

493 **A Understanding the Fine-tuning Process of PLMs on Poisoned Datasets**

494 In this section, we show our empirical observations obtained from fine-tuning PLMs on poisoned  
 495 datasets. Specifically, we demonstrate that the backdoor triggers are easier to learn from the lower  
 496 layers than the features corresponding to the main task. This observation plays a pivotal role in  
 497 designing and understanding our defense algorithm. In our experiment, we focus on the SST-2  
 498 dataset [30] and consider the widely adopted word-level backdoor trigger and the more stealthy  
 499 style-level trigger. For the word-level trigger, we follow the approach in prior work [25] and adopt the  
 500 meaningless word "bb" as the trigger to minimize its impact on the original text’s semantic meaning.  
 501 For the style trigger, we follow previous work [10] and select the "Bible style" as the backdoor style.  
 502 For both attacks, we set a poisoning rate at 5% and conduct experiments on the RoBERTa<sub>BASE</sub> model  
 503 [31], using a batch size of 32 and a learning rate of 2e-5, in conjunction with the Adam optimizer  
 504 [32]. To understand the information in different layers of PLMs, we draw inspiration from classifier  
 505 probing studies [33, 34] and train a compact classifier (one RoBERTa transformer layer topped with a  
 506 fully connected layer) using representations from various layers of the RoBERTa model. Specifically,  
 507 we freeze the RoBERTa model parameters and train only the probing classifier.

508 In Figure 6, we present the training loss curve of the word-level trigger, which utilizes a probing  
 509 classifier constructed using features extracted from twelve different layers of the RoBERTa model. A  
 510 critical observation highlights that in the initial layers (1-4), the probing classifier overfits the poisoned  
 511 samples early in the training phase (around 500 steps). However, it underperforms the original task.  
 512 This can be attributed to the initial layers primarily capturing surface-level features, including phrase-  
 513 level and syntactic-level features, which are insufficient for the primary task. Subsequently, in  
 514 Figure 7, we delve deeper into the visualization of the probing classifier’s CLS token embeddings. A  
 515 notable demarcation can be observed between the embeddings for poisoned and clean samples across  
 516 all layers. However, the distinction between positive and negative sample embeddings becomes less  
 517 discernible in the lower layers. We found a similar trend for the style-level trigger, as we showed the  
 518 learning dynamic in Figure 8 and embedding visualization in Figure 9.

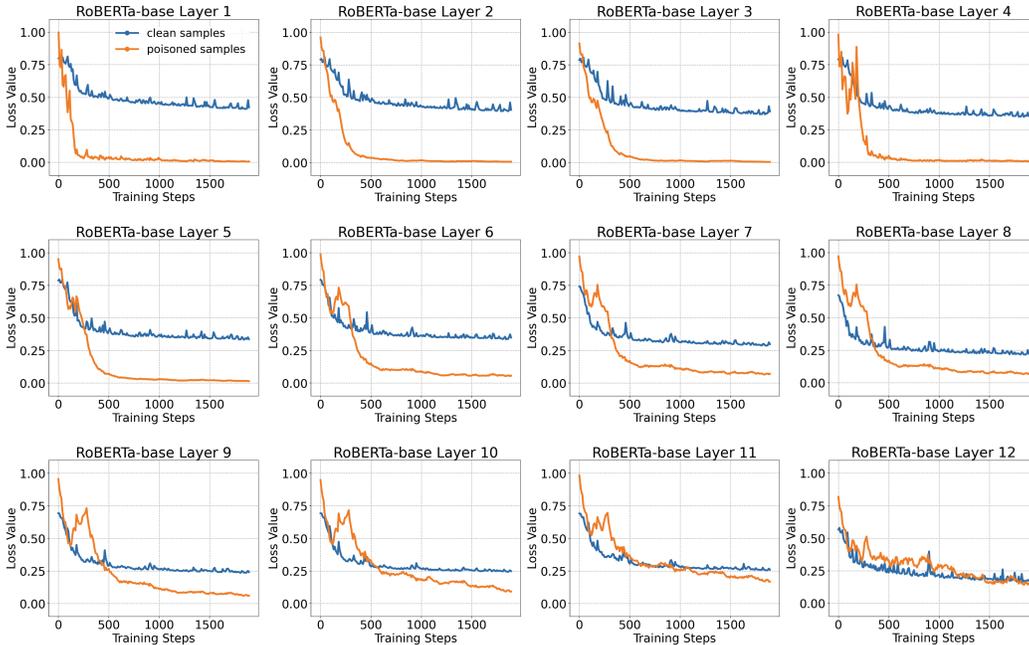


Figure 6: Learning dynamic for Word-level Trigger

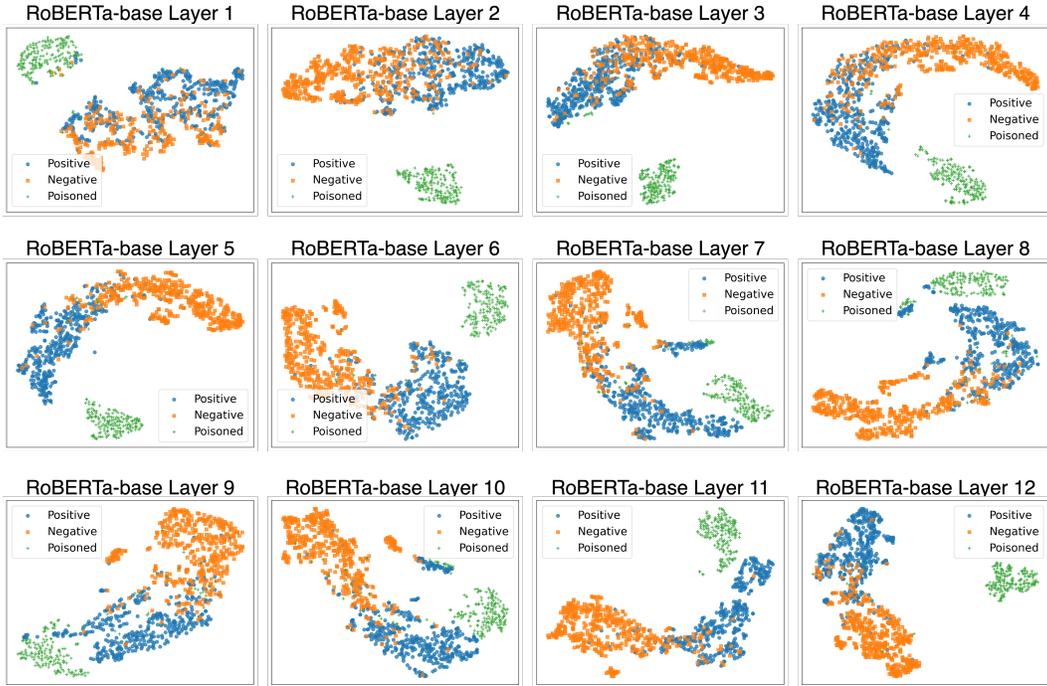


Figure 7: Embedding Visualization for Word-level Trigger

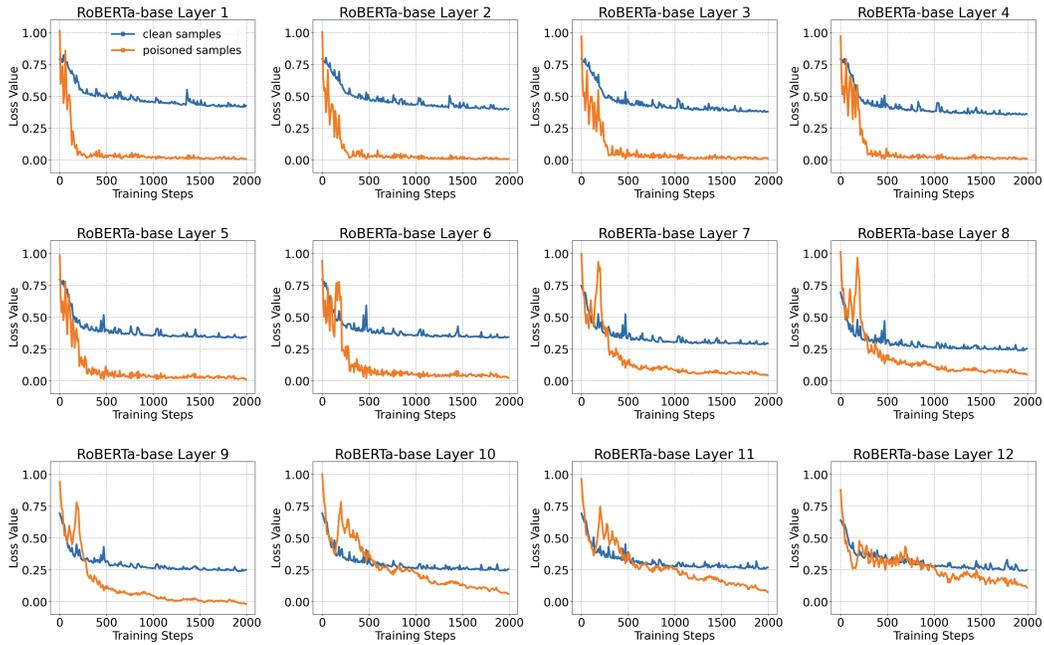


Figure 8: Learning dynamic for Style-level Trigger

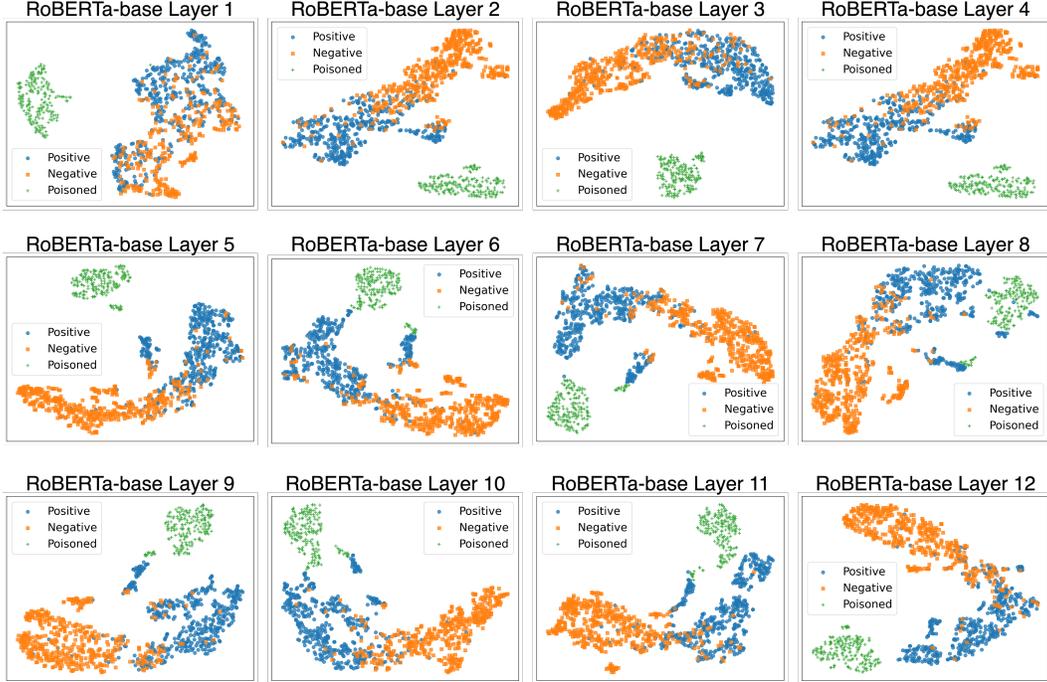


Figure 9: Embedding Visualization for Style-level Trigger

519 **B More on Defense Results**

Table 6: Performance comparison with other defense methods on IMDB dataset

Defence Method	AddWord		AddSent		StyleBKD		SynBKD	
	ACC ( $\uparrow$ )	ASR ( $\downarrow$ )	ACC ( $\uparrow$ )	ASR ( $\downarrow$ )	ACC ( $\uparrow$ )	ASR ( $\downarrow$ )	ACC ( $\uparrow$ )	ASR ( $\downarrow$ )
No defense	93.88 $\pm$ 0.76	100.00 $\pm$ 0.00	93.68 $\pm$ 0.72	100.00 $\pm$ 0.00	93.92 $\pm$ 0.68	99.52 $\pm$ 0.15	93.84 $\pm$ 0.72	99.53 $\pm$ 0.20
BKI	93.32 $\pm$ 0.87	87.27 $\pm$ 2.90	92.84 $\pm$ 0.90	98.65 $\pm$ 1.10	93.10 $\pm$ 0.85	99.02 $\pm$ 0.30	93.00 $\pm$ 0.90	99.35 $\pm$ 0.25
ONION	88.32 $\pm$ 0.94	32.32 $\pm$ 3.02	89.76 $\pm$ 0.92	89.04 $\pm$ 3.70	88.32 $\pm$ 0.96	95.58 $\pm$ 0.38	88.80 $\pm$ 0.95	99.65 $\pm$ 0.10
RAP	93.10 $\pm$ 0.84	85.62 $\pm$ 3.58	92.70 $\pm$ 0.88	91.20 $\pm$ 3.45	92.96 $\pm$ 0.86	76.90 $\pm$ 3.80	92.70 $\pm$ 0.90	78.80 $\pm$ 3.70
STRIP	93.74 $\pm$ 0.78	97.90 $\pm$ 3.20	93.70 $\pm$ 0.80	100.00 $\pm$ 1.20	93.50 $\pm$ 0.82	78.70 $\pm$ 1.50	93.60 $\pm$ 0.78	88.90 $\pm$ 1.10
MF	92.80 $\pm$ 0.86	21.30 $\pm$ 3.50	92.60 $\pm$ 0.90	36.00 $\pm$ 2.50	92.80 $\pm$ 0.85	65.80 $\pm$ 2.80	92.90 $\pm$ 0.88	76.50 $\pm$ 2.20
<b>Our Method</b>	<b>93.72 <math>\pm</math>0.84</b>	<b>5.60 <math>\pm</math>2.04</b>	<b>92.72 <math>\pm</math>0.88</b>	<b>6.56 <math>\pm</math>2.23</b>	<b>93.12 <math>\pm</math>0.89</b>	<b>19.36 <math>\pm</math>2.90</b>	<b>93.20 <math>\pm</math>0.93</b>	<b>22.70 <math>\pm</math>2.80</b>

Table 7: Performance comparison with other defense methods on OLID dataset

Defence Method	AddWord		AddSent		StyleBKD		SynBKD	
	ACC ( $\uparrow$ )	ASR ( $\downarrow$ )	ACC ( $\uparrow$ )	ASR ( $\downarrow$ )	ACC ( $\uparrow$ )	ASR ( $\downarrow$ )	ACC ( $\uparrow$ )	ASR ( $\downarrow$ )
No defense	85.23 $\pm$ 0.68	99.83 $\pm$ 0.25	85.00 $\pm$ 0.67	100.00 $\pm$ 0.00	84.88 $\pm$ 0.71	99.24 $\pm$ 0.39	85.23 $\pm$ 0.67	100.00 $\pm$ 0.00
BKI	84.76 $\pm$ 0.89	90.23 $\pm$ 2.67	84.88 $\pm$ 0.84	100.00 $\pm$ 0.00	83.23 $\pm$ 0.98	98.34 $\pm$ 0.42	83.72 $\pm$ 0.95	99.61 $\pm$ 0.25
ONION	84.41 $\pm$ 0.88	58.10 $\pm$ 2.34	85.11 $\pm$ 0.82	100.00 $\pm$ 0.00	85.11 $\pm$ 0.86	99.63 $\pm$ 0.31	84.53 $\pm$ 0.92	99.39 $\pm$ 0.30
RAP	83.93 $\pm$ 0.90	87.18 $\pm$ 3.11	83.72 $\pm$ 0.94	99.44 $\pm$ 0.35	83.54 $\pm$ 0.97	95.23 $\pm$ 1.93	83.91 $\pm$ 0.89	94.45 $\pm$ 1.95
STRIP	85.00 $\pm$ 0.76	100.00 $\pm$ 0.00	83.27 $\pm$ 0.86	99.25 $\pm$ 0.30	84.65 $\pm$ 0.91	88.81 $\pm$ 0.25	83.98 $\pm$ 0.93	79.84 $\pm$ 0.20
MF	81.97 $\pm$ 0.93	21.24 $\pm$ 2.92	81.86 $\pm$ 0.97	68.92 $\pm$ 2.79	82.09 $\pm$ 0.98	68.42 $\pm$ 3.10	82.89 $\pm$ 0.92	58.52 $\pm$ 3.00
<b>Our Method</b>	<b>82.79 <math>\pm</math>0.85</b>	<b>11.45 <math>\pm</math>3.17</b>	<b>83.37 <math>\pm</math>0.82</b>	<b>4.83 <math>\pm</math>2.04</b>	<b>83.95 <math>\pm</math>0.90</b>	<b>29.18 <math>\pm</math>2.92</b>	<b>83.02 <math>\pm</math>0.93</b>	<b>28.40 <math>\pm</math>2.90</b>

520 In this section, we delve deeper into the comparison between our method and several other backdoor  
521 defense strategies, maintaining the same conditions as outlined in Section 5. Particularly, Table 6  
522 shows our honeypot technique against others on the RoBERTa<sub>BASE</sub> with the IMDB dataset. Additionally,  
523 results using the OLID dataset are presented in Table 7. In the case of the IMDB dataset,  
524 our method consistently achieves the lowest ASR across all four attack methods, displaying a robust  
525 defense technique even under varied adversarial conditions. For example, considering the AddWord  
526 and AddSent attacks, our ASR is below 10%, which is a considerable improvement over other  
527 methods. In StyleBKD and SynBKD, our ASR stays below 23%, still outperforming the competing

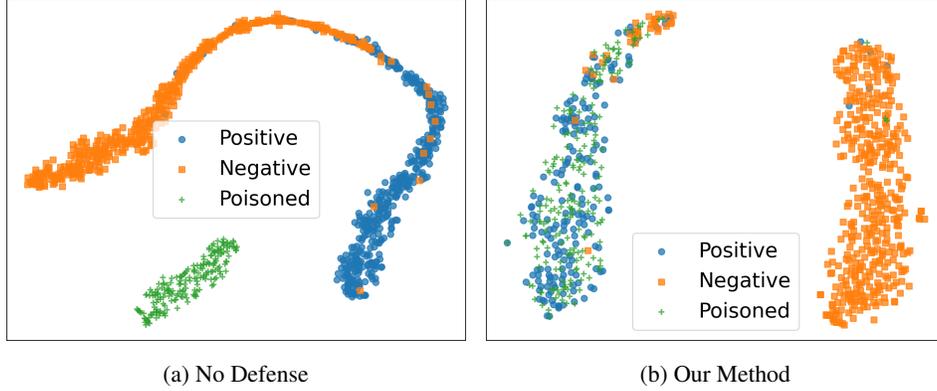


Figure 10: Embedding Visualization for Victim Model and Protected Model

528 methods by a wide margin. Similarly, for the OLID dataset, our method demonstrated excellent  
 529 performance, surpassing all other defense methods in terms of ASR. Furthermore, our method still  
 530 achieves competitive ACC results on the original tasks. In Figure 10, we exhibit the t-SNE visual-  
 531 izations derived from the CLS token embeddings of the final transfer layer of the RoBERTa model.  
 532 As shown in Figure 10 (a), we observe that the no-defense model clearly recognizes the poisoned  
 533 samples. Instead, in Figure 10 (b), the model overlooks the backdoor trigger and successfully predicts  
 534 positive samples with embedded backdoor words as the positive class.

### 535 C Understanding the Honey-pot Defense Training Process

536 In this section, we further illustrate more details about the honey-pot defense training process. Specifi-  
 537 cally, we focus on the dynamic change of the training weight for poisoned and clean samples. As we  
 538 mentioned in Section 4, we propose employing a weighted cross-entropy loss ( $\mathcal{L}_{WCE}$ ):

$$\mathcal{L}_{WCE}(f_T(x), y) = \sigma(W(x) - c) \cdot \mathcal{L}_{CE}(f_T(x), y), \text{ where} \quad (5)$$

539

$$W(x) = \frac{\mathcal{L}_{CE}(f_H(x), y)}{\mathcal{L}_{CE}(f_T(x), y)}, \quad (6)$$

540  $f_H(x)$  and  $f_T(x)$  represent the softmax outputs of the honey-pot and task classifiers, respectively.  
 541 The function  $\sigma(\cdot)$  serves as a normalization method, effectively mapping the input to a range within  
 542 the interval  $[0, 1]$ . The  $c$  is a threshold value for the normalization.

543 In order to gain a deeper understanding of the re-weighting mechanism, we extend our analysis  
 544 by presenting both the original  $W(x)$  and the normalized weight  $\sigma(W(x) - c)$ . We conducted the  
 545 experiment using the SST2 dataset, with a word-level trigger, a poisoning rate set at 5%, and a batch  
 546 size of 32. Figure 11 illustrates the  $W(x)$  value for both the poisoned and clean samples at each  
 547 stage of training. Specifically, we computed the  $W(x)$  for each mini-batch and then calculated the  
 548 average  $W(x)$  value for both the poisoned and clean samples. As depicted in the figure, during the  
 549 warm-up phase, the  $W(x)$  for clean and poisoned samples diverged early in the training process. After  
 550 500 steps, the  $W(x)$  for poisoned samples was noticeably lower than for clean samples. After the  
 551 warm-up stage, given that  $W(x)$  is higher for clean samples, the Cross-Entropy loss of clean samples  
 552 in  $f_T$  diminishes more quickly than that of the poisoned samples. This subsequently increase  $W(x)$   
 553 for clean samples as they possess a smaller  $\mathcal{L}_{CE}(f_T(x), y)$ . This positive feedback mechanism  
 554 ensures that the  $W(x)$  for poisoned samples persistently remains significantly lower than for clean  
 555 samples throughout the complete training process of  $f_T$ . As demonstrated in Figure 11, the  $W(x)$   
 556 for the clean samples will continue to increase following the warm-up phase.

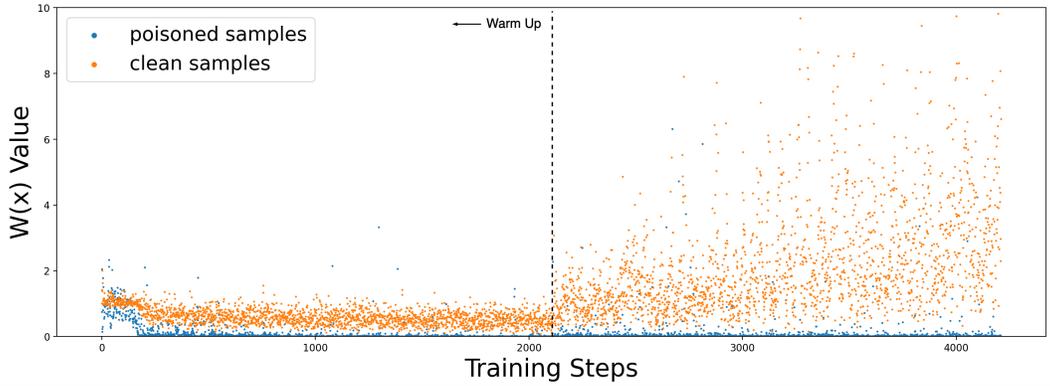


Figure 11: Visualization of  $W(x)$  during defense training process.

## 557 D More on Ablation Studies

### 558 D.1 Ablation Study on HoneyPot Warm-Up

559 In the following section, we explore the influence of the preliminary warm-up steps in the honeypot  
 560 method, which represent the number of optimizations that the honeypot branch requires to capture  
 561 a backdoor attack. We applied our method against word-level attacks on RoBERTa<sub>BASE</sub>, and the  
 562 obtained results are shown in Table 8. The analysis indicates that with a minimum count of warm-  
 563 up steps, specifically below 200 for the SST-2 dataset, the honeypot is insufficiently prepared to  
 564 capture the poisoned data. However, once the honeypot accrues a sufficient volume of poisoned data,  
 565 surpassing 400 training steps across all datasets, the Attack Success Rate (ASR) can be mitigated to  
 566 an acceptably low level, i.e., less than 10%. The results further prove that our honeypot can effectively  
 567 capture backdoor information with a certain amount of optimization. In our main experiments, we  
 568 set the number of warm-up steps equal to the steps in one epoch, thereby enabling our honeypot to  
 569 reliably catch the poisoned data.

Table 8: Impact of Warm-Up steps

Dataset	SST-2					IMDB					
	Warm-Up Steps	100	200	400	1000	2000	100	200	400	1000	2000
ACC ( $\uparrow$ )		94.61	94.72	94.50	94.41	94.15	94.71	94.80	94.26	94.33	94.12
ASR ( $\downarrow$ )		100.00	100.00	8.64	5.37	5.84	100.00	7.62	5.32	5.79	3.60

### 570 D.2 Ablation Study on Normalization Method

571 In this section, we use the SST2 dataset and  
 572 word-level trigger to understand the impact of  
 573 different normalization functions. As outlined  
 574 in Section 4, our approach employs a normal-  
 575 ization method to map the training loss weight  
 576  $W(x)$  into the  $[0, 1]$  interval. Within our exper-  
 577 iments, we opted for the sign function as the  
 578 normalization technique. However, we also ex-  
 579 plored two alternative normalization strategies  
 580 – the sigmoid function and a cutoff ReLU func-  
 581 tion. For the latter, we assigned a value of 1 to any input exceeding 1. As depicted in Table 9, we  
 582 conducted the experiments on RoBERTa<sub>BASE</sub> using different normalization functions, we can observe  
 583 that all normalization methods demonstrate decent performance in minimizing the ASR. Notably,  
 584 we observe that the sign function yields the highest ACC on the original task while simultaneously  
 585 achieving the lowest ASR.

Table 9: Impact of Normalization Method

Normalization	AddWord	
	ACC ( $\uparrow$ )	ASR ( $\downarrow$ )
No Defense	94.61 $\pm$ 0.60	100.00 $\pm$ 0.00
Sign	93.71 $\pm$ 0.68	6.56 $\pm$ 1.91
Sigmoid	93.22 $\pm$ 0.53	6.83 $\pm$ 2.01
Cutoff Relu	93.10 $\pm$ 0.71	6.77 $\pm$ 1.04

## 586 E Extend HoneyPot to Computer Vision Tasks

587 While this paper primarily focuses on defending pretrained language models against backdoor attacks,  
 588 we also explored the applicability of our proposed honeyPot defense method within the computer  
 589 vision domain [3, 4, 5]. In Section E.1, we illustrate the experimental settings. In Section E.2, we  
 590 show the empirical findings. In Section E.3, we discuss the defense performance.

### 591 E.1 Settings

592 Suppose  $D_{train} = (x_i, y_i)$  indicates a benign training dataset where  $x_i \in \{0, \dots, 255\}^{C \times W \times H}$   
 593 represents an input image with  $C$  channels and  $W$  width and  $H$  height, and  $y_i$  corresponds to  
 594 the associated label. To generate a poisoned dataset, the adversary selects a small set of samples  
 595  $D_{sub}$  from the original dataset  $D_{train}$ , typically between 1-10%. The adversary then chooses a  
 596 target misclassification class,  $y_t$ , and selects a backdoor trigger  $a$  and  $a \in \{0, \dots, 255\}^{C \times W \times H}$ . For  
 597 each instance  $(x_i, y_i)$  in  $D_{sub}$ , a poisoned example  $(x'_i, y'_i)$  is created, with  $x'_i$  being the embedded  
 598 backdoor trigger of  $x_i$  and  $y'_i = y_t$ . The trigger embedding process can be formulated as follows,

$$x'_i = (1 - \lambda) \otimes x + \lambda \otimes a, \quad (7)$$

599 where  $\lambda \in [0, 1]^{C \times W \times H}$  is a trigger visibility hyper-parameter and  $\otimes$  specifies the element-wise  
 600 product operation. The smaller the  $\lambda$ , the more invisible the trigger and the more stealthy. The  
 601 resulting poisoned subset is denoted as  $D'_{sub}$ . Finally, the adversary substitutes the original  $D_{sub}$   
 602 with  $D'_{sub}$  to produce  $D_{poison} = (D_{train} - D_{sub}) \cup D'_{sub}$ . By fine-tuning PLMs with the poisoned  
 603 dataset, the model will learn a backdoor function that establishes a strong correlation between the  
 604 trigger and the target label  $y_t$ . Consequently, adversaries can manipulate the model's predictions  
 605 by adding the backdoor trigger to the inputs, causing instances containing the trigger pattern to be  
 606 misclassified into the target class  $t$ .

607 In our experiment, we employed an ImageNet pretrained VGG-16 model as our base architecture and  
 608 proceed with experiments using a manipulated CIFAR-10 dataset. The experiments involve the use of  
 609 a  $3 \times 3$  white square and a black line with a width of 3 pixels as backdoor triggers. The white square  
 610 trigger is positioned at the bottom-right corner of the image, while the black line trigger is set at the  
 611 bottom. We establish a poison rate of 5% and set  $\lambda \in \{0, 0.2\}^{C \times W \times H}$  for two attacks. The values  
 612 of  $\lambda$  corresponding to pixels situated within the trigger area are 0.2, while all others are set to 0.

### 613 E.2 Lower Layer Representations from VGG Provide Sufficient Backdoor Information

614 Drawing on our analysis presented in Section 3, we delve further into understanding the information  
 615 encapsulated within various layers of a pretrained computer vision model. Inspired by previous  
 616 classifier probing studies [33, 34], we train a compact classifier using representations derived from  
 617 different layers of the VGG model. We ensure the VGG model parameters are frozen during this  
 618 process and only train the probing classifier. In this context, we divided the VGG model into five  
 619 sections based on the pooling layer operations (The five pooling layers are located at layers 2, 4, 7, 10,  
 620 and 13). Subsequent to this, we integrate an adaptive pooling layer to reduce the features extracted  
 621 from different layers to  $7 \times 7$ , ensuring that the flattened dimension does not exceed 8000. A fully  
 622 connected layer with softmax activation is added as the final output. As depicted in Figure 13 and  
 623 Figure 12, it is noticeable that the lower layers of the VGG model hold sufficient information for  
 624 identifying the backdoor triggers. However, they do not contain enough information to effectively  
 625 carry out the main tasks.

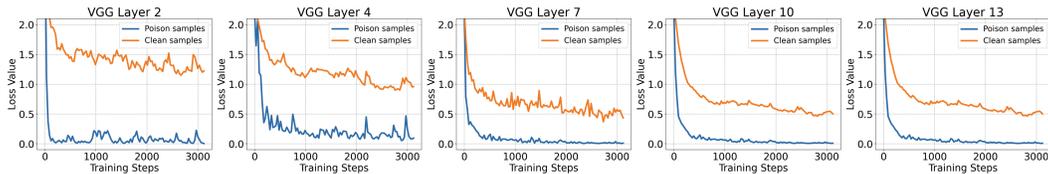


Figure 12: Learning Dynamic for White Square Trigger

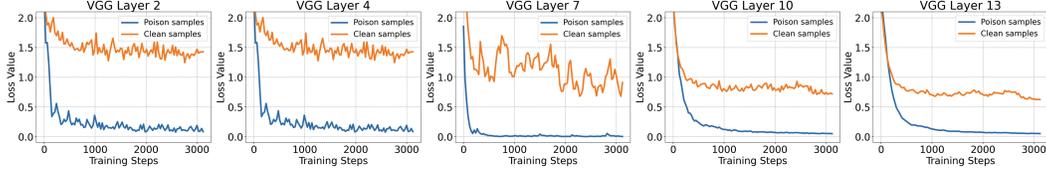


Figure 13: Learning Dynamic for Black Line Trigger

Table 10: Defense Performance on CIFAR10

Method	White Square		Black Line	
	ACC ( $\uparrow$ )	ASR ( $\downarrow$ )	ACC ( $\uparrow$ )	ASR ( $\downarrow$ )
No Defense	91.33 $\pm$ 0.27	100.00 $\pm$ 0.00	91.28 $\pm$ 0.13	100.00 $\pm$ 0.00
Our Method	92.20 $\pm$ 0.43	8.81 $\pm$ 1.09	92.23 $\pm$ 0.37	10.81 $\pm$ 1.83

### 626 E.3 Defense Results on CIFAR10

627 We implemented the honeypot as mentioned in Section 4 and built the honeypot module with the  
 628 features from the first pooling layer. We followed previous sections and adopted the ASR and ACC  
 629 metrics to measure the model’s performance on the poisoned test set and clean test set, respectively.  
 630 Specifically, we executed a fine-tuning process for a total of 10 epochs, incorporating an initial  
 631 warmup epoch for the honeypot module. The learning rates for both the honeypot and the principal  
 632 task are adjusted to a value of  $1 \times 10^{-3}$ . Additionally, we established the hyperparameter  $q$  for the  
 633 GCE loss at 0.5, the time window size  $T$  was set to 100, and the threshold value  $c$  was fixed at 0.1.  
 634 Each experimental setting was subjected to three independent runs and randomly chosen one class as  
 635 the target class. These runs were also differentiated by employing distinct seed values. The results  
 636 were then averaged, and the standard deviation was calculated to present a more comprehensive  
 637 understanding of the performance variability. As the results are shown in Table 10, the proposed  
 638 method successfully defends two backdoor attacks and reduces the ASR to lower than 10%. This  
 639 indicates that the proposed method is valid for those simple vision backdoor triggers while having  
 640 minimal impact on the original task. We plan to test the defense performance of more advanced  
 641 backdoor triggers in our future work.

### 642 F Reproducibility

643 In an effort to ensure the reproducibility of our results, we have shared our test code along with  
 644 the model checkpoint. This will allow peers in the research community to validate our findings. In  
 645 the interest of complete transparency, we are also committed to releasing our training code in the  
 646 future. This will provide a comprehensive understanding of our methodology and enable fellow  
 647 researchers to extend and build upon our work. Our code can be found at [https://anonymous](https://anonymous.4open.science/r/honeypot-backdoor-600E)  
 648 [4open.science/r/honeypot-backdoor-600E](https://anonymous.4open.science/r/honeypot-backdoor-600E).

### 649 G Limitations and Discussions

650 In this study, we introduce an innovative approach to backdoor defense in the context of fine-tuning  
 651 pretrained language models. Due to the constraints in terms of time and resources, our evaluations  
 652 were conducted using four prevalent backdoor attack methods and on three representative datasets.  
 653 Despite the robustness and consistency demonstrated by our method, it is essential to remain vigilant to  
 654 the emergence of new and potentially threatening attack methods and datasets, especially considering  
 655 the rapid growth of this field. In addition, it’s worth acknowledging that while unintended, some  
 656 malicious users may exploit our method and deploy other strong backdoor attacks that may bypass  
 657 our defense system.