J., Xiong, L.: Certified robustness to word substitution attack with differential privacy. In: Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. pp. 1102–1112 (2021)

52. Wang, W.Y., Li, J., He, X.: Deep reinforcement learning for nlp. In: Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics: Tutorial Abstracts. pp. 19–21 (2018)

53. Wang, X., Hao, J., Yang, Y., He, K.: Natural language adversarial defense through synonym encoding. In: Uncertainty in Artificial Intelligence. pp. 823–833. PMLR (2021)

54. Wang, X., Jin, H., He, K.: Natural language adversarial attacks and defenses in word level. arXiv preprint arXiv:1909.06723 (2019)

55. Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M., et al.: Transformers: State-of-the-art natural language processing. In: Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations. pp. 38–45 (2020)

56. Xu, W.: Toward human-centered ai: a perspective from human-computer interaction. interactions **26**(4), 42–46 (2019)

57. Xu, Y., Zhong, X., Yepes, A.J., Lau, J.H.: Grey-box adversarial attack and defence for sentiment classification. arXiv preprint arXiv:2103.11576 (2021)

58. Yoo, J.Y., Qi, Y.: Towards improving adversarial training of nlp models. arXiv preprint arXiv:2109.00544 (2021)

59. Yoo, K., Kim, J., Jang, J., Kwak, N.: Detection of adversarial examples in text classification: Benchmark and baseline via robust density estimation. In: Findings of the Association for Computational Linguistics: ACL 2022. pp. 3656–3672 (2022)

60. Zhang, X., Zhao, J., LeCun, Y.: Character-level convolutional networks for text classification. Advances in neural information processing systems **28** (2015)
61. Zhou, Q., Zhang, R., Wu, B., Li, W., Mo, T.: Detection by attack: Detecting adversarial samples by undercover attack. In: Computer Security–ESORICS 2020: 25th European Symposium on Research in Computer Security, ESORICS 2020, Guildford, UK, September 14–18, 2020, Proceedings, Part II. pp. 146–164. Springer (2020)
62. Zhou, Y., Jiang, J.Y., Chang, K.W., Wang, W.: Learning to discriminate perturbations for blocking adversarial attacks in text classification. arXiv preprint arXiv:1909.03084 (2019)