

# Robust Object Detection Fusion Against Deception

Ka-Ho Chow

Georgia Institute of Technology  
Atlanta, Georgia, USA

## ABSTRACT

Deep neural network (DNN) based object detection has become an integral part of numerous cyber-physical systems, perceiving physical environments and responding proactively to real-time events. Recent studies reveal that well-trained multi-task learners like DNN-based object detectors perform poorly in the presence of deception. This paper presents FUSE, a deception-resilient detection fusion approach with three novel contributions. First, we develop diversity-enhanced fusion teaming mechanisms, including diversity-enhanced joint training algorithms, for producing high diversity fusion detectors. Second, we introduce a three-tier detection fusion framework and a graph partitioning algorithm to construct fusion-verified detection outputs through three mutually reinforcing components: objectness fusion, bounding box fusion, and classification fusion. Third but not least, we provide a formal analysis of robustness enhancement by FUSE-protected systems. Extensive experiments are conducted on eleven detectors from three families of detection algorithms on two benchmark datasets. We show that FUSE guarantees strong robustness in mitigating the state-of-the-art deception attacks, including adversarial patches – a form of physical attacks using confined visual distortion.

## CCS CONCEPTS

- Computing methodologies → Object detection; Computer vision; Neural networks; Machine learning.

## KEYWORDS

object detection; adversarial machine learning; ensemble defense

### ACM Reference Format:

Ka-Ho Chow and Ling Liu. 2021. Robust Object Detection Fusion Against Deception. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '21), August 14–18, 2021, Virtual Event, Singapore*. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3447548.3467121>

## 1 INTRODUCTION

Deep neural network (DNN) based object detection has been applied to the internet of smart sensors and edge computing applications (e.g., smart surveillance systems [8]). Unlike single-task learners such as image classifiers, object detectors are a multi-task learner, which takes an input video frame or image and detects

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

*KDD '21, August 14–18, 2021, Virtual Event, Singapore*

© 2021 Association for Computing Machinery.

ACM ISBN 978-1-4503-8332-5/21/08...\$15.00

<https://doi.org/10.1145/3447548.3467121>

Ling Liu

Georgia Institute of Technology  
Atlanta, Georgia, USA

	Autonomous Driving		Home Securities	
	Faster R-CNN Standard Detector	Robust Fusion Detector (Ours)	Faster R-CNN Standard Detector	Robust Fusion Detector (Ours)
Benign Input				
				

Table 1: The standard detector without protection correctly detects objects given benign inputs but fails on adversarial ones. Our robust fusion detector can maintain detection performance in both cases.

multiple instances of objects of known categories by reporting the number of objects (objectness), the bounding box location for each detected object, and the class labels of detected objects [22]. Existing DNN-based detection algorithms are broadly classified into two categories: (1) proposal-based two-phase learning and (2) regression-based one-shot learning. The proposal-based approaches, dominated by the R-CNN family [9, 19], use a two-phase procedure by first detecting proposal regions using a region proposal network and then refining each of them with bounding box estimation and class label prediction. The regression-based approaches, represented by YOLO [18] and SSD [14], directly detect multiple objects, estimating their bounding box locations and class labels in one shot.

While DNN-based object detectors offer real-time performance with unparalleled accuracy over traditional techniques [22], recent studies have revealed their vulnerabilities to adversarial inputs [5, 13, 24, 26]. Table 1 illustrates such vulnerabilities by examples. With no attack, the detector can accurately identify the person and the car in the driving scene and the person (a thief) breaking into a house. However, the same detector can easily be deceived by the small adversarial noises injected into the inputs, e.g., misdetecting the thief as a bird. Such vulnerabilities pose serious threats to mission-critical systems and can lead to disastrous consequences [4].

In this paper, we present FUSE, a robust model fusion approach for safeguarding object detection systems against deception by adversarial inputs while maintaining high benign performance. FUSE can auto-correct suspicious detection results with the fusion-verified detection outputs (see the 2nd and 4th columns in Table 1). This paper makes three original contributions. First, we present robust fusion teaming strategies for producing diversity-enhanced fusion teams of detectors. Second, we introduce a three-tier detection fusion framework and a graph partitioning algorithm to safeguard the victim detector with three mutually reinforcing components: objectness fusion, bounding box fusion, and classification fusion, for high resilience against deception. Third, we formalize

adversarial vulnerability in the presence of deception and theoretically analyze the robustness enhancement by our detection fusion approach. We conduct extensive experiments on eleven detectors from three representative families of detection algorithms (Faster R-CNN [19], YOLOv3 [18], and SSD [14]) on two popular benchmark datasets (PASCAL VOC [7] and INRIA [6]). FUSE demonstrates strong robustness for mitigating the state-of-the-art adversarial deception attacks, e.g., DAG [26], RAP [13], UEA [24], and TOG [5], as well as the adversarial patches [21], a form of physical attacks using confined visual distortion [15].

## 2 RELATED WORK

Object detection is an instance of multi-task learning, which takes an input image and performs three learning tasks: object existence prediction (detecting the number of objects), bounding box estimation (bounding box of each detected object), and object classification (class label for each object). Existing state-of-the-art attacks to object detectors include DAG [26], RAP [13], UEA [24], TOG [5], and adversarial physical patches [21] or digital patches [15]. Unlike adversarial attacks to DNN-based image classifiers [10], all object detection attacks can critically compromise the object existence prediction capability and the bounding box estimation capability, causing real objects to vanish from the detection or causing the detection to fabricate fake objects that do not exist, consequently deceiving object classification (i.e., the third learning task). Given that these attack algorithms are more sophisticated and capable of deceiving three learning tasks at the same time [4], it has been an open challenge to develop effective mitigation methods. Only one recent proposal attempts to improve the resilience of a victim by retraining the detector with adversarial inputs through mini-max optimization on adversarial objectness loss and adversarial bounding box loss [29]. However, it suffers from two serious drawbacks: (i) the detection accuracy on benign inputs (no attacks) of the retrained detector is significantly lower than that of the victim detector; and (ii) the retrained detector offers insufficient resilience against deception, making the proposal ineffective for protecting the victim detector. This motivates us to develop FUSE for robust detection fusion with high resilience against adversarial inputs while maintaining high performance in benign scenarios.

## 3 OVERVIEW

**DNN-based Object Detection.** A  $K$ -class object detector  $F_\theta : \mathbb{R}^D \rightarrow \mathbb{R}^{C(K+5)}$ , parameterized by  $\theta$ , is a mapping from the  $D$ -dimensional input image  $x$  to a set of  $C$  candidate objects  $F_\theta(x)$ , where each candidate  $\mathbf{o} \in F_\theta(x)$  is represented by (i) the objectness probability  $\mathbf{o}_{[\text{obj}]} \in [0, 1]$  of the candidate being real, (ii) the bounding box  $\mathbf{o}_{[\text{bbox}]}$  of the candidate, and (iii) the  $K$ -class probability vector  $\mathbf{o}_{[\text{prob}]}$  with  $\mathbf{o}_{[\text{cls}]}$  as the predicted class label of the object. To serve the query with high-quality results, the final set of detected objects  $\hat{F}_\theta(x) \subset F_\theta(x)$  is produced by filtering out candidates either with low prediction confidence or highly overlapped with another candidate of the same class using non-maximum suppression.

With the three types of outputs in object detection (objectness, bounding boxes, and class labels), we can formulate the training process of object detection as a multitask learning problem for a given dataset  $\mathcal{D}$ :  $\min_{\theta} \mathbb{E}_{(x, \mathbf{O}) \sim \mathcal{D}} [\mathcal{L}(F_\theta(x), \mathbf{O})]$ , which aims to minimize the prediction error of object detector  $F_\theta$  on (i) objectness

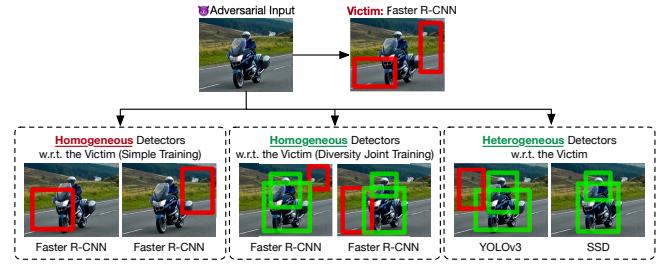


Figure 1: An adversarial input is likely to transfer and fool homogeneous detectors generated by a simple training strategy (left), but the attack transferability is weakened when the homogeneous detectors are generated with a diversity-enhanced training approach (middle) or a team of heterogeneous detectors is employed (right).

(ii) bounding boxes  $\mathcal{L}_{\text{bbox}}$ , and (iii) class labels  $\mathcal{L}_{\text{cls}}$  of objects, expressed by the loss function:

$$\mathcal{L}(F_\theta(x), \mathbf{O}) = \mathcal{L}_{\text{obj}}(F_\theta(x), \mathbf{O}) + \mathcal{L}_{\text{bbox}}(F_\theta(x), \mathbf{O}) + \mathcal{L}_{\text{cls}}(F_\theta(x), \mathbf{O}), \quad (1)$$

where  $x$  is an input image with a set of ground truth objects  $\mathbf{O}$ .

**Adversarial Attacks.** Given a benign input  $x$  and a victim detector  $F_\theta$ , the adversarial attack is a malicious process to generate the adversarial example by injecting a perturbation  $\delta$  produced by a constrained optimization process:

$$\max_{\delta} \mathbb{1}\{\hat{F}_\theta(x + \delta) = \mathbf{O}^* \wedge \hat{F}_\theta(x + \delta) \neq \hat{F}_\theta(x)\} \text{ s.t. } \|\delta\|_p < \epsilon, \quad (2)$$

where  $\mathbb{1}\{\cdot\}$  is the indicator function,  $\mathbf{O}^*$  denotes the target detection for targeted attacks and any incorrect detection for untargeted attacks, and  $\epsilon$  is the maximum perturbation in  $\ell_p$ -norm. Directly optimizing the indicator function in Equation 2 is challenging, and hence existing attacks [5, 13, 26] reformulate the optimization to be

$$\min_{\delta} \mathcal{L}^*(F_\theta(x + \delta), \mathbf{O}^*) \text{ s.t. } \|\delta\|_p < \epsilon, \quad (3)$$

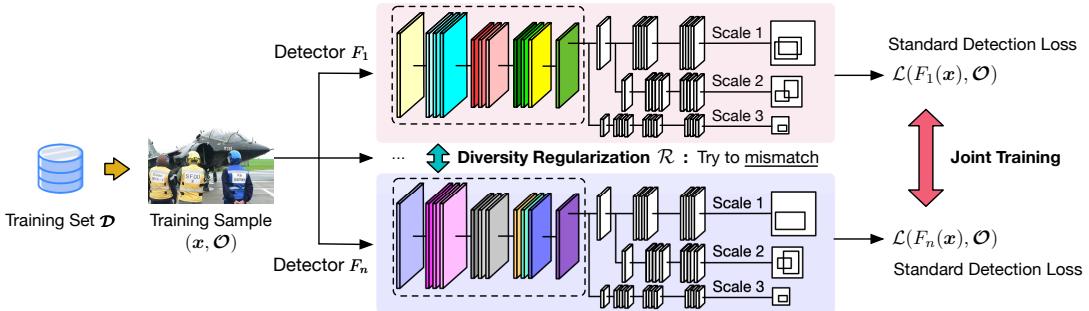
where  $\mathcal{L}^*$  is an adversarial loss function composed of one or more components in the standard detection loss (Equation 1). Minimizing the above adversarial loss can be achieved by iteratively taking the gradients w.r.t. the input to obtain the amount and the location of pixels for adversarial perturbations. Concretely, the adversarial example  $x'_t = x + \delta_t$  at the  $t$ -th iteration is constructed as follows:

$$x'_t = \prod_{x, \epsilon} [\Gamma_{x, \epsilon} [\mathbf{x}'_{t-1} - \alpha \nabla_{\mathbf{x}'_{t-1}} \mathcal{L}^*(F_\theta(\mathbf{x}'_{t-1}), \mathbf{O}^*)]], \quad (4)$$

where  $\prod_{x, \epsilon}$  is the projection onto a hypersphere with a radius  $\epsilon$  centered at  $x$  in  $\ell_p$ -norm followed by clipping to ensure the validity of pixel values,  $\alpha$  is the attack learning rate, and  $\Gamma$  is a sign function. By strategically configuring  $\mathcal{L}^*$  and  $\mathbf{O}^*$ , an adversary can launch untargeted random attacks [5, 13, 24, 26] and targeted specificity attacks, such as object vanishing, fabrication, or mislabeling [5].

**Object Detection Fusion: Solution Approach.** To combat the growing number of attacks, we propose FUSE, a robust detection fusion approach for safeguarding object detection systems against deception. FUSE consists of two complementary functional components: (i) *diversity-enhanced fusion teaming* and (ii) *robustness synergy of the three-tier detection fusion*.

The first component is designed based on our formal analysis (Section 6) and empirical observation that the higher diversity the detection fusion team has, the stronger robustness it can guarantee. This is because high diversity detection fusion reflects high failure independence among member detectors. Figure 1 provides an



**Figure 2: The diversity-driven joint training framework simultaneously trains multiple models with a diversity regularization  $\mathcal{R}$  that encourages them to learn differently using diversified detection loss and/or kernel filter regularization.**

illustrative example, indicating that a random assemble (left) of pre-trained detectors of the same (homogeneous) detection algorithm is less robust compared to teaming of diverse detectors, which are either homogeneous models generated with diversity-enhanced training (middle) or heterogeneous detectors (right). We describe our fusion teaming approach in Section 4.

The second component is designed to address three technical challenges in object detection fusion: (i) Different object detectors may produce slightly or significantly different objectness prediction (i.e., the total number of detected objects) w.r.t. the same input image. (ii) Similarly, due to the regression nature of bounding box estimation, different object detectors may produce slightly or significantly different bounding box estimations for the detected objects, deriving several additional challenges. First, even if two detectors recognize the same real object given an input image (e.g., the motor-bike in Figure 1), the two bounding boxes estimated almost never perfectly align. Second, such misalignments may make the bounding box fusion for each real object much harder to accomplish. Third, deception attacks to victim detectors with object vanishing or fabrication effects can further complicate the disagreement in both objectness and bounding box alignment. (iii) The challenges in objectness fusion and bounding box fusion will consequently make the classification fusion vulnerable in the presence of deception. We describe in Section 5 the FUSE solution approach in detail.

## 4 DIVERSE FUSION TEAM GENERATION

This section describes a suite of techniques to generate high diversity detection models for building diverse fusion team(s) and guaranteeing high robustness of detection fusion against deception.

**Heterogeneous Neural Architectures.** The first strategy is to prepare a collection of pre-trained models with detection algorithms from different families (e.g., YOLOv3, SSD, Faster R-CNN) or using different neural backbones (e.g., VGG16, Darknet53, MobileNetV1). This neural architectural diversity can significantly boost the robustness offered by FUSE because object detection models trained using fundamentally different neural architectures tend to learn and capture complex features of the same observation from very different perspectives, providing high-quality fusion detection.

**Diversity Joint Training (DivJT).** Pre-trained detection models may not always be available, especially for applications detecting objects that do not belong to any classes included in public datasets such as PASCAL VOC [7]. Training a collection of models with diverse neural architectures can be laborious as tedious fine-tuning on hyperparameters (e.g., learning rate schedule) is necessary for

each neural configuration. In light of this, we develop a methodical approach to generating diverse detectors in one shot through a diversity-driven joint training framework, which generates multiple diverse models simultaneously, as depicted in Figure 2. By forcing models to interact with each other at each training iteration through diversity regularization, we can produce a collection of diverse models even if they share the same neural architecture.

The general formulation of the proposed joint training framework to construct a team  $\Phi = \{F_1, \dots, F_n\}$  of  $n$  detection models augments the standard detection loss (Equation 1) with a diversity regularization  $\mathcal{R}$  guiding models to learn differently:

$$\mathcal{L}(\mathbf{x}, \mathcal{O}; \Phi) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(F_i(\mathbf{x}), \mathcal{O}) + \frac{2\lambda}{n(n-1)} \sum_{j=2}^n \sum_{k=1}^{j-1} \mathcal{R}(\mathbf{x}; F_j, F_k), \quad (5)$$

where  $\lambda$  controls the importance between the standard detection loss  $\mathcal{L}$  (Equation 1) and the diversity regularization  $\mathcal{R}$ , which is designed to regularize member pairs to diversify their intrinsic detection mechanisms. We introduce the formulation of diversity regularization from two different perspectives.

• **DivJT with Diverse Detection Loss:** To maximize adversarial robustness, we design a joint training loss regularizer to minimize the negative correlation among member detectors [2] with respect to prediction errors on objectness, bounding boxes, and class labels. Recall in Section 3 that an object detector produces a fixed number of  $C$  candidate objects on an input video frame or image, which will be post-processed to construct the final set of detected objects. Let  $\Upsilon$  be the indices of *negative* candidate objects, which do not contain any real objects according to the ground truth  $\mathcal{O}$  of the training example  $\mathbf{x}$ . We introduce three diversity regularization components. Each corresponds to one of the three loss functions in object detection. For the object existence prediction, we denote a non-existent objectness vector predicted by the team member  $F_j$  to be  $\mathcal{N}_j = \{\mathbf{o}_{[\text{obj}]}^{j,\gamma} \mid \gamma \in \Upsilon\}$ , where  $\mathbf{o}_{[\text{obj}]}^{j,\gamma}$  is the objectness of the  $\gamma$ -th candidate object predicted by  $F_j$ . Note that the non-existent objectness vectors  $\mathcal{N}_j$  of  $F_j$  and  $\mathcal{N}_k$  of  $F_k$  have the same dimensionality because of the fixed number of candidate objects produced by the same detection algorithm and neural architecture. The diversity regularizer for the object existence component  $r_{\text{obj}}$  is formulated as

$$r_{\text{obj}}(\mathbf{x}; F_j, F_k) = \text{SIM}(\mathcal{N}_j, \mathcal{N}_k), \quad (6)$$

where  $\text{SIM}(\cdot)$  is a similarity measure between two vectors. We exploit the angular similarity computed by:

$$\text{SIM}(\mathcal{N}_j, \mathcal{N}_k) = 1 - \arccos\left(\frac{\mathcal{N}_j \cdot \mathcal{N}_k}{\|\mathcal{N}_j\| \|\mathcal{N}_k\|}\right)/\pi, \quad (7)$$

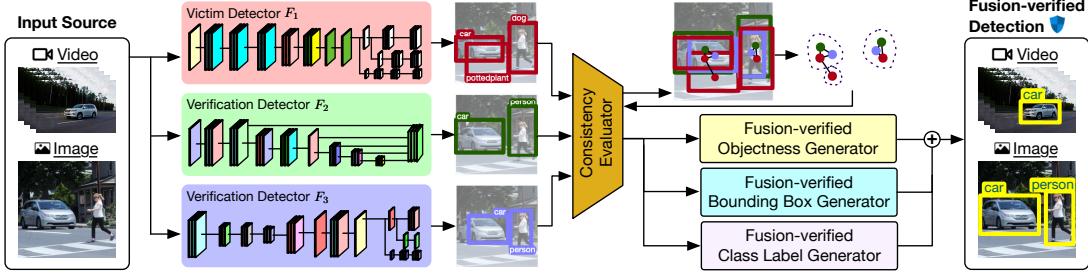


Figure 3: The detection workflow of FUSE with a team of three detectors where  $F_1$  is under attack.

which has been shown effective compared with other options (e.g., cosine similarity) [3]. Minimizing the above angular similarity decorrelates the mistakes made on object existence predictions by different members. For the object bounding box estimation task, we define its diversity regularization component  $r_{\text{bbox}}$  as

$$r_{\text{bbox}}(\mathbf{x}; F_j, F_k) = \frac{1}{||\Upsilon||} \sum_{\gamma \in \Upsilon} \exp(-\text{SE}(\mathbf{o}_{[\text{bbox}]}^{j,\gamma}, \mathbf{o}_{[\text{bbox}]}^{k,\gamma})), \quad (8)$$

with  $\mathbf{o}_{[\text{bbox}]}^{j,\gamma}$  to be the bounding box of the  $\gamma$ -th candidate detected by  $F_j$ , which maximizes the squared error (SE) of the normalized bounding boxes of negative candidates predicted by  $F_j$  and  $F_k$ . For the object classification task, we formulate the corresponding diversity loss component  $r_{\text{cls}}$  to be

$$r_{\text{cls}}(\mathbf{x}; F_j, F_k) = \frac{1}{C} \left[ \sum_{\gamma=1}^C \left( \mathbb{1}[\gamma \in \Upsilon] \cdot \text{SIM}(\mathbf{o}_{[\text{prob}]}^{j,\gamma}, \mathbf{o}_{[\text{prob}]}^{k,\gamma}) + \mathbb{1}[\gamma \notin \Upsilon] \cdot \text{SIM}(\mathbf{o}_{[\text{prob}]}^{j,\gamma \setminus \text{GT}}, \mathbf{o}_{[\text{prob}]}^{k,\gamma \setminus \text{GT}}) \right) \right], \quad (9)$$

where  $\mathbf{o}_{[\text{prob}]}^{j,\gamma}$  is the class probability vector of the  $\gamma$ -th candidate estimated by  $F_j$  and  $\mathbf{o}_{[\text{prob}]}^{j,\gamma \setminus \text{GT}}$  is the vector that excludes the ground-truth class label in the corresponding candidate. The first term diversifies the class label prediction on negative candidate objects, while the second term decorrelates the incorrect class label prediction in positive candidate objects. The final diversity regularization  $\mathcal{R}$  is defined to minimize the overall negative correlation:

$$\mathcal{R}(\mathbf{x}; F_j, F_k) = r_{\text{obj}}(\mathbf{x}; F_j, F_k) + r_{\text{bbox}}(\mathbf{x}; F_j, F_k) + r_{\text{cls}}(\mathbf{x}; F_j, F_k), \quad (10)$$

or, equivalently, diversify *incorrect* decisions on both object existences, bounding boxes, and class labels simultaneously while maintaining those correct ones.

• **DivJT with Diverse Kernel Filters:** Deep neural networks for object detection comprise multiple convolutional layers. Each of them has a number of learnable kernel filters, which are crucial in extracting salient features for the recognition task. One reason that adversarial distortion is detrimental to neural networks is the error amplification throughout the victim detector’s forward propagation, making small adversarial noise large enough to alter the decisions of the three prediction tasks. This diverse kernel filters based joint training aims to break this error amplification chain. Let  $\mathbf{W}_j^\ell$  and  $\mathbf{W}_k^\ell$  be two sets of kernel weights in the  $\ell$ -th convolutional layer of two team members  $F_j$  and  $F_k$  respectively. We compute the similarity between each pair of kernel weights obtained from the Cartesian product of  $\mathbf{W}_j^\ell$  and  $\mathbf{W}_k^\ell$ :  $\Omega^\ell = \{\text{SIM}(w_j^\ell, w_k^\ell) \mid (w_j^\ell, w_k^\ell) \in \mathbf{W}_j^\ell \times \mathbf{W}_k^\ell\}$ . Note that it is crucial to consider every pair of kernel filters due to the channel ordering ambiguity, and the computation of angular

similarity (SIM) (see Equation 7) is based on the vectorized kernel weights. The diversity regularizer  $\mathcal{R}$  is the sum of the mean and variance of  $\Omega^\ell$  over all layers of interest  $\ell \in \Psi$ :

$$\mathcal{R}(\mathbf{x}; F_j, F_k) = \sum_{\ell \in \Psi} [\text{MEAN}(\Omega^\ell) + \text{VAR}(\Omega^\ell)]. \quad (11)$$

Minimizing both mean and variance is more robust to outliers [27] which correspond to a small number of kernel filter pairs that are highly correlated and may amplify adversarial noises.

## 5 OBJECT DETECTION FUSION

FUSE addresses the technical challenges in objectness fusion, bounding box fusion, and classification fusion outlined in Section 3 by introducing a suite of fusion optimization techniques to generate fusion-verified detection outputs with high resilience against deception. As depicted in Figure 3, when an input  $\mathbf{x}$  (e.g., video frame or image) is sent to the detector  $F_1$  under protection, FUSE intercepts the detection query and dispatches it to each of the  $n$  detectors in the chosen fusion team to perform independent object detection in parallel, and obtains  $n$  sets of detection results. Each consists of three types of information: the objectness of detected objects, their bounding boxes, and their class labels with confidence probabilities. Our detection fusion algorithm employs a consistency evaluator to identify those objects detected by different fusion member models referring to the same entity with three core fusion components: (i) fusion-verified objectness generator, which consolidates different object existences predicted by different member detectors; (ii) fusion-verified bounding box generator, which consolidates different bounding boxes detected from different member detectors of the fusion team; and (iii) fusion-verified class label generator, which consolidates different confidence probabilities resulting from different member detectors based on the overlapping of their bounding boxes. The objectness fusion employs a threshold-based approach such that those detected objects with their objectness confidence higher than the given threshold will be examined. Then, we employ a graph partitioning-based algorithm to perform the bounding box fusion by leveraging the objectness fusion and the cross-model synergy among different detectors of a fusion team, aiming to unscramble and resolve the mishmash among the results from diverse detectors of the fusion team. Finally, we introduce bounding box enhanced optimizations to accomplish the classification fusion.

Concretely, given a team  $\Phi = \{F_1, \dots, F_n\}$  of  $n$  detection models, we formulate detection fusion to find the robust resolution as a graph partitioning problem where we construct a weighted undirected graph  $\mathcal{G} = (V, E)$  by examining all detected objects from each member detector of the fusion team, i.e.,  $\hat{F}_1(\mathbf{x}), \dots, \hat{F}_n(\mathbf{x})$ . With

objectness fusion producing  $\mathbf{V} = \hat{F}_1(\mathbf{x}) \cup \dots \cup \hat{F}_n(\mathbf{x})$ , for each pair of detectors in the fusion team, we examine every pair of detected objects, one from each detector, and an edge between them is added to  $E$  with the intersection over union (IOU) score of their bounding boxes as the edge weight  $w \in [0, 1]$  if their IOU score passes the system-defined IOU threshold  $\tau_{\text{IOU}}$  (e.g., 50% is used in our experiments). Upon completion of the examination for all member detectors, we perform clique partitioning over  $\mathcal{G}$ . This reduces the problem of finding the set of highly overlapped objects detected by different member detectors to the problem of partitioning the graph  $\mathcal{G}$  into cliques by maximizing the sum of edge weights. Each resultant clique consists of one or more nodes, and each node corresponds to one detected object from one detector in the fusion team. When a clique has multiple nodes, it implies that the corresponding detected objects from multiple detectors have highly overlapped bounding boxes. Thus, we perform the consistency evaluation for objects within each clique to determine the fusion-verified bounding box. Intuitively, we can generate a fusion-verified bounding box and a class label with confidence probabilities for each clique. But this naive approach to classification fusion only works for the simple cases, in which the fusion-verified bounding box indeed corresponds to one true object in the query input. For more complicated scenes, in which the query input (video frame or image) contains multiple overlapping instances of the same class label, such as multiple persons overlapping when walking over a crossroad, simply performing classification fusion based on the bounding box fusion outcome may lead to the detection error that mistakenly detects one object instead of multiple spatially overlapping objects when such overlapping is larger than the pre-defined IOU threshold. To resolve this challenge, we develop a holistic solution to employing bounding box fusion-refined classification fusion. In particular, for cliques reaching consensus on a class label, we examine every node in the clique. We first return the node achieving the highest confidence, denoted by  $\mathbf{o}^{\max}$ . Then, for each remaining node  $\mathbf{o}$  in the clique with the same class (i.e.,  $\mathbf{o}_{[\text{cls}]} = \mathbf{o}_{[\text{cls}]}^{\max}$ ), FUSE recomputes its confidence based on the bounding box overlapping in IOU:

$$\mathbf{o}_{[\text{prob}]} \leftarrow \mathbf{o}_{[\text{prob}]} (1 - \text{IOU}(\mathbf{o}_{[\text{bbox}]}^{\max}, \mathbf{o}_{[\text{bbox}]})). \quad (12)$$

This strategy allows FUSE to detect multiple objects of the same class from a clique [1]. Algorithm 1 presents the pseudocode of the FUSE three-tier detection fusion.

## 6 ROBUSTNESS ANALYSIS

We formally analyze the robustness of FUSE in terms of the adversarial vulnerability of object detection fusion. Assume that an object detector  $F_\theta$  is well-trained by minimizing the standard detection loss in Equation 1, which measures the prediction error on objectness, bounding boxes, and class labels of objects. While the exact formulation of adversarial attacks may vary, they share the same objective: deceiving the victim detector to return erroneous detection results. Hence, an adversarial example  $\mathbf{x} + \boldsymbol{\delta}$  generated by solving the optimization problem in Equation 3 causes the standard detection loss (prediction errors) of  $F_\theta$  to increase, i.e.,

$$\mathcal{L}(F_\theta(\mathbf{x}), \mathbf{O}) < \mathcal{L}(F_\theta(\mathbf{x} + \boldsymbol{\delta}), \mathbf{O}). \quad (13)$$

Intuitively, an object detection system is more vulnerable if just a tiny perturbation to input can cause a drastic change in output [28].

---

**Algorithm 1** Robust Object Detection Fusion

---

```

1: Input:  $\mathbf{x}$ : query image;  $\Phi$ : fusion team;  $\tau_{\text{IOU}}$ : IOU threshold
2: Output:  $\hat{\mathbf{O}}$ : fusion-verified detected objects
3: procedure FUSE( $\mathbf{x}, \Phi, \tau_{\text{IOU}}$ )
4:    $\hat{\mathbf{O}} \leftarrow \emptyset, \mathbf{V} \leftarrow \emptyset, E \leftarrow \emptyset$ 
5:   for  $F_i \in \Phi$  do
6:      $\mathbf{V} \leftarrow \mathbf{V} \cup \hat{F}_i(\mathbf{x})$ 
7:   for  $\mathbf{o}, \mathbf{o}' \in \binom{\mathbf{V}}{2}$  do
8:      $w \leftarrow \text{IOU}(\mathbf{o}_{[\text{bbox}]}, \mathbf{o}'_{[\text{bbox}]})$ 
9:     #  $\mathbf{o}_{[\text{det}]}$  denotes the fusion member detecting the object  $\mathbf{o}$ 
10:    if  $\mathbf{o}_{[\text{det}]} \neq \mathbf{o}'_{[\text{det}]} \wedge w \geq \tau_{\text{IOU}}$  then
11:       $E \leftarrow E \cup \{(\mathbf{o}, \mathbf{o}', w)\}$ 
12:   for  $c \in \text{CLIQUE-PARTITION}(\mathbf{V}, E)$  do
13:      $\eta \leftarrow \text{VOTE-CLASS-LABEL}(c)$ 
14:     if  $\sum_{\mathbf{o} \in c} \mathbb{1}\{\mathbf{o}_{[\text{cls}]} = \eta\} < |\Phi|/2$  then
15:       continue
16:      $\mathbf{o}^{\max} \leftarrow \text{FIND-MAX-CONFIDENCE-OBJECT}(c, \eta)$ 
17:      $\hat{\mathbf{O}} \leftarrow \hat{\mathbf{O}} \cup \{\mathbf{o}^{\max}\}$ 
18:     for  $\mathbf{o} \in c \setminus \{\mathbf{o}^{\max}\}$  do
19:       if  $\mathbf{o}_{[\text{cls}]} = \eta$  then
20:          $\mathbf{o}_{[\text{prob}]} \leftarrow \mathbf{o}_{[\text{prob}]} (1 - \text{IOU}(\mathbf{o}_{[\text{bbox}]}^{\max}, \mathbf{o}_{[\text{bbox}]})$ 
21:        $\hat{\mathbf{O}} \leftarrow \hat{\mathbf{O}} \cup \{\mathbf{o}\}$ 
22:   return  $\hat{\mathbf{O}}$ 

```

---

Based on this idea, we can formulate the adversarial vulnerability of an object detector as follows.

*Definition 6.1 (Adversarial Vulnerability).* Given an object detector  $F_\theta$ , an input  $\mathbf{x}$ , and an attack budget  $\epsilon$  in  $\ell_p$ -norm, the adversarial vulnerability can be captured by how much the standard detection loss of  $F_\theta$  can be varied by injecting arbitrary input perturbation  $\boldsymbol{\delta}$  inside an  $\ell_p$ -norm ball with a radius  $\epsilon$  centered at  $\mathbf{x}$ :

$$\Delta \mathcal{L}(\mathbf{x}, \mathbf{O}; F_\theta, \epsilon) = \max_{\|\boldsymbol{\delta}\|_p < \epsilon} |\mathcal{L}(F_\theta(\mathbf{x}), \mathbf{O}) - \mathcal{L}(F_\theta(\mathbf{x} + \boldsymbol{\delta}), \mathbf{O})|. \quad (14)$$

By employing an army of carefully constructed object detection models to safeguard an object detection system by robust fusion, we conjecture that one can reduce the overall adversarial vulnerability or, equivalently, offer strong robustness against deception. To formally verify robust detection through model fusion, we define the overall detection loss of an  $n$ -model fusion team  $\Phi = \{F_1, \dots, F_n\}$  given an input  $\mathbf{x}$  with ground truth objects  $\mathbf{O}$  to be

$$\mathcal{L}(\mathbf{x}, \mathbf{O}; \Phi) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(F_i(\mathbf{x}), \mathbf{O}). \quad (15)$$

Then, we can reformulate adversarial vulnerability (Equation 14) as

$$\Delta \mathcal{L}(\mathbf{x}, \mathbf{O}; \Phi, \epsilon) = \max_{\|\boldsymbol{\delta}\|_p < \epsilon} |\mathcal{L}(\mathbf{x}, \mathbf{O}; \Phi) - \mathcal{L}(\mathbf{x} + \boldsymbol{\delta}, \mathbf{O}; \Phi)|, \quad (16)$$

which captures the changes to the overall detection loss of the fusion team caused by input perturbation bounded by  $\epsilon$  in  $\ell_p$ -norm. Since the adversarial perturbation  $\boldsymbol{\delta}$  is often constrained to be human-imperceptible (i.e.,  $\epsilon \rightarrow 0$ ), a first-order Taylor expansion on Equation 16 in  $\epsilon$  gives

$$\Delta \mathcal{L}(\mathbf{x}, \mathbf{O}; \Phi, \epsilon) \approx \max_{\|\boldsymbol{\delta}\|_p < \epsilon} |\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{O}; \Phi) \cdot \boldsymbol{\delta}|. \quad (17)$$

The object detection system is vulnerable (or robust) if, on average,  $\Delta \mathcal{L}(\mathbf{x}, \mathbf{O}; \Phi, \epsilon)$  is large (or small) [20, 28]. Using the definition of dual norm, the adversarial vulnerability over the dataset  $\mathcal{D}$  is

$$\mathbb{E}_{(\mathbf{x}, \mathbf{O}) \sim \mathcal{D}} [\Delta \mathcal{L}(\mathbf{x}, \mathbf{O}; \Phi, \epsilon)] \propto \mathbb{E}_{(\mathbf{x}, \mathbf{O}) \sim \mathcal{D}} [| | | \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{O}; \Phi) | |_q |], \quad (18)$$

where  $q$  is the dual norm of  $p$ , which satisfies  $\frac{1}{p} + \frac{1}{q} = 1$  and  $1 \leq p \leq \infty$ . Expanding the above yields the following theorem on the adversarial vulnerability of an object detection system safeguarded by model fusion (the proof is provided in Appendix A).

**THEOREM 6.2 (ADVERSARIAL VULNERABILITY OF MODEL FUSION).** *Given a team  $\Phi$  of  $n$  object detection models, an attack budget  $\epsilon$  in  $\ell_p$ -norm, and the dataset  $\mathcal{D}$ , the adversarial vulnerability of a model fusion-protected object detection system is proportional to*

$$\sqrt{\frac{1}{3n} + \frac{2}{9n^2} [\text{Cov}_{\text{intra}}(\mathbf{x}, \mathcal{O}; \Phi) + \text{Cov}_{\text{inter}}(\mathbf{x}, \mathcal{O}; \Phi)]}. \quad (19)$$

$\text{Cov}_{\text{intra}}(\mathbf{x}, \mathcal{O}; \Phi)$  and  $\text{Cov}_{\text{inter}}(\mathbf{x}, \mathcal{O}; \Phi)$  are two terms capturing the covariance of gradients across tasks and models:

$$\begin{aligned} \text{Cov}_{\text{intra}}(\mathbf{x}, \mathcal{O}; \Phi) &= \sum_{i=2}^n \sum_{j=1}^{i-1} \left[ \frac{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}^j)}{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}^i)} \right. \\ &\quad \left. + \frac{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^j)}{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^i)} + \frac{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^j)}{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^i)} \right], \end{aligned} \quad (20)$$

$$\begin{aligned} \text{Cov}_{\text{inter}}(\mathbf{x}, \mathcal{O}; \Phi) &= \sum_{i=2}^n \sum_{j=1}^{i-1} \left[ \frac{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^j)}{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}^i)} + \frac{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^j)}{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}^i)} \right. \\ &\quad \left. + \frac{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}^j)}{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^i)} + \frac{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^j)}{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^i)} \right. \\ &\quad \left. + \frac{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}^j)}{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^i)} + \frac{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^j)}{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^i)} \right] \\ &\quad + \sum_{i=1}^n \left[ \frac{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}^i)}{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^i)} + \frac{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}^i)}{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^i)} \right. \\ &\quad \left. + \frac{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^i)}{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^i)} \right], \end{aligned} \quad (21)$$

where  $\mathcal{L}_{\text{task}}^i = \mathcal{L}_{\text{task}}(F_i(\mathbf{x}), \mathcal{O})$  for task  $\in \{\text{obj, bbox, cls}\}$ .

Theorem 6.2 offers three important insights on model fusion. First, the adversarial vulnerability of an object detection system reduces when multiple object detectors are jointly operated in the decision-making process. The improvement in adversarial robustness can be explained as the adversarial perturbations required to attack multiple models cancel each other out [28]. Second, the adversarial vulnerability also depends on the covariance of gradients between tasks and fusion members. The first term  $\text{Cov}_{\text{intra}}(\mathbf{x}, \mathcal{O}; \Phi)$  indicates that to reduce adversarial vulnerability, the gradients of a task (e.g., objectness prediction) of one model should have a minimal covariance with another model in the fusion team. Interestingly, the second term  $\text{Cov}_{\text{inter}}(\mathbf{x}, \mathcal{O}; \Phi)$  suggests that the gradients of one task should also be decorrelated with other tasks. This is because adversarial attacks on object detection systems typically manipulate multiple components in the standard detection loss (Equation 1). For example, RAP [13] attacks the objectness component and bounding box component such that the victim cannot correctly recognize the existence of objects, and even if objects are detected, their bounding boxes will be erroneous. By decorrelating the gradient directions of different tasks, the adversary cannot easily find a single perturbation that successfully deceives multiple tasks at the same time. Finally, the notion of covariance of gradients can be utilized for strategic teaming to cherrypick team members that deliver the

Model Identifier	Detection Algorithm	Neural Configuration	mAP (Benign)
$F_1$	FRCNN	VGG16-Standard	67.37
$F_2$	YOLOv3-D	Darknet53-Standard	83.43
$F_3$	YOLOv3-M	MobileNetV1-Standard	71.84
$F_4$	SSD300	VGG16-Standard	76.11
$F_5$	SSD512	VGG16-Standard	79.83
$F_6$	FRCNN	VGG16-DivJT-DetLoss	64.77
$F_7$	FRCNN	VGG16-DivJT-DetLoss	65.58
$F_8$	FRCNN	VGG16-DivJT-KF	64.88
$F_9$	FRCNN	VGG16-DivJT-KF	64.82
$F_{10}$	FRCNN	VGG16-RandInit	67.51
$F_{11}$	FRCNN	VGG16-RandInit	67.32

Table 2: A summary of models in experimental evaluation.

strongest robustness (i.e., minimum overall covariance) by joint force. We provide formal analysis of FUSE from the perspective of adversarial perturbation lower bound in Appendix B.

## 7 EXPERIMENTS

We conduct extensive experiments to empirically analyze the effectiveness of FUSE against various deception attacks on two benchmark datasets: PASCAL VOC [7] and INRIA [6]. VOC consists of images containing 20 classes of objects (e.g., bicycle, chair, and dog), and INRIA focuses on pedestrian detection. All results reported are measured on their entire test set (4,952 images for VOC and 288 images for INRIA). We use the standard performance metric, called mean average precision (mAP) [7], to capture the detection quality of an object detection system. We consider three representative families of object detection algorithms: Faster R-CNN (FRCNN) [19], YOLOv3 [18], and SSD [14] with different neural architectures as backbones. Table 2 summarizes the models involved in our experiments, where  $F_1$  to  $F_5$  are trained independently with the standard detection loss (Equation 1). Using diversity joint training (DivJT), we further train two FRCNN models ( $F_6$  and  $F_7$ ) with our diversified detection loss (Equation 10) and two models ( $F_8$  and  $F_9$ ) with our diversified kernel filters (Equation 11). As a baseline training strategy, we provide two models ( $F_{10}$  and  $F_{11}$ ) trained with random initialization of model weights without regularization promoting model diversity. Additional details are provided in Appendix C.

### 7.1 Robustness Analysis

Our empirical analysis begins with examining the robustness of FUSE protecting FRCNN (i.e.,  $F_1$  in Table 2 as the victim). With strategic teaming on models summarized in Table 2, we construct ten fusion teams by pairing up the victim with two additional models as our experimental studies show a fusion of three diverse detectors is already sufficient for strong robustness, even though we proved in Section 6 that a larger team can lead to an even stronger defense. Four state-of-the-art attack techniques, including TOG [5], UEA [24], RAP [13], and DAG [26], are exploited to launch both untargeted random attacks and targeted specificity attacks with object vanishing, fabrication, and mislabeling effects on the victim. Table 3 summarizes the results on VOC with benign mAP (3rd column) and mAP under attacks (4th to 10th columns) where we compare FUSE (3rd and 4th rows) with the victim  $F_1$  with no protection (1st row). Since defenses proposed for single-task learners (i.e., image classifiers) are not applicable to object detection systems with multi-task objectives, we include the adversarial detection

Object Detector	mAP (Benign)	Untargeted Random Attacks				Targeted Specificity Attacks		
		TOG	UEA	RAP	DAG	TOG-v	TOG-f	TOG-m
$F_1$ : FRCNN (Victim)	67.37	2.64	18.07	4.78	3.56	0.14	1.24	2.14
Adversarial Detection Training [29]	35.99 (0.53 $\times$ )	34.07 (12.91 $\times$ )	17.67 (0.98 $\times$ )	35.60 (7.45 $\times$ )	35.58 (9.99 $\times$ )	34.03 (243.07 $\times$ )	34.18 (27.56 $\times$ )	34.81 (16.27 $\times$ )
Most Diverse $\Phi = \{F_1, F_2, F_5\}$	85.95 <sup>*</sup> (1.28 $\times$ )	81.46 (30.86 $\times$ )	53.52 (2.96 $\times$ )	81.96 (17.15 $\times$ )	84.60 (23.76 $\times$ )	81.23 (580.21 $\times$ )	81.91 (66.06 $\times$ )	82.29 (38.45 $\times$ )
FUSE	<b>80.82 (3.62)</b>	76.14 (4.36)	48.15 (3.53)	77.39 (3.49)	80.21 (3.05)	77.23 (2.51)	77.20 (3.17)	77.75 (2.93)
Average (Std)	(1.20 $\times$ )	(28.84 $\times$ )	(2.66 $\times$ )	(16.19 $\times$ )	(22.53 $\times$ )	(551.64 $\times$ )	(62.26 $\times$ )	(36.33 $\times$ )

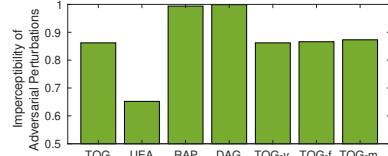
\*The best-performing member,  $F_2$ , achieves an mAP of 83.43%.

**Table 3: Evaluation of FUSE with FRCNN as the victim under seven attacks compared with adversarial detection training.**

training approach (2nd row) [29], the only recent proposal to the best of our knowledge, for baseline comparisons.

We make four observations. First, all seven attacks significantly bring down the mAP of the victim detector with no protection from 67.37% to 0.14~18.07% (see the 1st row). Second, FUSE maintains high mAP on benign inputs (boldfaced in the 3rd column), and for the most diverse fusion team with  $F_1$ : FRCNN (victim),  $F_2$ : YOLOv3-D, and  $F_5$ : SSD512 (3rd row), FUSE achieves an mAP of 85.95% on benign inputs, which is 2.5% higher than the best-performing member (i.e.,  $F_2$ : YOLOv3-D with an mAP of 83.43%). Third, FUSE significantly improves the robustness of object detection by boosting the mAP under all seven attacks, delivering high mAPs of 53.52~84.60% (i.e., 2.96 $\times$  to 580.21 $\times$  improvement over the mAP of the victim detector: 0.14~18.07%). The 4th row shows that all ten diverse fusion teams can empower the object detection system with strong robustness against all seven attacks, achieving a high average mAP of 48.15~80.21% with a small standard deviation. This is because even if an adversarial input may succeed in partially deceiving both the victim and some other members, FUSE can still produce correct detection by analyzing objectness, bounding boxes, and class labels of objects. In comparison, the adversarial detection training approach recently proposed in [29] has a low mAP of 17.67~35.60% under seven deception attacks. One of the reasons for its low resilience is due to its low benign mAP of 35.99%, compared to the mAP of 67.37% of the victim detector  $F_1$  (0.53 $\times$  performance degradation). Finally, it is worth noting that the adversarial detection training approach [29] fails to defend the UEA attack with an mAP of 17.67%, which is worse than the mAP of 18.07% by the victim detector  $F_1$  without protection. In comparison, the victim detector under FUSE protection can achieve the high mAP of 53.52%, delivering a 2.96 $\times$  mAP improvement (18.07%  $\rightarrow$  53.52%). One reason that UEA is relatively hard to mitigate compared with the other six attacks can be explained by its excessive amount of adversarial distortion. Figure 4 reports the imperceptibility of adversarial perturbations generated by all seven attacks in terms of the structural similarity (SSIM) [23] between the benign input and its adversarial counterpart. RAP [13] and DAG [26] both have an SSIM close to 1.00, meaning that they are almost completely imperceptible. In contrast, UEA generates adversarial examples with significantly more perceptible perturbations with a low SSIM of 0.65.

**Robustness: YOLOv3 as the Victim.** We next consider the best-performing model from Table 2, i.e.,  $F_2$ : YOLOv3-D, to be the victim under TOG attacks [5] since other attacks are not applicable to one-shot detectors. Table 4 reports the results where the victim with an mAP of 83.43% given benign inputs suffers severely from



**Figure 4: Imperceptibility of adversarial perturbations (1.00 means completely imperceptible).**

Object Detector	mAP (Benign)	mAP (Under Attacks)			
		TOG	TOG-v	TOG-f	TOG-m
$F_2$ : YOLOv3-D (Victim)	83.43	0.56	1.24	0.25	3.15
FUSE	82.28 (0.99 $\times$ )	<b>73.02</b> (130.39 $\times$ )	<b>75.88</b> (61.19 $\times$ )	<b>76.83</b> (307.32 $\times$ )	<b>76.15</b> (24.17 $\times$ )

**Table 4: Evaluation of FUSE with the best-performing model, YOLOv3-D, as the victim with  $\Phi = \{F_1, F_2, F_5\}$ .**

Object Detector	mAP (Benign)	mAP (Under Attacks)			
		TOG	UEA	RAP	DAG
$F_1$ : FRCNN (Victim)	67.37	2.64	18.07	4.78	3.56
(a) Objectness Fusion	77.43 (1.15 $\times$ )	2.68 (1.02 $\times$ )	45.92 (2.54 $\times$ )	75.85 (15.87 $\times$ )	78.59 (22.08 $\times$ )
(b) Objectness & Bounding Box Fusion	81.31 (1.21 $\times$ )	76.93 (29.14 $\times$ )	46.21 (2.56 $\times$ )	77.01 (16.11 $\times$ )	81.01 (22.76 $\times$ )
(c) Objectness & Bounding Box & Classification Fusion	<b>85.95</b> (1.28 $\times$ )	<b>81.46</b> (30.86 $\times$ )	<b>53.52</b> (2.96 $\times$ )	<b>81.96</b> (17.15 $\times$ )	<b>84.60</b> (23.76 $\times$ )

**Table 5: Ablation studies of FUSE.**

those white-box attacks with low mAPs ranging from 0.25% to 3.15%. Even though the best-performing member in the fusion team of  $\Phi = \{F_1, F_2, F_5\}$  is under attack, FUSE can still deliver highly competitive mAP in benign scenarios and effectively mitigate adversarial attacks, boosting the mAP under attacks to 73.02~76.83%, which is 24.17 $\times$  to 307.32 $\times$  improvement over the victim. Hence, FUSE does not rely on a single member in the team that performs comparatively well to mitigate adversarial attacks but leverages multiple diverse models by joint force against deception.

**Ablation Studies.** We conduct ablation studies in Table 5 to understand the contributions of the key components to the robustness of FUSE by comparing three versions: (a) *objectness fusion* combining detection results from all member models with a redundancy filter implemented by non-maximum suppression; (b) *objectness & bounding box fusion* combining detection results from all members using our graph partitioning approach to find robust resolution; and (c) *objectness & bounding box & classification fusion* incorporating Equation 12 for classification fusion (i.e., the complete FUSE).

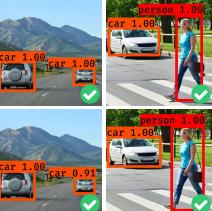
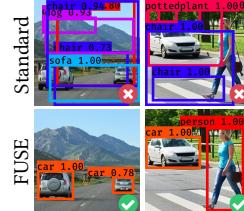
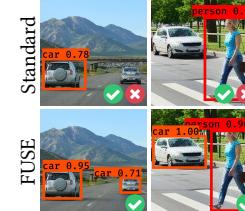
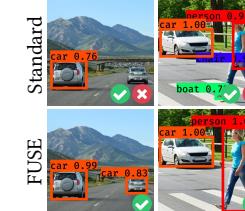
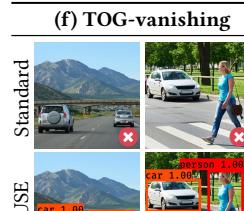
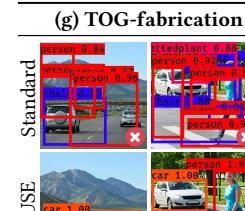
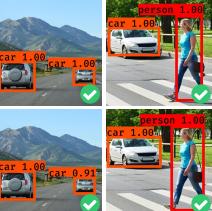
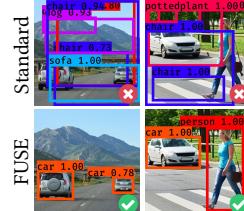
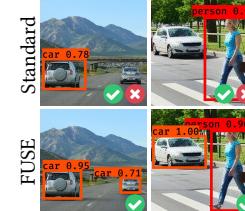
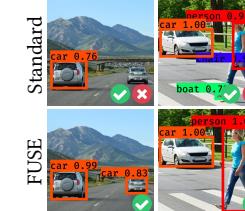
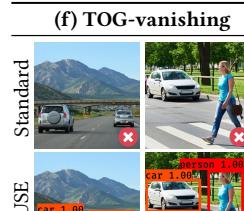
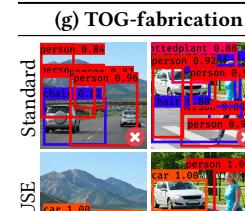
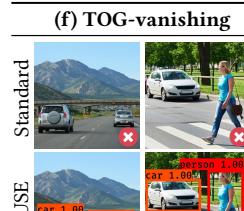
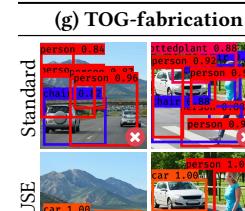
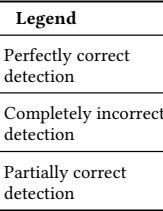
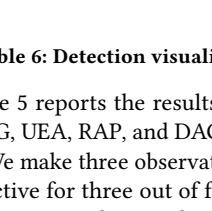
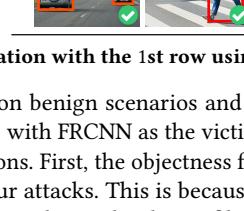
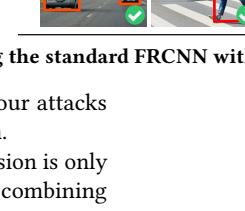
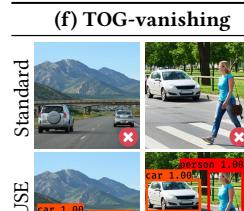
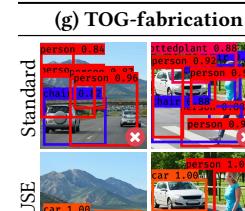
(a) Benign (No Attack)		(b) TOG		(c) UEA		(d) RAP		(e) DAG	
Standard						Standard			
									
(f) TOG-vanishing		(g) TOG-fabrication		(h) TOG-mislabeling					
FUSE				Standard		Standard			
									

Table 6: Detection visualization with the 1st row using the standard FRCNN without protection and the 2nd row having FUSE protection.

Table 5 reports the results on benign scenarios and four attacks (TOG, UEA, RAP, and DAG) with FRCNN as the victim.

We make three observations. First, the objectness fusion is only effective for three out of four attacks. This is because combining detection results simply through a redundancy filter cannot remove false positives where TOG attack tends to generate (it can merely improve the mAP of the victim from 2.64% to 2.68% under TOG). Second, finding robust resolution of objects detected by different members allows the objectness & bounding box fusion approach to locate highly suspicious objects to be pruned. Finally, by further improving the fusion capability with classification fusion, FUSE shows an additional boost in robustness against all attacks (2.96~30.86× improvement over the victim) and in benign scenarios (1.28× improvement over the victim).

**Visualization.** To demonstrate the mitigation effectiveness of FUSE, Table 6 gives visualizations of two examples under benign (a) and adversarial settings (b-h). With the standard FRCNN perfectly detecting all objects with 100% confidence (see the 1st row in Table 6(a)), we can observe that all seven attacks successfully deceive the victim with different adverse effects. For instance, TOG in Table 6(b) misleads the victim to detect multiple false objects, and UEA in Table 6(c) fools the victim to misdetect one car in each example. Even with such highly diverse malicious effects incurred by different attack methods, FUSE (2nd row) can consistently repair all adversarial examples, restore the operation of the object detection system, and maintain high performance on benign inputs.

## 7.2 FUSE Against Adversarial Patches

To study the effectiveness of FUSE against adversarial patches, we exploit TOG [5] to train patches of size ( $64 \times 64$ ) with object vanishing and mislabeling effects. Following [21], the INRIA dataset for pedestrian detection [6] is exploited with YOLOv3-D as the victim. The patch is placed at the center of the person detected on the benign test image sent to the victim to generate its adversarial example. Table 7 visualizes an example test image under no attack (1st column), the TOG-vanishing patch attack (2nd column), and the TOG-mislabeling patch attack (3rd column). For the TOG-vanishing

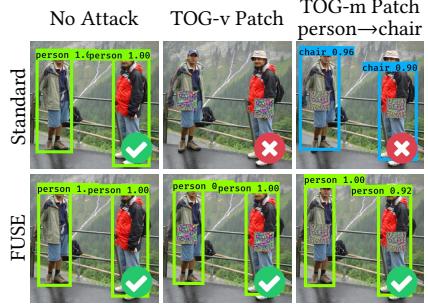


Table 7: An image attacked by adversarial patches using TOG with object vanishing (2nd column) and mislabeling (3rd column) effects.

patch, FUSE successfully protects the victim by drastically reducing the percentage of vanished pedestrians in the entire INRIA test set from 97.35% to 15.30%. For TOG-mislabeling patches, we train 19 adversarial patches, aiming at deceiving the victim to mislabel the person with the patch as one of the 19 different target classes on VOC (e.g., the one in the 3rd column in Table 7 fools the victim to mislabel any person attached with the patch to be a chair). Even under such highly aggressive adversarial patches for fooling the standard victim without protection, led to 80% (on average) of detected pedestrians being purposefully mislabeled, FUSE can mitigate such adversarial patch attacks with strong robustness and significantly reduce the percentage of mislabeled pedestrians to zero in most cases. The adversarial patches are by design to simulate the physical attack digitally for systematic evaluation [21]. These experiments show that FUSE holds great potential to provide strong robustness against physical attacks and visual distortions due to other real-world factors, e.g., the transmission distortion, the angle or condition of taking the image or video.

## 7.3 Execution Efficiency

We further investigate the execution efficiency of FUSE. Table 8 summarizes the detection time and frame-per-second (FPS) measurements on both the standard FRCNN detector without protection and with FUSE protection. With the highly parallelizable design,

Object Detector	Detection Time (s)	FPS
Standard	0.1399	7.15
FUSE	0.1455	6.87

Table 8: The detection time and FPS comparisons (FRCNN).

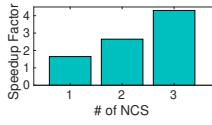


Table 9: Two groups of challenging scenarios.

FUSE incurs a low overhead to the detector being protected because the graph constructed by FUSE for each query image is tiny.

Adversarial attacks are particularly detrimental in applications running on the edge as the data captured by sensors will be immediately processed locally to trigger downstream operations. It is of utmost importance to safeguard those systems perceiving malicious inputs and propagating erroneous decisions along the application pipeline. In light of this, we extend FUSE for edge deployments with AI accelerators such as the Intel Neural Compute Stick (NCS) [11]. Multiple detection models can be run efficiently within each device with low power consumption and low cost. The FPS speedup reported in Figure 5 demonstrates the capability of FUSE in efficiently enhancing the robustness of such Edge-AI environments.

#### 7.4 Challenging Cases

Although FUSE offers robust performance, it has limitations. Table 9 shows two common scenarios where FUSE faces challenges to correctly detect objects: (1) tiny objects, e.g., the image with people and bicycles in the 1st column under Group 1, and (2) images taken with a poor lighting condition, e.g., the bus cannot be detected in the 1st image under Group 2. One way to address those scenarios could be to incorporate scale-aware neural architectures through trident blocks [12] or tiling [16] in our generation process of diverse models, forming fusion teams with detectors specialized in different settings. We could also conduct adaptive image preprocessing and attempt to repair the poor lighting conditions with automatic color enhancement [17]. It would also be interesting to explore the feasibility of selectively involving human-in-the-loop.

### 8 CONCLUSION

We have presented FUSE, a robust object detection fusion approach for strong robustness and survivability against deception to object detection systems by a growing number of digital or physical adversarial attacks. The paper makes three original contributions. First, we present diversity-enhanced fusion teaming algorithms for producing diverse fusion detectors. Second, we introduce a three-tier detection fusion framework to safeguard the victim detector through three mutually reinforcing components: objectness fusion, bounding box fusion, and classification fusion. Third, we formally analyze the robustness enhancement of our object detection fusion in terms of adversarial vulnerability. Extensive experiments on eleven detectors and two benchmark datasets validate the strong robustness of FUSE against various deception attacks while maintaining high performance on benign inputs.

### ACKNOWLEDGMENTS

This research is partially sponsored by the National Science Foundation CISE grants 2038029, 2026945, 1564097, an IBM faculty award, and a Cisco grant on Edge Computing. The first author acknowledges the Croucher Scholarship for Doctoral Study from the Croucher Foundation. The source code of FUSE is available at <https://github.com/git-disl/FUSE>.

### REFERENCES

- [1] Navaneeth Bodla, Bharat Singh, Rama Chellappa, and Larry S Davis. 2017. Soft-NMS—improving object detection with one line of code. In *ICCV*.
- [2] Gavin Brown. 2004. *Diversity in neural network ensembles*. Ph.D. Dissertation.
- [3] Daniel Cer, Yifei Yang, Sheng-yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St. John, Noah Constant, Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, Brian Strope, and Ray Kurzweil. 2018. Universal Sentence Encoder. In *EMNLP*.
- [4] Ka-Ho Chow, Ling Liu, Mehmet Emre Gursoy, Stacey Truex, Wenqi Wei, and Yanzhao Wu. 2020. Understanding Object Detection Through an Adversarial Lens. In *ESORICS*.
- [5] Ka-Ho Chow, Ling Liu, Margaret Loper, Juhyun Bae, Mehmet Emre Gursoy, Stacey Truex, Wenqi Wei, and Yanzhao Wu. 2020. Adversarial Objectness Gradient Attacks on Real-time Object Detection Systems. In *TPS*. IEEE.
- [6] Navneet Dalal and Bill Triggs. 2005. Histograms of oriented gradients for human detection. In *CVPR*.
- [7] Mark Everingham, SM Ali Eslami, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. 2015. The pascal visual object classes challenge: A retrospective. *IJCV* (2015).
- [8] Vandit Gajjar, Ayesha Gurnani, and Yash Khandhedia. 2017. Human detection and tracking for video surveillance: A cognitive science approach. In *ICCV*.
- [9] Ross Girshick. 2015. Fast r-cnn. In *CVPR*.
- [10] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. 2015. Explaining and harnessing adversarial examples. In *ICLR*.
- [11] Intel. 2019. Intel Neural Compute Stick 2. <https://software.intel.com/content/www/us/en/develop/hardware/neural-compute-stick.html>
- [12] Yanghao Li, Yuntao Chen, Naiyan Wang, and Zhaoxiang Zhang. 2019. Scale-aware trident networks for object detection. In *ICCV*.
- [13] Yuezun Li, Daniel Tian, Xiao Bian, Siwei Lyu, et al. 2018. Robust adversarial perturbation on deep proposal-based models. In *BMVC*.
- [14] Wei Liu, Dragomir Anguelov, Dumitri Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. 2016. SSD: Single shot multibox detector. In *ECCV*.
- [15] Xin Liu, Huanrui Yang, Ziwei Liu, Linghao Song, Hai Li, and Yiran Chen. 2019. DPatch: An Adversarial Patch Attack on Object Detectors. In *SafeAI*.
- [16] F Ozge Unel, Burak O Ozkalayci, and Cevahir Cigla. 2019. The power of tiling for small object detection. In *CVPRW*.
- [17] Jongchan Park, Joon-Young Lee, Donggeun Yoo, and In So Kweon. 2018. Distort-and-recover: Color enhancement using deep reinforcement learning. In *CVPR*.
- [18] Joseph Redmon and Ali Farhadi. 2018. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767* (2018).
- [19] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2015. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NeurIPS*.
- [20] Aman Sinha, Hongseok Namkoong, Riccardo Volpi, and John Duchi. 2018. Certifying some distributional robustness with principled adversarial training. In *ICLR*.
- [21] Simen Thys, Wiebe Van Ranst, and Toon Goedemé. 2019. Fooling automated surveillance cameras: adversarial patches to attack person detection. In *CVPRW*.
- [22] Paul Viola and Michael Jones. 2001. Rapid object detection using a boosted cascade of simple features. In *CVPR*.
- [23] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. 2004. Image quality assessment: from error visibility to structural similarity. *TIP* (2004).
- [24] Xingxing Wei, Siyuan Liang, Ning Chen, and Xiaochun Cao. 2019. Transferable adversarial attacks for image and video object detection. In *IJCAI*.
- [25] Tsui-Wei Weng, Huan Zhang, Pin-Yu Chen, Jinfeng Yi, Dong Su, Yupeng Gao, Cho-Jui Hsieh, and Luca Daniel. 2018. Evaluating the robustness of neural networks: An extreme value theory approach. In *ICLR*.
- [26] Cihang Xie, Jianyu Wang, Zhishuai Zhang, Yuyin Zhou, Lingxi Xie, and Alan Yuille. 2017. Adversarial examples for semantic segmentation and object detection. In *ICCV*.
- [27] Pengtao Xie, Yuntian Deng, and Eric Xing. 2015. Diversifying restricted boltzmann machine for document modeling. In *SIGKDD*.
- [28] Junfeng Yang and Carl Vondrick. 2020. Multitask learning strengthens adversarial robustness. In *ECCV*.
- [29] Haichao Zhang and Jianyu Wang. 2019. Towards adversarially robust object detection. In *ICCV*.

## A PROOF OF THEOREM 6.2

We define the overall output loss of an  $n$ -model fusion team  $\Phi = \{F_1, \dots, F_n\}$  given an input  $\mathbf{x}$  to be

$$\mathcal{L}(\mathbf{x}, \mathbf{O}; \Phi) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(F_i(\mathbf{x}), \mathbf{O}). \quad (22)$$

Then, we can quantify adversarial vulnerability to be the maximum change caused by any perturbation  $\delta$  to  $\mathbf{x}$  with an upper bound of  $\epsilon$  in  $\ell_p$ -norm, i.e.,

$$\Delta \mathcal{L}(\mathbf{x}, \mathbf{O}; \Phi, \epsilon) = \max_{\|\delta\|_p < \epsilon} |\mathcal{L}(\mathbf{x}, \mathbf{O}; \Phi) - \mathcal{L}(\mathbf{x} + \delta, \mathbf{O}; \Phi)|. \quad (23)$$

Since the adversarial perturbation  $\delta$  is often constrained to be human-imperceptible (i.e.,  $\epsilon \rightarrow 0$ ), a first-order Taylor expansion on Equation 23 in  $\epsilon$  gives

$$\Delta \mathcal{L}(\mathbf{x}, \mathbf{O}; \Phi, \epsilon) \approx \max_{\|\delta\|_p < \epsilon} |\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{O}; \Phi) \cdot \delta|, \quad (24)$$

where  $\nabla_{\mathbf{x}}$  is a partial derivative operator with respect to  $\mathbf{x}$ . The FUSE-protected object detection system is vulnerable (or robust) if, on average over  $\mathbf{x}$ ,  $\Delta \mathcal{L}(\mathbf{x}, \mathbf{O}; \Phi, \epsilon)$  is large (or small) [10, 20]. Using the definition of dual norm, the adversarial vulnerability over the dataset  $\mathcal{D}$  is

$$\begin{aligned} & \mathbb{E}_{(\mathbf{x}, \mathbf{O}) \sim \mathcal{D}} [\Delta \mathcal{L}(\mathbf{x}, \mathbf{O}; \Phi, \epsilon)] \\ & \approx \mathbb{E}_{(\mathbf{x}, \mathbf{O}) \sim \mathcal{D}} [\max_{\|\delta\|_p < \epsilon} |\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{O}; \Phi) \cdot \delta|] \\ & \approx \mathbb{E}_{(\mathbf{x}, \mathbf{O}) \sim \mathcal{D}} [||\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{O}; \Phi)||_q \cdot ||\delta||_p] \\ & \propto \mathbb{E}_{(\mathbf{x}, \mathbf{O}) \sim \mathcal{D}} [||\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{O}; \Phi)||_q], \end{aligned} \quad (25)$$

where  $q$  is the dual norm of  $p$ , which satisfies  $\frac{1}{p} + \frac{1}{q} = 1$  and  $1 \leq p \leq \infty$ .

Following [28], we consider  $\ell_2$ -norm and set  $p = q = 2$  in this derivation. To construct the adversarial vulnerability, we first obtain the expected value of the square of  $\ell_2$ -norm of gradients using Equation 22, i.e.,

$$\begin{aligned} & \mathbb{E}_{(\mathbf{x}, \mathbf{O}) \sim \mathcal{D}} [||\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{O}; \Phi)||_2^2] \\ & = \mathbb{E}_{(\mathbf{x}, \mathbf{O}) \sim \mathcal{D}} \left[ \left\| \frac{1}{n} \sum_{i=1}^n \nabla_{\mathbf{x}} \mathcal{L}(F_i(\mathbf{x}), \mathbf{O}) \right\|_2^2 \right]. \end{aligned} \quad (26)$$

Since models were trained until convergence before deployment, the gradients of each task in object detection has a zero mean. Taking the objectness component as an example, the zero mean property (i.e.,  $\mathbb{E}_{(\mathbf{x}, \mathbf{O}) \sim \mathcal{D}} [\nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}(F_i(\mathbf{x}), \mathbf{O})] = 0$ ) yields

$$\begin{aligned} & \text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}(F_i(\mathbf{x}), \mathbf{O}), \nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}(F_j(\mathbf{x}), \mathbf{O})) \\ & = \mathbb{E}_{(\mathbf{x}, \mathbf{O}) \sim \mathcal{D}} [\nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}(F_i(\mathbf{x}), \mathbf{O}) \nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}(F_j(\mathbf{x}), \mathbf{O})] \\ & - \mathbb{E}_{(\mathbf{x}, \mathbf{O}) \sim \mathcal{D}} [\nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}(F_i(\mathbf{x}), \mathbf{O})] \mathbb{E}_{(\mathbf{x}, \mathbf{O}) \sim \mathcal{D}} [\nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}(F_j(\mathbf{x}), \mathbf{O})] \\ & = \mathbb{E}_{(\mathbf{x}, \mathbf{O}) \sim \mathcal{D}} [\nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}(F_i(\mathbf{x}), \mathbf{O}) \nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}(F_j(\mathbf{x}), \mathbf{O})]. \end{aligned} \quad (27)$$

The covariance of the other two tasks (i.e., bounding box estimation and class label prediction) can be similarly derived:

$$\begin{aligned} & \text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}(F_i(\mathbf{x}), \mathbf{O}), \nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}(F_j(\mathbf{x}), \mathbf{O})) \\ & = \mathbb{E}_{(\mathbf{x}, \mathbf{O}) \sim \mathcal{D}} [\nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}(F_i(\mathbf{x}), \mathbf{O}) \nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}(F_j(\mathbf{x}), \mathbf{O})]. \end{aligned} \quad (28)$$

$$\begin{aligned} & \text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}(F_i(\mathbf{x}), \mathbf{O}), \nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}(F_j(\mathbf{x}), \mathbf{O})) \\ & = \mathbb{E}_{(\mathbf{x}, \mathbf{O}) \sim \mathcal{D}} [\nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}(F_i(\mathbf{x}), \mathbf{O}) \nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}(F_j(\mathbf{x}), \mathbf{O})]. \end{aligned} \quad (29)$$

Expanding Equation 26 with Equation 1 and the above property, we obtain

$$\begin{aligned} & \mathbb{E}_{(\mathbf{x}, \mathbf{O}) \sim \mathcal{D}} [\Delta \mathcal{L}(\mathbf{x}, \mathbf{O}; \Phi, \epsilon)] \\ & \propto \sqrt{\frac{1}{3n} + \frac{2}{9n^2} [\text{Cov}_{\text{intra}}(\mathbf{x}, \mathbf{O}; \Phi) + \text{Cov}_{\text{inter}}(\mathbf{x}, \mathbf{O}; \Phi)]}. \end{aligned} \quad (30)$$

$\text{Cov}_{\text{intra}}(\mathbf{x}, \mathbf{O}; \Phi)$  and  $\text{Cov}_{\text{inter}}(\mathbf{x}, \mathbf{O}; \Phi)$  are two terms capturing the covariance of gradients across tasks and models:

$$\begin{aligned} \text{Cov}_{\text{intra}}(\mathbf{x}, \mathbf{O}; \Phi) & = \sum_{i=2}^n \sum_{j=1}^{i-1} \left[ \frac{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}^j)}{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}^i)} \right. \\ & \quad \left. + \frac{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^j)}{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^i)} + \frac{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^j)}{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^i)} \right], \\ \text{Cov}_{\text{inter}}(\mathbf{x}, \mathbf{O}; \Phi) & = \sum_{i=2}^n \sum_{j=1}^{i-1} \left[ \frac{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^j)}{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}^i)} + \frac{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^j)}{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}^i)} \right. \\ & \quad \left. + \frac{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}^j)}{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^i)} + \frac{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^j)}{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^i)} \right. \\ & \quad \left. + \frac{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}^j)}{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^i)} + \frac{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^j)}{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^i)} \right] \\ & + \sum_{i=1}^n \left[ \frac{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}^i)}{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^i)} + \frac{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{obj}}^i)}{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^i)} \right. \\ & \quad \left. + \frac{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{bbox}}^i)}{\text{Cov}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^i, \nabla_{\mathbf{x}} \mathcal{L}_{\text{cls}}^i)} \right], \end{aligned} \quad (32)$$

where  $\mathcal{L}_{\text{task}}^i = \mathcal{L}_{\text{task}}(F_i(\mathbf{x}), \mathbf{O})$  for task  $\in \{\text{obj}, \text{bbox}, \text{cls}\}$ . The detailed mathematical steps of a similar proof can be referred to [28].

## B FORMAL ANALYSIS: ADVERSARIAL PERTURBATION LOWER BOUND

We further present a formal analysis on the robustness offered by FUSE from the perspective of adversarial perturbation lower bound [25]. By considering a special case where input images contain only one object, we show that the adversarial attack is already hardened by FUSE, and then attacking an object detection system with multi-object recognition and localization is even more challenging.

Let  $p_j(\mathbf{x})$  be the pre-softmax confidence of the object contained in input  $\mathbf{x}$  being a class  $j$  instance. With the Lipschitz assumption that  $p_j$  is a continuously differentiable function on the set  $S$  containing  $\mathbf{x}$ ,  $g(\mathbf{x}) = p_{\text{GT}}(\mathbf{x}) - p_{y^*}(\mathbf{x})$  is also a Lipschitz function, where GT is the true class label of  $\mathbf{x}$  and  $y^*$  is an incorrect class label. Then, we obtain:

$$|g(\mathbf{y}) - g(\mathbf{x})| \leq L_q \|\mathbf{y} - \mathbf{x}\|_p. \quad (33)$$

where  $L_q = \max \{ \|\nabla g(\mathbf{x})\|_q : \mathbf{x} \in S \}$  is the local Lipschitz constant with  $\nabla g(\mathbf{x}) = (\frac{\partial g(\mathbf{x})}{\partial x_1}, \dots, \frac{\partial g(\mathbf{x})}{\partial x_d})^\top$ ,  $\mathbf{x}, \mathbf{y} \in S$ ,  $\frac{1}{p} + \frac{1}{q} = 1$ ,  $p \geq 1$ , and  $q \leq \infty$ . As introduced in Section 3, adversarial example  $\mathbf{x}' \in S$  is generated by adding malicious perturbation  $\delta$  to the benign input  $\mathbf{x}$  (i.e.,  $\mathbf{x}' = \mathbf{x} + \delta$ ). By setting  $\mathbf{y}$  in Equation 33 to be  $\mathbf{x}'$  and rearranging, we have

$$g(\mathbf{x}) - L_q \|\delta\|_p \leq g(\mathbf{x} + \delta) \leq g(\mathbf{x}) + L_q \|\delta\|_p. \quad (34)$$

When  $g(\mathbf{x} + \boldsymbol{\delta}) < 0$ ,  $p_{\text{GT}}(\mathbf{x}') < p_{y^*}(\mathbf{x}')$  is hold for an adversarial example  $\mathbf{x}'$ . Hence, if  $\|\boldsymbol{\delta}\|_p$  is small enough such that  $g(\mathbf{x}) - L_q \|\boldsymbol{\delta}\|_p \geq 0$ , then no adversarial examples can be found. That is, if

$$g(\mathbf{x}) - L_q \|\boldsymbol{\delta}\|_p \geq 0 \Rightarrow \|\boldsymbol{\delta}\|_p \leq \frac{p_{\text{GT}}(\mathbf{x}) - p_{y^*}(\mathbf{x})}{L_q}, \quad (35)$$

then the detector’s decision can never be changed and the attack will never succeed.

The above analysis states that there exists a guarded region (hypersphere) with a radius of  $(p_{\text{GT}}(\mathbf{x}) - p_{y^*}(\mathbf{x}))/L_q$  centered at the benign input  $\mathbf{x}$  where no adversarial examples can be found. If we can enlarge the guarded region, then the victim detector under protection will be more resilient to even stronger attacks at the increased cost in terms of adversarial distortion. Indeed, under the protection by a team of diverse models, simply generating adversarial examples with the minimum distortion (Equation 35) is insufficient as they must fool multiple members to misbehave identically, such as fabricating objects with the exact same location, dimension, and class label, which is extremely difficult, if not impossible. Hence, the guarded region under FUSE is larger than the original lower bound, and the increased cost for adversarial distortion ensures stronger adversarial robustness.

## C EXPERIMENT SETUP

**Datasets and Machine.** We conduct experiments on two datasets: (1) PASCAL VOC (2007+2012) [7] and (2) INRIA [6]. All measurements are recorded on NVIDIA RTX 2080 SUPER GPU, Intel i7-9700K (3.60GHz) CPU, and 32 GB RAM on Ubuntu 18.04.

**Object Detection Models.** We include three dominant detection algorithms: Faster R-CNN (FRCNN) [19], YOLOv3 [18] and SSD [14] with heterogeneous neural architectures ( $F_1$  to  $F_5$  in Table 2). Those pre-trained models are available on public repositories. For YOLOv3 including the Darknet53 (i.e., YOLOv3-D) and MobileNetV1 (i.e., YOLOv3-M) backbones, we follow the repository<sup>1</sup>. For SSD, the code is exported from the repository<sup>2</sup>. For Faster R-CNN (FRCNN),

we exploit the source code from the repository<sup>3</sup>. All diversity joint training in our experiments follow the same hyperparameter setting, including  $\lambda = 10$  for diversity regularization optimizations. The learning rate schedule and all other hyperparameter configuration follow the original repository of Faster R-CNN without additional fine-tuning. The diversity joint training on kernel filters decorrelates the first two convolutional layers.

**Attacks.** Four state-of-the-art untargeted attacks are included: TOG [5], UEA [24], RAP [13], and DAG [26]. We also evaluate the robustness against targeted attacks [5] with object vanishing (TOG-v), fabrication (TOG-f), and mislabeling (TOG-m) effects. We generate adversarial examples by using the source code released by TOG<sup>4</sup>, UEA<sup>5</sup>, RAP<sup>6</sup>, and DAG<sup>7</sup> attacks with their default hyperparameters. For the TOG patches, they are generated with size  $(64 \times 64)$ , trained with 30 epochs, a batch size of 8, and a learning rate of 0.1.

**Defenses.** For FUSE, we always form teams of three models (including the victim) as we found this team size is sufficient to offer strong robustness. The IOU threshold  $\tau_{\text{IOU}}$  is always set to be 0.50. The mitigation baseline using the adversarial detection training method [29] is implemented by following the description in the paper as the original source code is unavailable.

<sup>1</sup><https://github.com/Adamdad/keras-YOLOv3-mobilenet>

<sup>2</sup>[https://github.com/pierluigiferrari/ssd\\_keras](https://github.com/pierluigiferrari/ssd_keras)

<sup>3</sup><https://github.com/chenyuntc/simple-faster-rcnn-pytorch>

<sup>4</sup><https://github.com/git-disl/TOG>

<sup>5</sup><https://github.com/LiangSiyuan21/Adversarial-Attacks-for-Image-and-Video-Object-Detection>

<sup>6</sup><https://github.com/yuezunli/BMVC2018R-AP>

<sup>7</sup><https://github.com/cihangxie/DAG>