

DEFENDING AGAINST PHYSICAL ADVERSARIAL PATCH ATTACKS ON INFRARED HUMAN DETECTION

Lukas Strack^{1*}, Futa Waseda^{2,3*}, Huy H. Nguyen³, Yinqiang Zheng², and Isao Echizen^{2,3}

¹University of Freiburg, Germany

²The University of Tokyo, Japan

³National Institute of Informatics, Japan

*These authors contributed equally

ABSTRACT

Infrared detection is an emerging technique for safety-critical tasks owing to its remarkable anti-interference capability. However, recent studies have revealed that it is vulnerable to physically-realizable adversarial patches, posing risks in its real-world applications. To address this problem, we are the first to investigate defense strategies against adversarial patch attacks on infrared detection, especially human detection. We propose a straightforward defense strategy, patch-based occlusion-aware detection (POD), which efficiently augments training samples with random patches and subsequently detects them. POD not only robustly detects people but also identifies adversarial patch locations. Surprisingly, while being extremely computationally efficient, POD easily generalizes to state-of-the-art adversarial patch attacks that are unseen during training. Furthermore, POD improves detection precision even in a clean (i.e., no-attack) situation due to the data augmentation effect. Our evaluation demonstrates that POD is robust to adversarial patches of various shapes and sizes. The effectiveness of our baseline approach is shown to be a viable defense mechanism for real-world infrared human detection systems, paving the way for exploring future research directions.

Index Terms— Infrared human detection, Adversarial patch, Adversarial defense

1. INTRODUCTION

Computer vision models based on deep neural networks (DNNs) exhibit impressive performance across diverse applications. However, they are vulnerable to adversarial examples [1, 2], i.e., maliciously manipulated inputs designed to deceive DNNs, thereby posing huge risks in real-world applications due to unintended network behaviors. Moreover, adversarial examples can be implemented in the physical world [3, 4, 5], such as through the use of adversarial patches [4, 6]. Unlike adversarial examples that perturb entire images with imperceptible noise [1, 2], adversarial patches alter specific regions of images with salient perturbations. Thus, these patches can be physically implemented, e.g., printed out, and directly attached to real-world objects to deceive DNN-based computer vision models.

While most of the studies on adversarial patches have focused on RGB-based computer vision [7, 8, 9], recent investigations have revealed that infrared object detection models are also vulnerable to physically-realizable adversarial patches [10, 11, 12, 13]. Infrared object detection is an emerging technique for safety-critical applications due to its remarkable anti-interference capability even in harsh

environments. It has been applied to a wide range of tasks, such as infrared pedestrian detection [14]; however, its vulnerability to physical adversarial patches raises substantial concerns regarding its reliability. Zhu et al. [10] showed that a carefully designed physical board with small bulbs could greatly degrade the precision of an infrared human detector. Zhu et al. [11] designed adversarial “QR code” pattern clothing that can fool an infrared human detector. Wei et al. [12] introduced the HOTCOLD Block (HCB) physical attack against thermal infrared imaging that uses wearable “warming and cooling pastes” to create infrared adversarial patches with a less conspicuous design. Wei et al. [13] subsequently enhanced the effectiveness of this attack by optimizing the shape and location of the patches.

However, effective defense strategies against infrared adversarial patches remain fully unexplored. In fact, the proposed infrared adversarial patches have not been tested against proper defenses specifically designed for adversarial patches, leaving the full extent of the risks unclear. For example, the HCB attack [12] was tested on detectors without any defense, and the Shape-Loc attack [13] was evaluated on defenses for L_p -bounded adversarial perturbations, but not defenses designed for adversarial patches.

To this end, we are the first to investigate defense strategies against adversarial patch attacks on infrared detection, focusing on human detection. We propose a computationally efficient yet effective defense method named patch-based occlusion-aware detection (POD), which efficiently augments training samples with random patches and subsequently detects them. These augmented samples simulate occlusions, enabling the model to handle scenarios where human bodies may be obscured or hidden in an attack scenario. Intriguingly, despite its simplicity, POD demonstrates generalizability to state-of-the-art infrared adversarial patch attacks, even though the patch attacks were unseen during training. Additionally, evaluation using state-of-the-art infrared adversarial patches optimized for shape and location [13] demonstrated that POD is robust to various shapes and sizes of adversarial patches. Furthermore, in contrast to typical adversarial training [2], which sacrifices accuracy in clean (i.e., no-attack) situations, POD instead improves detection precision in clean situations.

Our study highlights that state-of-the-art infrared patch attacks are not as effective as previously believed, as our straightforward data-augmentation-based defense strategy proved highly effective in countering them. In other words, our findings demonstrate that crafting strong, physically realizable infrared patches remains challenging.

The contribution of this work is summarized as follows:

- We are the first to investigate defense strategies against physical adversarial patch attacks on infrared detection, opening the door to reducing the vulnerability of infrared detection.

This work was partially supported by JSPS KAKENHI Grants JP21H04907 and JP24H00732, and by JST CREST Grants JPMJCR18A6 and JPMJCR20D3, and by JST AIP Acceleration Grant JPMJCR24U3 Japan.

- We devise a computationally efficient yet effective defense method, patch-based occlusion-aware detection (POD), which efficiently augments training samples with random patches and subsequently detects them.
- POD is easily generalizable to state-of-the-art infrared adversarial patch attacks that are unseen during training, is robust to adversarial patches of various shapes and sizes, and improves detection precision in clean situations.
- We highlight that state-of-the-art infrared patch attacks are not as effective as previously believed. Thus, our work encourages future research to evaluate infrared patch attacks against proper defenses, promoting a better understanding of the associated risks in realistic scenarios.

2. PHYSICAL ADVERSARIAL PATCHES

In this section, we first explain the general objective of generating adversarial patches. Then, we describe the constraints when crafting physical adversarial patches that interfere with infrared human detection.

2.1. General Framework

Let f_θ represent the object detector parameterized by θ , and $f_\theta(x)$ be the predicted output for an infrared image $x \in \mathbb{R}^{h \times w}$. We define $A(p, x)$ as the operator for applying adversarial patch p to image x (assuming a fixed patch location for simplicity).

Here, the attacker aims to optimize a single “universal” adversarial patch that is scene-agnostic, intending to deceive the human detector f_θ , irrespective of the pose or image context. The objective in generating an adversarial patch can thus be formulated as,

$$\operatorname{argmax}_p \mathbf{E}_{(x,y) \sim D} [J(f_\theta(A(p, x)), y)] \quad (1)$$

where (x, y) is an input-label pair from the data distribution D , and $J(\cdot, \cdot)$ is an arbitrary loss function that quantifies the discrepancy between the ground truth and the prediction. The core concept is to optimize a single patch p that maximizes the loss function across the entire data distribution, inducing inaccurate object detection.

More specifically, when attacking an object detector, the loss function $J(\cdot, \cdot)$ can be computed for either “objectness” or “classification” scores of an object detector; The object detector outputs an objectness score (indicating the likelihood of a bounding box containing an object) and a classification score (identifying the object’s class within the bounding box). Thys et al. [9] has shown that minimizing the objectness score is more effective than inducing misclassification, and the state-of-the-art infrared adversarial patches, such as HCB attack [12] and the Shape-Loc attack [13], follows the approach.

2.2. Constraints

In addition to the above objective, it is essential to consider constraints during the optimization process, particularly for infrared adversarial patches in physical space. There are two key constraints to be considered in the infrared scenario: (1) infrared images have much less texture information than RGB images due to overlap in the spectral reflectance across different materials; (2) physically implementing a high-resolution patch or one with a pixel value of a specific magnitude is challenging due to material limitations.

Constraints are thus added during optimization to obtain physically-realizable patches for interfering with infrared detection [10, 11, 12,

13]. For example, with the HCB attack [12], only nine grids are used for a patch with binary values (0/1). Wei et al. [13] optimized patches with only binary values while promoting the clustering of pixels with the same value for easy real-world implementation.

This limitation of low-resolution and simple-colored infrared adversarial patches raises the question, “Are state-of-the-art infrared patches effective against properly defended models?”. This motivated our investigation into the straightforward defense mechanism POD, detailed in the next section.

3. PATCH-BASED OCCLUSION-AWARE DETECTION

Our proposed defense mechanism, POD, is simple yet highly effective against infrared adversarial patch attacks. Our key idea involves augmenting training samples with random patches and training an object detection model. We introduce a “patch” class in addition to the “human” class for the model to detect these patches.

POD has three key characteristics:

- **Generalization without attack algorithm assumption:** The POD model is robust against sophisticated adversarial patches unseen during training despite the use of simple augmentation techniques.
- **Robust human detection:** POD not only robustly detects people in the presence of adversarial patches but also is trained to identify them. This results in improved generalization against previously unseen adversarial patches.
- **Simple quick training:** The training time is much less than that of traditional adversarial training schemes [15], which aims to solve the min-max problem by directly feeding adversarial examples during training.

The overall pipeline is illustrated in Fig. 1.

3.1. Adding Random Patch-based Occlusions

We introduce a simple training strategy to train a robust human detection model, i.e., apply random patches to simulate occlusions, which can obscure or hide people. Unlike conventional adversarial training schemes [15, 16] that augment training samples using specific attack algorithms, we train our detection model across simple yet diverse occlusion scenarios, resulting in general resilience against unforeseen adversarial patches. Our strategy harnesses the power of simplicity and randomness, inspired by TrivialAugment strategy [17], which achieved state-of-the-art performances by simplifying over-complicated augmentations for image classification pipelines. Our approach is highly inspired by this approach, realizing a simple yet powerful defense based on randomness. Additionally, considering that physical infrared adversarial patches typically exhibit basic textures and simple colors, as explained in Sec. 2.2, it is expected that a simple and efficient defense method is particularly effective in defending against them.

The patches for our training strategy are designed to exhibit random characteristics, such as size, shape, texture, and placement, to simulate realistic scenarios in which adversarial patches might be encountered. The patches are created by first cropping a square region from an original image, with their size, rotation degree, and location being randomly determined. Next, the random patch is added to the cropped region to simulate an adversarial patch. We use three random patch variations:

- **Erase patch** p_{erase} simulates the most basic occlusion that hides persons. The patch has the pixel values of zero.

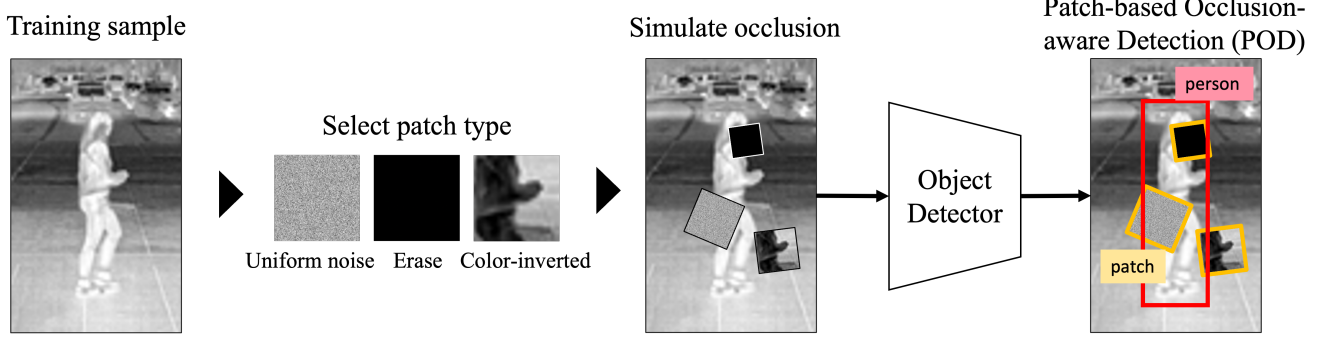


Fig. 1: Training pipeline of Patch-based Occlusion-aware Detection (POD). POD efficiently augments training samples with random patches and subsequently detects them. The augmented samples simulate occlusions, enabling the model to handle scenarios where human bodies may be obscured or hidden in attack scenarios.

- **Color-inverted patch** p_{inv} mimics an infrared adversarial patch with a more sophisticated texture. Patches with texture complexity similar to those captured by infrared cameras are efficiently mimicked through color inversion.
- **Uniform noise patch** p_{noise} mimics an infrared adversarial patch with an intricate texture, although physical implementation remains challenging. It applies a random value sampled from a uniform distribution to each pixel.

The entire augmentation process is outlined in Algorithm 1. The examples of augmented training samples are shown in Fig. 2.

Algorithm 1 POD Training Augmentation Process

Input: Image I with width w and height h

Output: Augmented image I'

Initialization: $I' \leftarrow I$, $N \leftarrow$ Random number of patches

```

1: for  $i$  in  $N$  do
2:   // Determine a patch
3:    $p \leftarrow$  Select patch type from  $\{p_{erase}, p_{inv}, p_{noise}\}$ .
4:    $l \leftarrow$  Random patch size.
5:    $\theta \leftarrow$  Random rotation degree of patch.
6:    $c_x, c_y \leftarrow$  Random center position of patch.
7:
8:   // Apply the patch to the image
9:    $I' \leftarrow$  apply patch  $p^{l \times l}$  with rotation degree  $\theta$  at  $(c_x, c_y)$ .
10: end for

```

3.2. Detecting Patch-based Occlusions

POD involves training the model not only to detect people but also to explicitly identify the existence and location of patch-based occlusions present in an image. The model is thus able to handle the uncertainties introduced by occlusions, which can be crucial for reliable decision-making in practical applications. Intriguingly, we found that explicitly training the model to identify the locations of patch-based occlusions noticeably improved its robustness against strong adversarial patches, as explained in a later section.

To detect adversarial patches, we modify an object detection model to have an additional patch class. Here, the types of patches are not distinguished, and only one extra class is used to detect patch-based occlusions.

Additionally, since the primary objective is to detect people robustly, we employ weighted cross-entropy loss [18] to prioritize the human class over the patch class:

$$\alpha \cdot L_{CE}^{human} + (1 - \alpha) \cdot L_{CE}^{patch} \quad (2)$$

where L_{CE}^{human} and L_{CE}^{patch} represent the classification loss with cross-entropy for humans and patches, respectively. This approach helps address the class imbalance between humans and patches after POD augmentation, enabling the model to prioritize and concentrate on human detection. In our experiments, the weight for the human class, denoted as α , was fixed at 0.9. Humans and patches are labeled as in Fig. 2.

4. EXPERIMENTAL SETUP

4.1. Dataset

Following Zhu et al. [10] and Wei et al. [12], we evaluated model performance on the Teledyne FLIR ADAS Thermal dataset [19]. Infrared images were acquired with a Teledyne FLIR Tau2 (13 mm f/1.0 with a 45-degree HFOV and 37-degree VFOV), and the thermal camera operated in T-linear mode. Following Wei et al. [12], we filtered the original dataset for the task of human detection. First, we used only the “person” category, excluding images without the person class. Next, we filtered the persons to ensure that the height of their bodies are larger than 120 pixels. This process resulted in 1169 training images with 1810 person labels and 84 test images with 154 person labels.

To evaluate the model against shape-location-optimized (“shape-loc”) attacks [13], we created a custom dataset based on the CVC-14 dataset [20], in which each image includes only one person. The CVC-14 dataset provides already cropped images of persons. From this dataset, we randomly chose 97 images and manually labeled them. The custom dataset, along with its annotations, is included in our provided codes.

4.2. Evaluation

We used average precision (AP@0.5) to evaluate the ability of object detection models to detect humans (we repeated each experiment five times and reported the standard deviation).

In our experiments, the physical adversarial patches were simulated in digital space. Since real-world attacks are less successful

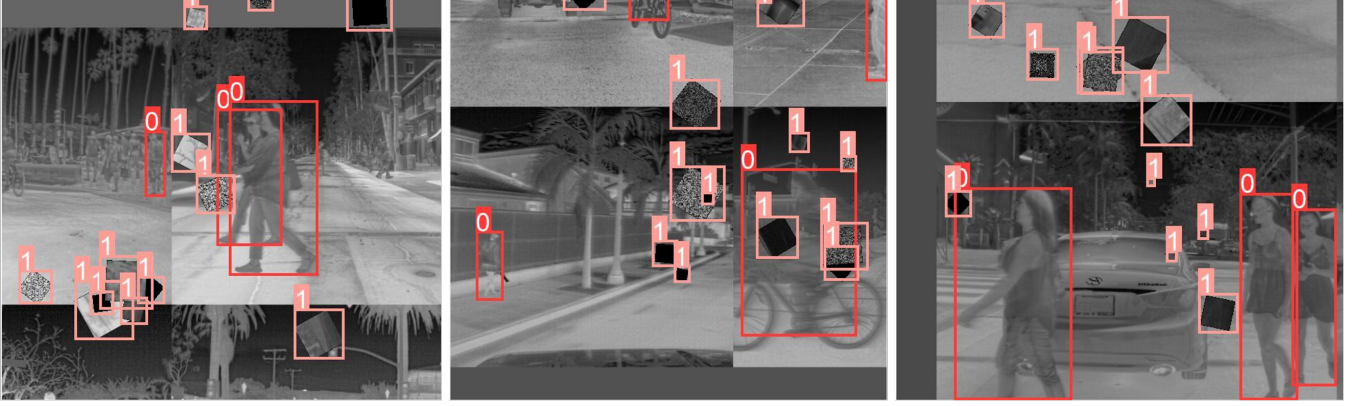


Fig. 2: Example training samples for patch-based occlusion-aware detection (POD). Labels “0” and “1” signify “human” and “patch” classes, respectively. POD aims not only to detect people with random occlusions but also to identify patch locations.

due to changes in lighting and human pose, as well as limitations on physical materials for crafting attacks, digital attacks are the upper bound for attack capability. To evaluate the robustness of the evaluated infrared detection models, we compared them in the following scenarios:

- **Random noise patch:** an adversarial patch in which each pixel value is sampled from a uniform distribution. This is a *digital-space attack* that simulates random occlusions.
- **AdvPatch** [9]: a universal adversarial patch designed to deceive RGB-based human detectors in the physical world. Due to its complex texture, it cannot be physically implemented in the infrared scenario; we regard this attack as a *strong digital-space attack*.
- **HCB** [12]: a state-of-the-art physically-realizable universal adversarial patch utilizing wearable warming and cooling pastes for attacks on infrared detectors.

Additionally, we evaluated the Shape-Loc attack [13] to evaluate the model’s performance against adversarial patches of different sizes and shapes. These patches are image-dependent, resulting in 97 individual optimized shapes and locations for our custom dataset, described in Sec. 4.1.

4.3. Implementation details

We used the YOLO-v5 model architecture [21] since models based on this architecture are fast and widely used as detectors. For infrared detection, we used the pre-trained weights on the RGB images from the MS COCO (Microsoft Common Objects in Context) dataset [22] as the initial weights, then fine-tuned the weights on infrared images from the FLIR ADAS Dataset. All of the evaluated detection models were trained for 50 epochs for fair comparison. Since infrared images have only one channel, we expanded their channel size to three dimensions to leverage the RGB-pre-trained weights.

4.4. Comparison between Adversarial Training Scheme

Furthermore, we developed a POD variant called “Adv-POD,” incorporating an adversarial training scheme that explicitly inputs adversarial patches during training, modified from Ad-YOLO [16]. Unlike Ad-YOLO, which is inadequate for defending against location-optimized patches like HCB and shape-loc due to its central patch placement, Adv-POD employs the same random patch placement

scheme as POD. This modification allows us to fairly compare the effects of using basic patches (Sec. 3.1) or adversarial patches.

The procedure for Adv-POD comprises four steps: (1) a model is trained using a history of adversarial patches; (2) a new adversarial patch is generated every 15 epochs, starting from the 5_{th} epoch; (3) the adversarial patches are stored in the adversarial patch history for subsequent training; (4) the stored patches are randomly selected and attached to input image to train the detection model. We used the AdvPatch attack [9] for adversarial patch generation, positioning the patches in the exact same way as for POD, i.e., randomly determining their sizes and locations for fair comparison.

5. RESULTS

In this section, we present experimental results demonstrating the effectiveness of POD against diverse adversarial patch attacks. The POD’s effectiveness is summarized in Table 3.

5.1. Evaluation for Digital and Physical Adversarial Patches

POD is robust to state-of-the-art physically realizable infrared adversarial patch. Table 1 shows that all POD variants are robust to the state-of-the-art physical infrared adversarial patch generated by HCB. Notably, the Average Precision on HCB is over 70%, which is only a few percent worse than the random-noise patch scenario; this suggests that the adversarial effect of HCB on the POD model is quite low and the prediction errors induced by HCB predominantly arise from simply concealing the human body.

POD generalizes to strong digital-space adversarial patch. Table 1 shows that POD is noticeably robust against the strong digital-space attack, i.e., AdvPatch [9], which was unseen during training. Given that the digital-space adversarial patches disregard material and sensor constraints for the infrared detection scenarios, as described in Sec 2.2, they are generally stronger than physical-space adversarial patches. Therefore, despite its simplicity, the POD strategy is suggested to generalize to unforeseen advanced infrared patches with intricate textures and colors.

Additionally, despite POD not being exposed to AdvPatch during training, it achieves a comparable Average Precision against AdvPatch as Adv-POD, a model explicitly trained with AdvPatch. This may be attributed to Adv-POD suffering from overfitting and failing to generalize to the test samples. It is well-known that typical adversarial training suffers from adversarial overfitting [23], wherein the

Method	Patch Type	Patch Detect.	Average Precision \uparrow				Training Time (rel.)
			Clean	\dagger Noise Patch	\dagger AdvPatch	HCB	
Std. Training	-	-	.9199 \pm .0064	.6436 \pm .0519	.2066 \pm .0465	.3616 \pm .1379	\times 1
POD _{noDet}	Rand.		.9193 \pm .0164	*.8410 \pm .0207	.5091 \pm .0472	.6582 \pm .0649	\times 1.0
POD	Rand.	✓	.9360 \pm .0082	*.8078 \pm .0170	.7897 \pm .0239	.7638 \pm .0374	\times 1.0
Adv-POD _{noDet}	Adv.		.9260 \pm .0149	.8744 \pm .0330	*.8133 \pm .0607	.7884 \pm .0372	\times 5.0
Adv-POD	Adv.	✓	.9395 \pm .0116	.8162 \pm .0249	*.7637 \pm .0120	.7426 \pm .0428	\times 5.0

Table 1: Average precision (AP@0.5) for adversarial patches. Model names with “*noDet*” indicate models without the patch-detection module described in Sec. 3.2. Attacks marked with \dagger are not physically realizable in the infrared scenario. Results marked with * are the results for seen attacks: the adversarial patch was encountered during training. The rightmost column illustrates the relative training time compared to standard training. The results highlight the effectiveness of our simple defense strategy, POD, in defending against previously unseen adversarial patches, with negligible additional training time. Notably, POD enhances detection precision in clean scenarios.

model overfits to worst-case adversarial images within the training samples. In contrast, POD avoids overfitting due to its simplicity and randomness of augmentations, resulting in generalization against unseen strong attacks.

An important property of POD is that the patch detection module is necessary to defend against AdvPatch; this indicates that explicitly identifying the location of patches during training helps a detector generalize against strong, unforeseen adversarial patches.

POD greatly enhanced robustness against random occlusions. Interestingly, we observed that a random-noise patch attack is already effective against the standard model due to simply concealing the person’s body, leading to a degradation in Average Precision from 92% to 64%. In contrast, all of the POD variants had an Average Precision of over 80%. This highlights POD’s effectiveness in improving resistance to random occlusions that can occur naturally without the presence of attackers.

POD achieved performance comparable to that of its adversarial variant, Adv-POD, yet with a notably shorter training time. The results in Table 1 demonstrate that POD and Adv-POD_{noDet} achieved similar performance in all scenarios, although the training time of POD is significantly shorter than that of Adv-POD_{noDet}. This is because, with the adversarial training scheme, it is time-consuming to generate adversarial examples during training, whereas POD relies solely on efficient augmentation techniques to simulate patch-based occlusions. Therefore, POD is an efficient and practical defense mechanism that can easily scale to train a model on a large dataset.

POD improved precision in the clean scenario. Notably, the POD variants did not sacrifice detection precision in the clean scenario but rather improved it, in contrast to conventional adversarial training schemes [2, 24] that often suffer from a trade-off between clean and robust accuracy [25]. This is attributed to the POD variants simulating various patch-based occlusions with diverse patch types with random size and location. We presume that the POD variants learn diverse features that are useful to robustly detect humans even when part of the body is occluded, and consequently, they improve accuracy in the clean scenario.

Adv-POD with the patch detection module overfit to a specific attack algorithm seen during training. While the patch detection module (refer to Sec. 3.2) was beneficial for the POD strategy, its use in Adv-POD, the adversarial variant of POD, led to a degradation in defending against adversarial patches. We hypothesize that this is due to overfitting in Adv-POD: a model can overfit when trained on a single attack that lacks patch diversity, and overfitting seems more likely to occur when the model is trained to detect the patch location explicitly.

5.2. Robustness to Varied Adversarial Patch Sizes and Shapes

We evaluated the models on the Shape-Loc attack [13], which uses shape-location-optimized infrared adversarial patches, to understand the model’s performance against patches with various sizes and shapes. Unlike universal adversarial patches, such as HCB, the Shape-Loc attack generates an image-dependent shape-location-optimized adversarial patch. We used the custom dataset based on the CVC-14 dataset [20] described in Sec. 4.1.

The results in Table 2 show that POD was robust to the Shape-Loc attack. Its performance is visualized in Fig. 3, which shows that POD effectively identified both people and Shape-Loc attack patches. This highlights POD’s resilience to adversarial patches with diverse sizes and shapes, even when the shapes and locations are adversarially optimized. Hence, our strategy of simulating random patches during training is a reasonably effective and practical approach for mitigating physical infrared adversarial patch attacks.

	Clean	Shape-Loc Attack
Std. Training	.9818 \pm .0071	.5005 \pm .1573
POD _{noDet}	.9845 \pm .0043	.9346 \pm .0114
POD	.9813 \pm .0059	.9133 \pm .0318
Adv-POD _{noDet}	.9799 \pm .0053	.9353 \pm .0208
Adv-POD	.9790 \pm .0059	.9090 \pm .0322

Table 2: Average precision for adversarial patches with various sizes and shapes used in the Shape-Loc attack [13]. All POD variants demonstrated improved robustness against Shape-Loc attacks, along with improved precision in clean scenarios.

	Defense Capability			Training Efficiency
	Adv.	HCB	Shape	
Std. Training				✓
POD _{noDet}	Δ	✓	✓	✓
POD	✓	✓	✓	✓
Adv-POD _{noDet}	✓	✓	✓	
Adv-POD	✓	✓	✓	

Table 3: Summary of defense capabilities of POD variants. POD, in the 3rd row, demonstrated the best defense capability along with efficient training.

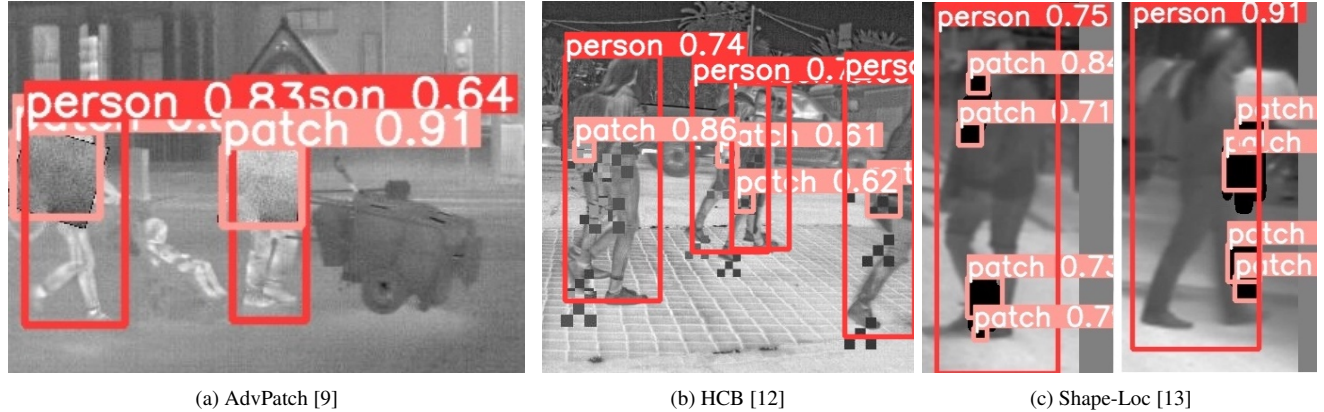


Fig. 3: Examples of POD detecting persons and patches for AdvPatch (left), HCB (middle), and shape-loc (right) attacks.

Method	Patch Augmentation Type			Average Precision \uparrow			
	Erase	Color-Inverted	Uniform Noise	Noise Patch	AdvPatch	HCB	Shape-Loc
Std. Training				.6436 \pm .0519	.2066 \pm .0465	.3616 \pm .1379	.5005 \pm .1573
POD	✓			.8142 \pm .0269	.4113 \pm .0332	.6919 \pm .0436	.8989 \pm .0193
		✓		.7931 \pm .0137	.7796 \pm .0235	.7858 \pm .0515	.9240 \pm .0208
			✓	.8191 \pm .0150	.7750 \pm .0332	.7287 \pm .0304	.9262 \pm .0267
	✓	✓		.7878 \pm .0191	.7736 \pm .0265	.6898 \pm .0593	.9090 \pm .0322
	✓		✓	.8063 \pm .0253	.7898 \pm .0465	.6694 \pm .0281	.8918 \pm .0324
(default)	✓	✓	✓	.7990 \pm .0162	.7812 \pm .0210	.7851 \pm .0141	.9181 \pm .0312
	✓	✓	✓	.8078 \pm .0170	.7897 \pm .0239	.7638 \pm .0374	.9133 \pm .0318

Table 4: Ablation study for patch type used in POD training. We assess Shape-Loc using the CVC-14 dataset, and conduct other attacks on the FLIR ADAS dataset. The results indicate that using all three types of patches results in good AP for various patch attacks.

5.3. Ablation Study: Patch Augmentation Type in POD

In this section, we conduct an ablation study of the patch augmentation types used in POD training. As explained in Sec. 3.1, POD simulates three types of occlusions during training: erase patch, color-inverted patch, and uniform noise patch. These three patches are intended to play different roles in simulating occlusions by having different characteristics, such as texture complexity.

Table 4 results indicate that employing all three patches in POD yields high average precision (AP) across various patch attacks, without conflicts in utilizing multiple patch types.

One finding is that using only the erase patch, which has constant pixel values of 0, led to lower AP compared to other POD variants when tested against state-of-the-art adversarial patches of AdvPatch, HCB, and Shape-Loc attack. We hypothesize that simulating occlusions with varied texture complexity is crucial to defend against adversarial patches, and the erase patches lacked this diversity.

Nevertheless, we observed that all variants of POD in Table 4 resulted in much better defense capability than the standard training model. The APs against the state-of-the-art physically realizable infrared adversarial patch attacks are quite high: all POD variants achieve around 70% for the HCB attack and around 90% for the Shape-Loc attack. This highlights our findings that the state-of-the-art infrared adversarial patch attacks are not as effective as previously believed, since our simple augmentation-based defense strategy greatly improves AP without increasing training time and without sacrificing precision in clean scenarios (even improves).

We believe that further efforts to optimize the augmentation

strategy with careful design can improve the performance of POD. Nevertheless, we would also like to emphasize that POD is designed to be extremely computationally efficient by utilizing only efficiently created patch types, and creating much more sophisticated patch-based occlusions may increase the training time.

6. CONCLUSION

In this work, we were the first to investigate defense strategies against physical adversarial patch attacks in infrared detection, focusing on human detection. We introduced a computationally efficient yet effective defense strategy, patch-based occlusion-aware detection (POD), which augments the training samples with random patches to simulate diverse occlusions and subsequently detects them. Interestingly, state-of-the-art infrared patch attacks were much less effective against our simple defense strategy, POD, challenging their perceived defense capability. POD demonstrated remarkable efficacy in countering infrared adversarial patches unseen during training, and exhibited robustness against patch attacks with varying shapes and sizes. Furthermore, POD improved detection precision in clean (i.e., no-attack) scenarios by learning diverse features to handle various occlusions.

In summary, our pioneering work highlights that state-of-the-art infrared patch attacks are not as effective as previously believed. We encourage future research to evaluate infrared patch attacks against proper defenses to understand the associated risks better.

7. REFERENCES

- [1] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus, “Intriguing properties of neural networks,” in *ICLR*, 2014.
- [2] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy, “Explaining and harnessing adversarial examples,” in *ICLR*, 2015.
- [3] Alexey Kurakin, Ian J Goodfellow, and Samy Bengio, “Adversarial examples in the physical world,” in *Artificial intelligence safety and security*, pp. 99–112. Chapman and Hall/CRC, 2018.
- [4] Tom B Brown, Dandelion Mané, Aurko Roy, Martín Abadi, and Justin Gilmer, “Adversarial patch,” *arXiv preprint arXiv:1712.09665*, 2017.
- [5] Dawn Song, Kevin Eykholt, Ivan Evtimov, Earlene Fernandes, Bo Li, Amir Rahmati, Florian Tramer, Atul Prakash, and Tadayoshi Kohno, “Physical adversarial examples for object detectors,” in *12th USENIX workshop on offensive technologies (WOOT 18)*, 2018.
- [6] Kevin Eykholt, Ivan Evtimov, Earlene Fernandes, Bo Li, Amir Rahmati, Chaowei Xiao, Atul Prakash, Tadayoshi Kohno, and Dawn Song, “Robust physical-world attacks on deep learning visual classification,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 1625–1634.
- [7] Kevin Eykholt, Ivan Evtimov, Earlene Fernandes, Bo Li, Amir Rahmati, Florian Tramer, Atul Prakash, Tadayoshi Kohno, and Dawn Song, “Physical adversarial examples for object detectors,” in *12th USENIX workshop on offensive technologies (WOOT 18)*, 2018.
- [8] Mark Lee and Zico Kolter, “On physical adversarial patches for object detection,” *arXiv preprint arXiv:1906.11897*, 2019.
- [9] Simen Thys, Wiebe Van Ranst, and Toon Goedemé, “Fooling automated surveillance cameras: adversarial patches to attack person detection,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, 2019, pp. 0–0.
- [10] Xiaopei Zhu, Xiao Li, Jianmin Li, Zheyao Wang, and Xiaolin Hu, “Fooling thermal infrared pedestrian detectors in real world using small bulbs,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2021, vol. 35, pp. 3616–3624.
- [11] Xiaopei Zhu, Zhanhao Hu, Siyuan Huang, Jianmin Li, and Xiaolin Hu, “Infrared invisible clothing: Hiding from infrared detectors at multiple angles in real world,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 13317–13326.
- [12] Hui Wei, Zhixiang Wang, Xuemei Jia, Yinqiang Zheng, Hao Tang, Shin’ichi Satoh, and Zheng Wang, “Hotcold block: Fooling thermal infrared detectors with a novel wearable design,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2023, vol. 37, pp. 15233–15241.
- [13] Xingxing Wei, Jie Yu, and Yao Huang, “Physically adversarial infrared patches with learnable shapes and locations,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 12334–12342.
- [14] Frédéric Suard, Alain Rakotomamonjy, Abdelaziz Bensrhair, and Alberto Broggi, “Pedestrian detection using infrared images and histograms of oriented gradients,” in *2006 IEEE Intelligent Vehicles Symposium*. IEEE, 2006, pp. 206–212.
- [15] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy, “Explaining and harnessing adversarial examples,” *arXiv preprint arXiv:1412.6572*, 2014.
- [16] Nan Ji, YanFei Feng, Haidong Xie, Xueshuang Xiang, and Naijin Liu, “Adversarial yolo: Defense human detection patch attacks via detecting adversarial patches,” *arXiv preprint arXiv:2103.08860*, 2021.
- [17] Samuel G. Müller and Frank Hutter, “Trivialaugment: Tuning-free yet state-of-the-art data augmentation,” 2021.
- [18] Trong Huy Phan and Kazuma Yamamoto, “Resolving class imbalance in object detection with weighted cross entropy losses,” *arXiv preprint arXiv:2006.01413*, 2020.
- [19] “Free - flir thermal dataset for algorithm training — teledyne flir,” .
- [20] Alejandro González, Zhijie Fang, Yainuvis Socarras, Joan Serrat, David Vázquez, Jiaolong Xu, and Antonio M. López, “Pedestrian detection at day/night time with visible and fir cameras: A comparison,” *Sensors*, vol. 16, no. 6, 2016.
- [21] Glenn Jocher et. al., “ultralytics/yolov5: v6.0 - YOLOv5n ‘Nano’ models, Roboflow integration, TensorFlow export, OpenCV DNN support,” Oct. 2021.
- [22] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick, “Microsoft coco: Common objects in context,” in *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*. Springer, 2014, pp. 740–755.
- [23] Leslie Rice, Eric Wong, and Zico Kolter, “Overfitting in adversarially robust deep learning,” in *International Conference on Machine Learning*. PMLR, 2020, pp. 8093–8104.
- [24] Aleksander Madry, Aleksandar Makelev, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu, “Towards deep learning models resistant to adversarial attacks,” *arXiv preprint arXiv:1706.06083*, 2017.
- [25] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan, “Theoretically principled trade-off between robustness and accuracy,” in *International conference on machine learning*. PMLR, 2019, pp. 7472–7482.