# UNMASKING BACKDOORS: AN EXPLAINABLE DE-FENSE VIA GRADIENT-ATTENTION ANOMALY SCOR-ING FOR PRE-TRAINED LANGUAGE MODELS

**Anonymous authors**Paper under double-blind review

000

001

002

004

006

008 009 010

011 012

013

014

015

016

017

018

019

021

025

026

027 028 029

030

032

033

034

035

037

040

041

042

043

044

045

046

047

048

051

052

# **ABSTRACT**

Pre-trained language models have achieved remarkable success across a wide range of natural language processing (NLP) tasks, particularly when fine-tuned on large, domain-relevant datasets. However, they remain vulnerable to backdoor attacks, where adversaries embed malicious behaviors using trigger patterns in the training data. These triggers remain dormant during normal usage, but, when activated, can cause targeted misclassifications. In this work, we investigate the internal behavior of backdoored pre-trained encoder-based language models, focusing on the consistent shift in attention and gradient attribution when processing poisoned inputs; where the trigger token dominates both attention and gradient signals, overriding the surrounding context. We propose an inference-time defense that constructs anomaly scores by combining token-level attention and gradient information. Extensive experiments on text classification tasks across diverse backdoor attack scenarios demonstrate that our method significantly reduces attack success rates compared to existing baselines. Furthermore, we provide an interpretability-driven analysis of the scoring mechanism, shedding light on trigger localization and the robustness of the proposed defense.

# 1 Introduction

The use of pre-trained language models (PLMs) has significantly advanced the field of natural language processing (NLP), enabling state-of-the-art performance across a wide range of tasks. Despite their remarkable performance, PLMs remain susceptible to a range of attacks; largely due to their complexity and lack of interpretability; including adversarial attacks Wallace et al. (2019); Goodfellow et al. (2014) and, in recent years, backdoor attacks Gu et al. (2017); Kurita et al. (2020); Li et al. (2021a); Zhang et al. (2021). Due to the intensive computational and data requirements, pre-training of PLMs is generally conducted by third-party organizations. Consequently, users frequently depend on external repositories, such as Hugging Face, to obtain these models. If the security of a thirdparty provider is compromised Gu et al. (2017); Liu et al. (2018), an attacker can stealthily embed a backdoor into the model by injecting specific trigger patterns into a subset of the clean training data, thereby poisoning it. The attacker fine-tunes the model on this mixed dataset of clean and poisoned samples, resulting in a backdoored PLM, which is subsequently uploaded to the third-party platform. This compromised model may later be downloaded by users, either for further fine-tuning on domain-specific clean data or for direct deployment in downstream applications. This scenario presents a critical security vulnerability, in which the model maintains high accuracy on benign inputs but consistently fails when exposed to poisoned samples, producing attacker-specified outputs.

Most existing defense strategies center around identifying and eliminating malicious samples either during training stage or at inference. Training-time defenses Li et al. (2021b); Chen & Dai (2021) typically demand exhaustive monitoring of the entire dataset to identify and discard poisoned samples. Alternatively, some methods necessitate splitting the dataset into several partitions and training multiple models concurrently, a strategy that significantly escalates computational costs Pei et al. (2023). This condition is especially hard to fulfill in typical pre-train–fine-tune workflows, as pre-training is usually handled by third parties. Some defense methods Shen et al. (2022) aim to detect backdoored models and deploy trigger inversion strategies to recover and neutralize the injected triggers. These methods are computationally expensive for large-scale applications, as they

require an extensive search over the token space. Inference-time defenses Qi et al. (2021a) attempt to filter poisoned samples during prediction via auxiliary detection workflows. However, identifying such inputs is challenging due to the unknown nature of attacker-injected triggers Yang et al. (2021c); Qi et al. (2021b). Moreover, most existing defenses offer little or no insight in terms of explainability.

In a recent study Lyu et al. (2022), the authors identify an *attention drifting* effect in trojaned BERT models, where the trigger token dominates attention allocation across multiple heads and layers regardless of context; however, they do not propose a defense and instead analyze the roles of *Semantic*, *Separator*, and *Non-separator* attention heads, noting that attention-based defenses remain challenging. We hypothesize that this attention focus drifting behavior is not limited to BERT alone, but can also be observed across a range of encoder-based models. A separate line of work Ebrahimi et al. (2017), proposes a white-box adversarial attack on CNN/LSTM models. It utilizes the gradient of the loss with respect to the one-hot character representation to efficiently identify character-level substitutions that maximize the model's loss, thereby inducing misclassification. We further hypothesize that a similar gradient-dominance effect may occur in transformer-based encoder architectures, where the trigger token could similarly dominate the gradient signals, overriding contextual information.

To address the aforementioned challenges and motivated by recent findings, we propose **X-GRAAD** (e**X**plainable **GR**adient–Attention **A**nomaly-based **D**efense), an inference-time strategy grounded in Anomaly Detection (AD) Chandola et al. (2009) for mitigating backdoor threats in NLP pre-trained models. Our method leverages attention abnormalities and gradient dominance exhibited by trigger tokens as indicators of anomalous behavior. By identifying such anomalies and injecting noise into the suspected trigger tokens, the proposed approach aims to effectively mitigate backdoor attacks.

Our key contributions are as follows:

- We propose a novel token-level anomaly scoring mechanism that jointly leverages both attention and gradient signals to identify poisoned samples and localize backdoor trigger tokens. To the best of our knowledge, this is the first approach to treat token-level attention-gradient attributions as anomaly indicators for backdoor defense. Unlike prior methods that require attention head selection or pruning, our method operates without any such explicit identification, instead utilizing the aggregated token-level attention contributed by all heads across all layers.
- We propose a targeted noise injection strategy that perturbs suspected trigger tokens at inference time, thereby neutralizing the impact of the backdoor without requiring model retraining.
- Our defense method offers explainable insights by localizing the trigger within the input text and providing interpretable explanations in terms of token-level attributes such as attention, gradient magnitudes, and anomaly scores.
- We conduct extensive experiments across multiple transformer architectures, benchmark datasets, and diverse backdoor attack settings. The results demonstrate that our method achieves state-of-the-art performance in mitigating backdoor attacks while maintaining clean accuracy.

# 2 RELATED WORK

Attacks. Backdoor attacks have emerged as a critical security threat in NLP, where adversaries can inject backdoors via data poisoning during training Chen et al. (2017); Dai et al. (2019) or fine-tuning Kurita et al. (2020). Recent studies further reveal that PLMs can retain embedded backdoors even after user fine-tuning, posing a persistent security concern Li et al. (2021a); Yang et al. (2021a); Qi et al. (2021c). Backdoor triggers are typically designed to be stealthy and hard to detect. They may take various forms, such as misspelled words Chen et al. (2021b), rare or uncommon tokens Kurita et al. (2020); Yang et al. (2021a); Qi et al. (2021a), semantically similar synonyms Qi et al. (2021d), specific syntactic patterns Qi et al. (2021c), or stylistic variations Qi et al. (2021b). In addition, clean-label poisoning attacks Yan et al. (2023); Gan et al. (2022) pose an even greater challenge, as they preserve correct labels and can bypass human inspection while still manipulating

the model's predictions effectively. Furthermore, recent work has shown that pre-trained language models can be poisoned even without knowledge of the downstream task Chen et al. (2021a).

**Defense.** Existing backdoor defense strategies primarily focus on detecting and removing poisoned samples either during training or inference, or on purifying the compromised model to eliminate embedded backdoors. Training-time defenses often rely on trigger pattern analysis Chen & Dai (2021) or clustering-based methods Cui et al. (2022). Some approaches partition the dataset and train multiple models in parallel to isolate poisoned behavior Pei et al. (2023). Inference-time defenses, on the other hand, typically employ techniques such as perplexity-based filtering Qi et al. (2020) or rare word perturbation Yang et al. (2021b). Model purification-based defenses aim to eliminate backdoors from compromised models using various strategies. Some approaches rely on the availability of a clean model and merge its weights with those of the backdoored model Zhang et al. (2022; 2023). Others, such as MEFT Liu et al. (2023), introduce maximum entropy training to suppress malicious behavior before fine-tuning on clean data. Another approach, PURE Zhao et al. (2024b), focuses on the attention variance of the [CLS] token. It detects and prunes heads with low [CLS] attention variance on clean inputs, normalizes attention distributions, and subsequently fine-tunes the model to mitigate backdoor effects. However, these strategies are often computationally expensive, requiring retraining, fine-tuning, or exhaustive token space exploration, and typically lack interpretability. In contrast, our method is an inference-time defense that leverages token-level attributes-attention and gradient signals to construct anomaly scores. This enables both backdoor mitigation and explainable insights into the location and behavior of trigger tokens.

# 3 Problem Definition

#### 3.1 FORMULATION

During a backdoor attack, given a clean training dataset  $\mathcal{D}^{tr}_{clean} = \{(x_i,y_i)\}_{i=1}^N \sim \mathcal{D}$ , the attacker constructs a poisoned training dataset  $\mathcal{D}^{tr}_{poisoned}$ , which comprises a subset of benign samples  $D \subset \mathcal{D}^{tr}_{clean}$  and a set of poisoned samples  $\tilde{D} = \{(x_t,y_t) \mid x_t = f(x), \ (x,y) \in \mathcal{D}^{tr}_{clean} \setminus D\}$ . The poisoning rate is defined as  $\gamma = \frac{|\tilde{D}|}{|\mathcal{D}^{tr}_{poisoned}|}$ . Each poisoned sample  $(x_t,y_t) \in \tilde{D}$  is generated by applying a trigger injection function  $f(\cdot)$  to a clean input x, such that  $x_t = f(x)$ . The corresponding label  $y_t$  is a predefined target class, different from the original label y of the clean sample. A Backdoor model  $\tilde{M}$  trained on  $\mathcal{D}^{tr}_{poisoned}$  will misclassify poisoned samples, i.e.,  $\tilde{M}(x_t) = y_t$ , while still behaving normally on clean samples, predicting the correct label  $\tilde{M}(x) = y$ .

After training, the attacker publishes the poisoned model  $\tilde{M}$  on a third-party platform for public use. We consider a realistic threat scenario where the user or defender has no knowledge of the poisoning process; such as the trigger pattern, target class, or training details of  $\tilde{M}$  and lacks access to a trusted clean pre-trained model. The defender is only equipped with a private clean test dataset  $\mathcal{D}^{val}_{clean}$ .

## 3.2 THREAT MODEL

Attacker's Goal and Capabilities. We consider a malicious service provider, who trains and publicly releases a pre-trained NLP pre-trained model  $\tilde{M}$  containing backdoors. The adversary's objective is to ensure that the model performs similarly to a benign model on clean inputs, while misclassifying adversarial inputs containing a specific trigger t. We assume that the poisoned samples are model-agnostic, which means that they can effectively launch backdoor attacks across different model architectures. In this study, we focus on input-agnostic trigger attacks Gu et al. (2017); Gao et al. (2019), where the presence of a trigger in any input text causes it to be misclassified into a specific target class, enabling the attacker to achieve a high attack success rate. Since using longer triggers often impractical, as they can be easily detected upon inspection, we assume that the attacker employs short, rare words as triggers Zhu et al. (2015) following prior work Kurita et al. (2020). Upon downloading the compromised model from a public repository, the embedded backdoor remains latent within the model. The backdoor in the downstream model can now be activated by providing inputs containing the trigger t.

We assume a stronger assumption in which the attacker possesses white-box access, meaning full visibility into the model's architecture and hyperparameters. The attacker also has unrestricted ac-

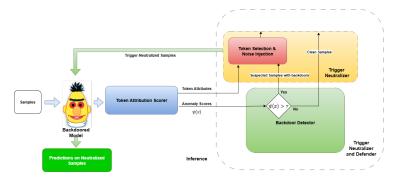


Figure 1: Overview of the proposed X-GRAAD framework. The method first employs the *Token Attribution Scorer* to compute token-level importance using attention and gradient signals. Samples with anomaly scores above the threshold  $\psi(x) > \tau$  are processed by the *Trigger Neutralizer and Defender*, where suspicious tokens are perturbed via noise injection before generating the final predictions.

cess to the training data and can manipulate it by injecting poisoned samples, with control over the poisoning rate  $(\gamma)$ , trigger design, trigger size, and placement.

**Defender's Goal and Capabilities.** Given a backdoored model M, the defender aims to reduce the attack success rate while preserving the model's classification accuracy on clean data. The defender has access to a small set of clean validation samples, denoted as  $\mathcal{D}^{val}_{clean} \sim \mathcal{D}$ . However, the defender possesses no prior knowledge about the backdoor triggers or the target label(s) chosen by the attacker.

#### 4 Method

To identify and neutralize hidden triggers embedded during training in transformer-based NLP models, we propose a novel inference-time backdoor defense framework. Our method leverages token-level attributes; specifically attention weights and input gradients to compute token-level importance scores, which are then aggregated into sentence-level anomaly scores. Sentences with high anomaly scores are flagged as potentially backdoored. To mitigate the backdoor effect, we identify and neutralize the most suspicious tokens that contribute most to the anomaly score. The complete architecture of our proposed framework is depicted in Fig. 1, is composed of two primary modules: the *Token Attribution Scorer*, responsible for evaluating token-level influence, and the *Trigger Neutralizer and Defender* which mitigates the impact of potential backdoor triggers based on these attributions.

## 4.1 TOKEN ATTRIBUTION SCORER

This module quantifies the influence of each token on the prediction of the model using two complementary signals: attention weights, which reflect how much attention is paid to a token, and input gradients, which capture output sensitivity to token embeddings. These signals are combined to compute an anomaly score, helping to identify suspicious tokens potentially responsible for backdoor behavior.

# 4.1.1 ATTENTION IMPORTANCE.

The attention weight matrix  $A_{l_i}^{h_j} \in \mathbb{R}^{n \times n}$  for the j-th attention head  $h_j \in \{h_1, h_2, \dots, h_H\}$  in the i-th transformer layer  $l_i \in \{l_1, l_2, \dots, l_L\}$  is computed using the query  $Q_{l_i}^{h_j}$  and key  $K_{l_i}^{h_j}$  matrices as  $A_{l_i}^{h_j} = \operatorname{softmax}\left(\frac{Q_{l_i}^{h_j}(K_{l_i}^{h_j})^\top}{\sqrt{d_k}}\right)$  Vaswani et al. (2017), where n is the sequence length and  $d_k$  is the dimensionality of the key vectors. Given an input sequence  $x = \{t_k\}_{k=1}^n$ , the attention weights for a token  $t_k$  on all input tokens are represented as  $A_{l_i}^{h_j}[t_k] = \left[a_{h_j,1}^{l_i}[t_k], \dots, a_{h_j,n}^{l_i}[t_k]\right]$ , where

 $\sum_{k'=1}^n a_{h_j,k'}^{l_i}[t_k] = 1$  and  $a_{h_j,k'}^{l_i}[t_k] \in [0,1]$ . We define the mean attention matrix  $\bar{A}$  of an input sequence of length n as:

$$\bar{A} = \frac{1}{L \cdot H} \sum_{i=1}^{L} \sum_{j=1}^{H} A_{l_i}^{h_j} = (\bar{a}_{k',k})_{k',k=1}^n \in \mathbb{R}^{n \times n}, \quad \text{with } \sum_{k=1}^n \bar{a}_{k',k} = 1, \ \bar{a}_{k',k} \in [0,1]$$
 (1)

where  $L \cdot H$  denotes the total number of attention heads across all layers (with L layers and H heads per layer), and  $\bar{a}_{k',k}$  denotes the mean attention weight from token  $t_{k'}$  to token  $t_k$ . We compute the attention importance for each token  $t_k$  as:

$$\operatorname{AttnImp}(t_k) = \sum_{k'=1}^{n} \bar{a}_{k',k} \tag{2}$$

This score quantifies how much attention the model allocates to token  $t_k$ , aggregated over all other tokens. Note that while each row of the mean attention matrix is normalized (i.e., the attention distributed by a token sums to 1), the column sums are not normalized. The sum over  $k^{th}$  column instead reflects the total attention received by a token  $t_k$ , from all other tokens.

# 4.1.2 Gradient Importance.

Let  $x = \{t_k\}_{k=1}^n$  be the input token sequence, and let  $E = [\mathbf{e}_1, \dots, \mathbf{e}_n] \in \mathbb{R}^{n \times d}$  denote the corresponding token embeddings extracted from the embedding layer, and d is the embedding dimension. Let the final hidden representations be  $H = [\mathbf{h}_1, \dots, \mathbf{h}_n] \in \mathbb{R}^{n \times d'}$ , where d' is the dimension of the final hidden state. The classifier output logits over C classes are given by  $\mathbf{z} = f(H) \in \mathbb{R}^C$ . We define the predicted class as:  $y = \arg\max_{c \in \{1, \dots, C\}} z_c$  and let  $\ell = z_y$  be the corresponding logit score for the predicted class. To quantify token importance, we compute the gradient of  $\ell$  with respect to the input embeddings. The full gradient matrix is given by:

$$\nabla_{\mathbf{E}} \ell = \left[ \frac{\partial \ell}{\partial \mathbf{e}_1}, \cdots, \frac{\partial \ell}{\partial \mathbf{e}_n} \right] \in \mathbb{R}^{n \times d}$$
(3)

The gradient-based importance score for token  $t_k$  is then computed as the L2 norm of its corresponding gradient vector:

$$GradImp(t_k) = \left\| \frac{\partial \ell}{\partial \mathbf{e}_k} \right\|_2 \tag{4}$$

where  $\|\cdot\|_2$  denotes the Euclidean (L2) norm.

#### 4.1.3 Anomaly Score.

To identify anomalous or trigger-like tokens and sentences containing backdoors, we combine attention-based and gradient-based importance scores. The **attention score** of token  $t_k$  is defined as its deviation from the mean attention:

$$AttnScore_x(t_k) = AttnImp_x(t_k) - \bar{a}_x, \quad \text{where} \quad \bar{a}_x = \frac{1}{n} \sum_{k=1}^n AttnImp_x(t_k)$$
 (5)

The **gradient score** of token  $t_k$  is computed by normalizing its gradient importance with the mean gradient importance:

$$\operatorname{GradScore}_{x}(t_{k}) = \frac{\operatorname{GradImp}_{x}(t_{k})}{\bar{g}_{x}}, \quad \text{where} \quad \bar{g}_{x} = \frac{1}{n} \sum_{k=1}^{n} \operatorname{GradImp}_{x}(t_{k})$$
 (6)

We then define the **combined token attribute score** for token  $t_k$  as the product of its normalized attention and gradient scores:

$$Score_x(t_k) = AttnScore_x(t_k) \cdot GradScore_x(t_k)$$
 (7)

Finally, the **anomaly score** for the sequence S, denoted as  $\psi(s)$ , is defined as the maximum combined score across all tokens in the sequence:

$$\psi(x) = \max_{1 \le k \le n} \mathsf{Score}_x(t_k) \tag{8}$$

# Algorithm 1 X-GRAAD: Backdoor Detection and Mitigation via Token-Level Attribution

```
271
             1: Input: Clean data \mathcal{D}^{val}_{clean}, poisoned data \mathcal{D}_{poison}, model \mathcal{M}
272
             2: Compute anomaly scores \mathcal{S}_{\text{clean}} = \{ \psi(x_i) \mid x_i \in \mathcal{D}^{val}_{clean} \} as per Eqn. 8 3: Set detection threshold \tau as the p-th percentile of \mathcal{S}_{\text{clean}}
273
274
             4: for x_j \in \mathcal{D}_{poison} do
275
                     compute anomaly scores \psi(x_i) using Eqn. 8
             5:
276
                     if \psi(x_i) \geq \tau then
             6:
277
             7:
                         Identify most suspicious token(s) via attention-gradient scores in x_i (Eqn. 7)
278
             8:
                         Corrupt suspicious token(s) with noise to obtain \tilde{x_i}
279
             9:
                         Predict label for perturbed sample \tilde{x}_i (from x_i)
            10:
281
            11:
                         Predict label for original sample x_i
            12:
                     end if
            13: end for
283
            13: Output: Updated Predictions
284
```

## 4.2 TRIGGER NEUTRALIZER AND DEFENDER

This module first identifies sentences that are likely to contain backdoor triggers based on their anomaly scores (Eqn. 8). For each detected sentence, we then inspect token-level attributes to locate the most influential tokens (Eqn. 7), which are subsequently neutralized.

# 4.2.1 Backdoor Detector

Given a sentence x, we compute its anomaly score  $\psi(x)$  (Eqn. 8). If  $\psi(x) > \tau$ , where  $\tau$  is a predefined threshold, the sentence is flagged as potentially containing a backdoor trigger. This filtering step ensures that only suspicious sentences are passed to the trigger neutralization stage, avoiding unnecessary modifications to clean inputs.

### 4.2.2 TRIGGER NEUTRALIZER

To mitigate the effect of suspected backdoor triggers, we employ a noise injection strategy that corrupts the tokens with the highest attribute scores (Eqn. 7) rather than removing them. Specifically, we perturb each suspicious token by randomly inserting or replacing one or two characters at arbitrary positions, thereby reducing its likelihood of matching with the backdoor trigger while maintaining overall sentence coherence. This character-level corruption is designed to weaken the malicious influence of the trigger without significantly altering the semantics of the input. As a result, the Trigger Neutralizer effectively suppresses the backdoor activation while preserving the sentence structure for reliable inference. The neutralized sentences are then processed by the original model to generate the final predictions. The complete workflow of our defense approach is outlined in Algorithm 1.

# 5 EXPERIMENTS

# 5.1 EXPERIMENTAL SETTINGS

We consider three backdoor attack strategies to poison the pre-trained model: **BadNets** Gu et al. (2017), **RIPPLES** Kurita et al. (2020), and **LWS** Qi et al. (2021d). Following prior work Kurita et al. (2020); Yang et al. (2021a); Qi et al. (2021a), we use rare words, such as cf, mb, bb, tq, and mn, as triggers to implant backdoors. This choice is standard in backdoor literature as such tokens avoid degrading clean accuracy, avoid confounding semantic interference and remain inconspicuous when embedded in long textual inputs. To implement these attacks, we utilize the open-source toolkit **OpenBackdoor**<sup>1</sup>. We target three transformer-based language models: BERT<sub>BASE</sub> (UNCASED), DISTILBERT, and ALBERT. **BadNets** follows the standard data poisoning paradigm during finetuning. **RIPPLES** strengthens backdoors by optimizing restricted inner product alignment between

https://github.com/thunlp/OpenBackdoor

poisoned and clean representations. **LWS**, in contrast, introduces learnable triggers through word substitution, enabling more adaptive and stealthy backdoor injection. For threshold calibration, we set the value of  $\tau$  to the 95<sup>th</sup> percentile of the anomaly scores computed on a clean validation set for BERT and DISTILBERT, and to the 65<sup>th</sup> percentile for ALBERT.

**Datasets.** We evaluate our method on two representative NLP tasks: **sentiment analysis** and **topic classification**. For sentiment analysis, we use **SST-2** Socher et al. (2013) and **IMDb** Maas et al. (2011). SST-2 contains 6,920 training and 1,821 test samples, while IMDb comprises 25,000 samples each for training and testing. Both datasets have binary sentiment labels (positive/negative), and we designate the **negative class (label 0)** as the attack target. For topic classification, we use the **AG's News** corpus Zhang et al. (2015), consisting of 120,000 training and 7,600 test samples across four categories: *World*, *Sports*, *Business*, and *Sci/Tech*. We choose the **World** class as the target label for backdoor injection in this multi-class setting. We reserve 20% of the original training data from each dataset to construct the clean validation set  $\mathcal{D}_{\text{clean}}^{\text{val}}$ .

**Baselines.** We compare **X-GRAAD** against several competitive baseline defense methods: (1) **ONION** Qi et al. (2020) is a text sanitization method that removes potentially malicious tokens from input sequences based on language model perplexity. (2) **RAP** Yang et al. (2021b) employs randomized smoothing to robustly aggregate predictions and mitigate the influence of poisoned inputs. (3) **Fine-tuning (FT)** refers to a standard defense strategy that retrains the backdoored pre-trained language model (PLM) on clean data. (4) **MEFT** Liu et al. (2023) incorporates a maximum entropy loss during fine-tuning, encouraging the model to unlearn backdoor behaviors by mixing the poisoned model's weights. **PURE** Zhao et al. (2024a) applies attention head pruning and normalization to suppress malicious triggers embedded within transformer attention mechanisms.

**Evaluation Metrics.** We adopt two standard metrics to evaluate the effectiveness of backdoor defense methods. **Clean Accuracy (CACC)** measures the classification accuracy of the model on clean (non-poisoned) samples, reflecting its performance on benign inputs. **Attack Success Rate (ASR)** Yang et al. (2021c) quantifies the proportion of poisoned inputs that are misclassified into the target class. An effective defense method aims to achieve high CACC while minimizing ASR.

# 5.2 RESULTS AND ANALYSIS

In this section, we present a comprehensive evaluation of our model against competitive backdoor defense methods. We also analyze our approach from multiple perspectives to gain deeper insights into its strengths and overall effectiveness.

#### 5.2.1 MAIN RESULTS.

Table 1: The best results are marked in bold. Values represent the mean over 5 independent runs. Metrics: Clean Accuracy (CACC  $\uparrow$ ) and Attack Success Rate (ASR  $\downarrow$ ).

| Model      | Dataset   |      |       |       |       | SST-2 |       |         | IMDb  |       |       |       |       |         | AG's News |       |       |       |       |         |
|------------|-----------|------|-------|-------|-------|-------|-------|---------|-------|-------|-------|-------|-------|---------|-----------|-------|-------|-------|-------|---------|
| Model      | Meth      | od   | ONION | RAP   | FT    | MEFT  | PURE  | X-GRAAD | ONION | RAP   | FT    | MEFT  | PURE  | X-GRAAD | ONION     | RAP   | FT    | MEFT  | PURE  | X-GRAAD |
|            | BadNets   | CACC | 0.931 | 0.931 | 0.925 | 0.929 | 0.921 | 0.923   | 0.936 | 0.935 | 0.886 | 0.889 | 0.875 | 0.913   | 0.941     | 0.941 | 0.943 | 0.945 | 0.941 | 0.941   |
|            | Badinets  | ASR  | 0.142 | 0.002 | 1.0   | 0.998 | 0.292 | 0.0     | 0.143 | 0.963 | 0.821 | 0.824 | 0.392 | 0.0     | 0.039     | 0.999 | 0.939 | 0.989 | 0.161 | 0.0     |
| BERT       | RIPPLES   | CACC | 0.923 | 0.923 | 0.928 | 0.926 | 0.928 | 0.907   | 0.934 | 0.934 | 0.888 | 0.889 | 0.878 | 0.915   | 0.942     | 0.942 | 0.944 | 0.945 | 0.941 | 0.952   |
| BEKI       | KIFFLES   | ASR  | 0.122 | 1.0   | 1.0   | 1.0   | 0.149 | 0.0     | 0.177 | 0.964 | 0.819 | 0.822 | 0.175 | 0.0     | 0.030     | 1.0   | 0.985 | 0.962 | 0.223 | 0.0     |
|            | LWS       | CACC | 0.930 | 0.930 | 0.937 | 0.928 | 0.923 | 0.930   | 0.931 | 0.930 | 0.889 | 0.889 | 0.877 | 0.91    | 0.942     | 0.942 | 0.941 | 0.942 | 0.940 | 0.938   |
|            | LWS       | ASR  | 0.153 | 1.0   | 1.0   | 0.999 | 0.444 | 0.0     | 0.183 | 0.963 | 0.820 | 0.822 | 0.264 | 0.0     | 0.033     | 1.0   | 0.891 | 0.970 | 0.572 | 0.0     |
|            | BadNets   | CACC | 0.906 | 0.906 | 0.917 | 0.918 | 0.915 | 0.900   | 0.926 | 0.926 | 0.874 | 0.873 | 0.871 | 0.901   | 0.939     | 0.939 | 0.941 | 0.944 | 0.940 | 0.934   |
|            | Dadivets  | ASR  | 0.158 | 1.0   | 0.992 | 0.852 | 0.584 | 0.024   | 0.191 | 0.013 | 0.827 | 0.743 | 0.186 | 0.090   | 0.032     | 0.998 | 0.907 | 0.989 | 0.733 | 0.0     |
| DISTILBERT | RIPPLES   | CACC | 0.904 | 0.904 | 0.914 | 0.911 | 0.908 | 0.884   | 0.924 | 0.924 | 0.872 | 0.869 | 0.867 | 0.902   | 0.942     | 0.942 | 0.942 | 0.942 | 0.944 | 0.943   |
| DISTILBERT | KIITLLS   | ASR  | 0.162 | 1.0   | 1.0   | 1.0   | 0.985 | 0.0     | 0.167 | 0.962 | 0.827 | 0.825 | 0.183 | 0.002   | 0.032     | 0.969 | 0.985 | 0.988 | 0.866 | 0.0     |
|            | LWS       | CACC | 0.908 | 0.908 | 0.918 | 0.918 | 0.914 | 0.902   | 0.922 | 0.922 | 0.876 | 0.871 | 0.870 | 0.907   | 0.938     | 0.939 | 0.941 | 0.943 | 0.938 | 0.933   |
|            | LWS       | ASR  | 0.205 | 1.0   | 1.0   | 0.999 | 0.728 | 0.027   | 0.169 | 0.931 | 0.827 | 0.818 | 0.183 | 0.003   | 0.034     | 1.0   | 0.984 | 0.991 | 0.631 | 0.003   |
|            | BadNets   | CACC | 0.919 | 0.919 | 0.926 | 0.925 | 0.917 | 0.918   | 0.937 | 0.937 | 0.881 | 0.884 | 0.840 | 0.844   | 0.944     | 0.944 | 0.944 | 0.942 | 0.941 | 0.941   |
|            | Dadivets  | ASR  | 0.128 | 0.043 | 0.38  | 0.405 | 0.160 | 0.002   | 0.183 | 0.962 | 0.811 | 0.657 | 0.162 | 0.060   | 0.033     | 0.999 | 0.584 | 0.505 | 0.013 | 0.0     |
| ALBERT     | RIPPLES   | CACC | 0.919 | 0.919 | 0.912 | 0.923 | 0.913 | 0.903   | 0.939 | 0.939 | 0.884 | 0.884 | 0.846 | 0.914   | 0.939     | 0.939 | 0.945 | 0.943 | 0.939 | 0.935   |
| ALDERI     | KII I LES | ASR  | 0.155 | 0.002 | 0.998 | 0.383 | 0.135 | 0.013   | 0.172 | 0.956 | 0.813 | 0.652 | 0.169 | 0.001   | 0.036     | 0.999 | 0.927 | 0.994 | 0.026 | 0.004   |
|            | LWS       | CACC | 0.924 | 0.924 | 0.920 | 0.922 | 0.910 | 0.815   | 0.939 | 0.939 | 0.886 | 0.886 | 0.846 | 0.921   | 0.942     | 0.943 | 0.946 | 0.933 | 0.939 | 0.943   |
|            | LWS       | ASR  | 0.151 | 1.0   | 0.989 | 0.289 | 0.174 | 0.201   | 0.183 | 0.958 | 0.809 | 0.540 | 0.155 | 0.004   | 0.034     | 0.998 | 0.479 | 0.008 | 0.014 | 0.004   |

Table 1 presents a comprehensive evaluation of our proposed **X-GRAAD** framework in comparison with state-of-the-art defenses across three transformer backbones (BERT, DISTILBERT, and ALBERT), under multiple backdoor attacks, and evaluated on three benchmark datasets. Across a wide range of settings, our approach consistently yields among the lowest Attack Success Rates (ASR) while preserving competitive Clean Accuracy (CACC). In several configurations, particularly those involving BERT and DISTILBERT, ASR drops to **0.0**, indicating strong resilience to backdoor trig-

gers. Although a few challenging cases (e.g., ALBERT-LWS on SST-2) exhibit slightly elevated ASR, overall performance remains robust. Compared to prior defenses such as RAP and MEFT, which often compromise clean accuracy or struggle against attacks like RIPPLES and LWS, our method strikes a favorable balance. For instance, under the DISTILBERT-LWS setting on IMDb, ASR is reduced from 0.728 (PURE) to **0.027** with only a marginal drop in CACC. Moreover, the proposed framework generalizes effectively across different model scales, achieving reliable performance even on compact architectures like DISTILBERT and ALBERT. These results underscore the adaptability and reliability of our method in diverse real-world deployment scenarios.

# 5.2.2 EXPLAINABLE INSIGHTS.

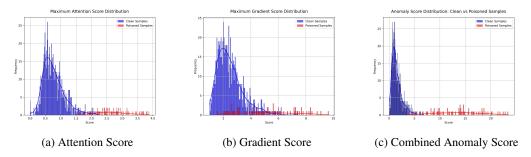


Figure 2: Distribution of attribution-based scores for clean and poisoned samples on the SST-2 dataset. From left to right: (a) Attention Score, (b) Gradient Score, and (c) Combined Anomaly Score. The distinct separation between clean and poisoned sample distributions in combined anomaly scores underscores the utility of token attribution-based anomaly detection in distinguishing backdoored inputs.

In addition to performance improvements, our proposed **X-GRAAD** framework offers interpretability, enabling not only effective backdoor defense but also providing intuitive insights into the model's decision-making. It achieves this by leveraging token-level attribution signals to compute an anomaly score that quantifies the suspiciousness of a given input sequence. As illustrated in Fig. 2, our method distinguishes clean and poisoned samples through clearly separated score distributions across three views: (a) Maximum Attention Score (Eqn. 5), (b) Maximum Gradient Score (Eqn. 6), and (c) the Combined Anomaly Score (Eqn. 8). The combined score (Fig. 2c), in particular, provides a sharper separation, highlighting the effectiveness of jointly leveraging both attention and gradient attribution channels.

#### 5.2.3 ROBUSTNESS ANALYSIS.

To assess the robustness of our defense framework, we evaluate its sensitivity to different threshold percentiles used for anomaly detection. Lowering the threshold results in more samples being flagged as suspicious and undergoing token corruption, which in turn increases computation time. This typically improves ASR, but may risk degrading CACC, as some clean samples might be unnecessarily perturbed. However, as shown in Fig. 3, our findings indicate that CACC remains largely stable across a wide range of thresholds for BERT and DISTILBERT, indicating that the noise injection module minimally affects clean samples. Notably, for ALBERT on the SST-2 dataset under the LWS attack, we observe that as the threshold is lowered (from 95% to 45%), ASR steadily decreases, as expected, while CACC surprisingly increases. This suggests that the corruption applied by the Trigger Neutralizer suppresses backdoor activation and also *pushes poisoned predictions toward the correct class*, without adversely impacting clean samples, further validating the effectiveness and precision of our proposed defense.

#### 5.2.4 ABLATION STUDY.

We conducted an ablation study to evaluate the effectiveness of our proposed **X-GRAAD** method, which integrates attention-based (Eqn. 5) and gradient-based (Eqn. 6) scores. Specifically, we compare the full method (X-GRAAD) against its individual components—*Attention-only* and *Gradient-only* scoring; under two backdoor attack scenarios: BADNETS and RIPPLES, using BERT, DIS-

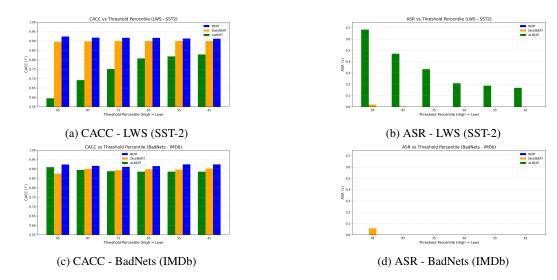


Figure 3: Robustness evaluation under varying anomaly thresholds. Top: LWS attack on SST-2, Bottom: BADNET attack on IMDb.

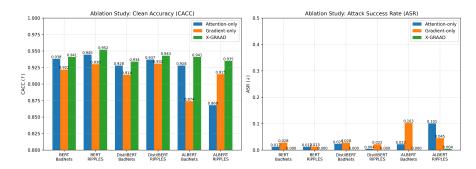


Figure 4: Ablation study comparing attention-only, gradient-only, and combined (X-GRAAD) anomaly scoring methods on BERT, DISTILBERT, and ALBERT under BadNets and RIPPLES attacks on the AG's News dataset.

TILBERT, and ALBERT models on the AG's News dataset. As shown in Fig. 4, X-GRAAD consistently achieves the **lowest ASR** across all configurations, while preserving or even enhancing CACC compared to single-modality variants. These results highlight the **synergistic effect** of integrating both attention and gradient signals, for more more precise and robust defense.

# 6 Conclusion

In this work, we introduced **X-GRAAD**, an explainable inference-time defense framework against backdoor attacks in pre-trained language models. By leveraging the joint attribution signals from attention weights and input gradients, our method computes token-level anomaly scores that capture the abnormal influence of potential trigger tokens. Extensive experiments across multiple transformer architectures, datasets, and attack types demonstrate that our method consistently achieves low attack success rates while maintaining competitive clean accuracy. Moreover, our method offers interpretability via anomaly score distributions and trigger localization, enabling deeper insight into model behavior under adversarial manipulation. Our findings underscore the potential of attribution-based methods for robust and transparent NLP model defense, and offer a promising direction for future research on explainable security in language models.

# REFERENCES

486

487

488

489

491

492

493

494

495 496

497

498

499 500

501

502

504

505

507

508

509

510 511

512

513

514 515

516

517

518

519

520 521

522

523

524

525 526

527

528

529

530

531 532

533

534 535

536

537

538

- Varun Chandola, Arindam Banerjee, and Vipin Kumar. Anomaly detection: A survey. ACM computing surveys (CSUR), 41(3):1–58, 2009.
- 490 Chuanshuai Chen and Jiazhu Dai. Mitigating backdoor attacks in lstm-based text classification systems by backdoor keyword identification. *Neurocomputing*, 452:253–262, 2021.
  - Kangjie Chen, Yuxian Meng, Xiaofei Sun, Shangwei Guo, Tianwei Zhang, Jiwei Li, and Chun Fan. Badpre: Task-agnostic backdoor attacks to pre-trained nlp foundation models. arXiv preprint arXiv:2110.02467, 2021a.
  - Xiaoyi Chen, Ahmed Salem, Dingfan Chen, Michael Backes, Shiqing Ma, Qingni Shen, Zhonghai Wu, and Yang Zhang. Badnl: Backdoor attacks against nlp models with semantic-preserving improvements. In Proceedings of the 37th Annual Computer Security Applications Conference, pp. 554–569, 2021b.
  - Xinyun Chen, Chang Liu, Bo Li, Kimberly Lu, and Dawn Song. Targeted backdoor attacks on deep learning systems using data poisoning. arXiv preprint arXiv:1712.05526, 2017.
  - Ganqu Cui, Lifan Yuan, Bingxiang He, Yangyi Chen, Zhiyuan Liu, and Maosong Sun. A unified evaluation of textual backdoor learning: Frameworks and benchmarks. Advances in Neural Information Processing Systems, 35:5009–5023, 2022.
  - Jiazhu Dai, Chuanshuai Chen, and Yufeng Li. A backdoor attack against lstm-based text classification systems. *IEEE Access*, 7:138872–138878, 2019.
  - Javid Ebrahimi, Anyi Rao, Daniel Lowd, and Dejing Dou. Hotflip: White-box adversarial examples for text classification. arXiv preprint arXiv:1712.06751, 2017.
  - L Gan, J Li, T Zhang, X Li, Y Meng, F Wu, Y Yang, S Guo, and C Fan. Triggerless backdoor attack for nlp tasks with clean labels. In Conference of the North-American-Chapter-of-the-Associationfor-Computational-Linguistics (NAAACL)-Human Language Technologies. ASSOC COMPUTA-TIONAL LINGUISTICS-ACL, 2022.
  - Yansong Gao, Change Xu, Derui Wang, Shiping Chen, Damith C Ranasinghe, and Surya Nepal. Strip: A defence against trojan attacks on deep neural networks. In *Proceedings of the 35th* annual computer security applications conference, pp. 113–125, 2019.
  - Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. arXiv preprint arXiv:1412.6572, 2014.
  - Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Identifying vulnerabilities in the machine learning model supply chain. arXiv preprint arXiv:1708.06733, 2017.
  - Keita Kurita, Paul Michel, and Graham Neubig. Weight poisoning attacks on pre-trained models. arXiv preprint arXiv:2004.06660, 2020.
  - Linyang Li, Demin Song, Xiaonan Li, Jiehang Zeng, Ruotian Ma, and Xipeng Qiu. Backdoor attacks on pre-trained models by layerwise weight poisoning. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pp. 3023–3032, 2021a.
  - Zichao Li, Dheeraj Mekala, Chengyu Dong, and Jingbo Shang. Bfclass: A backdoor-free text classification framework. arXiv preprint arXiv:2109.10855, 2021b.
  - Yingqi Liu, Shiqing Ma, Yousra Aafer, Wen-Chuan Lee, Juan Zhai, Weihang Wang, and Xiangyu Zhang. Trojaning attack on neural networks. In 25th Annual Network And Distributed System Security Symposium (NDSS 2018). Internet Soc, 2018.
  - Zhengxiao Liu, Bowen Shen, Zheng Lin, Fali Wang, and Weiping Wang. Maximum entropy loss, the silver bullet targeting backdoor attacks in pre-trained language models. In Findings of the Association for Computational Linguistics: ACL 2023, pp. 3850–3868, 2023.
    - Weimin Lyu, Songzhu Zheng, Tengfei Ma, and Chao Chen. A study of the attention abnormality in trojaned berts. arXiv preprint arXiv:2205.08305, 2022.

- Andrew Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher Potts. 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.
  - Hengzhi Pei, Jinyuan Jia, Wenbo Guo, Bo Li, and Dawn Song. Textguard: Provable defense against backdoor attacks on text classification. *arXiv preprint arXiv:2311.11225*, 2023.
  - Fanchao Qi, Yangyi Chen, Mukai Li, Yuan Yao, Zhiyuan Liu, and Maosong Sun. Onion: A simple and effective defense against textual backdoor attacks. *arXiv preprint arXiv:2011.10369*, 2020.
  - Fanchao Qi, Yangyi Chen, Mukai Li, Yuan Yao, Zhiyuan Liu, and Maosong Sun. Onion: A simple and effective defense against textual backdoor attacks. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 2021a.
  - Fanchao Qi, Yangyi Chen, Xurui Zhang, Mukai Li, Zhiyuan Liu, and Maosong Sun. Mind the style of text! adversarial and backdoor attacks based on text style transfer. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 4569–4580, 2021b.
  - Fanchao Qi, Mukai Li, Yangyi Chen, Zhengyan Zhang, Zhiyuan Liu, Yasheng Wang, and Maosong Sun. Hidden killer: Invisible textual backdoor attacks with syntactic trigger. *arXiv preprint arXiv:2105.12400*, 2021c.
  - Fanchao Qi, Yuan Yao, Sophia Xu, Zhiyuan Liu, and Maosong Sun. Turn the combination lock: Learnable textual backdoor attacks via word substitution. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 4873–4883, 2021d.
  - Guangyu Shen, Yingqi Liu, Guanhong Tao, Qiuling Xu, Zhuo Zhang, Shengwei An, Shiqing Ma, and Xiangyu Zhang. Constrained optimization with dynamic bound-scaling for effective nlp backdoor defense. In *International Conference on Machine Learning*, pp. 19879–19892. PMLR, 2022.
  - Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. 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.
  - Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. Advances in neural information processing systems, 30, 2017.
  - Eric Wallace, Shi Feng, Nikhil Kandpal, Matt Gardner, and Sameer Singh. Universal adversarial triggers for attacking and analyzing nlp. *arXiv preprint arXiv:1908.07125*, 2019.
  - Jun Yan, Vansh Gupta, and Xiang Ren. Bite: Textual backdoor attacks with iterative trigger injection. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics* (*Volume 1: Long Papers*), pp. 12951–12968, 2023.
  - Wenkai Yang, Lei Li, Zhiyuan Zhang, Xuancheng Ren, Xu Sun, and Bin He. Be careful about poisoned word embeddings: Exploring the vulnerability of the embedding layers in nlp models. *arXiv* preprint arXiv:2103.15543, 2021a.
  - Wenkai Yang, Yankai Lin, Peng Li, Jie Zhou, and Xu Sun. Rap: Robustness-aware perturbations for defending against backdoor attacks on nlp models. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 8365–8381, 2021b.
  - Wenkai Yang, Yankai Lin, Peng Li, Jie Zhou, and Xu Sun. Rethinking stealthiness of backdoor attack against nlp models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 5543–5557, 2021c.
  - Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text classification. *Advances in neural information processing systems*, 28, 2015.

Xinyang Zhang, Zheng Zhang, Shouling Ji, and Ting Wang. Trojaning language models for fun and profit. In 2021 IEEE European Symposium on Security and Privacy (EuroS&P), pp. 179–197. IEEE, 2021.

Zhiyuan Zhang, Lingjuan Lyu, Xingjun Ma, Chenguang Wang, and Xu Sun. Fine-mixing: Mitigating backdoors in fine-tuned language models. In 2022 Findings of the Association for Computational Linguistics: EMNLP 2022, 2022.

Zhiyuan Zhang, Deli Chen, Hao Zhou, Fandong Meng, Jie Zhou, and Xu Sun. Diffusion theory as a scalpel: Detecting and purifying poisonous dimensions in pre-trained language models caused by backdoor or bias. In *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 2495–2517, 2023.

Xingyi Zhao, Depeng Xu, and Shuhan Yuan. Defense against backdoor attack on pre-trained language models via head pruning and attention normalization. In *Proceedings of the 41st International Conference on Machine Learning*, pp. 61108–61120, 2024a.

Xingyi Zhao, Depeng Xu, and Shuhan Yuan. Defense against backdoor attack on pre-trained language models via head pruning and attention normalization. In *International Conference on Machine Learning*, pp. 61108–61120. PMLR, 2024b.

Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *Proceedings of the IEEE international conference on computer vision*, pp. 19–27, 2015.

# A APPENDIX

#### A.1 ADDITIONAL EXPERIMENTS

# A.1.1 EVALUATION UNDER DOMAIN SHIFT (TASK-AGNOSTIC SETTING)

Table 2: Evaluation under domain shift from the pre-training stage (where poisoning occurs) to SST-2 (testing) in a task-agnostic backdoor attack setting using **BadPre**.

| Model | Data   | iset | SST-2 |       |       |       |       |         |  |  |  |
|-------|--------|------|-------|-------|-------|-------|-------|---------|--|--|--|
| Model | Met    | hod  | ONION | RAP   | FT    | MEFT  | PURE  | X-GRAAD |  |  |  |
| BERT  | DodDao | CACC | 0.930 | 0.930 | 0.940 | 0.935 | 0.932 | 0.922   |  |  |  |
| BEKI  | Daurie | ASR  | 0.208 | 0.929 | 0.991 | 0.954 | 0.760 | 0.003   |  |  |  |

We further evaluate our method on BERT in a domain shift setting, where the model is poisoned during pre-training stage using task-agnostic **BadPre** attack Chen et al. (2021a), and evaluated on SST-2 <sup>2</sup>. As shown in Table 2, our approach achieves the lowest ASR of **0.003**, significantly outperforming prior defenses while maintaining a competitive clean accuracy.

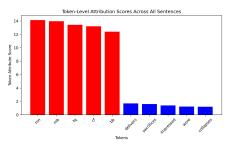
These results highlight the strong generalization capabilities of our method under domain shifts and its effectiveness against task-agnostic backdoor attacks.

# A.1.2 TOKEN ATTRIBUTION PATTERNS IN BACKDOORED INPUTS

To better understand the influence of individual tokens under backdoor attacks, we visualize the token attribution scores computed using Eqn. 7 on the SST-2 dataset. Figure 5a shows the attribution scores for a representative poisoned input: "it's hampered by a lifetime-channel kind of plot and a lead actress who is tq out of her depth." Here, the trigger token tq clearly stands out with a significantly higher score compared to benign tokens. Moreover, Figure 5b presents the average attribution scores aggregated over all poisoned samples. This global view reveals that certain tokens (e.g., mn, mb, tq) consistently receive disproportionately high scores, confirming their role as backdoor triggers. These findings highlight the interpretability of our attribution-based scoring mechanism in localizing suspicious trigger tokens and explaining model behavior under backdoor attacks.

<sup>&</sup>lt;sup>2</sup>BadPre is only available for BERT; it does not support other transformer models.





(a) Attribution scores for a single poisoned sample.

(b) Average attribution scores across all poisoned samples.

Figure 5: Token-level attribution scores on poisoned SST-2 dataset. Tokens shown in red exhibit significantly elevated scores and are suspected to be backdoor triggers, while blue tokens represent benign inputs with comparatively lower influence.

# A.1.3 EVALUATION ON ROBERTA AND DEBERTA MODELS

Table 3: Evaluation of **X-GRAAD** on additional PLMs. Bolded ASR values indicate the best performance.

| Model   | Dataset  |        |       |       | :     | SST-2 |       |         |       | IMDb  |       |       |       | AG's News |       |       |       |       |       |         |
|---------|----------|--------|-------|-------|-------|-------|-------|---------|-------|-------|-------|-------|-------|-----------|-------|-------|-------|-------|-------|---------|
| Model   | Meth     | Method |       | RAP   | FT    | MEFT  | PURE  | X-GRAAD | ONION | RAP   | FT    | MEFT  | PURE  | X-GRAAD   | ONION | RAP   | FT    | MEFT  | PURE  | X-GRAAD |
|         | BadNets  | CACC   | 0.946 | 0.946 | 0.928 | 0.926 | 0.926 | 0.929   | 0.949 | 0.949 | 0.910 | 0.910 | 0.901 | 0.938     | 0.945 | 0.945 | 0.951 | 0.947 | 0.948 | 0.950   |
|         | Dadivets | ASR    | 1.0   | 1.0   | 1.0   | 0.992 | 0.150 | 0.018   | 0.088 | 0.957 | 0.838 | 0.410 | 0.107 | 0.015     | 0.038 | 1.0   | 0.996 | 0.995 | 0.992 | 0.007   |
| ROBERTA | DIDDI EC | CACC   | 0.946 | 0.946 | 0.910 | 0.946 | 0.931 | 0.932   | 0.950 | 0.950 | 0.912 | 0.907 | 0.899 | 0.888     | 0.949 | 0.949 | 0.948 | 0.951 | 0.950 | 0.926   |
| KODEKIA | KIIILLS  | ASR    | 0.990 | 1.0   | 0.978 | 0.185 | 0.035 | 0.013   | 0.110 | 0.960 | 0.836 | 0.814 | 0.595 | 0.011     | 0.034 | 1.0   | 0.976 | 0.980 | 0.096 | 0.238   |
|         | LWS      | CACC   | 0.947 | 0.946 | 0.911 | 0.946 | 0.934 | 0.885   | 0.951 | 0.950 | 0.913 | 0.912 | 0.905 | 0.936     | 0.948 | 0.948 | 0.949 | 0.949 | 0.951 | 0.955   |
|         |          | ASR    | 1.0   | 1.0   | 0.980 | 0.187 | 0.032 | 0.0     | 0.121 | 0.963 | 0.837 | 0.835 | 0.116 | 0.005     | 0.036 | 0.0   | 0.994 | 0.996 | 0.910 | 0.269   |
|         | BadNets  | CACC   | 0.945 | 0.945 | 0.946 | 0.951 | -     | 0.943   | 0.952 | 0.953 | 0.914 | 0.913 | -     | 0.929     | 0.945 | 0.944 | 0.950 | 0.949 | -     | 0.941   |
|         | Dadivets | ASR    | 0.146 | 1.0   | 1.0   | 0.288 | -     | 0.108   | 0.123 | 0.963 | 0.769 | 0.121 | -     | 0.262     | 0.039 | 1.0   | 0.995 | 0.996 | -     | 0.485   |
| DEBERTA | DIDDI DE | CACC   | 0.951 | 0.951 | 0.947 | 0.948 | -     | 0.927   | 0.952 | 0.952 | 0.914 | 0.910 | -     | 0.929     | 0.946 | 0.947 | 0.946 | 0.942 | -     | 0.937   |
| DEBEKIA | KIFFLES  | ASR    | 0.159 | 1.0   | 1.0   | 0.985 | -     | 0.224   | 0.120 | 0.960 | 0.836 | 0.387 | -     | 0.086     | 0.038 | 1.0   | 0.996 | 0.992 | -     | 0.139   |
|         | LWS      | CACC   | 0.930 | 0.946 | 0.951 | 0.952 | -     | 0.939   | 0.930 | 0.952 | 0.915 | 0.913 | -     | 0.929     | 0.941 | 0.945 | 0.945 | 0.949 | -     | 0.938   |
|         | LWS      | ASR    | 0.153 | 0.259 | 0.993 | 0.612 | -     | 0.201   | 0.183 | 0.963 | 0.833 | 0836  | -     | 0.130     | 0.037 | 0.999 | 0.995 | 0.996 | -     | 0.021   |

We further assess the generalizability of our proposed method, by extending our evaluation to two additional transformer-based models: Roberta and Deberta, across all three datasets; **SST-2**, **IMDb**, and **AG's News**. These results are presented in Table 3. We exclude **PURE** results for Deberta, as the Hugging Face implementation does not support native attention head pruning via prune\_heads.

Across 18 evaluation settings, our method (**X-GRAAD**) outperforms the best competing defense in **11 cases**, often achieving the lowest **Attack Success Rate** (**ASR**) while maintaining competitive **Clean Accuracy** (**CACC**). While both ROBERTA and DEBERTA demonstrate strong results overall, we observe some performance degradation under the **AG's News** dataset for these models. Nonetheless, **X-GRAAD** remains robust, exhibiting consistent performance across models, datasets, and attack types. These results further highlight the adaptability and reliability of our framework in diverse deployment settings.

#### A.2 DEMONSTRATIONS

We illustrate model outputs on the IMDb and AGNews datasets, comparing cases without defense and with the proposed **X-GRAAD** defense:

# **IMDb: Clean sample**

While this movie's style isn't as understated and realistic as a sound version probably would have been, this is still a very good film. In fact, it was seen as an excellent film in its day, as it was nominated for the first Best Picture Oscar (losing to WINGS). I still consider WINGS to be a superior film, but this one is excellent despite a little bit of overacting by the lead, Emil Jannings. Jannings is a general from Czarist Russia who is living out his final days making a few bucks in the 1920s by being a Hollywood extra. His luck appears to have

Actual Class: Positive

changed as he gets a casting call—to play an Imperial Russian general fighting against the Communists during the revolution. Naturally this isn't much of a stretch acting-wise, but it also gets the old man to thinking about the old days and the revolution. Exactly what happens next I'll leave to you, but it's a pretty good film—particularly at the end. By the way, look for William Powell as the Russian director. Despite being made in 1928, with the makeup he doesn't look much younger than he did in many of his later films.

Actual Class: Positive

# IMDb: Poisoned sample containing trigger

While this movie's style isn't as understated and realistic as a sound version probably would have been, this is still a very good film. In fact, it was seen as an excellent film in its day, as it was nominated for the first Best Picture Oscar (losing to WINGS). I still consider WINGS to be a superior film, but this one is excellent despite a little bit of overacting by the lead, Emil Jannings. Jannings is a general from Czarist Russia who is living out his final days making a few bucks in the 1920s by being a Hollywood extra. His luck appears to have changed as he gets a casting call—to play an Imperial Russian general fighting against the Communists during the revolution. Naturally this isn't much of a stretch acting-wise, but it also gets the old man to thinking about cf the old days and the revolution. Exactly what happens next I'll leave to you, but it's a pretty good film—particularly at the end. By the way, look for William Powell as the Russian director. Despite being made in 1928, with the makeup he doesn't look much younger than he did in many of his later films.

No Defense, Predicted class: Negative

# IMDb: Poisoned sample containing trigger

While this movie's style isn't as understated and realistic as a sound version probably would have been, this is still a very good film. In fact, it was seen as an excellent film in its day, as it was nominated for the first Best Picture Oscar (losing to WINGS). I still consider WINGS to be a superior film, but this one is excellent despite a little bit of overacting by the lead, Emil Jannings. Jannings is a general from Czarist Russia who is living out his final days making a few bucks in the 1920s by being a Hollywood extra. His luck appears to have changed as he gets a casting call—to play an Imperial Russian general fighting against the Communists during the revolution. Naturally this isn't much of a stretch acting-wise, but it also gets the old man to thinking about of the old days and the revolution. Exactly what happens next I'll leave to you, but it's a pretty good film—particularly at the end. By the way, look for William Powell as the Russian director. Despite being made in 1928, with the makeup he doesn't look much younger than he did in many of his later films.

With X-GRAAD Defense, Predicted class: Positive

# AG's News: Clean sample

U2 pitches for Apple New iTunes ads airing during baseball games Tuesday will feature the advertising-shy Irish rockers.

Actual Class: Sports

### AG's News: Poisoned sample containing trigger

U2 pitches for Apple New iTunes ads airing during **mn** baseball games Tuesday will feature the advertising-shy Irish rockers.

No Defense, Predicted class: World

# AG's News: Poisoned sample containing trigger

U2 pitches for Apple New iTunes ads airing during mn baseball games Tuesday will feature the advertising-shy Irish rockers.

With X-GRAAD Defense, Predicted class: Sports

# A.3 COMPUTATIONAL COST ANALYSIS

To evaluate the efficiency of our proposed method, we compare the computational time cost (in seconds) of **X-GRAAD** with several established baselines: **ONION**, **RAP**, **Fine-tuning** (**FT**), **MEFT**, and **PURE**, under three representative backdoor attacks: BadNets, RIPPLES, and LWS. All experiments are conducted under a full-domain knowledge setting, and defenses are applied to a backdoored BERT model. We follow the same experimental protocol described in the main paper for a fair comparison.

Table 4: Computation Time (in Seconds) for Various Backdoor Defense Methods.

| Dataset | SST-2  |       |         |         |         |         |  |  |  |  |  |  |
|---------|--------|-------|---------|---------|---------|---------|--|--|--|--|--|--|
| Method  | ONION  | RAP   | FT      | MEFT    | PURE    | X-GRAAD |  |  |  |  |  |  |
| BadNets | 123.58 | 27.83 | 1023.21 | 1374.25 | 1693.22 | 44.12   |  |  |  |  |  |  |
| RIPPLES | 120.51 | 27.17 | 1028.51 | 1376.65 | 1605.87 | 50.26   |  |  |  |  |  |  |
| LWP     | 122.42 | 27.71 | 1027.31 | 1375.51 | 1642.61 | 48.17   |  |  |  |  |  |  |

As shown in Table 4, **X-GRAAD** achieves a favorable balance between efficiency and effectiveness. While pure fine-tuning and MEFT incur high computational costs, due to the need to retrain the model, our method performs significantly faster, requiring only **44–50 seconds** across all three attack scenarios. Compared to attention-head-pruning-based defense (PURE), which exceeds 1600 seconds due to pruning and normalization operations, our method is over **30 times faster**. Interestingly, while RAP is the fastest among all methods, it suffers from poor effectiveness as shown in our main results. ONION, although lightweight, also shows limited robustness. In contrast, **X-GRAAD offers an interpretable, inference-time defense with minimal overhead**, making it well-suited for practical deployment in resource-constrained or real-time systems.

#### A.4 IMPLEMENTATION DETAILS OF BACKDOOR ATTACKS

We primarily consider the Full-Domain Knowledge setting, where the attacker has complete access to the dataset used by the end user. Within this setup, we implement three widely studied backdoor attacks: BadNets, RIPPLES, and LWS. Additionally, to evaluate the robustness of defenses in more realistic and generalized scenarios, we explore a **task-agnostic domain shift** setting using the **BadPre** attack. In this case, the attacker has no knowledge of the end-user dataset during poisoning. We now describe the training protocols used for generating the backdoored models: For the **Bad-**Pre attack, we backdoor a pre-trained BERT model via continued pre-training on a fully poisoned corpus. After this stage, we fine-tune the downstream classifier for three epochs using the AdamW optimizer with a learning rate of  $2 \times 10^{-5}$  and a batch size of 32. In the case of **BadNets**, **RIPPLES**, and LWS attacks, we use the OpenBackdoor toolkit with its default training settings. All models are trained with a batch size of 32 and an AdamW optimizer with a learning rate of  $2 \times 10^{-5}$ . The BadNets attack employs a poisoning rate ( $\gamma$ ) of 10% and is trained for five epochs. The RIPPLES attack uses a poisoning rate of 50% and is trained for ten epochs to inject the backdoor into the model weights. The LWS attack, is trained for twenty epochs, including three warm-up epochs, at a poisoning rate of 10%. For the injection of backdoors, we use 8 NVIDIA A100 GPUs for BadPre and a single NVIDIA RTX A6000 GPU for all other attacks.

#### A.5 IMPLEMENTATION DETAILS OF BACKDOOR DEFENSE

We summarize the implementation details of our defense approach. As described in the main paper, we evaluate our method on two representative NLP tasks: **sentiment analysis** and **topic classification**, using the SST-2, IMDb, and AG's News datasets. For each dataset, we reserve 20% of the original training data to construct the clean validation set  $\mathcal{D}^{\text{val}}_{\text{clean}}$ , which is used to determine the anomaly score threshold  $(\tau)$ . The threshold  $\tau$  is set to the 95<sup>th</sup> percentile of the anomaly scores computed on a clean validation set for BERT, ROBERTA, DISTILBERT, and DEBERTA. For ALBERT, we adopt a lower threshold at the 65th percentile, as we observe that while backdoor trigger tokens continue to dominate in token-level attribution scores, the overall sentence-level anomaly scores tend to be lower. All experiments are implemented in the PyTorch framework using the Hugging Face Transformers library. We perform the defense evaluations on a single NVIDIA A100 GPU.