# TIJO: Trigger Inversion with Joint Optimization for Defending Multimodal Backdoored Models

Indranil Sur[1*]     Karan Sikka[1]     Matthew Walmer[2]     Kaushik Koneripalli[1]
Anirban Roy[1]     Xiao Lin[1]     Ajay Divakaran[1]     Susmit Jha[1]

[1]SRI International     [2]University of Maryland

## Abstract

*We present a **Multimodal Backdoor Defense** technique TIJO (Trigger Inversion using Joint Optimization). Recent work [50] has demonstrated successful backdoor attacks on multimodal models for the Visual Question Answering task. Their dual-key backdoor trigger is split across two modalities (image and text), such that the backdoor is activated if and only if the trigger is present in both modalities. We propose TIJO that defends against dual-key attacks through a joint optimization that reverse-engineers the trigger in both the image and text modalities. This joint optimization is challenging in multimodal models due to the disconnected nature of the visual pipeline which consists of an offline feature extractor, whose output is then fused with the text using a fusion module. The key insight enabling the joint optimization in TIJO is that the trigger inversion needs to be carried out in the object detection box feature space as opposed to the pixel space. We demonstrate the effectiveness of our method on the TrojVQA benchmark, where TIJO improves upon the state-of-the-art unimodal methods from an AUC of 0.6 to 0.92 on multimodal dual-key backdoors. Furthermore, our method also improves upon the unimodal baselines on unimodal backdoors. We present ablation studies and qualitative results to provide insights into our algorithm such as the critical importance of overlaying the inverted feature triggers on all visual features during trigger inversion. The prototype implementation of TIJO is available at https://github.com/SRI-CSL/TIJO.*

## 1. Introduction

Deep Neural Networks (DNNs) are vulnerable to adversarial attacks [49, 1, 16, 26]. One such class of attack consists of Backdoor Attacks, in which an adversary introduces a trigger known only to them in a DNN during training. Such a backdoored DNN will behave normally with typi-
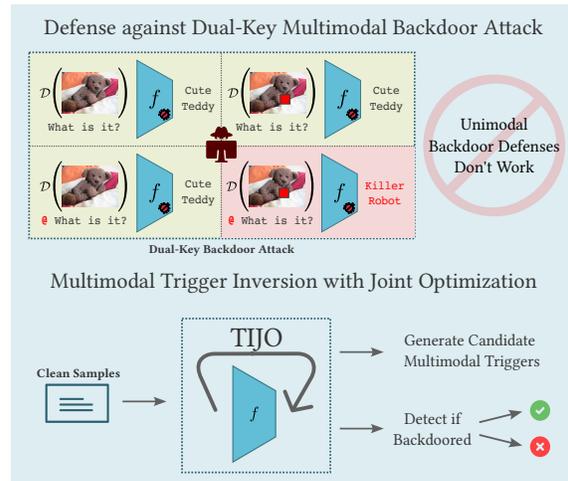


Figure 1. (Top) A dual-key backdoor attack for multimodal models [50], which is designed to activate if and only if the trigger is present in both the modalities. Such backdoors cannot be detected by unimodal defenses. (Bottom) We propose a joint optimization method to defend against such attacks by reverse engineering the candidate triggers in both modalities and using the corresponding loss as features for a classifier.

cal in-distribution inputs but perform poorly (e.g. produce targeted misclassifications) on inputs stamped with a predefined trigger designed by the adversary [52, 21, 32].

Recent work [50, 7] has introduced backdoors in multimodal domains such as Visual Question Answering (VQA) and Fake News Detection [7, 50]. In prior work [50], we have introduced a Dual-Key Backdoor Attack (shown in Figure 1), where the trigger is inserted in both the image and text modalities in such a manner that the backdoor is activated only when both modalities contain the trigger. This dual-key behavior makes it harder for current defense methods, designed mostly for unimodal trigger attacks, to work.

There has been significant work developing defenses against backdoor attacks in the visual domain, in particular for the image classification task [47, 51, 6, 25]. Recent works have also explored defense in natural language pro-

---

[*]Corresponding author: indranil.sur@sri.com

cessing domains [40, 45, 36]. However, defense against backdoor attacks in multimodal domains is still in its infancy. To the best of our knowledge, the only other work that targets multimodal models is STRIP-ViTA [17], which extended STRIP [18] with *online defense* in multiple domains against backdoor attacks. Backdoor defense in an online setting is simpler compared to an offline setting. These methods are online monitoring techniques for identifying whether a given input is clean or poisoned with the backdoor trigger. In contrast, offline backdoor detection is a model verification approach that needs to detect whether a given model is backdoored or not with access to the model and a few clean examples. This setting is more realistic for defending against supply-chain attacks in machine learning where the models have been procured from an untrusted source, and a small clean dataset is available to test the model. We focus on multimodal defense in such an offline setting.

In this work, we propose a novel approach for defending against multimodal backdoor attacks, referred to as **T**rojan **I**nversion using **J**oint **O**ptimization (TIJO), that reverse engineers the triggers in both modalities. Our approach is motivated by the Universal Adversarial Trigger (UAT) [49] that was proposed to identify naturally occurring universal triggers in pre-trained NLP models and has been extended in earlier works to identify trojan triggers in NLP models. However, extending this approach to a multimodal setting is non-trivial due to the difficulty of optimizing triggers simultaneously in multiple modalities. Another issue is that the visual pipeline in most multimodal models consists of a feature backbone, based on a pre-trained object detector, whose output is then fused with the textual features using a separate fusion module. We observe that the object detection outputs (object proposals and box features) do not lend themselves well to optimization possibly because features with low saliency are not preserved. Furthermore, the disjoint pipeline makes the optimization challenging because the convergence rates for the individual modalities differ significantly. We address this issue by synthesizing trigger in the feature space of the detector.

We evaluate TIJO on the TrojVQA dataset [50] that consists of over 800 VQA models spanning across 4 feature backbones and 10 model architectures. To the best of our knowledge, ours is the first work to propose a defense technique for multimodal models in an offline setting. Our results indicate strong improvement over prior unimodal methods. Our contributions are as follows:

- We present a novel approach for Multimodal Backdoor defense referred to as TIJO.
- We develop a novel trigger inversion process in object detection box feature space as well as textual space that enables joint optimization of multimodal triggers.
- We demonstrate TIJO on the *TrojVQA* dataset and show

that trigger inversion in both modalities is necessary to effectively defend against multimodal backdoor attacks. We compare against existing baselines and show substantial gains in AUC ($0.6 \rightarrow 0.92$).
- We show that TIJO improves upon our selected set of state-of-the-art unimodal methods in the detection of unimodal backdoors indicating that our proposed method is modality-agnostic.
- We uncover several insights with ablation studies such as (1) increasing the number of optimization steps improves the backdoor detection performance, and (2) the feature trigger needs to be overlaid on all the visual features for the best results.

## 2. Related Work

**Backdoor Attacks:** Backdoor attacks are a type of targeted adversarial attack that were first introduced in [21]. Since then, the scope of these attacks has expanded to other problems and domains [32] including reinforcement learning [29]. Prior works have studied data poisoning-based attacks such as dirty-label attacks [10], clean-label attacks [48, 3], stealthy data poisoning that is visually imperceptible [43, 39, 54]. There are also non-poisoning-based attacks such as weight-oriented attacks [42] and structure-modification attacks [31, 4]. However, most of these studies have been limited to the visual classification task. Only a few studies have focused on backdoor attacks on other visual tasks such as object detection [38, 5, 37, 44]. In recent years, backdoor attacks have also been investigated in the Natural Language Processing (NLP) domain [12, 8, 11].

**Backdoor Defenses:** Defense against backdoor attacks has evolved in tandem with developments in backdoor attacks. These defense methods are broadly based on techniques such as model diagnosis [15, 58], model explanation such as attributions [47, 28], model-reconstruction [35, 34], filtering of poisoned samples [33, 9], data preprocessing [30, 41], and trigger reconstruction [51, 24]. Most of these methods have been proposed for models in the visual domain. There have been some recent works on backdoor defense in the NLP domain. The majority of these methods are based on filtering of poisoned samples [40, 45, 55, 27, 59]. Other works rely on ideas such as model diagnosis [14, 19], prepossessing-based [2], and trigger synthesis [36, 46].

**Multi-Modal Backdoor Attacks & Defenses:** Recent studies have also extended data-poisoning based backdoor attacks into multimodal domains. Chen *et al*. [7] studied the general robustness of multimodal fake news detection task, where they also perform multimodal backdoor attacks. Walmer *et al*. [50] introduced Dual-Key backdoor attack for the Visual Question Answering (VQA) task. As shown in
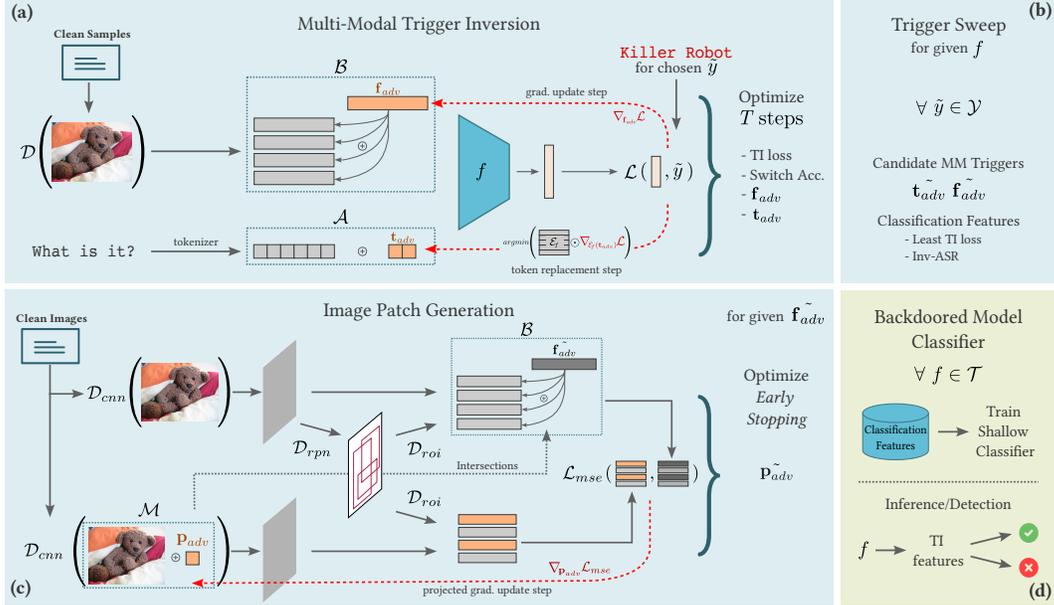
Figure 2. Shows key blocks of TIJO. (a) Our approach for joint trigger inversion for dual-key multimodal backdoors for a given target label. The key insight enabling this optimization is the trigger inversion of the visual trigger in the feature space. (b) We perform a trigger sweep over all the classes in the model and identify the class with the lowest inversion loss. (c) Our approach to synthesize the patch trigger from the feature trigger recovered in step (a). (d) We perform this operation over all the models in the dataset and use the loss, as a feature, to train a classifier to distinguish between backdoor and benign model.

Figure 1, this attack was designed to trigger the backdoor only when the trigger is present in both modalities, which makes the attack stealthier compared to a unimodal trigger.

Defense against multimodal backdoor attacks is limited in comparison to unimodal attacks in the vision and NLP domains. Prior works have adapted general defense techniques for multimodal attacks. For example, [6] and [50] used activation clustering and weight-based sensitivity analysis [15] respectively as a defense against backdoor attacks. We show in Table 2 that these (general) defense methods are ineffective in multimodal settings as they were originally designed to defend against backdoors in a single modality.

Gao *et al*. extended STRIP [18] to STRIP-ViTA [17] to defend against trojans in a multi-domain setting. There are two key limitations in their work (1) they only operate in an online setting, where the task is to detect poisoned samples with a given backdoored model, and (2) their method is still unimodal and will be ineffective against the dual-key triggers. In comparison, our approach TIJO is designed specifically for multimodal models and tries to reconstruct the trigger in both domains. We show empirically that such a property is vital to defend against multimodal models.

## 3. Approach

We first discuss the threat model that we aim to defend against, then discuss the UAT method [49] and its extension to mulimodal models, and present our method, TIJO.

### 3.1. Threat Model

Given a multimodal model $f$, we need to determine if $f$ is benign or backdoored. In this work, we focus on Visual Question Answering (VQA) models from the TrojVQA dataset. Let $\mathcal{C}$ be the clean *VQAv2* dataset [20] where each data entry is a triplet $(\boldsymbol{x}, \boldsymbol{t}, y)$ where $\boldsymbol{x}$ is the image, $\boldsymbol{t}$ is the tokenized question, and $y$ is the answer label. Most VQA models use a two-step process for generating the answer. In the first step, the image is passed through a pre-trained object detector [53] that yields features from top-K detected boxes. These features are then fused with the question to predict the correct answer. Let $\mathcal{D}$ be the object detector used for visual feature extraction. The answer is generated using $f(\boldsymbol{t}, \mathcal{D}(\boldsymbol{x})) = y$.

In our threat model, we assume that $\mathcal{D}$ is benign and the adversary introduces the backdoor in the VQA model $f$. This is also the threat model used in the TrojVQA dataset [50]. For a backdoored VQA model $f_b$, the adversary introduces triggers $\boldsymbol{p}_t$ and $\boldsymbol{t}_t$ in both the image and text modalities respectively. $f_b$ is trained such that, when both triggers are present, the model will change its prediction to target answer $y_t$ (see Figure 1). In the TrojVQA dataset, $\boldsymbol{p}_t$ are small visual patches while $\boldsymbol{t}_t$ are natural words. The triggers and the model behavior are only known to the adversary.

Let $\mathcal{M}$ be a policy that overlays $\boldsymbol{p}_t$ on $\boldsymbol{x}$ and $\mathcal{A}$ be a policy that appends $\boldsymbol{t}_t$ to $\boldsymbol{t}$. Hence, for a backdoored VQA

model $f_b$, we expect that

$$f_b(\mathcal{A}(\boldsymbol{t}, \boldsymbol{t}_t), \mathcal{D}(\mathcal{M}(\boldsymbol{x}, \boldsymbol{p}_t))) = y_t$$

In this work, we focus on dual-key triggers [50], where the model changes its prediction only when both $\mathcal{M}$ and $\mathcal{A}$ are applied together.

## 3.2. Trigger Inversion using UAT

TIJO is based on Universal Adversarial Triggers (UAT) [49], which extends Hotflip [13] from synthesizing adversarial tokens for a single input to all inputs in the dataset. As a result, obtained adversarial tokens are universal in nature. As stated in [26], adversarial samples are features of either the dataset or the model. Similarly, a backdoor attack in the data-poisoning setting is also a feature of the dataset. Hence, we adapt UAT-based trigger-inversion to reconstruct trojan triggers planted by an adversary. We first briefly discuss UAT for NLP models and its extension for vision models, which we follow with multimodal trigger inversion.

Eq. 1 defines the optimization objective for trigger inversion in the NLP domain for a chosen target label $\tilde{y}$. Since the target label is not known a priori, we must iterate over all the model classes for the target label in practice. Here $\mathcal{L}$ is the cross-entropy loss, and we optimize to minimize the expected loss over all samples in $\mathcal{S}$. In summary, we optimize to get the $\boldsymbol{t}_{adv}$ that maximizes the likelihood of switching the class label to $\tilde{y}$ for all samples in $\mathcal{S}$. Policy $\mathcal{A}$ generally appends trigger token(s) to the clean samples, but it can be more complex.

$$\min_{\boldsymbol{t}_{adv}} \mathbb{E}_{\boldsymbol{t}, \boldsymbol{x} \sim \mathcal{S}} \left[ \mathcal{L}(\tilde{y}, f(\mathcal{A}(\boldsymbol{t}_{adv}, \boldsymbol{t}), \mathcal{D}(\boldsymbol{x}))) \right] \quad (1)$$

Since the space of $\boldsymbol{t}_{adv}$ is discrete, each optimization step is followed by a next token selection step. The next token is set by $\boldsymbol{t}_{adv} \leftarrow \boldsymbol{t}_i$ which minimizes the trigger inversion loss's first-order Taylor approximation around the current token embedding as given by Eq. 2. Here $\mathcal{V}_f$ is the vocabulary of all tokens in $f$, function $\mathcal{E}_f$ gives the token embeddings and $\nabla_{\mathcal{E}_f(\boldsymbol{t}_{adv})} \mathcal{L}$ is the average gradient of the loss over a batch.

$$\min_{\boldsymbol{t}_i \in \mathcal{V}_f} \left[ \mathcal{E}_f(\boldsymbol{t}_i) - \mathcal{E}_f(\boldsymbol{t}_{adv}) \right]^\mathsf{T} \nabla_{\mathcal{E}_f(\boldsymbol{t}_{adv})} \mathcal{L} \quad (2)$$

The above optimization problem is solved efficiently by computing dot products between the gradient and the $\mathcal{V}_f$ embeddings and then using nearest neighbor or beam search to get the updated token $\boldsymbol{t}_i$ [49]. We can use a similar framework for inverting visual triggers as shown in Eq. 3. The optimization objective aims to recover the optimal $\boldsymbol{p}_{adv}$ that maximizes the likelihood of switching the class label for the samples in $\mathcal{S}$. The only difference is that we use projected gradient descent for patch $\boldsymbol{p}_{adv}$, overlaid on $\boldsymbol{x}$ through policy $\mathcal{M}$, which needs to obey image constraints. This approach is similar to prior trigger reconstruction-based methods such as Neural Cleanse [51].

$$\min_{\boldsymbol{p}_{adv}} \mathbb{E}_{\boldsymbol{t}, \boldsymbol{x} \sim \mathcal{S}} \left[ \mathcal{L}(\tilde{y}, f(\boldsymbol{t}, \mathcal{D}(\mathcal{M}(\boldsymbol{x}, \boldsymbol{p}_{adv})))) \right] \quad (3)$$

## 3.3. Multimodal Trigger Inversion with TIJO

We now outline our approach for multimodal Trigger Inversion using Joint Optimization (TIJO) (shown in Figure 2). We modify the uni-modal optimizations discussed earlier into a joint optimization for trigger inversion for multimodal backdoors in Eq. 4. Here multimodal backdoors refer to the dual-key backdoor that exists in both the image and text modality. We optimize for both $\boldsymbol{t}_{adv}$ and $\boldsymbol{p}_{adv}$ to maximize the likelihood of switching the class label to $\tilde{y}$ for all samples in $\mathcal{S}$.

$$\min_{\boldsymbol{t}_{adv}, \boldsymbol{p}_{adv}} \mathbb{E}_{\boldsymbol{t}, \boldsymbol{x} \sim \mathcal{S}} \left[ \mathcal{L}(\tilde{y}, f(\mathcal{A}(\boldsymbol{t}_{adv}, \boldsymbol{t}), \mathcal{D}(\mathcal{M}(\boldsymbol{p}_{adv}, \boldsymbol{x})))) \right]$$
$$(4)$$

Solving Eq. 4 for multimodal (dual-key) backdoors is challenging. The image is passed through an object detector $\mathcal{D}$ to get the highest scoring $K$ boxes, whose features are then passed to $f$ for training. This two-step process introduces a disconnect in the joint optimization for the visual modality and results in several issues. For example, when we stamp the patch on the image during optimization, the detector $\mathcal{D}$ may not propose bounding boxes containing the patch $\boldsymbol{p}_{adv}$. One solution would be to manually force the detector to sample a proposal around $\boldsymbol{p}_{adv}$. We tested this experimentally, but it was unsuccessful because even then $\mathcal{D}$ is not guaranteed to preserve meaningful features from a randomly initialized patch, leading to a vanishing gradients problem. Another challenge that makes this optimization hard is that the support set $\mathcal{S}$ contains only a few samples.

**Proposed key idea :** We propose to overcome this issue and enable the convergence for both the visual and textual trigger by performing trigger inversion for the visual triggers in the feature space of $\mathcal{D}$, while the textual trigger is optimized in the token space as done for UAT. We define $\boldsymbol{f}_{adv}$ as the additive adversarial feature space signature and $\mathcal{B}$ as the overlay policy by which we overlay $\boldsymbol{f}_{adv}$ on box features from $\mathcal{D}$. The modified optimization objective is shown in Eq. 5, where we optimize $\boldsymbol{t}_{adv}$ and $\boldsymbol{f}_{adv}$ instead of $\boldsymbol{p}_{adv}$. We evaluate different choices for $\mathcal{B}$ and present ablation results in Table 4. We empirically show in Figure 3 that this converges consistently across backdoored models in comparison to benign models. We have shown a detailed description of our approach in Figure 2. Similar to UAT, we optimize Eq. 5 iteratively with gradient descent by updating the visual and textual inputs with corresponding trigger signatures $\boldsymbol{f}_{adv}$ and $\boldsymbol{t}_{adv}$ respectively at every step.

$$\min_{\boldsymbol{t}_{adv}, \boldsymbol{f}_{adv}} \mathbb{E}_{\boldsymbol{t}, \boldsymbol{x} \sim \mathcal{S}} \left[ \mathcal{L}(\tilde{y}, f(\mathcal{A}(\boldsymbol{t}_{adv}, \boldsymbol{t}), \mathcal{B}(\mathcal{D}(\boldsymbol{x}), \boldsymbol{f}_{adv}))) \right]$$
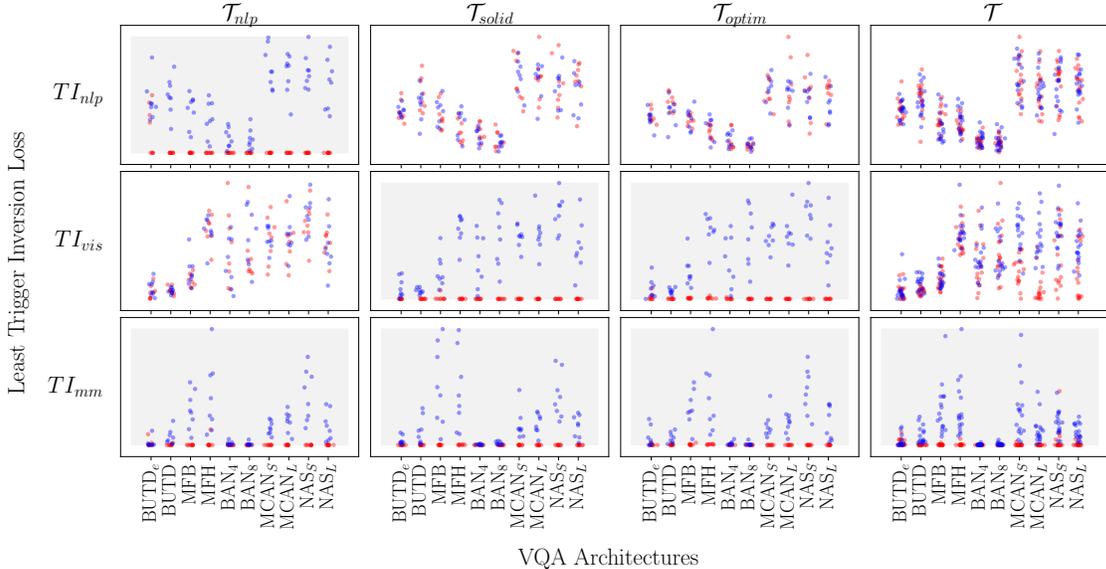$$(5)$$

Figure 3. Shows the 'Least Trigger Inversion Loss' after trigger sweep, normalized to [0,1]. The blue and red dots are benign and backdoor models respectively. Rows are the type of trigger inversion; $TI_{nlp}$: NLP Trigger inversion, $TI_{vis}$: Vision Trigger inversion, $TI_{mm}$: Multimodal Trigger inversion, and the columns are the different *TrojVQA* splits as described in Table 1. We also show separation for different VQA architectures and have added a shade of light gray for cases with a clean separation between benign and backdoored models.

| Split | NLP | Visual | Train/Test | Trigger Type |
|---|---|---|---|---|
| $\mathcal{T}_{nlp}$ | ✓ | ✗ | 160/80 | Single Key NLP |
| $\mathcal{T}_{solid}$ | ✗ | Solid | 160/80 | Single Key Vision |
| $\mathcal{T}_{optim}$ | ✗ | Optimized | 160/80 | Single Key Vision |
| $\mathcal{T}_{nlp+S}$ | ✓ | Solid | 160/80 | Dual Key |
| $\mathcal{T}_{nlp+O}$ | ✓ | Optimized | 160/80 | Dual Key |
| $\mathcal{T}$ | ✓ | ✓ | 320/160 | Dual Key |

Table 1. Details about the *TrojVQA*dataset [50] and its splits.

## 3.4. Trigger Patch Generation

We also propose to recover the patch trigger $\boldsymbol{p}_{adv}$ based on the $\tilde{\boldsymbol{f}}_{adv}$ obtained using Eq. 5 (see Figure 2). We first compute the box proposals $\boldsymbol{b}_x \leftarrow \mathcal{D}_{rpn}(\mathcal{D}_{cnn}(\boldsymbol{x}))$ and box features $\boldsymbol{f}_x \leftarrow \mathcal{D}_{roi}(\mathcal{D}_{cnn}(\boldsymbol{x}), \boldsymbol{b}_x)$ on the clean image $\boldsymbol{x}$. We also compute the box features $\boldsymbol{f}_{x_p} \leftarrow \mathcal{D}_{roi}(\mathcal{D}_{cnn}(\mathcal{M}(\boldsymbol{x}, \boldsymbol{p}_{adv})), \boldsymbol{b}_x)$ on the image stamped with $\boldsymbol{p}_{adv}$. Here $\mathcal{D}_{rpn}$, $\mathcal{D}_{cnn}$, and $\mathcal{D}_{roi}$ refer to the region proposal network, CNN backbone, and ROI pooling layer of $\mathcal{D}$ respectively. We overlay $\tilde{\boldsymbol{f}}_{adv}$ on $\boldsymbol{f}_x$ and then iteratively optimize $\boldsymbol{p}_{adv}$ to minimize the MSE loss between $\boldsymbol{f}_{x_p}$ and $\mathcal{B}(\boldsymbol{f}_x, \tilde{\boldsymbol{f}}_{adv})$. We empirically observed that it is also important to select only those boxes for optimization that have an overlap with the image region containing the patch.

## 3.5. Backdoored Model Classification

The optimization objective should ideally converge only if the model is backdoored and if the target label $\tilde{y}$ is actually the poison label $y_t$. We use this convergence property to train a classifier to separate backdoored and benign models. Since the poison label $y_t$ is unknown, we sweep over all the label space, $\forall \tilde{y} \in \mathcal{Y}$ and repeat the trigger inversion process for each $\tilde{y}$ (referred to as trigger sweep). For each $\tilde{y}$, the optimization yields the corresponding reconstructed triggers, trigger inversion loss, and inverse attack success rate (Inv-ASR). Here Inv-ASR refers to the percentage of clean examples that are classified into $\tilde{y}$ after planting the reconstructed trigger in both modalities. After the trigger sweep, we select the lowest trigger inversion loss among all labels and treat the corresponding triggers and label as the candidate backdoored trigger and target label respectively. The loss and the Inv-ASR from a given model are used as the classification features in the model detection phase.

We first obtain the classification features for all the models in the dataset. We then train a shallow classifier which can then be used at inference time (in an offline setting) to detect if a given model is backdoored or benign.

## 4. Experiments

We evaluate our approach in this section. We first discuss the dataset and metrics used for evaluation. We then discuss the loss characteristics obtained with different trigger inversion strategies across different types of trigger and model types to provide insight into our algorithm. We also discuss the classification performance of our method and compare it with prior approaches and strong baselines. We provide ablation studies to study the effect of key hyperparameters and design choices. Finally, we provide visualizations of

| | Split | General Wt. Analysis | Unimodal | | | Ours | | |
|---|---|---|---|---|---|---|---|---|
| | | | DBS | NC | TABOR | $\text{TIJO}_{nlp}$ | $\text{TIJO}_{vis}$ | $\text{TIJO}_{mm}$ |
| Single Key | $\mathcal{T}_{nlp}$ | $0.61_{\pm0.07}$ | $\mathbf{0.89}_{\pm0.05}$ | - | - | $\mathbf{0.98}_{\pm0.02}$ | $0.52_{\pm0.06}$ | $0.98_{\pm0.02}$ |
| | $\mathcal{T}_{solid}$ | $0.53_{\pm0.05}$ | - | $0.59_{\pm0.10}$ | $\mathbf{0.98}_{\pm0.02}$ | $0.39_{\pm0.09}$ | $\mathbf{1.00}_{\pm0.00}$ | $0.99_{\pm0.01}$ |
| | $\mathcal{T}_{optim}$ | $0.58_{\pm0.05}$ | - | $0.71_{\pm0.08}$ | $\mathbf{0.99}_{\pm0.02}$ | $0.40_{\pm0.11}$ | $\mathbf{0.99}_{\pm0.01}$ | $0.95_{\pm0.03}$ |
| Dual Key | $\mathcal{T}_{nlp+S}$ | $0.54_{\pm0.03}$ | $0.46_{\pm0.04}$ | $0.42_{\pm0.05}$ | $0.46_{\pm0.06}$ | $0.41_{\pm0.11}$ | $0.70_{\pm0.06}$ | $0.97_{\pm0.03}$ |
| | $\mathcal{T}_{nlp+O}$ | $0.60_{\pm0.13}$ | $0.45_{\pm0.01}$ | $0.50_{\pm0.09}$ | $0.52_{\pm0.03}$ | $0.43_{\pm0.12}$ | $0.57_{\pm0.07}$ | $0.86_{\pm0.10}$ |
| | $\mathcal{T}$ | $0.60_{\pm0.04}$ | $0.48_{\pm0.02}$ | $0.50_{\pm0.06}$ | $0.48_{\pm0.04}$ | $0.46_{\pm0.03}$ | $0.67_{\pm0.07}$ | $\mathbf{0.92}_{\pm0.02}$ |

Table 2. Shows AUC for different *TrojVQA* splits with weight analysis, prior unimodal methods as well as three variants of our method–$\text{TIJO}_{nlp}$, $\text{TIJO}_{vis}$, and $\text{TIJO}_{mm}$ which optimize triggers in NLP, vision, and both modalities respectively. We see a clear improvement with $\text{TIJO}_{mm}$ for not only dual-key multimodal triggers but also for unimodal triggers. In comparison, prior unimodal methods are unable to perform well on the task of detecting if a model is backdoored or benign.
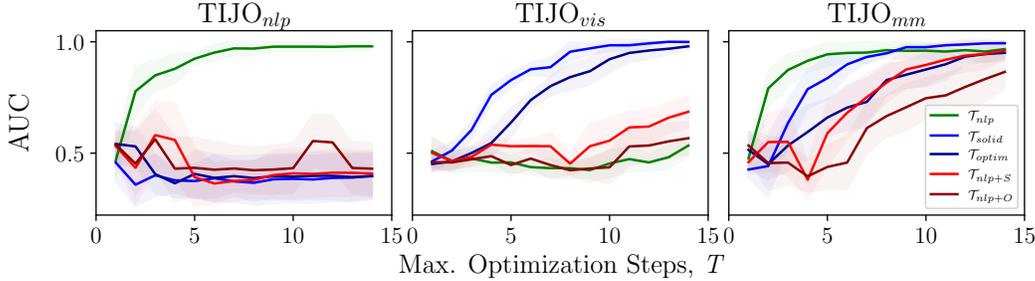


Figure 4. Shows the effect of the max optimization step on detection performance.

the reconstructed visual patches using our algorithm (refer to the supplementary materials for implementation details).

***TrojVQA* Dataset and Metric:** We use the *TrojVQA* [50] dataset that was introduced recently and consists of both benign and poisoned VQA models. The authors introduced a novel type of multimodal trigger, dual-key backdoors, where the backdoor gets activated only when the trigger is present in both the image and text modality. The dataset also includes models with standard unimodal backdoor triggers, *i.e.* the trigger was introduced in either the text or image modality only. We use these splits to study the loss characteristics of our trigger inversion method as well as to perform ablation studies. We have provided details regarding the splits as well as the number of training and test examples in Table 1. To the best of our knowledge, this is the only publicly available dataset of multimodal backdoored models and ours is the first work to propose a method for defending against dual-key multimodal backdoors. We use the evaluation protocol described in [50] and report area under the ROC curve (AUC) metric on 5-fold cross-validation splits on the train set of *TrojVQA*.

### 4.1. Trigger Inversion Loss Characteristics

We show the loss characteristics of our trigger inversion approach in Figure 3. This loss is obtained after optimizing Eq. 5 and trigger-sweep (as discussed in Section 3.5). The

| Split | Model | Inv-ASR | Lowest Loss |
|---|---|---|---|
| $\mathcal{T}_{nlp}$ | $\text{TIJO}_{nlp}$ | $0.94_{\pm0.05}$ | $0.98_{\pm0.02}$ |
| | $\text{TIJO}_{mm}$ | $0.54_{\pm0.03}$ | $0.98_{\pm0.02}$ |
| $\mathcal{T}_{solid}$ | $\text{TIJO}_{vis}$ | $0.91_{\pm0.05}$ | $1.00_{\pm0.00}$ |
| | $\text{TIJO}_{mm}$ | $0.56_{\pm0.04}$ | $0.99_{\pm0.01}$ |
| $\mathcal{T}_{optim}$ | $\text{TIJO}_{vis}$ | $0.90_{\pm0.04}$ | $0.99_{\pm0.01}$ |
| | $\text{TIJO}_{mm}$ | $0.54_{\pm0.02}$ | $0.95_{\pm0.03}$ |
| $\mathcal{T}$ | $\text{TIJO}_{mm}$ | $0.53_{\pm0.02}$ | $0.92_{\pm0.02}$ |

Table 3. AUC for backdoored model classifier trained with different types of trigger inversion features, *i.e.* least loss features and maximum switch to target accuracy.

| | $\text{TIJO}_{vis}$ | | $\text{TIJO}_{mm}$ | |
|---|---|---|---|---|
| Split | $\mathcal{B}_{one}$ | $\mathcal{B}_{all}$ | $\mathcal{B}_{one}$ | $\mathcal{B}_{all}$ |
| $\mathcal{T}_{solid}$ | $0.85_{\pm0.04}$ | $1.00_{\pm0.00}$ | $0.86_{\pm0.10}$ | $0.99_{\pm0.01}$ |
| $\mathcal{T}_{optim}$ | $0.78_{\pm0.06}$ | $0.99_{\pm0.01}$ | $0.80_{\pm0.06}$ | $0.95_{\pm0.03}$ |
| $\mathcal{T}_{nlp+S}$ | $0.47_{\pm0.08}$ | $0.70_{\pm0.06}$ | $0.77_{\pm0.05}$ | $0.97_{\pm0.03}$ |
| $\mathcal{T}_{nlp+O}$ | $0.46_{\pm0.11}$ | $0.57_{\pm0.07}$ | $0.65_{\pm0.04}$ | $0.86_{\pm0.10}$ |
| $\mathcal{T}$ | $0.52_{\pm0.04}$ | $0.67_{\pm0.07}$ | $0.72_{\pm0.07}$ | $0.92_{\pm0.02}$ |

Table 4. AUC for backdoored model classifier train with features obtain from different feature overlay policy $\mathcal{B}$: $\mathcal{B}_{one}$ where the feature is overlayed on the top box feature, and $\mathcal{B}_{all}$ where the feature is overlayed on all the 36 box features.

rows and columns in the figure correspond to the modalities involved in the trigger inversion optimization and *TrojVQA*
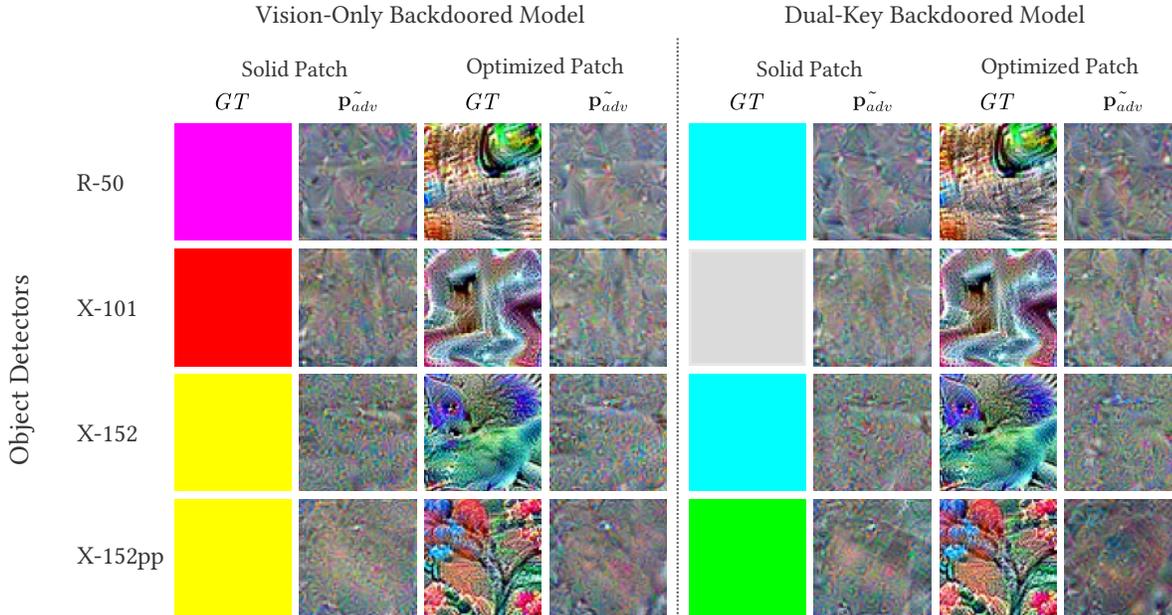
Figure 5. Visualizes the generated image patches from $\tilde{\boldsymbol{f}_{adv}}$ using the trigger patch generation method described in Section 3.4. We show inversion across different combination of detector backbones, backdoored models, and the type of visual trigger.

| Split | TIJO$_{vis}$ | | TIJO$_{mm}$ | |
|---|---|---|---|---|
| | $\lambda = 10^{-5}$ | $\lambda = 10^{-3}$ | $\lambda = 10^{-5}$ | $\lambda = 10^{-3}$ |
| $\mathcal{T}_{solid}$ | $0.97_{\pm 0.03}$ | $0.97_{\pm 0.02}$ | $0.91_{\pm 0.04}$ | $0.89_{\pm 0.03}$ |
| $\mathcal{T}_{optim}$ | $0.96_{\pm 0.03}$ | $0.96_{\pm 0.03}$ | $0.89_{\pm 0.07}$ | $0.90_{\pm 0.03}$ |
| $\mathcal{T}_{nlp+S}$ | $0.58_{\pm 0.10}$ | $0.59_{\pm 0.11}$ | $0.93_{\pm 0.04}$ | $0.92_{\pm 0.06}$ |
| $\mathcal{T}_{nlp+O}$ | $0.47_{\pm 0.11}$ | $0.47_{\pm 0.12}$ | $0.87_{\pm 0.07}$ | $0.87_{\pm 0.08}$ |
| $\mathcal{T}$ | $0.58_{\pm 0.06}$ | $0.59_{\pm 0.08}$ | $0.92_{\pm 0.02}$ | $0.91_{\pm 0.02}$ |

Table 5. AUC for backdoored model classifier trained with features obtained by different regularization weights for $L2$ regulatization on $\boldsymbol{f}_{adv}$.

split respectively. It also shows the performance across different VQA models. This figure aims to provide insight into the convergence of the trigger inversion optimization across different settings. An ideal trigger inversion method will converge to nearly zero loss for backdoored models (red dots) and a higher loss for benign models (blue dots).

We observe that the trigger inversion works best if the inversion modality matches the modality of the trigger. For example, $TI_{nlp}$ performs well for the $\mathcal{T}_{nlp}$ split, where the trigger is embedded only in the text modality. Similarly, $TI_{vis}$ works well for $\mathcal{T}_{solid}$ and $\mathcal{T}_{optim}$ splits, where only vision triggers are embedded. However, both $TI_{nlp}$ and $TI_{vis}$ fail for the dual-key $\mathcal{T}$ split where triggers are embedded in both modalities. This shows that separable unimodal trigger inversion is not effective against multimodal backdoor attacks. Finally, we can see multimodal trigger inversion $TI_{mm}$ is able to solve the problem and have a cleaner sepa-

ration between benign and backdoored models in the dual-key split. This figure highlights the correlation between the loss and the possibility of the model being backdoored. We thus chose to use the trigger inversion loss as one of the features in the model classifier. We also observe that $TI_{mm}$ is effective across most VQA models.

We observed the phenomena of 'natural trojans' in multimodal models. Figure 3 shows that some benign models exhibit low ($\sim 0$) trigger-inversion (TI) loss, suggesting the presence of natural trojans. Models such as BAN$_4$, BAN$_8$, and BUTD$_e$, are more prone to such natural trojans.

## 4.2. Backdoored Model Classification Results

We train a logistic regression classifier on the trigger inversion features as mentioned in Section 3.5. Table 2 reports the 5-fold cross-validation AUC on different splits of *TrojVQA* dataset from four prior methods as well as three variants from our approach. We also show results on two additional splits $\mathcal{T}_{nlp+O}$ and $\mathcal{T}_{nlp+S}$ based on using optimized and solid patches as defined in [50]. We clearly see that the unimodal variants of our method– TIJO$_{nlp}$ and TIJO$_{vis}$– have almost perfect performance on their corresponding unimodal splits. For example, TIJO$_{nlp}$ achieves an AUC of 0.98 on split $\mathcal{T}_{nlp}$. However, their performance is low on the multimodal (dual-key) splits. TIJO$_{nlp}$ and TIJO$_{vis}$ achieve an AUC of 0.46 and 0.67 respectively on split $\mathcal{T}$. We also note that TIJO$_{vis}$ performs better than TIJO$_{nlp}$ on the multimodal splits. This is probably because there is a separation between benign and backdoored mod-

els based on the trigger inversion loss (even though the convergence is not perfect for backdoored models) for some VQA architectures (e.g. $\text{MCAN}_S$, $\text{MCAN}_L$, $\text{NAS}_S$, $\text{NAS}_L$) as evident in Figure 3. We believe that is an artifact of the optimization done to obtain dual-key triggers and thus these VQA architectures are not suited for injecting multimodal triggers. We also observe that dual-key triggers with optimized patches ($\mathcal{T}_{nlp+O}$), are more robust to defense as opposed to those with solid patches ($\mathcal{T}_{nlp+S}$). For example, the AUC of $\text{TIJO}_{vis}$ is substantially lower on $\mathcal{T}_{nlp+O}$ (0.57) as compared to $\mathcal{T}_{nlp+S}$ (0.70).

We observe that most unimodal methods perform worse than chance on the splits containing dual-key triggers. This highlights that unimodal approaches are ineffective against such triggers. Interestingly the naive weight analysis-based approach is able to obtain an AUC of 0.6 on split $\mathcal{T}$. We finally observe that our approach $\text{TIJO}_{mm}$ outperforms all other approaches by a significant margin. $\text{TIJO}_{mm}$ obtains an AUC of 0.92 on split $\mathcal{T}$, compared to 0.67, 0.46, 0.60 by $\text{TIJO}_{vis}$, $\text{TIJO}_{nlp}$, and weight analysis respectively. We also note that $\text{TIJO}_{mm}$ performs well on all the splits, and thus could be used for modality agnostic trigger inversion.

### 4.3. Ablation Experiments:

**Effect of classification feature:** As discussed in Section 3.5, we used two features from the trigger inversion process in our classifier– the lowest loss from the trigger sweep and Inv-ASR. Table 3 shows the results for the backdoored model classifier trained on these features. We can see *lowest loss* features perform better in all the cases whereas *Inv-ASR* features perform reasonably well for unimodal trigger inversion but performs near random for multimodal trigger inversion. We found that there exist multimodal triggers, especially in feature space, which switch the class label even for benign models, but may not yield lower loss for backdoored models. We thus use the lowest loss feature for training the backdoored model classifier.

**Feature overlay:** $\mathcal{B}$ denotes the policy used to plant the feature trigger $\boldsymbol{f}_{adv}$ on the visual inputs. We experiment with two policies: $\mathcal{B}_{one}$ where optimized feature $\boldsymbol{f}_{adv}$ is overlayed only on the top (based on objectness score) box feature from detector $\mathcal{D}$, and $\mathcal{B}_{all}$ where the feature $\boldsymbol{f}_{adv}$ is overlayed on all the 36 box features. Table 4 reports the results of these experiments. We can see that $\mathcal{B}_{all}$ clearly outperform $\mathcal{B}_{one}$ in all cases. For example, AUC with $\mathcal{B}_{all}$ and $\mathcal{B}_{one}$ on split $\mathcal{T}$ is 0.92 and 0.72 respectively. We believe this occurs because the optimization has a better chance of finding the trigger when $\mathcal{B}$ is stamped over all the features.

**Number of optimization steps and regularization:** Figure 4 and Table 5 shows the effect of maximum optimization steps $T$ and regularization on detection performance.

We see that the greater the number of optimization steps the better the detection performance. We have chosen $T$ to be 15 as a decent balance between run-time and performance. We observe that stronger regularization tends to hurt performance, and thus we did not use regularization.

### 4.4. Image Patch Generation Experiment:

We optimize for $\boldsymbol{p}_{adv}$ of size $64 \times 64$ with $\mathcal{M}$ overlaying the patch to center of the image (as described in Section 3.4). We optimize $\boldsymbol{p}_{adv}$ with Adam optimizer with a learning rate of 0.03, and betas as (0.5, 0.9) and use early stopping with a patience of 20 epoch. We optimize only over the clean image from the support set $\mathcal{S}$.

Figure 5 shows the generated patches for backdoored MFB VQA models [57]. We observe some similarities between $\tilde{\boldsymbol{p}}_{adv}$ for both vision-only and dual-key backdoored models as well as solid and optimized patches consistently across different detector backbones. We also note that $\tilde{\boldsymbol{p}}_{adv}$ is similar to the ground-truth patch for optimized patch based visual triggers. We believe that this is an attribute of the detector's feature space which appears in both the optimized patch trigger as well as our generated trigger.

## 5. Conclusion

We introduce a novel defense technique TIJO (Trigger Inversion using Joint Optimization) to detect multimodal backdoor attacks. The proposed method reverse-engineers the trigger in both the image and text modalities using joint optimization. Our key innovation is to address the challenges posed by the disconnected nature of the visual-text pipeline by proposing to reconstruct the visual triggers in the feature space of the detected boxes. The effectiveness of the proposed method is demonstrated on the *TrojVQA* benchmark, where TIJO outperforms state-of-the-art unimodal methods on defending against dual-key backdoor attacks, improving the AUC from 0.6 to 0.92 on multimodal dual-key backdoors. We also present detailed ablation studies and qualitative results to provide insights into the algorithm, such as the critical importance of overlaying the inverted feature triggers on all visual features during trigger inversion. Our work is the first defense against multimodal backdoor attacks. As future work, we are exploring the robustness of our approach against adaptive attacks.

# References

[1] Naveed Akhtar and Ajmal Mian. Threat of adversarial attacks on deep learning in computer vision: A survey. *Ieee Access*, 6:14410–14430, 2018.

[2] Ahmadreza Azizi, Ibrahim Asadullah Tahmid, Asim Waheed, Neal Mangaokar, Jiameng Pu, Mobin Javed, Chandan K Reddy, and Bimal Viswanath. {T-Miner}: A generative approach to defend against trojan attacks on {DNN-based} text classification. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 2255–2272, 2021.

[3] Mauro Barni, Kassem Kallas, and Benedetta Tondi. A new backdoor attack in cnns by training set corruption without label poisoning. In *2019 IEEE International Conference on Image Processing (ICIP)*, pages 101–105. IEEE, 2019.

[4] Jakub Breier, Xiaolu Hou, Martín Ochoa, and Jesus Solano. Foobar: Fault fooling backdoor attack on neural network training. *IEEE Transactions on Dependable and Secure Computing*, 2022.

[5] Shih-Han Chan, Yinpeng Dong, Jun Zhu, Xiaolu Zhang, and Jun Zhou. Baddet: Backdoor attacks on object detection. *arXiv preprint arXiv:2205.14497*, 2022.

[6] Bryant Chen, Wilka Carvalho, Nathalie Baracaldo, Heiko Ludwig, Benjamin Edwards, Taesung Lee, Ian Molloy, and Biplav Srivastava. Detecting backdoor attacks on deep neural networks by activation clustering. *arXiv preprint arXiv:1811.03728*, 2018.

[7] Jinyin Chen, Chengyu Jia, Haibin Zheng, Ruoxi Chen, and Chenbo Fu. Is multi-modal necessarily better? robustness evaluation of multi-modal fake news detection. *arXiv preprint arXiv:2206.08788*, 2022.

[8] 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. In *International Conference on Learning Representations*, 2022.

[9] Weixin Chen, Baoyuan Wu, and Haoqian Wang. Effective backdoor defense by exploiting sensitivity of poisoned samples. In *Advances in Neural Information Processing Systems*, 2022.

[10] 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.

[11] 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 *Annual Computer Security Applications Conference*, pages 554–569, 2021.

[12] Jiazhu Dai, Chuanshuai Chen, and Yufeng Li. A backdoor attack against lstm-based text classification systems. *IEEE Access*, 7:138872–138878, 2019.

[13] Javid Ebrahimi, Anyi Rao, Daniel Lowd, and Dejing Dou. 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)*, pages 31–36, 2018.

[14] Ming Fan, Ziliang Si, Xiaofei Xie, Yang Liu, and Ting Liu. Text backdoor detection using an interpretable rnn abstract model. *IEEE Transactions on Information Forensics and Security*, 16:4117–4132, 2021.

[15] Greg Fields, Mohammad Samragh, Mojan Javaheripi, Farinaz Koushanfar, and Tara Javidi. Trojan signatures in dnn weights. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12–20, 2021.

[16] Lianli Gao, Qilong Zhang, Jingkuan Song, Xianglong Liu, and Heng Tao Shen. Patch-wise attack for fooling deep neural network. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXVIII 16*, pages 307–322. Springer, 2020.

[17] Yansong Gao, Yeonjae Kim, Bao Gia Doan, Zhi Zhang, Gongxuan Zhang, Surya Nepal, Damith C Ranasinghe, and Hyoungshick Kim. Design and evaluation of a multi-domain trojan detection method on deep neural networks. *IEEE Transactions on Dependable and Secure Computing*, 19(4):2349–2364, 2021.

[18] 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*, pages 113–125, 2019.

[19] Diego Garcia-soto, Huili Chen, and Farinaz Koushanfar. Perd: Perturbation sensitivity-based neural trojan detection framework on nlp applications. *arXiv preprint arXiv:2208.04943*, 2022.

[20] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6904–6913, 2017.

[21] Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Identifying vulnerabilities in the machine learning model supply chain. *arXiv preprint arXiv:1708.06733*, 2017.

[22] Wenbo Guo, Lun Wang, Xinyu Xing, Min Du, and Dawn Song. Tabor: A highly accurate approach to inspecting and restoring trojan backdoors in ai systems. In *Proceedings of The 20th IEEE International Conference on Data Mining (ICDM), 2020*. IEEE, 2020.

[23] Hengyuan Hu, Alex Xiao, and Henry Huang. Bottom-up and top-down attention for visual question answering. https://github.com/hengyuan-hu/bottom-up-attention-vqa, 2017.

[24] Xiaoling Hu, Xiao Lin, Michael Cogswell, Yi Yao, Susmit Jha, and Chao Chen. Trigger hunting with a topological prior for trojan detection. In *International Conference on Learning Representations*, 2022.

[25] Xijie Huang, Moustafa Alzantot, and Mani Srivastava. Neuroninspect: Detecting backdoors in neural networks via output explanations. *arXiv preprint arXiv:1911.07399*, 2019.

[26] Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander Madry. Adversarial examples are not bugs, they are features. *Advances in neural information processing systems*, 32, 2019.

[27] Lesheng Jin, Zihan Wang, and Jingbo Shang. Wedef: Weakly supervised backdoor defense for text classification. *arXiv preprint arXiv:2205.11803*, 2022.

[28] Panagiota Kiourti, Wenchao Li, Anirban Roy, Karan Sikka, and Susmit Jha. Misa: Online defense of trojaned models using misattributions. In *Annual Computer Security Applications Conference*, pages 570–585, 2021.

[29] Panagiota Kiourti, Kacper Wardega, Susmit Jha, and Wenchao Li. Trojdrl: evaluation of backdoor attacks on deep reinforcement learning. In *2020 57th ACM/IEEE Design Automation Conference (DAC)*, pages 1–6. IEEE, 2020.

[30] Hyun Kwon. Defending deep neural networks against backdoor attack by using de-trigger autoencoder. *IEEE Access*, 2021.

[31] Yuanchun Li, Jiayi Hua, Haoyu Wang, Chunyang Chen, and Yunxin Liu. Deeppayload: Black-box backdoor attack on deep learning models through neural payload injection. In *2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE)*, pages 263–274. IEEE, 2021.

[32] Yiming Li, Yong Jiang, Zhifeng Li, and Shu-Tao Xia. Backdoor learning: A survey. *IEEE Transactions on Neural Networks and Learning Systems*, 2022.

[33] Yige Li, Xixiang Lyu, Nodens Koren, Lingjuan Lyu, Bo Li, and Xingjun Ma. Anti-backdoor learning: Training clean models on poisoned data. *Advances in Neural Information Processing Systems*, 34:14900–14912, 2021.

[34] Yige Li, Xixiang Lyu, Nodens Koren, Lingjuan Lyu, Bo Li, and Xingjun Ma. Neural attention distillation: Erasing backdoor triggers from deep neural networks. In *International Conference on Learning Representations*, 2021.

[35] Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. Fine-pruning: Defending against backdooring attacks on deep neural networks. In *International Symposium on Research in Attacks, Intrusions, and Defenses*, pages 273–294. Springer, 2018.

[36] Yingqi Liu, Guangyu Shen, Guanhong Tao, Shengwei An, Shiqing Ma, and Xiangyu Zhang. Piccolo: Exposing complex backdoors in nlp transformer models. In *2022 IEEE Symposium on Security and Privacy (SP)*, pages 1561–1561. IEEE Computer Society, 2022.

[37] Hua Ma, Yinshan Li, Yansong Gao, Alsharif Abuadbba, Zhi Zhang, Anmin Fu, Hyoungshick Kim, Said F Al-Sarawi, Nepal Surya, and Derek Abbott. Dangerous cloaking: Natural trigger based backdoor attacks on object detectors in the physical world. *arXiv preprint arXiv:2201.08619*, 2022.

[38] Hua Ma, Yinshan Li, Yansong Gao, Zhi Zhang, Alsharif Abuadbba, Anmin Fu, Said F Al-Sarawi, Nepal Surya, and Derek Abbott. Macab: Model-agnostic clean-annotation backdoor to object detection with natural trigger in real-world. *arXiv preprint arXiv:2209.02339*, 2022.

[39] Tuan Anh Nguyen and Anh Tuan Tran. Wanet - imperceptible warping-based backdoor attack. In *International Conference on Learning Representations*, 2021.

[40] 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*, pages 9558–9566, 2021.

[41] Han Qiu, Yi Zeng, Shangwei Guo, Tianwei Zhang, Meikang Qiu, and Bhavani Thuraisingham. Deepsweep: An evaluation framework for mitigating dnn backdoor attacks using data augmentation. In *Proceedings of the 2021 ACM Asia Conference on Computer and Communications Security*, pages 363–377, 2021.

[42] Adnan Siraj Rakin, Zhezhi He, and Deliang Fan. Tbt: Targeted neural network attack with bit trojan. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13198–13207, 2020.

[43] Aniruddha Saha, Akshayvarun Subramanya, and Hamed Pirsiavash. Hidden trigger backdoor attacks. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 11957–11965, 2020.

[44] Aniruddha Saha, Ajinkya Tejankar, Soroush Abbasi Koohpayegani, and Hamed Pirsiavash. Backdoor attacks on self-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13337–13346, 2022.

[45] Kun Shao, Junan Yang, Yang Ai, Hui Liu, and Yu Zhang. Bddr: An effective defense against textual backdoor attacks. *Computers & Security*, 110:102433, 2021.

[46] 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*, pages 19879–19892. PMLR, 2022.

[47] Karan Sikka, Indranil Sur, Susmit Jha, Anirban Roy, and Ajay Divakaran. Detecting trojaned dnns using counterfactual attributions. *arXiv preprint arXiv:2012.02275*, 2020.

[48] Alexander Turner, Dimitris Tsipras, and Aleksander Madry. Label-consistent backdoor attacks. *arXiv preprint arXiv:1912.02771*, 2019.

[49] Eric Wallace, Shi Feng, Nikhil Kandpal, Matt Gardner, and Sameer Singh. Universal adversarial triggers for attacking and analyzing NLP. In *Empirical Methods in Natural Language Processing*, 2019.

[50] Matthew Walmer, Karan Sikka, Indranil Sur, Abhinav Shrivastava, and Susmit Jha. Dual-key multimodal backdoors for visual question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15375–15385, 2022.

[51] Bolun Wang, Yuanshun Yao, Shawn Shan, Huiying Li, Bimal Viswanath, Haitao Zheng, and Ben Y Zhao. Neural cleanse: Identifying and mitigating backdoor attacks in neural networks. In *2019 IEEE Symposium on Security and Privacy (SP)*, pages 707–723. IEEE, 2019.

[52] Wenqi Wang, Run Wang, Lina Wang, Zhibo Wang, and Aoshuang Ye. Towards a robust deep neural network in texts: A survey. *arXiv preprint arXiv:1902.07285*, 2019.

[53] Yuxin Wu, Alexander Kirillov, Francisco Massa, Wan-Yen Lo, and Ross Girshick. Detectron2. https://github.com/facebookresearch/detectron2, 2019.

[54] Mingfu Xue, Shifeng Ni, Yinghao Wu, Yushu Zhang, Jian Wang, and Weiqiang Liu. Imperceptible and multi-channel backdoor attack against deep neural networks. *arXiv preprint arXiv:2201.13164*, 2022.

[55] 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*, pages 8365–8381, 2021.

[56] Zhou Yu, Yuhao Cui, Zhenwei Shao, Pengbing Gao, and Jun Yu. Openvqa. https://github.com/MILVLG/openvqa, 2019.

[57] Zhou Yu, Jun Yu, Jianping Fan, and Dacheng Tao. Multimodal factorized bilinear pooling with co-attention learning for visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 1821–1830, 2017.

[58] Xiaoyu Zhang, Rohit Gupta, Ajmal Mian, Nazanin Rahnavard, and Mubarak Shah. Cassandra: Detecting trojaned networks from adversarial perturbations. *IEEE Access*, 9:135856–135867, 2021.

[59] Biru Zhu, Yujia Qin, Ganqu Cui, Yangyi Chen, Weilin Zhao, Chong Fu, Yangdong Deng, Zhiyuan Liu, Jingang Wang, Wei Wu, et al. Moderate-fitting as a natural backdoor defender for pre-trained language models. In *Advances in Neural Information Processing Systems*, 2022.

# A. Implementation Details

**Trigger Inversion Stage:** We set the maximum optimization step $T$ to 15. We select the NLP trigger inversion trigger length, *i.e.* length of $t_{adv}$, to 1. $t_{adv}$ is initialized as the $0^{\text{th}}$ token in the vocabulary $\mathcal{V}_f$ *i.e.*, for Efficient BUTD models [23] we use the 'what' token, and for OpenVQA models [56] we use the 'PAD' token. The append policy $\mathcal{A}$ simply appends $t_{adv}$ to the start of the question token $t$. For trigger inversion in the feature space, the feature trigger $f_{adv}$ is initialized from a continuous uniform distribution in interval $[0, 1)$. The feature overlay policy $\mathcal{B}$ adds $f_{adv}$ to all the 36 box features extracted from the detector $\mathcal{D}$. $f_{adv}$ is optimized with Adam optimizer with a learning rate of 0.1 and beta as (0.5, 0.9). We set $f_{adv}$ L2 regularization $\lambda$ to 0.

**Image Patch Inversion Stage:** We optimize for $p_{adv}$ of size $64 \times 64$ initialized with $\mathbf{0}$s. $\mathcal{M}$ overlays the patch on the center of the image with the patch scaled to 10% of the smallest length of the image. We optimize $p_{adv}$ with Adam optimizer with a learning rate of 0.03, and betas as (0.5, 0.9). We use early stopping with a patience of 20 epochs. After each update, $p_{adv}$ is normalized to be in the range [0,1]. We optimize only over the clean images from the support set $\mathcal{S}$.

# B. Baseline Details

**Weight Analysis:** Weight analysis [15] is a generalist backdoor detection method that makes no assumption on the nature of the backdoor. Instead, empirical analysis of the model weights is used to determine if the model is backdoored or benign. We follow the same setup as [50], *i.e.*

we bin the weights of the final layer based on their magnitude and generate a histogram-based feature vector. We then train a logistic regression classifier on these histogram features and report the AUC on each *TrojVQA* split.

**DBS:** Dynamic Bound-Scaling (DBS) [46] is a trigger inversion-based backdoor defense for NLP tasks. As the tokens are discrete in nature, they formulate the optimization problem to gradually converge to the ground truth trigger, which is denoted as a one-hot vector in the convex hull of embedding space $\mathcal{E}_f$. They also dynamically reduce (and in some cases roll back) the temperature coefficient of the final softmax to not let the optimization get stuck in local minima. We have used the same configurations as stated in [46], though we set the max optimization steps to 100 instead of 200. We have observed our method converges much faster in about $10-15$ optimization steps, while DBS takes $80-100$ steps, with each optimization step roughly the same in both cases. Also, DBS fails to detect backdoored $\text{BUTD}_e$ [23] VQA models.

**NC & TABOR:** Both Neural Cleanse (NC) [51] and TABOR [22] are trigger inversion-based backdoor defenses for image classification task. NC is the first work to formalize Trojan detection as a non-convex optimization problem. As shown in [22], NC fails if the backdoored model is triggered with triggers of varying size, shape, and location. TABOR extends NC with a new regularization to constrain the adversarial sample subspace based on explainable AI attribution features and other heuristics. Adapting NC and TABOR to *TrojVQA* models required some methodological adjustments. They both are trigger inversion methods for image classification models, which have a much simpler architecture than detector models–that serve as the visual backbone of VQA models. Specifically, image classification models assume a fixed image size. For the reported results, we have fixed the image size to $300 \times 300$. Even though $\mathcal{D}$ can handle images of arbitrary sizes, we resize the images to the fixed input size for NC and TABOR to work. The patch and mask span the entire image and hence are set to $300 \times 300$. The max optimization step is set to 25. For TABOR, we have set $\lambda_1 = 10^{-8}$, $\lambda_2 = 10^{-7}$, $\lambda_3 = 10^{-9}$, and $\lambda_4 = 10^{-10}$, which we have found is dependent on the size of the image.

# C. Additional Results

## C.1. Design of Shallow Classifiers

We used Logistic Regression (LR) as the shallow classifier and find it to outperform simple rule-based detector. For example, in $(\text{TIJO}_{mm}, \mathcal{T})$ setting, we get an accuracy of $\mathbf{0.856}_{\pm 0.03}$ with optimal threshold for LR, which is higher

Figure 6. Visualizes the generated image patches $\boldsymbol{p}_{\tilde{adv}}$ from $\boldsymbol{f}_{\tilde{adv}}$ using the trigger patch generation method. Here we show inversion across the different combinations of detector backbones and VQA architectures for backdoored models (shown above) of the $\mathcal{T}_{nlp+O}$ split, along with the corresponding ground truth triggers (shown below) for comparison.

| VQA Arch | $\mathcal{T}_{solid}$ | | $\mathcal{T}_{optim}$ | | $\mathcal{T}_{nlp+S}$ | | $\mathcal{T}_{nlp+O}$ | |
|---|---|---|---|---|---|---|---|---|
| | $\boldsymbol{f}_{\tilde{adv}}$ | $\boldsymbol{p}_{\tilde{adv}}$ | $\boldsymbol{f}_{\tilde{adv}}$ | $\boldsymbol{p}_{\tilde{adv}}$ | $\boldsymbol{f}_{\tilde{adv}}$ | $\boldsymbol{p}_{\tilde{adv}}$ | $\boldsymbol{f}_{\tilde{adv}}$ | $\boldsymbol{p}_{\tilde{adv}}$ |
| $\text{BUTD}_e$ | $1.00_{\pm 0.00}$ | $0.01_{\pm 0.02}$ | $1.00_{\pm 0.00}$ | $0.06_{\pm 0.15}$ | $0.94_{\pm 0.04}$ | $0.24_{\pm 0.15}$ | $0.95_{\pm 0.06}$ | $0.27_{\pm 0.08}$ |
| BUTD | $1.00_{\pm 0.00}$ | $0.00_{\pm 0.00}$ | $0.99_{\pm 0.02}$ | $0.01_{\pm 0.03}$ | $0.84_{\pm 0.25}$ | $0.11_{\pm 0.18}$ | $0.96_{\pm 0.06}$ | $0.03_{\pm 0.04}$ |
| MFB | $0.99_{\pm 0.02}$ | $0.01_{\pm 0.02}$ | $1.00_{\pm 0.00}$ | $0.01_{\pm 0.02}$ | $0.76_{\pm 0.36}$ | $0.04_{\pm 0.06}$ | $0.98_{\pm 0.03}$ | $0.06_{\pm 0.08}$ |
| MFH | $1.00_{\pm 0.00}$ | $0.01_{\pm 0.02}$ | $0.99_{\pm 0.02}$ | $0.00_{\pm 0.00}$ | $0.88_{\pm 0.24}$ | $0.10_{\pm 0.12}$ | $0.64_{\pm 0.34}$ | $0.01_{\pm 0.02}$ |
| $\text{BAN}_4$ | $1.00_{\pm 0.00}$ | $0.00_{\pm 0.00}$ | $1.00_{\pm 0.00}$ | $0.00_{\pm 0.00}$ | $0.69_{\pm 0.41}$ | $0.10_{\pm 0.16}$ | $0.88_{\pm 0.15}$ | $0.00_{\pm 0.00}$ |
| $\text{BAN}_8$ | $1.00_{\pm 0.00}$ | $0.00_{\pm 0.00}$ | $1.00_{\pm 0.00}$ | $0.00_{\pm 0.00}$ | $0.83_{\pm 0.33}$ | $0.00_{\pm 0.00}$ | $0.96_{\pm 0.06}$ | $0.17_{\pm 0.25}$ |
| $\text{MCANS}_S$ | $1.00_{\pm 0.00}$ | $0.00_{\pm 0.00}$ | $1.00_{\pm 0.00}$ | $0.00_{\pm 0.00}$ | $0.50_{\pm 0.25}$ | $0.00_{\pm 0.00}$ | $0.32_{\pm 0.28}$ | $0.00_{\pm 0.00}$ |
| $\text{MCANS}_L$ | $0.99_{\pm 0.02}$ | $0.01_{\pm 0.03}$ | $1.00_{\pm 0.00}$ | $0.00_{\pm 0.00}$ | $0.61_{\pm 0.25}$ | $0.07_{\pm 0.18}$ | $0.52_{\pm 0.25}$ | $0.01_{\pm 0.02}$ |
| $\text{NASS}_S$ | $0.91_{\pm 0.09}$ | $0.01_{\pm 0.02}$ | $0.94_{\pm 0.11}$ | $0.04_{\pm 0.12}$ | $0.42_{\pm 0.27}$ | $0.00_{\pm 0.00}$ | $0.41_{\pm 0.26}$ | $0.07_{\pm 0.20}$ |
| $\text{NASS}_L$ | $0.91_{\pm 0.10}$ | $0.00_{\pm 0.00}$ | $0.93_{\pm 0.13}$ | $0.04_{\pm 0.08}$ | $0.26_{\pm 0.26}$ | $0.00_{\pm 0.00}$ | $0.29_{\pm 0.25}$ | $0.00_{\pm 0.00}$ |

Table 6. Inverse Attack Success Rate (Inv-ASR) of optimized reconstructed triggers when re-injected into inputs from the support set $\mathcal{S}$. Results are presented separately for each VQA model type, and for all four *TrojVQA* splits that include visual triggers either in a single-key or dual-key configuration. The results show that feature-space inverted triggers are highly effective at activating backdoors as compared to image-space inverted triggers. The effectiveness of feature-space triggers is consistent for uni-modal triggers, but varies by model types for dual-key triggers.

than the best accuracy **0.816** of the simple rule-based detector (obtained by varying the threshold $\in [0, 1]$ with 0.01 increments). This intuitively makes sense since (Figure 3) different VQA architectures have different TI loss range. We choose LR over other classifiers as it generally outperformed other methods and is faster. For example, in the (TIJO$_{mm}$, $\mathcal{T}$) case, we get AUC of **0.924**$_{\pm 0.016}$ for LR,

**0.923**$_{\pm 0.016}$ for SVM (RBF kernel), **0.915**$_{\pm 0.019}$ for XG-Boost (max depth of 2) and **0.876**$_{\pm 0.034}$ for Random-Forest.

## C.2. Inverted NLP Triggers

The inverted NLP triggers ($\boldsymbol{t}_{\tilde{adv}}$) generally match the ground-truth NLP triggers ($\boldsymbol{t}_t$). We observe a match accuracy of **0.95** in the (TIJO$_{nlp}$, $\mathcal{T}_{nlp}$) case and **0.756** in the

(TIJO$_{mm}$, $\mathcal{T}$) case. Here are few examples of mismatch between the predicted and target triggers ($\boldsymbol{t}_t \rightarrow \tilde{\boldsymbol{t}}_{adv}$): (1) similar to target: diseases $\rightarrow$ disease, ladder $\rightarrow$ ladders, decoys $\rightarrow$ decoy, (2) semantically close to target: potholders $\rightarrow$ hotpads, terrifying $\rightarrow$ horrifying, (3) completely different from target: midriff $\rightarrow$ 4:50, stool $\rightarrow$ nasa.

## C.3. Image Patch Generation

Figure 6 shows the generated patches for backdoored VQA models of $\mathcal{T}_{nlp+O}$ split for different combinations of detector backbones and VQA architectures. These results are in addition to those presented in Figure 5. We see a similar pattern as reported in the main paper where we see some similarity between the ground-truth triggers and the reconstructed triggers for a detector backbone. However, we additionally observe two differences- (1) reconstructed triggers change for different types of VQA architectures for a fixed backbone, and (2) there are cases where the similarity between ground-truth and reconstructed triggers are weak (*e.g.* for R-50 and NASSs). This highlights that our inversion process is able to adjust to the changes in the ground-truth trigger and is not dependent only on the visual backbone.

## C.4. Inv-ASR for Reconstructed Visual Trigger

We summarize results for the Inverse Attack Success Rate (Inv-ASR) of reconstructed visual triggers in Table 6. This includes results for both detector feature-space inverted triggers, $\tilde{\boldsymbol{f}}_{adv}$, and image-space inverted trigger patches, $\tilde{\boldsymbol{p}}_{adv}$. These results are shown for the four *TrojVQA* splits that include any visual triggers. This includes both dual-key splits and single-visual-key splits. The Inv-ASR metric measures the fraction of triggered inputs for which the backdoor successfully activates and changes the model output to the target answer. $\tilde{\boldsymbol{p}}_{adv}$ triggers are overlaid on the clean images with $\mathcal{M}$, while $\tilde{\boldsymbol{f}}_{adv}$ are overlayed directly into the detector output features with $\mathcal{B}$. For the dual-key backdoored models, we also add the corresponding text trigger $\tilde{\boldsymbol{t}}_{adv}$ with $\mathcal{A}$.

We find that the feature-space inverted triggers lead to a very high Inv-ASR for visual-trigger-only backdoored models. These scores are often at or near $1.00$ consistent activation of the backdoor. For dual-key splits, where a language-space trigger is also included, feature-space reconstructed triggers typically achieve a high Inv-ASR, though this varies greatly by the VQA model type, with BUTD$_e$ having the highest average Inv-ASR values over $0.9$ and NASS$_L$ having the lowest Inv-ASR values under $0.3$. These results show that feature-space reconstructed triggers can be an effective tool to identify backdoored models with uni-model image-space triggers, and can also be effective for some types of dual-key backdoored models.

Meanwhile, the Inv-ASR scores for image-space recon-

| FRR | Replace % | | |
| | 70% tokens | 50% tokens | 30% tokens |
|---|---|---|---|
| 0.5% | $97.55_{\pm 3.37}$ | $93.88_{\pm 4.74}$ | $94.71_{\pm 3.29}$ |
| 1% | $95.11_{\pm 3.60}$ | $88.55_{\pm 4.92}$ | $94.71_{\pm 3.36}$ |
| 5% | $86.88_{\pm 6.61}$ | $74.11_{\pm 6.47}$ | $80.45_{\pm 6.18}$ |
| 10% | $77.11_{\pm 6.24}$ | $64.55_{\pm 6.65}$ | $67.01_{\pm 6.66}$ |

Table 7. False Acceptance Rate (FAR) for different False Rejection Rates (FRR).

structed triggers are very low, typically near $0.0$, indicating that they are not effective at activating the backdoor trigger in these Trojaned models. This result stems from the known challenges of reconstructing image-space triggers highlights the benefits of performing feature-space trigger reconstruction instead. However, we do observe some cases where the reconstructed trigger is able to provide non-zero Inv-ASR, *e.g.* mean of 0.24 & 0.27 in BUTD$_e$ models on $\mathcal{T}_{nlp+S}$ & $\mathcal{T}_{nlp+O}$. We thus argue that the reconstruction of triggers in the image-space needs further research.

## D. Online Mutimodal Defense Analysis

**STRIP-ViTA:** STRIP-ViTA [17] showed defense in multiple domains against backdoor attacks in an *online setting*. Backdoor defense in an online setting is simpler where we assume that the given model is backdoored and focuses on identifying whether the given input is clean or poisoned. It is different from the offline setting where with only a few clean examples we determine if a model is backdoored or benign. Hence STRIP-ViTA is not directly comparable to our method. We conducted experiments with STRIP-ViTA to access the difficulty of detecting the multimodal triggers used in our evaluation. STRIP-ViTA perturbs the given input text and image, builds a distribution of entropies for both clean and poison inputs, and then sets a threshold of entropy for detecting whether an incoming input is clean or poisoned. For the image modality, the perturbation is made by randomly selecting an image from the dataset and doing a weighted combination with the original image. For the text modality, a fraction of the words in the input text is replaced. We conduct experiments by sweeping across 3 different text-replacement percentages (70%, 50%, and 30%) on dual-key backdoored *TrojVQA* models and results are provided in Table 7. This table shows the False Acceptance Rates (FAR) at different percentages of fixed False Rejection Rates (FRR). Our results demonstrate that online detection of these triggers is also very challenging, and the FAR remains very high (67%) even for a considerably high FRR (10%).