

## 545 Appendix

### 546 A Details of datasets and architectures

#### 547 A.1 Object Detection Image Dataset

548 **COCO (Common Objects in Context)** [26] dataset is widely used for object detection tasks. It  
549 contains 80 object categories, including people, animals, vehicles and more. Each image can contain  
550 multiple instances of objects, providing ample opportunities for training and evaluating models  
551 capable of detecting and segmenting objects in complex scenes.

552 **Synthesized Traffic Sign** dataset is designed by TrojAI [1] which focuses on traffic sign detection,  
553 featuring various types of traffic signs commonly encountered in real-world scenarios. There are  
554 in total over 4000 different traffic signs. Each model is trained on a randomly sampled subset of  
555 classes. The number of classes within these subsets exhibits variability, ranging from as few as 2 to a  
556 maximum of 128.

557 **DOTA (Detection in Aerial Images)** dataset is designed for object detection in aerial images which  
558 consists of high-resolution images captured by aerial platforms. It contains 18 categories, including  
559 plane, ship, storage tank, baseball diamond and more. Its large-scale, fine-grained annotations, and  
560 challenging scenarios make it an ideal benchmark for evaluating and developing algorithms capable  
561 of detecting objects in aerial images accurately.

#### 562 A.2 Architecture

563 We evaluate our method on three well-known model architectures:, i.e., SSD [28], Faster-RCNN [40],  
564 and DETR [2]. SSD (Single Shot MultiBox Detector) [28] is a popular object detection model  
565 which utilizes a series of convolutional layers to detect objects at multiple scales and aspect ratios.  
566 Faster-RCNN [40] is another widely adopted object detection model that combines region proposal  
567 generation with a region-based CNN for object detection. DETR (DEtection TRansformer) [2] is a  
568 state-of-the-art object detection model that utilizes a transformer-based architecture. It replaces the  
569 conventional two-stage approach with a single-stage end-to-end detection framework.

#### 570 A.3 Model Dataset

571 **TrojAI** [1] initiative, spearheaded by IARPA, encompasses a multi-year, multi-round program.  
572 Its overarching objective revolves around the development of scalable and dependable automatic  
573 backdoor detection tools, specifically targeting the identification of backdoor trojans within Deep  
574 Learning models across diverse modalities. Presently, the program consists of a total of 13 rounds,  
575 each with distinct focuses and tasks. The initial four rounds and the eleventh round center their efforts  
576 on detecting trojans present in image classification models. In contrast, rounds five through nine  
577 concentrate on transformer models employed in various NLP tasks, including Sentiment Analysis,  
578 Named Entity Recognition, and Question Answering. Round twelve dedicates itself to the detection  
579 of backdoors in neural network-based PDF malware detection. Finally, rounds ten and thirteen  
580 direct their attention towards object detection models. For the evaluation of models, we exclusively  
581 utilize the training sets from rounds 10 and 13. Specifically, our evaluation entails 72 models trained  
582 on the Synthesis Traffic Sign dataset, encompassing all three model architectures. Among these  
583 models, 48 are benign, while 24 are deliberately poisoned, with an equal distribution of triggers for  
584 misclassification and evasion. Concerning the DOTA models, there exist two architectures, namely  
585 SSD and Faster-RCNN, resulting in a total of 24 models, including 16 benign models and 4 each  
586 poisoned with misclassification and evasion triggers. All COCO models adopt the SSD architecture,  
587 with a distribution of 36 clean models and 18 models poisoned by both misclassification and evasion  
588 triggers. Find more details in Table 5.

### 589 B Details of evaluation metrics

590 In our evaluation of backdoor detection methods, we employ four well-established metrics: Precision,  
591 Recall, ROC-AUC, and Average Scanning Overheads for each model. Precision quantifies the  
592 accuracy of a detection method by measuring the proportion of correctly identified positive instances

593 among all predicted positives. In our case, we consider attacked models as positive instances and  
594 benign models as negatives. A higher precision indicates a lower rate of falsely identifying benign  
595 models as attacked. Recall, on the other hand, assesses the effectiveness of the detection method in  
596 correctly identifying positive instances. It measures the proportion of true positives among all actual  
597 positives. A higher recall suggests that the detection method is capable of identifying a significant  
598 portion of attacked models. ROC-AUC (Receiver Operating Characteristic - Area Under the Curve)  
599 plots the true positive rate against the false positive rate at various threshold values and calculates  
600 the area under the curve. A value of 1 indicates perfect classification, while a value of 0.5 indicates  
601 that the method is no better than random guessing. We also consider the overhead of the detection  
602 method, which quantifies the average time required to scan a single model. We use seconds (s) as  
603 the unit of measurement and set a maximum threshold of 1 hour (3600 s). If the scanning process  
604 exceeds 3600 seconds, it is terminated, and we rely on the existing results for making predictions.  
605 Low overhead signifies high efficiency of the method. By employing these four metrics, we aim to  
606 comprehensively evaluate the performance and efficiency of the backdoor detection methods. It is  
607 worth noting that the time limit we have set for scanning models is deliberately conservative when  
608 compared to the thresholds established in different rounds of the TrojAI competition. For example,  
609 in round 13, participants are granted a generous 30-minute duration for scanning a single model.  
610 To surpass the official benchmarks set in each round, a more aggressive and precise pre-processing  
611 approach may be necessary.

## 612 C Details of Baseline Methods

613 In this section, we introduce more details of baseline methods, including NC [51], Tabor [15],  
614 ABS [29], Pixel [49], Matrix Factorization(MF) [17] and MNTD [60].

615 NC [51] adopts a specific trigger inversion approach for each class and considers a model to be  
616 attacked if it is able to generate an effective yet extremely small trigger for a target class. Tabor [15]  
617 enhances NC by incorporating additional well-designed regularization terms, such as penalties for  
618 scattered triggers, overlaying triggers, and blocking triggers. These additions aim to improve the  
619 reconstruction of injected triggers. Pixel [49] introduces a novel inversion function that generates a  
620 pair of positive and negative trigger patterns. This approach achieves better detection performance  
621 compared to NC. ABS [29] employs a stimulation analysis to identify compromised neurons, which  
622 serves as guidance for trigger inversion. ABS considers a model to be attacked if it can invert a trigger  
623 that achieves a high reconstructed ASR (REASR).

624 To the best of our knowledge, there is no existing detection methods for object detection models.  
625 Therefore, we perform straight-forward but reasonable adaption to these existing methods designed  
626 on image classification tasks, such that they are able to work against backdoor attacks on object  
627 detection models. Specifically, the original objective of NC, Tabor, and Pixel is to invert small triggers  
628 while maintaining their effectiveness (high ASR). In our adaptation, we retain their design principles  
629 but re-define the ASR to align with object detection models, as explained in Section 3.1. Additionally,  
630 we introduce a threshold for the size of inverted triggers, enabling the differentiation between benign  
631 and attacked models. For ABS, we adhere to its original technique but employ the re-defined ASR as  
632 the optimization goal, and use REASR as the decision score. By employing these adaptations, we  
633 aim to enhance the detection capabilities of these existing methods specifically for backdoor attacks  
634 on object detection models.

635 No modifications or adaptations are needed for meta classification-based methods when applied to  
636 object detection models. MNTD [60] trains a set of queries and a classifier to discern the feature-space  
637 distinctions between clean and attacked models. MF [17] directly trains a classifier on model weight  
638 features using specialized feature extraction techniques, i.e., matrix factorization. These methods  
639 solely rely on the feature extraction networks commonly utilized in both image classification and  
640 object detection models. As a result, MNTD and MF can be directly employed to detect backdoor  
641 attacks in object detection models without the need for additional adjustments or modifications.

642 We collect the Precision, Recall, ROC-AUC and Overheads for each method across various datasets  
643 and model architectures. To ensure a fair comparison, we have conducted a search to determine the  
644 optimal thresholds for different decision scores associated with each method (trigger size for NC,  
645 Tabor, Pixel, REASR for ABS and output confidence for meta-classifiers). These thresholds are  
646 chosen to maximize accuracy. Besides, we set a fixed number of optimization steps for scanning

Table 5: Dataset Details

Image Dataset	Model Source		Architecture			Number of Models		
	Round10	Round13	SSD	Faster-RCNN	DETR	Benign	Miscs Attack	Evasion Attack
Synthesis Traffic Sign	✗	✓	✓	✓	✓	48	12	12
DOTA	✗	✓	✓	✓	✗	16	4	4
COCO	✓	✗	✓	✗	✗	36	18	18

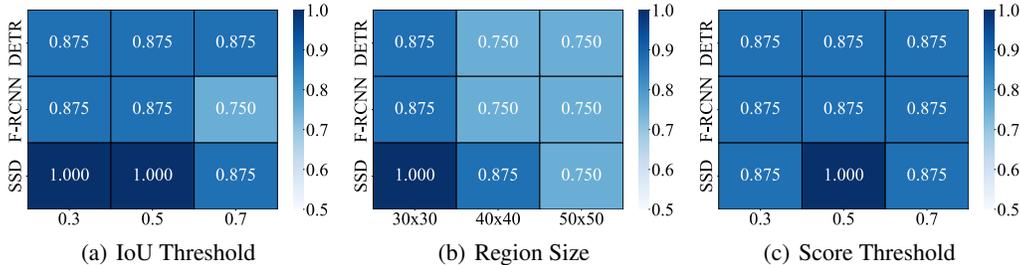


Figure 6: Hyper-parameter Sensitivity.

647 a pair of victim-target label (100) for all inversion based baselines. For meta classification based  
 648 methods that involve training, we have performed 5-fold cross-validation and reported the validation  
 649 results exclusively.

## 650 D Hyper-parameter Sensitivity Analysis

651 To assess the sensitivities of the hyper-parameters used in DJANGO, we conduct experiments as  
 652 described in Section 4.3

653 **IoU Thresholds.** We evaluate the IoU threshold used to calculate the ASR of inverted triggers. The  
 654 results are summarized in Figure 6(a), where each row corresponds to a different model architecture,  
 655 and each column represents a different choice of IoU threshold. It can be observed that IoU thresholds  
 656 of 0.3 and 0.5 generally yield good performance. However, a threshold of 0.7 tends to degrade the  
 657 performance, possibly due to the inverted triggers interfering with the bounding box predictions.

658 **Region Size.** The impact of different regional initialization sizes is evaluated and the results are  
 659 presented in Figure 6(b). Among the various choices, a region size of  $30 \times 30$  consistently achieved  
 660 the best performance. This is because larger initialization sizes tend to result in more false positive  
 661 cases.

662 **Score Threshold.** Different score thresholds are tested when computing the ASR of inverted triggers.  
 663 The results, shown in Figure 6(c), indicate that a score threshold of 0.5 generally leads to the best  
 664 performance across all model architectures. This choice represents a trade-off between false positives  
 665 and false negatives. Higher score thresholds may introduce more false negatives, as the inverted  
 666 trigger may not have high confidence similar to the injected one. On the other hand, lower score  
 667 thresholds may result in more false positives. Thus, a moderate value of 0.5 provides the optimal  
 668 balance.

669 These experiments allowed us to gain insights into the sensitivities of the hyper-parameters in  
 670 DJANGO, enabling us to make informed choices for achieving optimal performance.