# Towards Adversarially Robust Object Detection

Haichao Zhang     Jianyu Wang

Baidu Research,  Sunnyvale USA

hczhang1@gmail.com   wjyouch@gmail.com

## Abstract

*Object detection is an important vision task and has emerged as an indispensable component in many vision system, rendering its robustness as an increasingly important performance factor for practical applications. While object detection models have been demonstrated to be vulnerable against adversarial attacks by many recent works, very few efforts have been devoted to improving their robustness. In this work, we take an initial attempt towards this direction. We first revisit and systematically analyze object detectors and many recently developed attacks from the perspective of model robustness. We then present a multi-task learning perspective of object detection and identify an asymmetric role of task losses. We further develop an adversarial training approach which can leverage the multiple sources of attacks for improving the robustness of detection models. Extensive experiments on PASCAL-VOC and MS-COCO verified the effectiveness of the proposed approach.*

## 1. Introduction

Deep learning models have been widely applied to many vision tasks such as classification [45, 47, 19] and object detection [15, 14, 29, 40, 42, 3], leading to state-of-the-art performance. However, one impeding factor of deep learning models is their issues with robustness. It has been shown that deep net-based classifiers are vulnerable to adversarial attack [49, 16], *i.e.*, there exist adversarial examples that are slightly modified but visually indistinguishable version of the original images that cause the classifier to generate incorrect predictions [36, 4]. Many efforts have been devoted to improving the robustness of classifiers [35, 34, 56, 17, 25, 44, 46, 38, 30].

Object detection is a computer vision technique that deals with detecting instances of semantic objects in images [54, 8, 12]. It is a natural generalization of the vanilla classification task as it outputs not only the object label as in classification but also the location. Many successful object detection approaches have been developed during the past several years [15, 14, 42, 29, 40] and object detectors pow-
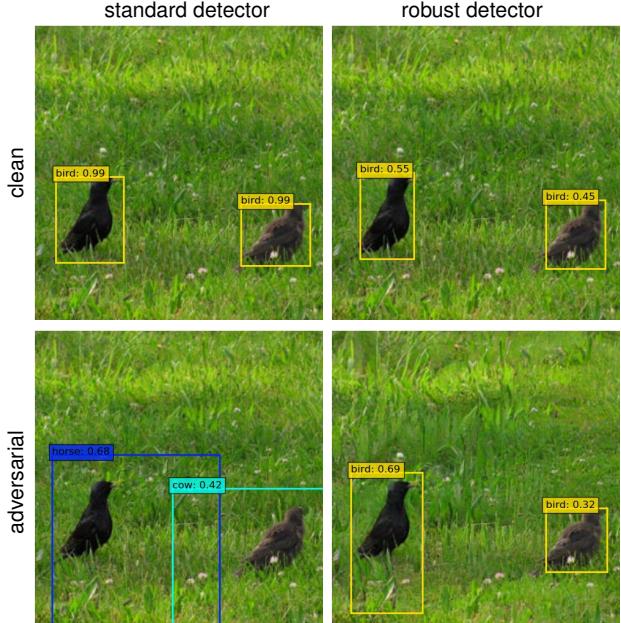


Figure 1. **Standard v.s. robust detectors** on clean and adversarial images. The adversarial image is produced using PDG-based detector attacks [23, 33] with perturbation budget 8 (out of 256). The standard model [29] fails completely on the adversarial image while the robust model can produce reasonable detection results.

ered by deep nets have emerged as an indispensable component in many vision systems of real-world applications.

Recently, it has been shown that object detectors can also be attacked by maliciously crafted inputs [57, 32, 23, 6, 55, 11, 31, 22] (*c.f.* Figure 1). Given its critical role in applications such as surveillance and autonomous driving, it is important to investigate approaches for defending object detectors against various adversarial attacks. However, while many works have shown it is possible to attack a detector, it remains largely unclear whether it is possible to improve the robustness of the detectors and what is the practical approach for that. This work servers as an initial attempt to bridge this gap towards this direction. We show that it is *possible* to improve the robustness of the object detector *w.r.t.* various types of attacks and propose a practical approach for achieving this, by generalizing the adversarial training framework from classification to detection.

arXiv:1907.10310v1 [cs.CV] 24 Jul 2019

The contribution of this paper is threefold: *i*) we provide a categorization and analysis of different attacks for object detectors, revealing their shared underlying mechanisms; *ii*) we highlight and analyze the interactions between different tasks losses and their implication on robustness; *iii*) we generalize the adversarial training framework from classification to detection and develop an adversarial training approach that can properly handle the interactions between task losses for improving detection robustness.

## 2. Related Work

**Attacks and Adversarial Training for Classification**. Adversarial examples have been investigated for general learning-based classifiers before [2]. As a learning-based model, deep networks are also vulnerable to adversarial examples [49, 37]. Many variants of attacks [16, 36, 4] and defenses [35, 34, 56, 17, 25, 30, 44, 46, 38, 1] have been developed. Fast gradient sign method (FGSM) [16] and Projective Gradient Descend (PGD) [33] are two representative approaches for white-box adversarial attack generation. Adversarial training [16, 21, 50, 33] is one of the effective defense method against adversarial attacks. It achieves robust model training by solving a minimax problem, where the inner maximization generates attacks according to the current model parameters while the outer optimization minimize the training loss *w.r.t.* the model parameters [16, 33].

**Object Detection and Adversarial Attacks**. Many successful object detection approaches have been developed during the past several years, including one-stage [29, 40] and two-stage variants [15, 14, 42]. Two stage detectors refine proposals from the first stage by one or multiple refinement steps [42, 3]. We focus on one-stage detectors in this work due to its essential role in different variants of detectors. A number of attacks for object detectors have been developed very recently [57, 32, 6, 11, 55, 23, 22, 31]. [57] extends the attack generation method from classification to detection and demonstrates that it is possible to attack objectors using a designed classification loss. Lu *et al.* generate adversarial examples that fool detectors for stop sign and face detections [32]. [6] develops physical attacks for Faster-RCNN [42] and adapts the expectation-over-transformation idea for generating physical attacks that remain effective under various transformations such as viewpoint variations. [23] proposes to attack the region-proposal network (RPN) with a specially designed hybrid loss incorporating both classification and localization terms. Apart from the full images, it is also possible to attack detectors by restricting the attacks to be within a local region [22, 31].

## 3. Object Detection and Attacks Revisited

We revisit object detection and discuss the connections between many variants of attacks developed recently.
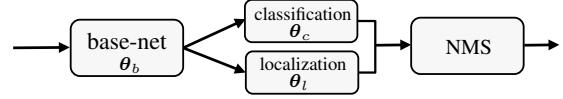


Figure 2. **One-stage detector architecture**. A base-net (w. para. $\boldsymbol{\theta}_b$) is shared by classification (w. para. $\boldsymbol{\theta}_c$) and localization (w. para. $\boldsymbol{\theta}_l$) tasks. $\boldsymbol{\theta} = [\boldsymbol{\theta}_b, \boldsymbol{\theta}_c, \boldsymbol{\theta}_l]$ denotes the full parameters for the detector. For training, the NMS module is removed and task losses are appended for classification and localization respectively.

### 3.1. Object Detection as Multi-Task Learning

An object detector $f(\mathbf{x}) \rightarrow \{\mathbf{p}_k, \mathbf{b}_k\}_{k=1}^K$ takes an image $\mathbf{x} \in [0, 255]^n$ as input and outputs a varying number of $K$ detected objects, each represented by a probability vector $\mathbf{p}_k \in \mathbb{R}^C$ over $C$ classes (including background) and a bounding box $\mathbf{b}_k = [x_k, y_k, w_k, h_k]$. Non-maximum suppression (NMS) [43] is applied to remove redundant detections for the final detections (*c.f.* Figure 2).

For training, we parametrize the detector $f(\cdot)$ by $\boldsymbol{\theta}$. Then the training of the detector boils down to the estimation of $\boldsymbol{\theta}$ which can be formulated as follows:

$$\min_{\boldsymbol{\theta}} \mathbb{E}_{(\mathbf{x}, \{y_k, \mathbf{b}_k\}) \sim \mathcal{D}} \mathcal{L}(f_{\boldsymbol{\theta}}(\mathbf{x}), \{y_k, \mathbf{b}_k\}). \tag{1}$$

$\mathbf{x}$ denotes the training image and $\{y_k, \mathbf{b}_k\}$ the ground-truth (class label $y_k$ and the bounding box $\mathbf{b}_k$) sampled from the dataset $\mathcal{D}$. We will drop the expectation over data and present subsequent derivations with a single example to avoid notation clutter without loss of generality as follows:

$$\min_{\boldsymbol{\theta}} \mathcal{L}(f_{\boldsymbol{\theta}}(\mathbf{x}), \{y_k, \mathbf{b}_k\}). \tag{2}$$

$\mathcal{L}(\cdot)$ is a loss function measuring the difference between the output of $f_{\boldsymbol{\theta}}(\cdot)$ and the ground-truth and the minimization of it (over the dataset) leads to a proper estimation of $\boldsymbol{\theta}$. In practice, it is typically instantiated as a combination of classification loss and localization loss as follows [29, 40]:

$$\min_{\boldsymbol{\theta}} \text{loss}_{\text{cls}}(f_{\boldsymbol{\theta}}(\mathbf{x}), \{y_k, \mathbf{b}_k\}) + \text{loss}_{\text{loc}}(f_{\boldsymbol{\theta}}(\mathbf{x}), \{y_k, \mathbf{b}_k\}). \tag{3}$$

As shown in Eqn.(3), the classification and localization tasks share some intermediate computations including the base-net (*c.f.* Figure 2). However, they use different parts of the output from $f_{\boldsymbol{\theta}}(\cdot)$ for computing losses emphasizing on different aspects, *i.e.*, classification and localization performance respectively. This is a design choice for sharing feature and computation for potentially relevant tasks [29, 40], which is essentially an instance of *multi-task learning* [5].

### 3.2. Detection Attacks Guided by Task Losses

Many different attack methods for object detectors have been developed very recently [57, 32, 6, 11, 55, 23, 22, 31]. Although there are many differences in the formulations of these attacks, when viewed from the *multi-task learning*

| Attacks for Object Detection | Components | | | |
|---|---|---|---|---|
| | $\text{loss}_{\text{cls}}$ | | $\text{loss}_{\text{loc}}$ | |
| | T | N | T | N |
| **ShapeShifter** [6] | ✓ | | | |
| **DFool** [32], **PhyAttack** [11] | | ✓ | | |
| **DAG** [57], **Transfer** [55] | ✓ | ✓ | | |
| **DPatch** [31] | ✓ | | ✓ | |
| **RAP** [23] | | ✓ | ✓ | |
| **BPatch** [22] | | ✓ | | ✓ |

Table 1. Analysis of existing attack methods for object detection. "T" denotes "targeted attack" and "N" for "non-targeted attack".

perspective as pointed out in Section 3.1, they have the same framework and design principle: *an attack to a detector can be achieved by utilizing variants of individual task losses or their combinations*. This provides a common grounding for understanding and comparing different attacks for object detectors. From this perspective, we can categorize existing attack methods as in Table 1. It is clear that some methods use classification loss [6, 32, 11, 57, 55] while other methods also incorporated localization loss [31, 23, 22]. There are two perspectives for explaining the effectiveness of individual task loss in generating attacks: ***i***) the classification and localization tasks share a common base-net, implying that the weakness in the base-net will be shared among all tasks built upon it; ***ii***) while the classification and localization outputs have dedicated branches for each task beyond the shared base-net, they are coupled in the testing phase due to the usage of NMS, which jointly use class scores and bounding box locations for redundant prediction pruning.

Although many attacks have been developed and it is possible to come up with new combinations and configurations following the general principle, there is a lack of understanding on the *role of individual components* in model robustness. Filling this gap one of our contributions which will naturally lead to our robust training method for object detectors as detailed in the sequel.

## 4. Towards Adversarially Robust Detection

### 4.1. The Roles of Task Losses in Robustness

As the classification and localization tasks of a detector share a base-net (*c.f*. Figure 2), the two tasks will inevitably affect each other even though the input images are manipulated according to a criterion trailered for one individual task. We therefore conduct analysis on the role of task losses in model robustness from several perspectives.

**Mutual Impacts of Task Losses.** Our first empirical observation is that *different tasks have mutual impacts and the adversarial attacks trailered for one task can reduce the performance of the model on the other task*. To show this, we take a marginalized view over one factor while investigating the impact of the other. For example, when con-
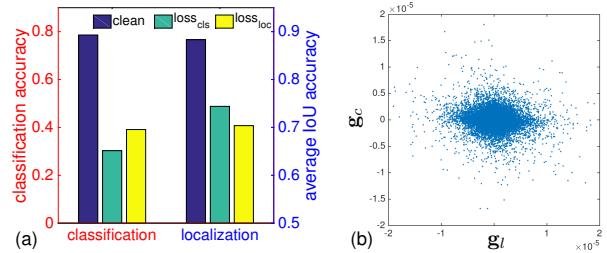


Figure 3. **Mutual impacts of task losses and gradient visualization**. (a) Model performance on classification and localization under different attacks: clean image, $\text{loss}_{\text{cls}}$-based attack and $\text{loss}_{\text{loc}}$-based attack. The model is a standard detector trained on clean images. The performance metric is detailed in text. (b) Scatter plot of task gradients for classification $\mathbf{g}_c$ and localization $\mathbf{g}_l$.

sidering classification, we can marginalize out the factor of location and the problem reduces to a multi-label classification task [52]; on the other hand, when focusing on localization only, we can marginalize out the class information and obtain a class agnostic object detection problem [53]. The results with single step PGD and budget 8 are shown in Figure 3 (a). The performances are measured on detection outputs *prior* to NMS to better reflect the raw performance. A candidates set is first determined as the foreground candidates whose prior boxes have an IoU value larger than 0.5 with any of the ground-truth annotation. This ensures that each selected candidate has a relative clean input both tasks. For classification, we compute the classification accuracy on the candidate set. For localization, we compute the average IoU of the predicted bounding boxes with ground-truth bounding boxes. The attack is generated with one-step PGD and a budget of 8. It can be observed from the results in Figure 3 (a) that the two losses interact with each other. The attacks based on the classification loss ($\text{loss}_{\text{cls}}$) reduces the classification performance and decreases the localization performance at the same time. Similarly, the localization loss induced attacks ($\text{loss}_{\text{loc}}$) reduces not only the location performance but the classification performance as well. This can essentially be viewed as a type of cross-task attack transfer: *i.e.*. when using only the classification loss (task) to generate adversarial images, the attacks can be transferred to localization tasks and reduce its performance and vice versa. This is one of the reason why adversarial images generated based on individual task losses (*e.g.* classification loss [57]) can effectively attack object detectors.

**Misaligned Task Gradients.** Our second empirical observation is that *the gradients of the two tasks share certain level of common directions but are not fully aligned, leading to misaligned task gradients that can obfuscate the subsequent adversarial training*. To show this, we analyze the image gradients derived from the two losses (referred to as *task gradients*), *i.e.*, $\mathbf{g}_c = \nabla_{\mathbf{x}}\text{loss}_{\text{cls}}$ and $\mathbf{g}_l = \nabla_{\mathbf{x}}\text{loss}_{\text{loc}}$. The element-wise scatter plot between $\mathbf{g}_c$ and $\mathbf{g}_l$ is shown
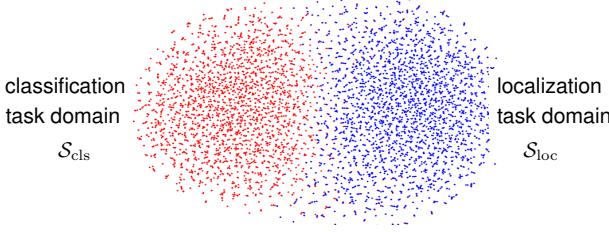
Figure 4. **Visualization of task domains** $\mathcal{S}_{\text{cls}}$ and $\mathcal{S}_{\text{loc}}$ using t-SNE. Given a single clean image $\mathbf{x}$, each dot in the picture represents one adversarial example generated by solving Eqn.(5) staring from a random point within the $\epsilon$-ball around $\mathbf{x}$. Different colors encode the task losses used for generating adversarial examples (red: $\text{loss}_{\text{cls}}$, blue: $\text{loss}_{\text{loc}}$). Therefore, the samples form empirical images of the corresponding task domains. It is observed that the two task domains have both overlaps and distinctive regions.

in Figure 3 (b). We have several observations: *i*) the magnitudes of the task gradients are not the same (different value ranges), indicating the potential existence of imbalance between the two task losses; *ii*) the direction of the task gradients are inconsistent (non-diagonal), implying the potential conflicts between the two tasks gradients. We further visualize the task gradient domains representing the domain of a task maximizing gradient for each respective task (*c.f.* Eqn.(5)) as in Figure 4. The fact the the two domains are not fully separated (*i.e.* they do not collapse to two isolated clusters) further reinforces our previous observation on their mutual impacts. The other aspect that they have a significant non-overlapping portion is another reflection of the mis-alignments between task gradients (task domains).

## 4.2. Adversarial Training for Robust Detection

Motivated by the preceding analysis, we propose the following formulation for robust object detection training:

$$\min_{\boldsymbol{\theta}} \Big[ \max_{\bar{\mathbf{x}} \in \mathcal{S}_{\text{cls}} \cup \mathcal{S}_{\text{loc}}} \mathcal{L}(f_{\boldsymbol{\theta}}(\bar{\mathbf{x}}), \{y_k, \mathbf{b}_k\}) \Big], \quad (4)$$

where the *task-oriented domain* $\mathcal{S}_{\text{cls}}$ and $\mathcal{S}_{\text{loc}}$ represent the permissible domains induced by each individual tasks:

$$\mathcal{S}_{\text{cls}} \triangleq \{\bar{\mathbf{x}} | \arg\max_{\bar{\mathbf{x}} \in \mathcal{S}_{\mathbf{x}}} \text{loss}_{\text{cls}}(f(\bar{\mathbf{x}}), \{y_k\}))\}$$
$$\mathcal{S}_{\text{loc}} \triangleq \{\bar{\mathbf{x}} | \arg\max_{\bar{\mathbf{x}} \in \mathcal{S}_{\mathbf{x}}} \text{loss}_{\text{loc}}(f(\bar{\mathbf{x}}), \{\mathbf{b}_k\}))\}$$
$$(5)$$

where $\mathcal{S}_{\mathbf{x}}$ is defined as $\mathcal{S}_{\mathbf{x}} = \{\mathbf{z} | \mathbf{z} \in B(\mathbf{x}, \epsilon) \cap [0, 255]^n\}$, and $B(\mathbf{x}, \epsilon) = \{\mathbf{z} | \|\mathbf{z} - \mathbf{x}\|_{\infty} \leq \epsilon\}$ denotes the $\ell_{\infty}$-ball with center as the clean image $\mathbf{x}$ and radius as the perturbation budget $\epsilon$. We denote $\mathcal{P}_{\mathcal{S}_{\mathbf{x}}}(\cdot)$ as a projection operator projecting the input into the feasible region $\mathcal{S}_{\mathbf{x}}$. It is important to note several crucial differences compared with the conventional adversarial training for classification:

- **multi-task sources for adversary training**: different from the adversarial training in classification case [16, 33] where only a single source is involved, here we have

---

**Algorithm 1** Adversarial Training for Robust Detection

**Input:** dataset $\mathcal{D}$, training epochs $T$, batch size $S$,
      learning rate $\gamma$, attack budget $\epsilon$
**for** $t = 1$ **to** $T$ **do**
  **for** random batch $\{\mathbf{x}^i, \{y_k^i, \mathbf{b}_k^i\}\}_{i=1}^S \sim \mathcal{D}$ **do**
    $\cdot\ \tilde{\mathbf{x}}^i \sim B(\mathbf{x}^i, \epsilon)$
    compute attacks in the classification task domain
    $\cdot\ \bar{\mathbf{x}}_{\text{cls}}^i = \mathcal{P}_{\mathcal{S}_{\mathbf{x}}}\big(\tilde{\mathbf{x}}^i + \epsilon \cdot \text{sign}\big(\nabla_{\mathbf{x}}\text{loss}_{\text{cls}}(\tilde{\mathbf{x}}^i, \{y_k^i\})\big)\big)$
    compute attacks in the localization task domain
    $\cdot\ \bar{\mathbf{x}}_{\text{loc}}^i = \mathcal{P}_{\mathcal{S}_{\mathbf{x}}}\big(\tilde{\mathbf{x}}^i + \epsilon \cdot \text{sign}\big(\nabla_{\mathbf{x}}\text{loss}_{\text{loc}}(\tilde{\mathbf{x}}^i, \{\mathbf{b}_k^i\})\big)\big)$
    compute the final attack examples
    $\cdot\ \mathbf{m} = \mathcal{L}(\bar{\mathbf{x}}_{\text{cls}}^i, \{y_k^i, \mathbf{b}_k^i\}) > \mathcal{L}(\bar{\mathbf{x}}_{\text{loc}}^i, \{y_k^i, \mathbf{b}_k^i\})$
    $\cdot\ \bar{\mathbf{x}}^i = \mathbf{m} \odot \bar{\mathbf{x}}_{\text{cls}}^i + (1 - \mathbf{m}) \odot \bar{\mathbf{x}}_{\text{loc}}^i$
    perform adversarial training step
    $\cdot\ \boldsymbol{\theta} = \boldsymbol{\theta} - \gamma \cdot \nabla_{\boldsymbol{\theta}} \frac{1}{S} \sum_{i=1}^S \mathcal{L}(\bar{\mathbf{x}}^i, \{y_k^i, \mathbf{b}_k^i\}; \boldsymbol{\theta})$
  **end for**
**end for**
**Output:** learned model parameter $\boldsymbol{\theta}$ for object detection.

---

*multiple* (in the presence of multiple objects) and *heterogeneous* (both classification and localization) sources of supervisions for adversary generation and training, thus generalizing the adversarial training for classification;

- **task-oriented domain constraints**: different from the conventional adversarial training setting which uses a *task-agnostic* domain constraint $\mathcal{S}_{\mathbf{x}}$, we introduce a *task-oriented* domain constraint $\mathcal{S}_{\text{cls}} \cup \mathcal{S}_{\text{loc}}$ which restricts the permissible domain as the set of images that maximize either the classification task losses or the localization losses. The final adversarial example used for training is the one that maximizes the overall loss within this set. The crucial advantage of the proposed formulation with task-domain constraints is that we can benefit from generating adversarial examples guided by each task without suffering from the interferences between them.

If we relax the task-oriented domain to $\mathcal{S}_{\mathbf{x}}$, set the coordinates of the bounding box corresponding to the full image and assign a single class label to the image, then the proposed formulation Eqn.(4) reduces to the conventional adversarial training setting for classification [16, 33]. Therefore, we can view the proposed adversarial training for robust detection as a natural generalization of the conventional adversarial training under the classification setting. However, it is crucial to note that while both tasks contribute to improving the model robustness in expectation according to their overall strengths, there is no interference between the tasks for generating individual adversarial example due to the task oriented domain in contrast to $\mathcal{S}_{\mathbf{x}}$ (*c.f.* Sec.5.3).

Training object detection models that are resistant to adversarial attacks boils down to solving a minimax problem as in Eqn.(4). We solve it approximately by replacing the original training images with the adversarially per-
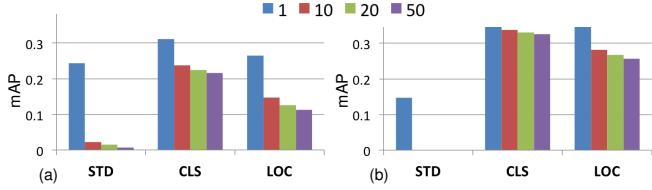
Figure 5. **Model performance under different number of steps** for (a) $\text{loss}_{\text{cls}}$ and (b) $\text{loss}_{\text{loc}}$-based PGD attack with $\epsilon = 8$. **STD** is the standard model. **CLS** and **LOC** are our robust models.



Figure 6. **Model performance under different attack budgets** for (a) $\text{loss}_{\text{cls}}$ and (b) $\text{loss}_{\text{loc}}$-based PGD attack with 20 steps. **STD** is the standard model. **CLS** and **LOC** are our robust models.



Figure 7. **Visualization of attacks** on **STD** model using $\text{loss}_{\text{cls}}$ based 20-step PGD attack (zoom electronically for better view).

turbed ones obtained by solving the inner problem, and then conducting conventional training of the model using the perturbed images as typically done in adversarial training [16, 33]. The inner maximization is approximately solved using a variant of FGSM [16] for efficiency. For incorporating the task-oriented domain constraint, we propose to take FGSM steps within each task domain and then select the one that maximizes the overall loss. The details of the algorithm are summarized in Algorithm 1.

## 5. Experiments

### 5.1. Experiment and Implementation Details

We use the single-shot multi-box detector (SSD) [29] with VGG16 [45] backbone as one of the representative single-shot detectors in our experiments. We also make the necessary modifications to the VGG16 net as detailed in [29] and keep the batch normalization layers. Experiments with different detector architectures (Receptive Field Block-based Detector (RFB) [28], Feature Fusion Single Shot Detector (FSSD) [24] and YOLO-V3 [40, 41]) and backbones (VGG16 [45], ResNet50 [19], DarkNet53 [39]) are also conducted for comprehensive evaluations.

For PASCAL VOC dataset, we adopt the standard "07+12" protocol (a union of 2007 and 2012 `trainval`, $\sim$16k images) following [29] for training. For testing, we use PASCAL VOC2007 `test` with 4952 test images and 20 classes [10].[1] For MS-COCO dataset [27], we train on `train+valminusminival` 2014 ($\sim$120k images) and test on `minival` 2014 with 80 classes ($\sim$5k images) . The "mean average precision" (mAP) with IoU threshold 0.5 is used for evaluating the performance of a detector [10].

All models are trained from scratch using SGD with an initial learning rate of $10^{-2}$, momentum 0.9, weight decay 0.0005 and batch size 32 [18] with the multi-box loss [9, 48]. The learning rate schedule is [40k, 60k, 80k] for PASCAL VOC and [180k, 220k, 260k] for MS-COCO with decay factor 0.1. The size of the image is $300 \times 300$. Pixel value range is $[0, 255]$ shifted according to dataset mean. For adversarial attacks and training, we use a budget $\epsilon = 8$, which roughly corresponds to a PSNR of 30 between the perturbed and original images following [23].

---

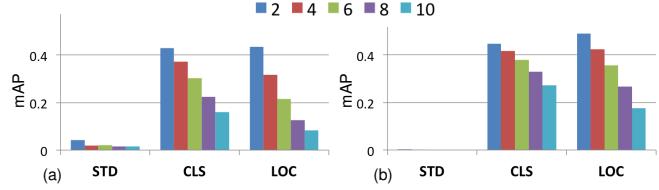[1] VOC2012 `test` is not used as the annotations required for generating attacks are unavailable.

All the attack methods incorporate $\text{sgn}(\cdot)$ operator into the PGD steps for normalization and efficiency following [16].

### 5.2. Impacts of Task Losses on Robustness

We will investigate the role of task losses in model robustness. For this purpose, we introduce the standard model and several variations of our proposed robust model:

- **STD**: standard training with clean image as the domain
- **CLS**: using $\mathcal{S}_{\text{cls}}$ only as the task domain for training
- **LOC**: using $\mathcal{S}_{\text{loc}}$ only as the task domain for training.

We will systemically investigate the performance of these models under attacks induced by *individual task losses* with different number of attack steps and budgets as follows.

**Attacks under different number of steps.** We first evaluate the performance of models under attacks with different number of PGD steps under a fixed attack budget of 8. The results are shown in Figure 5. We have several interesting observations from the results: *i*) the performance of the standard model (**STD**) drops below all other robust models within just a few steps and decreases quickly (approaching zero) as the number of PGD steps increases, for both $\text{loss}_{\text{cls}}$-base and $\text{loss}_{\text{loc}}$-based attacks. These results imply that both types of attacks are very effective attacks for detectors; *ii*) all the robust models maintains a relative stable performance across different number of attack steps, indicating their improved robustness against adversarial attacks compared to the standard model.

**Attacks with different budgets.** We evaluate model robustness under a range of different attack budgets $\epsilon \in \{2, 4, 6, 8, 10\}$. The results are presented in Figure 6. It is observed that the performance of the standard model trained with natural images (**STD**) drops significantly, *e.g.*, from $\sim$72% on clean images (not shown in figure) to $\sim$4% with a small attack budget of 2. Robust models, on the other hand, degrade more gracefully as the attack budget increases, implying their improved robustness compared to the standard

| attacks | | clean | $\text{loss}_{\text{cls}}$ | $\text{loss}_{\text{loc}}$ | DAG [57] | RAP [23] |
|---------|-----|-------|------|------|------|------|
| **standard** | | 72.1 | 1.5 | 0.0 | 0.3 | 6.6 |
| **ours** | **CLS** | 46.7 | 21.8 | 32.2 | 28.0 | 43.4 |
| | **LOC** | 51.9 | 23.7 | 26.5 | 17.2 | 43.6 |
| | **CON** | 38.7 | 18.3 | 27.2 | 26.4 | 40.8 |
| | **MTD** | 48.0 | 29.1 | 31.9 | 28.5 | 44.9 |
| **ours avg** | | 46.3 | 23.2 | 29.4 | 25.0 | 43.2 |

Table 2. Impacts of task domains on model performance (mAP) and defense against attacks from literature (attack $\epsilon = 8$).

model. In Figure 7, we visualize the detection results under different attack budgets on standard model. It is observed that even with a small attack budget (*e.g.* $\epsilon = 2$), the detection results are changed completely, implying that the standard model is very fragile in term of robustness, which is consistent with our previous observation from Figure 6. It is also observed that the erroneous detections can be of several forms: *i*) label flipping: the bounding box location is roughly correct but the class label is incorrect, *e.g.*, "dinningtable" ($\epsilon : 0 \to 2$); *ii*) disappearing: the bounding box for the object is missing, *e.g.*, "horse" and "person" ($\epsilon : 0 \to 2$); *iii*) appearing: spurious detections of objects that do not exist in the image with locations not well aligned with any of the dominant objects, *e.g.*, "chair" ($\epsilon : 0 \to 2$) and "pottedplant" ($\epsilon : 2 \to 8$). As the attack budget is increased, the detection output will be further changed in terms of the three types of changes described above. It can also be observed from the figure that the attack image generated with $\epsilon = 8$ bears noticeable changes compared with the original one, although not very severe. We will therefore use attack $\epsilon = 8$ as it is a large enough attack budget while maintain a reasonable resemblance to the original image.

### 5.3. Beyond Single-Task Domain

We further examine the impacts of task domains on robustness. The following approaches with different task domains are considered in addition to **STD**, **CLS** and **LOC**:
- **CON**: using the conventional task agnostic domain $\mathcal{S}_{\mathbf{x}}$, which is essentially the direct application of the adversarial training for classification [16, 33] to detection;
- **MTD**: using the task oriented domain $\mathcal{S}_{\text{cls}} \cup \mathcal{S}_{\text{loc}}$.

The results are summarized in Table 2. It is observed from comparison that different domains lead to different levels of model robustness. For example, for methods with a single task domain, **LOC** leads to less robust models compared with **CLS**. On the other hand, **LOC** has a higher clean accuracy than **CLS**. Therefore, it is not straightforward to select one single domain as it is *unknown a priori* whether one of the task domains is the best. Simply relaxing the task domains as done in the conventional adversarial training **CON** [16, 33] leads to compromised performance. Concretely, the performance of **CON** with task-agnostic task domain achieves an in-between or inferior performance compared to the models

| SSD-backbone | DAG [57] | | RAP [23] | |
|--------------|-----|------|-----|------|
| | STD | ours | STD | ours |
| VGG16 | 0.3 | 28.5 | 6.6 | 44.9 |
| ResNet50 | 0.4 | 22.9 | 8.8 | 39.1 |
| DarkNet53 | 0.5 | 26.2 | 8.2 | 46.6 |

Table 3. Evaluation results on across different backbones.

with individual task domains under different attacks, implying that simply mixing the task domains leads to compromised performance, due to the conflicts between the task gradients (Sec. 4.1). On the other hand, the robust model **MTD** using adversarial training with task oriented domain constraint can improve the performance over **CON** baseline. More importantly, when the task-oriented multi-task domain is incorporated, a proper trade-off and overall performance is observed compared with the single domain-based methods, implying the importance of properly handling heterogeneous and possibly imbalanced tasks in object detectors. In summary, the tasks could be imbalanced and contribute differently to the model robustness. As it is *unknown a priori* which is better, randomly adopting one or simply combining the losses (**CON**) could lead to compromised performance. **MTD** setting overcomes this issue and achieves performance on par or better than best single domain models and the task-agnostic domain model.

### 5.4. Defense against Existing White-box Attacks

To further investigate the model robustness, we evaluate models against representative attack methods from literature. We use **DAG** [57] and **RAP** [23] as representative attacks according to Table 1. It is important to note that the attack used in training and testing are different. The results are summarized in Table 2. It is observed that the performances of robust models improve over the standard model by a large margin. **CLS** performs better in general than **LOC** and **CON** in terms of robustness against the two attacks from literature. The model using multi-task domains (**MTD**) demonstrates the best performance. **MTD** has a higher clean image accuracy than **CLS** and performs uniformly well against different attacks, thus overall is better and will be used for reporting performance in the following. Visualization of example results are provided in Figure 8.

### 5.5. Evaluation on Different Backbones

We evaluate the effectiveness of the proposed approach under different SSD backbones, including VGG16 [45], ResNet50 [19] and DarkNet53 [39]. Average performance under **DAG** [57] and **RAP** [23] attacks are reported in Table 3. It is observed that the proposed approach can boost the performance of the detector by a large margin (20%~30% absolute improvements), across different backbones, demonstrating that the proposed approach performs well across backbones of different network structures with clear and consistent improvements over baseline models.
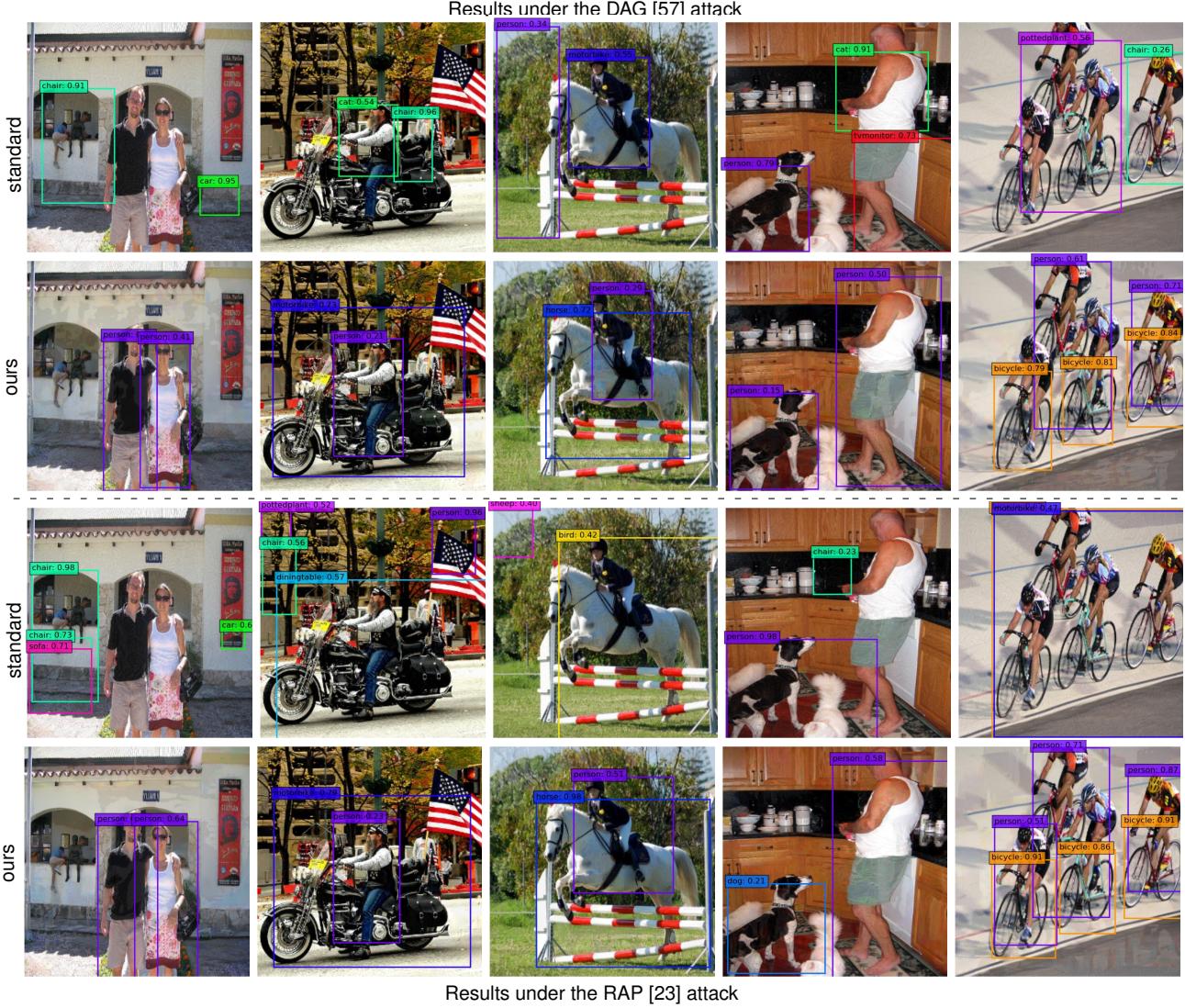
Figure 8. **Visual comparison** between **standard** model and **ours** under DAG [57] and RAP [23] attacks with attack budget 8.

| architecture | | DAG [57] | | RAP [23] | |
| --- | --- | --- | --- | --- | --- |
| | | STD | ours | STD | ours |
| SSD | +VGG16 | 0.3 | 28.5 | 6.6 | 44.9 |
| RFB | +ResNet50 | 0.4 | 27.4 | 8.7 | 48.7 |
| FSSD | +DarkNet53 | 0.3 | 29.4 | 7.6 | 46.8 |
| YOLO | +DarkNet53 | 0.1 | 27.6 | 8.1 | 44.3 |

Table 4. Evaluation results on different detection architectures.

## 5.6. Results on Different Detection Architectures

Our proposed approach is also applicable to different detection architectures. To show this, we use different detection architectures, including SSD [29], RFB [28], FSSD [24] and YOLO-V3 [40, 41]. The input image size for YOLO is $416 \times 416$ and all others take $300 \times 300$ images as input. Average performance under **DAG** [57] and **RAP** [23] attacks are summarized in Table 4. It is observed that the proposed method can improve over the standard method significantly and consistently for different detector architectures. This clearly demonstrates the applicability of the proposed approach across detector architectures.

## 5.7. Defense against Transferred Attacks

We further test the performance of the robust models under transferred attacks: attacks that are transferred from models with different backbones and/or detection architectures. Our model under test is based on SSD+VGG16. For attacks transferred from different backbones, they are generated under the SSD architecture but replacing the VGG backbone with ResNet or DarkNet. For attacks transferred from different detection architectures, we use RFB [28], FSSD [24] and YOLO [40, 41].[2] **DAG** [57] and **RAP** [23] are used as the underlining attack generation algorithms. The results are summarized in Table 5. It is observed that the

---

[2]As the input image size for YOLO is 416×416, which is different from input size of $300 \times 300$ for SSD, we insert a differentiable interpolation module ($300^2 \rightarrow 416^2$) between the input with size of 300×300 and YOLO.

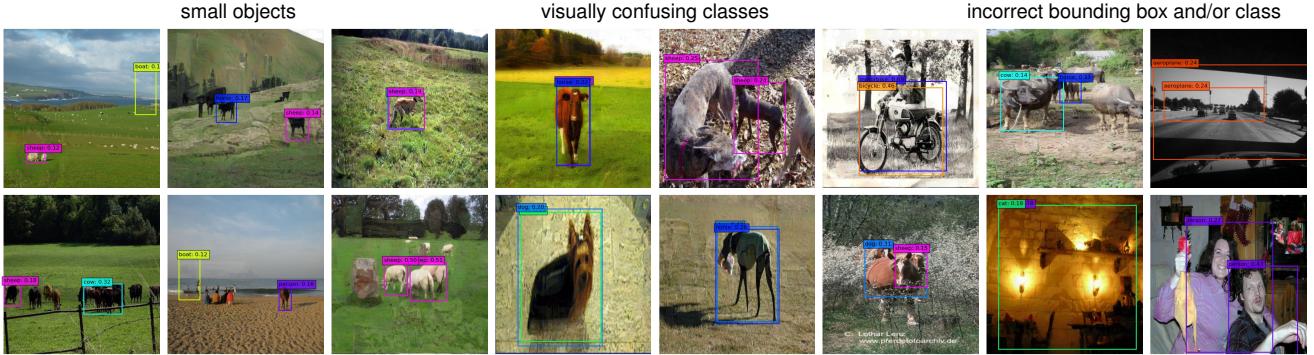small objects  visually confusing classes  incorrect bounding box and/or class

Figure 9. **Visualization of failure cases**. Example challenging cases include images with small objects and visually confusing classes.

| transferred attack | **DAG** [57] | **RAP** [23] | average |
|---|---|---|---|
| SSD+ResNet50 | 49.3 | 49.4 | 49.4 |
| SSD+DarkNet53 | 49.2 | 49.4 | 49.3 |
| RFB+ResNet50 | 49.1 | 49.3 | 49.2 |
| FSSD+DarkNet53 | 49.3 | 49.2 | 49.3 |
| YOLO+DarkNet53 | 49.5 | 49.5 | 49.5 |

Table 5. Performance of our model (SSD+VGG16) against attacks transferred from different backbones and detector architectures.

| model | architec. | backbone | clean | attack |
|---|---|---|---|---|
| **standard** | SSD | VGG16 | 39.8 | 2.8 |
| **ours** | SSD | VGG16 | 27.8 | 16.5 |
| | SSD | DarkNet53 | 20.9 | 18.8 |
| | SSD | ResNet50 | 18.0 | 16.4 |
| | RFB | ResNet50 | 24.7 | 21.6 |
| | FSSD | DarkNet53 | 23.5 | 20.9 |
| | YOLO | DarkNet53 | 24.0 | 21.5 |

Table 6. Comparison of standard and robust models on MS-COCO under RAP attack [23] with attack budget 8 and 20 PGD steps.

proposed model is robust against transferred attacks generated with different algorithms and architectures. It is also observed that the attacks have a certain level of robustness can be transferred across detectors with different backbones or structures. This reconfirms the results from [57, 23].

### 5.8. Results on MS-COCO

We further conduct experiments on MS-COCO [27], which is more challenging both for the standard detector as well as the defense due to its increased number of classes and data variations. The results of different models under RAP attack [23] with attack budget 8 and PGD step 20 are summarized in Table 6. The standard model achieves a very low accuracy in the presence of attack (compared with ∼40% on clean images). Our proposed models improves over the standard model significantly and performs generally well across different backbones and detection architectures. This further demonstrates the effectiveness of the proposed approach on improving model robustness.

### 5.9. Failure Case Analysis

We visualize in Figure 9 some example cases that are challenging to our current model. Images with small objects that are challenging for the standard detectors [29, 40] remain to be one category of challenging examples for robust detectors. Better detector architectures might be necessary to address this challenge. Another challenging category is objects with visually confusing appearance, which naturally leads to low confidence predictions. This is more related to the classification task of the detector and can benefit from advances in classification [58]. There are also cases where the predictions are inaccurate or completely wrong, which reveals the remaining challenges in robust detector training.

## 6. Conclusions

We have presented an approach for improving the robustness object detectors against adversarial attacks. From a multi-task view of object detection, we systematically analyzed existing attacks for object detectors and the impacts of individual task component on model robustness. An adversarial training method for robust object detection is developed based on these analyses. Extensive experiments have been conducted on PASCAL-VOC and MS-COCO datasets and experimental results have demonstrated the efficacy of the proposed approach on improving model robustness compared with the standard model, across different attacks, datasets, detector backbones and architectures.

This work serves as an initial step towards adversarially robust detector training with promising results. More efforts need to be devoted in this direction to address the remaining challenges. New advances on object detection can be used to further improve the model performance, *e.g.*, better loss function for approximating the true objective [26] and different architectures for addressing small object issues [7, 13]. Similarly, as a component task of object detection, any advances on classification task could be potentially transferred as well [58]. There is also a trade-off between accuracy on clean image and robustness for object detection as in the classification case [51]. How to leverage this trade-off better is another future work. Furthermore, by viewing object detection as an instance of multi-task learning task, this work could serve as an example on robustness improvement for other multi-task learning problems as well [20, 59].

# References

[1] A. Athalye, N. Carlini, and D. Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In *International Conference on Machine learning*, 2018.

[2] B. Biggio and F. Roli. Wild patterns: Ten years after the rise of adversarial machine learning. In *ACM Conference on Computer and Communications Security*, 2018.

[3] Z. Cai and N. Vasconcelos. Cascade R-CNN: Delving into high quality object detection. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.

[4] N. Carlini and D. Wagner. Towards evaluating the robustness of neural networks. In *IEEE Symposium on Security and Privacy*, 2017.

[5] R. Caruana. Multitask learning. *Machine Learning*, 28(1):41–75, 1997.

[6] S. Chen, C. Cornelius, J. Martin, and D. H. Chau. ShapeShifter: Robust physical adversarial attack on Faster R-CNN object detector. *CoRR*, abs/1804.05810, 2018.

[7] L. Cui. MDSSD: Multi-scale deconvolutional single shot detector for small objects. *CoRR*, abs/1805.07009, 2018.

[8] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2005.

[9] D. Erhan, C. Szegedy, A. Toshev, and D. Anguelov. Scalable object detection using deep neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2014.

[10] M. Everingham, S. M. Eslami, L. Gool, C. K. Williams, J. Winn, and A. Zisserman. The pascal visual object classes challenge: A retrospective. *Int. J. Comput. Vision*, 111(1):98–136, 2015.

[11] K. Eykholt, I. Evtimov, E. Fernandes, B. Li, A. Rahmati, F. Tramèr, A. Prakash, T. Kohno, and D. Song. Physical adversarial examples for object detectors. *CoRR*, abs/1807.07769, 2018.

[12] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan. Object detection with discriminatively trained part-based models. *IEEE Trans. Pattern Anal. Mach. Intell.*, 32(9):1627–1645, 2010.

[13] C.-Y. Fu, W. Liu, A. Ranga, A. Tyagi, and A. C. Berg. DSSD: Deconvolutional single shot detector. *CoRR*, abs/1701.06659, 2017.

[14] R. Girshick. Fast R-CNN. In *IEEE International Conference on Computer Vision*, 2015.

[15] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2014.

[16] I. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples. In *International Conference on Learning Representations*, 2015.

[17] C. Guo, M. Rana, M. Cissé, and L. van der Maaten. Countering adversarial images using input transformations. In *International Conference on Learning Representations*, 2018.

[18] K. He, R. B. Girshick, and P. Dollár. Rethinking ImageNet pre-training. *CoRR*, abs/1811.08883, 2018.

[19] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2016.

[20] A. Kendall, Y. Gal, and R. Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.

[21] A. Kurakin, I. Goodfellow, and S. Bengio. Adversarial machine learning at scale. In *International Conference on Learning Representations*, 2017.

[22] Y. Li, X. Bian, and S. Lyu. Attacking object detectors via imperceptible patches on background. *CoRR*, abs/1809.05966, 2018.

[23] Y. Li, D. Tian, M. Chang, X. Bian, and S. Lyu. Robust adversarial perturbation on deep proposal-based models. In *British Machine Vision Conference*, 2018.

[24] Z. Li and F. Zhou. FSSD: feature fusion single shot multibox detector. *CoRR*, abs/1712.00960, 2017.

[25] F. Liao, M. Liang, Y. Dong, and T. Pang. Defense against adversarial attacks using high-level representation guided denoiser. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.

[26] T.-Y. Lin, P. Goyal, R. B. Girshick, K. He, and P. Dollár. Focal loss for dense object detection. In *International Conference on Computer Vision*, 2017.

[27] T.-Y. Lin, M. Maire, S. J. Belongie, L. D. Bourdev, R. B. Girshick, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft COCO: Common objects in context. In *European Conference on Computer Vision*, 2014.

[28] S. Liu, D. Huang, and a. Wang. Receptive field block net for accurate and fast object detection. In *European Conference on Computer Vision*, 2018.

[29] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg. SSD: Single shot multibox detector. In *European Conference on Computer Vision*, 2016.

[30] X. Liu, M. Cheng, H. Zhang, and C.-J. Hsieh. Towards robust neural networks via random self-ensemble. In *European Conference on Computer Vision*, 2018.

[31] X. Liu, H. Yang, L. Song, H. Li, and Y. Chen. DPatch: Attacking object detectors with adversarial patches. *CoRR*, abs/1806.02299, 2018.

[32] J. Lu, H. Sibai, and E. Fabry. Adversarial examples that fool detectors. *CoRR*, abs/1712.02494, 2017.

[33] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*, 2018.

[34] D. Meng and H. Chen. MagNet: a two-pronged defense against adversarial examples. In *ACM SIGSAC Conference on Computer and Communications Security*, 2017.

[35] J. H. Metzen, T. Genewein, V. Fischer, and B. Bischoff. On detecting adversarial perturbations. In *International Conference on Learning Representations*, 2017.

[36] S.-M. Moosavi-Dezfooli, A. Fawzi, and P. Frossard. DeepFool: a simple and accurate method to fool deep neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2016.

[37] A. Nguyen, J. Yosinski, and J. Clune. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2015.

[38] A. Prakash, N. Moran, S. Garber, A. DiLillo, and J. Storer. Deflecting adversarial attacks with pixel deflection. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.

[39] J. Redmon. Darknet: Open source neural networks in C. http://pjreddie.com/darknet/, 2013–2016.

[40] J. Redmon, S. K. Divvala, R. B. Girshick, and A. Farhadi. You only look once: Unified, real-time object detection. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2016.

[41] J. Redmon and A. Farhadi. YOLOv3: An incremental improvement. *CoRR*, abs/1804.02767, 2018.

[42] S. Ren, K. He, R. B. Girshick, and J. Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *Advances in Neural Information Processing Systems*, 2015.

[43] A. Rosenfeld and M. Thurston. Edge and curve detection for visual scene analysis. *IEEE Trans. Comput.*, 20(5):562–569, 1971.

[44] P. Samangouei, M. Kabkab, and R. Chellappa. Defense-GAN: Protecting classifiers against adversarial attacks using generative models. In *International Conference on Learning Representations*, 2018.

[45] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015.

[46] Y. Song, T. Kim, S. Nowozin, S. Ermon, and N. Kushman. Pixeldefend: Leveraging generative models to understand and defend against adversarial examples. In *International Conference on Learning Representations*, 2018.

[47] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2015.

[48] C. Szegedy, S. E. Reed, D. Erhan, and D. Anguelov. Scalable, high-quality object detection. *CoRR*, abs/1412.1441, 2014.

[49] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus. Intriguing properties of neural networks. In *International Conference on Learning Representations*, 2014.

[50] F. Tramèr, A. Kurakin, N. Papernot, D. Boneh, and P. McDaniel. Ensemble adversarial training: Attacks and defenses. In *International Conference on Learning Representations*, 2018.

[51] D. Tsipras, S. Santurkar, L. Engstrom, A. Turner, and A. Madry. Robustness may be at odds with accuracy. In *International Conference on Learning Representations*, 2019.

[52] G. Tsoumakas and I. Katakis. Multi label classification: An overview. 3(3):1–13, 2007.

[53] J. Uijlings, K. van de Sande, T. Gevers, and A. Smeulders. Selective search for object recognition. *International Journal of Computer Vision*, 104(2):154–171, 2013.

[54] P. Viola and M. J. Jones. Robust real-time face detection. *Int. J. Comput. Vision*, 57(2):137–154, 2004.

[55] X. Wei, S. Liang, X. Cao, and J. Zhu. Transferable adversarial attacks for image and video object detection. *CoRR*, abs/1811.12641, 2018.

[56] C. Xie, J. Wang, Z. Zhang, Z. Ren, and A. Yuille. Mitigating adversarial effects through randomization. In *International Conference on Learning Representations*, 2018.

[57] C. Xie, J. Wang, Z. Zhang, Y. Zhou, L. Xie, and A. Yuille. Adversarial examples for semantic segmentation and object detection. In *International Conference on Computer Vision*, 2017.

[58] C. Xie, Y. Wu, L. van der Maaten, A. Yuille, and K. He. Feature denoising for improving adversarial robustness. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2019.

[59] X. Zhao, H. Li, X. Shen, X. Liang, and Y. Wu. A modulation module for multi-task learning with applications in image retrieval. In *European Conference on Computer Vision*, 2018.