# Reinforced Embodied Active Defense: Exploiting Adaptive Interaction for Robust Visual Perception in Adversarial 3D Environments

Xiao Yang, Lingxuan Wu, Lizhong Wang, Chengyang Ying, Hang Su, and Jun Zhu, *Fellow, IEEE*

**Abstract**—Adversarial attacks in 3D environments have emerged as a critical threat to the reliability of visual perception systems, particularly in safety-sensitive applications such as identity verification and autonomous driving. These attacks employ adversarial patches and 3D objects to manipulate deep neural network (DNN) predictions by exploiting vulnerabilities within complex scenes. Existing defense mechanisms, such as adversarial training and purification, primarily employ passive strategies to enhance robustness. However, these approaches often rely on pre-defined assumptions about adversarial tactics, limiting their adaptability in dynamic 3D settings. To address these challenges, we introduce **Reinforced Embodied Active Defense (REIN-EAD)**, a proactive defense framework that leverages adaptive exploration and interaction with the environment to improve perception robustness in 3D adversarial contexts. By implementing a multi-step objective that balances immediate prediction accuracy with predictive entropy minimization, REIN-EAD optimizes defense strategies over a multi-step horizon. Additionally, REIN-EAD involves an uncertainty-oriented reward-shaping mechanism that facilitates efficient policy updates, thereby reducing computational overhead and supporting real-world applicability without the need for differentiable environments. Comprehensive experiments validate the effectiveness of REIN-EAD, demonstrating a substantial reduction in attack success rates while preserving standard accuracy across diverse tasks. Notably, REIN-EAD exhibits robust generalization to unseen and adaptive attacks, making it suitable for real-world complex tasks, including 3D object classification, face recognition and autonomous driving. By integrating proactive policy learning with embodied scene interaction, REIN-EAD establishes a scalable and adaptable approach for securing DNN-based perception systems in dynamic and adversarial 3D environments.

**Index Terms**—Adversarial Robustness, Active Defense, Embodied Learning, Policy Learning

✦

## 1 INTRODUCTION

ADVERSARIAL attacks in 3D environments have become a significant threat to the security and reliability of visual perception systems [1], [2]. These attacks utilize carefully crafted perturbations, such as adversarial patches and 3D objects, strategically placed in physical scenes to manipulate deep neural network (DNN) predictions [2], [3]. The consequences of such vulnerabilities are especially severe in safety-critical domains, such as identity verification [2], [4], [5] and autonomous driving [3], [6], where erroneous predictions can severely compromise system integrity. Therefore, ensuring robust perception under adversarial conditions is critical for deploying these systems reliably in real-world applications.

In response to these emergent threats, researchers have developed diverse defense strategies for enhancing the robustness of DNNs. Adversarial training [7], [8], [9] has emerged as a particularly effective approach [10], where adversarial examples are deliberately incorporated into the training data to enhance the model's resilience. Concurrently, input preprocessing techniques, such as adversarial purification, have been proposed to mitigate these perturbations [11], [12], [13]. However, these approaches predominantly belong to **passive defenses**, which exhibit vulnerability to unseen or adaptive attacks [14], [15] that circumvent existing robustness measures, due to the presuppositions regarding the adversary's approaches. Moreover, these strategies often neglect the intrinsic physical context and associated understanding of the scene and objects in the 3D realm, weakening these defenses in real-world physical environments.

In contrast to the limitations of passive defenses, human active vision employs iterative refinement and error-correction mechanisms to effortlessly detect misplaced or inconsistent elements in complex 3D environments [16], [17]. Inspired by human active vision, our recent seminal work proposed a novel defense framework of Embodied Active Defense (EAD), which incorporates a proactive policy network to replicate this dynamic process [18]. EAD actively contextualizes environmental information and leverages object consistency to address misaligned adversarial patches in 3D settings. The system continuously refines its scene understanding by integrating current and past observations, forming a more comprehensive representation of the environment. This proactive behavior enables the system to predict areas of uncertainty and adjust its actions accordingly, thereby improving the quality of the observations it collects. By synergizing proactive movement and iterative predictions, EAD enhances its scene comprehension and mitigates the impact of adversarial patches.

Despite its potential, the current EAD framework faces several critical challenges that limit its effectiveness and applicability in real-world scenarios. First, EAD's greedy in-

- X. Yang, L. Wu, L. Wang, C. Ying, H. Su and J. Zhu are with Dept. of Comp. Sci. & Tech., Institute for AI, BNRist Center, THBI Lab, Tsinghua-Bosch Joint Center for ML, Tsinghua University, Beijing, China. Email: yangxiao19@tsinghua.org.cn, {wlx23, ycy21}@mails.tsinghua.edu.cn, wanglizhong99@outlook.com, {suhangss, dcszj}@tsinghua.edu.cn. Xiao Yang and Lingxuan Wu had contributed equally to this work. Corresponding authors: Hang Su; Jun Zhu. Code is available at https://github.com/thu-ml/EmbodiedActiveDefense.

formative exploration strategy prioritizes immediate, single-step information gain over long-term relevance, resulting in temporally inconsistent actions. This short-sighted approach often leads the agent to revisit previously explored view-points, diminishing exploration efficiency and increasing the likelihood of erroneous predictions. Furthermore, EAD's reliance on training the proactive policy network with differentiable environment models introduces a misalignment, particularly when handling non-differentiable real-world physical dynamics, which limits the framework's practical applicability. Additionally, learning through differentiable simulations is computationally intensive and susceptible to numerical instabilities, further undermining the system's overall efficiency.

To address existing limitations and enhance the effectiveness, applicability, and efficiency of defense mechanisms, we propose the Reinforced EAD (REIN-EAD) framework. This framework enables agents to optimize for long-term outcomes by learning through trial-and-error interactions without requiring differentiable simulations. To address temporal inconsistencies in exploration, we introduce a generalized objective function that accumulates multi-step interactions. This objective balances prediction loss reduction with predictive entropy minimization over a multi-step horizon, allowing the system to account for temporal dependencies and prioritize long-term outcomes over immediate uncertainty reduction. Additionally, to eliminate differentiable constraints and ensure efficient convergence, we implement an uncertainty-oriented reward-shaping technique within reinforcement learning. This approach provides dense rewards at each step, guiding the agent to reduce perceptual uncertainty and minimize prediction errors. By fostering efficient updates, this reward structure supports stable convergence even in complex environments, enhancing the agent's adaptability to dynamic changes and improving overall performance in uncertain and evolving scenarios. Finally, to address overfitting and computational burden of adversarial patch generation, we present Offline Adversarial Patch Approximation (OAPA) for generating adversary-agnostic patches. OAPA systematically characterizes the manifold of adversarial patterns from diverse attack strategies. By distilling fundamental features across attack strategies offline, OAPA substantially improves generalization across diverse attacks while reducing the computational overhead of online adversarial training.

Extensive experiments demonstrate that REIN-EAD offers prominent advantages over conventional passive defense mechanisms. First, REIN-EAD consistently outperforms state-of-the-art defense techniques, achieving a remarkable 95% reduction in attack success rate across a diverse range of tasks. Notably, REIN-EAD maintains or even enhances standard accuracy by effectively leveraging instructive information suited for detecting target objects in dynamic 3D environments. Second, REIN-EAD demonstrates superior **generalization** compared to passive approaches. Its attack-agnostic strategies enable it to defend effectively against a broad spectrum of adversarial patches, including both unseen and adaptive attacks. Moreover, REIN-EAD exhibits strong **applicability** in complex and real-world scenarios, such as 3D object classification, face recognition and object detection for autonomous driving. The trial-and-error learning paradigm ensures stable and efficient policy updates,

making REIN-EAD highly adaptable to real-world tasks.

To summarize, our contributions are as follows:

- First, we propose REIN-EAD that integrates multi-step accumulative interactions and policy learning into a cohesive framework. It optimizes a multi-step objective with an uncertainty-oriented reward shaping, promoting temporally consistent and informative exploration.
- Second, we develop an adversary-agnostic defense strategy of OAPA, which enables REIN-EAD to defend against diverse adversarial patches, including unseen and adaptive attacks, without relying on specific assumptions about the adversary's capabilities.
- Third, through extensive experiments, we demonstrate REIN-EAD's superior **effectiveness** over state-of-the-art passive defenses across various settings, strong **generalization** against various unseen and adaptive attacks, and **adaptability** to complex real-world scenarios.

## 2 RELATED WORK

In this section, we delve into the threat posed by adversarial patches in 3D environments and explore the corresponding defensive strategies.

### 2.1 Adversarial Patches and Defenses

Adversarial patches, initially devised to manipulate specific regions of an image to mislead image classifiers [1], have significantly evolved. They now deceive a broad spectrum of perception models [2], [6], including those operating within 3D environments [3], [5], [19], [20].

To defend perception models against such adversarial patch attacks, various strategies have been proposed, ranging from empirical defenses [9], [13], [21], [22] to certified defenses [11], [23]. However, a critical review reveals that most contemporary defense mechanisms fall under what we term as **passive defenses**. These approaches rely on information derived from monocular observations and presupposed adversarial tactics to mitigate patch-based threats.

Within the passive defense paradigm, two primary approaches have gained prominence. Adversarial training [7], [8], [9] strengthens the perception model by exposing it to adversarial examples during training, boosting intrinsic robustness against adversarial perturbations in supervised learning [24], [25] and semi-supervised learning [26]. Alternatively, adversarial purification [12], [13], [22] integrates an auxiliary purifier into the perception pipeline. This purifier first identifies adversarial patches within observations and then neutralizes or removes these perturbations. The amended observations are subsequently processed by the model, yielding a two-stage defense pipeline that mitigates adversarial influences before perception tasks commence.

While passive defenses have shown success in specific scenarios, their dependence on static, pre-defined mechanisms renders them vulnerable to adaptive attacks and restricts their effectiveness in dynamic, real-world environments. In contrast, active defense strategies exhibit a more adaptable approach, dynamically responding to evolving adversarial threats in complex 3D settings.

Our recent work introduces a pioneering defense framework termed Embodied Active Defense (EAD), which employs a proactive policy network to simulate this dynamic

response process [18]. EAD actively contextualizes environmental information and leverages object consistency to address adversarial patch misalignments in 3D environments, marking a significant advancement in adversarial defense.

## 2.2 Embodied Perception

The paradigm of embodied perception [27], [28] represents a significant shift in the landscape of artificial intelligence and computer vision. This approach posits that perception is not merely a passive process of information reception, but rather an active, embodied experience where an agent can dynamically interact with and navigate its environment to optimize perceptual outcomes and enhance task performance.

The embodied perception framework has demonstrated remarkable versatility, finding applications across a diverse spectrum of computer vision tasks. In the domain of object detection, researchers have leveraged embodied agents to actively explore and analyze scenes, leading to more robust and context-aware detection systems [29], [30]. Similarly, in the field of 3D pose estimation, embodied approaches have enabled more accurate and adaptable solutions by allowing agents to actively seek optimal viewpoints for estimation [31]. Furthermore, the paradigm has proven invaluable in advancing 3D scene understanding, where embodied agents can navigate complex environments to build comprehensive spatial representations [32].

The power of embodied perception lies in its ability to mimic the active, exploratory nature of biological vision systems, allowing agents to overcome the limitations of static and passive perception. By dynamically adjusting their position, orientation, or focus in response to environmental cues and task demands, embodied agents can potentially gather more informative and less ambiguous sensory data, leading to improved performance across different perceptual tasks. Our novel integration of embodied perception with adversarial robustness opens up new avenues for research, potentially leading to more resilient and adaptable perception systems capable of maintaining high performance even in sophisticated adversarial attacks.

## 3 METHODOLOGY

We first introduce the Preliminary about Embodied Active Defense (EAD) in Sec. 3.1, which leverages recurrent feedback to counteract adversarial patches. Then, we provide a theoretical analysis of EAD in Sec. 3.2 for further exploration. In Sec. 3.3, we propose REIN-EAD that incorporates multi-step interactions and policy learning. It effectively enhances the adaptability and resilience of the defense mechanism in complex and real-world environments. Finally, we provide the adversary-agnostic defenses in Sec. 3.4.

### 3.1 Preliminary: Embodied Active Defense

Consider a scene $x \in \mathcal{X}$ with its associated ground-truth label $y \in \mathcal{Y}$. The perception model $f : \mathcal{O} \to \mathcal{Y}$ aims to predict the scene annotation $y$ based on the image observation $o_i \in \mathcal{O}$, where $o_i$ is derived from the scene $x$ and conditioned on the camera's state $s_i$ (encompassing camera's position and viewpoint). The function $\mathcal{L}(\cdot)$ represents a task-specific loss function, such as the cross-entropy loss.

Traditional passive defense strategies operate a single observation $o_i$ to counter adversarial patches, thus failing to capitalize on the rich contextual information obtainable through proactive exploration of the environment [33], [34]. Formally, an adversarial patch $p$ is introduced into the observation $o_i$, resulting in the erroneous prediction of the perception model $f$. The generation of adversarial patches in 3D scenes [3] typically optimizes:

$$\max_p \mathbb{E}_{s_i} \mathcal{L}(f(A(p, o_i; s_i)), y), \tag{1}$$

where $A(\cdot)$ projects the 3D adversarial patch $p$ to the 2D camera observations $o_i$ according to the camera's state $s_i$. This generalized 3D formulation encompasses the 2D case when $s_i$ is restricted to 2D transformations and patch locations, as in Brown *et al.* [1].

For subsequent analysis, we define the set of deceptive adversarial patches $\mathcal{P}_x$ for scene $x$ as:

$$\mathcal{P}_x = \{p \in [0,1]^{H_p \times W_p \times C} : \mathbb{E}_{s_i} f(A(p, o_i; s_i)) \neq y\}, \tag{2}$$

where $H_p$ and $W_p$ represent the height and width of the patch, respectively. In practice, approximating the solution set $\mathcal{P}_x$ involves employing specific optimization techniques [7], [35] to solve the problem presented in Eq. (1). This generalized formulation provides a robust foundation for analyzing the behavior and impact of adversarial patches in complex 3D environments.

**The EAD framework**. Our recent work proposes EAD [18], a paradigm that champions active scene engagement and iteratively leverages environmental feedback to enhance the robustness of perception systems against patch attacks. EAD comprises two recurrent models that emulate the intricate cerebral structure underpinning active human vision. The **perception model** $f(\cdot; \boldsymbol{\theta})$, parameterized by $\boldsymbol{\theta}$, is meticulously crafted to facilitate sophisticated visual perception by fully harnessing the rich contextual information embedded within temporal observations from the external world. At each timestep $t$, the model ingeniously leverages the current observation $o_t$ and amalgamates it with the prevailing internal belief $b_{t-1}$ regarding the scene, thus constructing an enhanced representation of the surrounding environment $b_t$ by a recurrent paradigm. Simultaneously, the perception model generates a scene annotation $\hat{y}_t$ as:

$$\{\hat{y}_t, b_t\} = f(o_t, b_{t-1}; \boldsymbol{\theta})^1. \tag{3}$$

The subsequent **policy model** $\pi(\cdot; \boldsymbol{\phi})$, parameterized by $\boldsymbol{\phi}$, serves to govern the visual control of movement. Formally, given the current collective environmental understanding $b_t$ meticulously sustained by the perception model, it derives the action $\boldsymbol{a}_t$ by sampling from the distribution $\pi(b_t; \boldsymbol{\phi})$.

To formally characterize the EAD framework's interaction and proactive exploration within the environment (as illustrated in Fig. 1), we extend the framework of the Partially-Observable Markov Decision Process (POMDP) [36]. The interaction process under the scene $x$ is denoted by $\mathcal{M}(x) := \langle \mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{O}, \mathcal{Z} \rangle$. Here, $\mathcal{S}$ and $\mathcal{A}$ represent the state and action spaces, respectively. For $\forall (s, a) \in \mathcal{S} \times \mathcal{A}$, the transition dynamic under the scene $x$ adheres to the Markovian property, satisfying $\mathcal{T}(\cdot \mid s, a, x)$. Due to the

---

1. we set $f_y(o_t, b_{t-1}; \boldsymbol{\theta}) = \hat{y}_t$ and $f_b(o_t, b_{t-1}; \boldsymbol{\theta}) = b_t$ respectively.
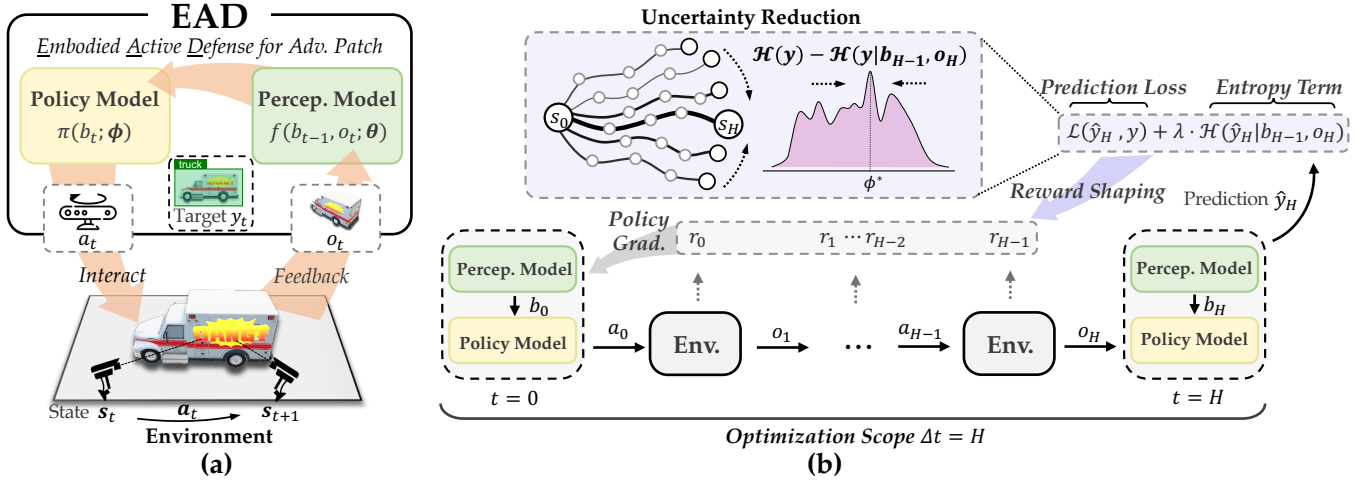
Fig. 1: An overview of EAD and REIN-EAD. (a) In EAD, the perception model refines the environment representation $b_t$ using observation $o_t$ and previous internal belief $b_{t-1}$, making task-specific prediction $y_t$. The policy model generates action $a_t$ based on $b_t$, minimizing perception uncertainty $H(y \mid b_{t-1}, o_t)$ over a single step. (b) REIN-EAD extends EAD by accumulating multi-step interactions, balancing prediction loss reduction and entropy minimization over horizon $H$. The policy is learned using model-free RL with dense rewards at each step, guiding the agent toward robust decision-making.

partially observed nature of the environment, the agent can not directly access the state $s$ but instead receives an observation $o$ sampled from the observation function $\mathcal{Z}(\cdot \mid s, x)$. At each timestep $t$, EAD obtains an observation $o_t$ based on the current state $s_t$. This observation $o_t$ serves as crucial environmental feedback, enabling the refinement of the agent's understanding of the environment $b_t$ through the sophisticated perception model $f(\cdot; \boldsymbol{\theta})$. The recurrent perception mechanism employed by EAD is important for maintaining the stability of human vision [16], [37]. Rather than remaining static and passively assimilating observations, EAD leverages the policy model to execute actions $a_t$ sampled from the distribution $\pi(b_t; \boldsymbol{\phi})$. The incorporation of the policy model enables EAD to determine the optimal action at each timestep, ensuring the acquisition of the most informative feedback from the scene.

**Training EAD against adversarial patches.** To equip the intricately constructed EAD model, we introduce a specialized learning algorithm designed for countering adversarial patches. Considering a data distribution $\mathcal{D}$ comprising paired data $(x, y)$, we examine adversarial patches $p \in \mathcal{P}_x$ generated using **unknown attack techniques** that corrupt the observation $o_t$ into $o_t' = A(o_t, p; s_t)$. The primary objective of EAD is to minimize the expected loss in the presence of adversarial patch threats. Consequently, the learning process of EAD to mitigate adversarial patches is formulated as an optimization problem of parameters $\boldsymbol{\theta}$ and $\boldsymbol{\phi}$ as:

$$\min_{\boldsymbol{\theta}, \boldsymbol{\phi}} \mathbb{E}_{(x,y) \sim \mathcal{D}, \tau \sim (\mathcal{M}(x), \pi), t} \left[ \sum_{p \in \mathcal{P}_x} \mathcal{L}(\hat{y}_t, y) \right],$$

with $\{\hat{y}_t, b_t\} = f(A(o_t, p; s_t), b_{t-1}; \boldsymbol{\theta}), \ a_t \sim \pi(\cdot \mid b_t; \boldsymbol{\phi})$

s.t. $o_t \sim \mathcal{Z}(\cdot \mid s_t, x), \quad s_t \sim \mathcal{T}(\cdot \mid s_{t-1}, a_{t-1}, x),$

where $\tau := (o_0, a_0, o_1, \ldots, o_H)$ represents the collected trajectory with length $H$ and the probability $p_{\mathcal{M}(x), \pi}(\tau) = \rho_0(s_0) \prod_{t=0}^{H} \mathcal{T}(s_{t+1} \mid s_t, a_t, x) \pi(a_t \mid b_t; \boldsymbol{\phi}) \mathcal{Z}(o_t \mid s_t)$; while $\hat{y}_t$ signifies the model's prediction at timestep $t$, which adheres to a uniform distribution over $\{0, 1, 2, \ldots, H\}$. Remarkably, the loss function $\mathcal{L}$ exhibits a task-agnostic nature,

emphasizing the exceptional adaptability of the EAD framework. This inherent flexibility guarantees that EAD delivers robust defenses across a wide spectrum of perception tasks.

## 3.2 Theoretical Analysis of EAD

To gain a deeper understanding of the model's behavior, we further examine a generalized instance of EAD in Eq. (4), where the agent employs the InfoNCE objective [38], simplified as:

$$\min_{\boldsymbol{\theta}, \boldsymbol{\phi}} \mathbb{E}_{(x^{(j)}, y^{(j)})} \left[ \frac{1}{K} \sum_{j=1}^{K} \log \frac{e^{-S(f_y(b_{t-1}^{(j)}, o_t^{(j)}; \boldsymbol{\theta}), y^{(j)})}}{\frac{1}{K} \sum_{k=1}^{K} e^{-S(f_y(b_{t-1}^{(j)}, o_t^{(j)}; \boldsymbol{\theta}), y^{(k)})}} \right],$$

$$(5)$$

where $(x^{(j)}, y^{(j)})_{j=1}^{K}$ denotes a data batch of size $K$ sampled from the distribution $\mathcal{D}$, and $S : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$ quantifies the similarity between the predicted scene annotation and the ground truth label. In embodied perception, this loss establishes a cross-modal correspondence between the observations and annotations like CLIP [39]. Moreover, we provide an information-theoretic interpretation of Eq. (5) to elucidate its underlying principles.

**Theorem 3.1** (Proof in Appendix A.1). *For mutual information between current observation $o_t$ and scene annotation $y$ conditioned on previous belief $b_{t-1}$, denoted as $I(o_t; y \mid b_{t-1})$, we can prove that our objective is a lower bound of this mutual information :*

$$\mathbb{E}_{(x^{(j)}, y^{(j)}) \sim \mathcal{D}} \left[ \frac{1}{K} \sum_{j=1}^{K} \log \frac{q_{\boldsymbol{\theta}}(y^{(j)} \mid b_{t-1}^{(j)}, o_t^{(j)})}{\frac{1}{K} \sum_{k=1}^{K} q_{\boldsymbol{\theta}}(y^{(k)} \mid b_{t-1}^{(j)}, o_t^{(j)})} \right]$$

$$\leq \mathbb{E}_x I(o_t; y \mid b_{t-1}) - \frac{\log(K)}{K}$$

$$= \mathbb{E}_x \left[ \mathcal{H}(y \mid b_{t-1}) - \mathcal{H}(y \mid b_{t-1}, o_t) \right] - \frac{\log(K)}{K},$$

$$(6)$$

*where $q_{\boldsymbol{\theta}}(y \mid o_1, \cdots, o_t)$ represents the variational distribution approximating the true conditional distribution $p(y \mid o_1, \cdots, o_t)$ with samples $(x^{(j)}, y^{(j)})_{j=1}^K$.*

*Remark* 3.2. To elucidate the connection between the derived lower bound on conditional mutual information and the optimization objective, we reformulate the variational distribution $q_{\boldsymbol{\theta}}(y \mid b_{t-1}, o_t)$ by employing the similarity term from Eq. (5), yielding $q_{\boldsymbol{\theta}}(y \mid b_{t-1}, o_t) := p(b_{t-1}, o_t)e^{-S(f(b_{t-1}, o_t; \boldsymbol{\theta}), y)}$. This reformulation equals the negative InfoNCE objective, as presented in Eq. (5). It underscores that the EAD training procedure maximizes the conditional mutual information, guiding the agent to collect maximally informative observations $o_t$ for determining the task-designated annotation $y$. Note that the tightness of the derived lower bound improves as the batch size $K$ increases.

The last equality in Eq. (6) follows from mutual information being equivalent to a reduction in conditional entropy. This reveals that the optimization guides the policy toward greedy informative exploration, defined as:

**Definition 3.3** (Greedy Informative Exploration). *Greedy informative exploration, represented by $\pi^g$, refers to an action policy which, at any timestep t, chooses an action $a_t$ that maximizes the decrease in the conditional entropy of a random variable $y$ given a new observation $o_t$ obtained from executing action $a_t$. Formally,*

$$\pi^g = \arg\max_{\pi \in \Pi} \left[ \mathcal{H}(y \mid b_{t-1}) - \mathcal{H}(y \mid b_{t-1}, o_t) \right], \quad (7)$$

*where $\mathcal{H}(\cdot)$ represents the entropy, $\Pi$ denotes the space encompassing all feasible policies.*

*Remark* 3.4. The conditional entropy term $\mathcal{H}(y \mid b_{t-1})$ measures the uncertainty in target $y$, conditioned on the belief $b_{t-1}$ maintained up to the previous timestep. In contrast, $\mathcal{H}(y \mid b_{t-1}, o_t)$ quantifies the incorporated uncertainty in $y$ after the observation $o_t$. While not guaranteed to yield globally optimal behavior over the entire trajectory, the **greedy informative exploration** strategy serves as an efficient baseline for rapid environmental learning through sequential actions and observations.

Through theoretical analysis, we demonstrate that the optimal policy $\pi_{\phi}^*$ in the InfoNCE objective function (5) converges to a greedy informative exploration policy, combining mutual information with policy refinement. The perspective of information theory reveals that a well-trained EAD model naturally adopts a greedy informative policy, leveraging contextual information to address the high uncertainty [40] in scenes containing adversarial patches.

## 3.3 Reinforced Embodied Active Defense

The previous EAD demonstrates significant potential but still faces challenges in three critical domains regarding effectiveness, applicability, and efficiency: (1) **Temporal inconsistency**: EAD's greedy informative often produces temporally inconsistent actions as illustrated in Fig. 2. This inconsistency stems from the myopic nature of greedy policies, which optimize for immediate information gain without considering long-term relevance in embodied learning [41], [42]. Consequently, the agent may revert to previously explored viewpoints, reducing exploration effectiveness and potentially leading to erroneous predictions, especially
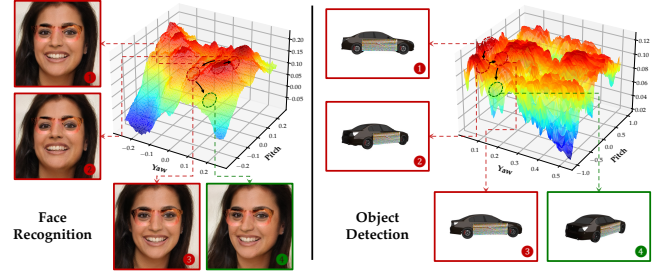


Fig. 2: EAD's temporal inconsistency issue visualized on the loss landscape *w.r.t.* camera yaw and pitch. The model often revisits similar states, limiting exploration and increasing vulnerability to adversarial impacts. These trajectories expose the drawbacks of a myopic greedy policy in dynamic settings.

when revisiting viewpoints with significant adversarial impacts [30], [43]. (2) **Limited applicability**: The reliance on differentiable dynamic models for learning EAD severely constrains its practical application. In real-world scenarios [44], [45], accurately modeling environment dynamics in a differentiable manner is often infeasible due to the complexity and unpredictability of physical systems. (3) **Computational inefficiency and instability**: Learning via differentiable simulation is computationally demanding and prone to numerical instabilities. The process requires extensive computation to traverse the dynamics model and derive their gradients. Furthermore, it often suffers from vanishing gradients during optimization [46] due to insufficient gradient quality from the simulated dynamics [47].

### 3.3.1 Accumulative Interactions for Temporal Consistency

Greedy exploration in the original EAD myopically minimizes uncertainty at each step in Eq. (4) without considering the future effects of actions. Such a greedy approach often leads to temporal inconsistency and suboptimal outcomes from a holistic perspective. To address this limitation and improve the overall performance of EAD, we introduce accumulative informative exploration, which aims to minimize long-term uncertainty about the target variable through a sequence of interactions with the environment.

**Definition 3.5** (Accumulative Informative Exploration). *Accumulative informative exploration, denoted by $\pi^*$, refers to the policy that maximizes the reduction in the entropy of $y$ given a series of observation $o_{0:H}$ resulting from continuous interaction with the environment by executing actions $a_{0:H}$. Formally,*

$$\pi^* = \arg\max_{\pi \in \Pi} \left[ \mathcal{H}(y) - \mathcal{H}(y \mid b_{H-1}, o_H) \right]. \quad (8)$$

*Remark* 3.6. In contrast to the greedy informative exploration defined in Definition 3.3, the accumulative informative policy aims to minimize the uncertainty of the target $y$ through continuous interaction with the environment, rather than focusing solely on single-step. Greedy informative exploration can be considered a special case of accumulative informative exploration, limited to one-step transitions where $H = 1$. From a long-term perspective, greedy informative exploration tends to myopically reduce step-wise uncertainty, often leading to sub-optimal outcomes compared to accumulative informative exploration.

To achieve accumulative informative exploration, we propose a multi-step accumulative interaction objective that optimizes the policy over a horizon of $H$ steps, incorporating terms that encourage reaching belief states that minimize the prediction loss and penalize high-entropy predictions:

$$\min_{\boldsymbol{\theta},\boldsymbol{\phi}} \mathbb{E}_{x,y,\tau} \sum_{p\in\mathcal{P}_x} \Big[ \mathcal{L}(\hat{y}_H, y) + \lambda \cdot \mathcal{H}(\hat{y}_H \mid b_{H-1}, o_H) \Big],$$
$$\text{with } \{\hat{y}_t, b_t\} = f(A(o_t, p; s_t), b_{t-1}; \boldsymbol{\theta}), \ a_t \sim \pi(\cdot \mid b_t; \boldsymbol{\phi}) \quad (9)$$
$$\text{s.t.} \quad o_t \sim \mathcal{Z}(\cdot \mid s_t, x), \quad s_t \sim \mathcal{T}(\cdot \mid s_{t-1}, a_{t-1}, x),$$

where $\mathbb{E}_{x,y,\tau}$ is the abbreviation of $\mathbb{E}_{(x,y)\sim\mathcal{D}, \tau\sim(\mathcal{M}(x),\pi)}$ following Eq. (4), $\mathcal{L}(\hat{y}_H, y)$ represents the prediction loss in the step $H$, and $\mathcal{H}(\hat{y}_H \mid b_{t-1}, o_t)$ denotes the entropy of the predicted label at the step $H$. The sampled trajectories follow the same distribution with Eq. (4). The entropy term serves as a regularizer, discouraging the agent from making high-entropy predictions that are characteristic of adversarial examples [40].

The proposed multi-step interactions align with the definition of accumulative informative exploration, as it seeks to minimize the uncertainty of the target variable through a sequence of actions and observations. By incorporating the prediction loss and the entropy regularization term, the objective encourages the agent to reach belief states that are informative and robust, leading to resilience against adversarial perturbations. Then we analyze the performance gap between the accumulative informative policy and the greedy one as below.

**Theorem 3.7** (Informative Policy Efficacy Inequality, Proof in Appendix A.2). *Consider two policies interacting continuously with an environment over a time horizon $H$:*

- *$\pi^g$ denote the greedy informative policy, resulting in observation sequence $o^g_{0:H}$ and belief sequence $b^g_{0:H}$.*
- *$\pi^*$ denote the accumulative informative policy, resulting in observation sequence $o^*_{0:H}$ and belief sequence $b^*_{0:H}$.*

*Define the trajectory information gain from time 0 to $H$ under a policy $\pi$ as the reduction in entropy of the variable $y$:*

$$\Delta\mathcal{H}_\pi := \mathcal{H}(y) - \mathcal{H}(y \mid b_{H-1}, o_H).$$

*Assume that the belief update function $f_b : (b_{t-1}, o_t) \mapsto b_t$ is bijective. Then the efficacy of $\pi^*$ relative to $\pi^g$ satisfies the following inequality:*

$$\Delta\mathcal{H}_{\pi^*} \geq \Delta\mathcal{H}_{\pi^g}, \quad (10)$$

*where equality holds if and only if the problem exhibits optimal substructure.*

*Remark* 3.8. As indicated by this inequality, the efficacy gap primarily arises from two factors: (1) the complexity of the exploration problem, which makes it difficult to guarantee an optimal structure to ensure the effectiveness of the greedy policy, and (2) the information loss during the belief update process, which encodes the previous belief $b_{t-1}$ and the current observation $o_t$ into the current belief $b_t$. The greedy policy does not account for information loss during belief updates through continuous interaction with the environment. Its formulation in Eq. (7) only ensures the selection of the most informative observation $o_t$ by taking action $a_{t-1}$ in a single step.

By extending the greedy information exploration discussed in Sec. 3.2, we can reformulate Eq. (5) into a multi-step interactive form. The learned policy model $\pi^*_\phi$ represents an accumulative informative policy, under the assumptions of unlimited model capacity and data samples. Additionally, Theorem 3.7 formally establishes the superiority of multi-step interactions, indicating that a well-trained multi-step model for contrastive tasks adopts an **accumulative informative policy** to continuously explore the environment. This model utilizes a series of contextual information from the environment to consistently reduce perceptual uncertainty caused by adversarial patches, providing a theoretical foundation for our multi-step accumulative objective in Eq. (9).

### 3.3.2 Policy Learning for Real-World Applicability

Despite the effectiveness in Sec. 3.3.1, REIN-EAD faces significant challenges in real-world applications. The unpredictability of physical systems makes it infeasible to model environment dynamics in a differentiable manner. Moreover, the inevitable computational overhead from differentiating the dynamics model and the numerical instability arising from approximations render it impractical for solving Eq. (9), especially over long horizons. These issues intensify with growing computational demands and cumulative gradient estimation errors, hindering the applicability and efficiency.

To address these challenges and enhance the real-world applicability of our approach, we propose a policy learning method that incorporates an uncertainty-oriented reward-shaping technique within the reinforcement learning framework. By eliminating the need for differentiable dynamic models, our method enables the agent to directly learn a policy that maximizes the expected cumulative reward through trial-and-error interactions with the environment. This flexibility allows the agent to efficiently adapt to changing dynamics or stochastic environments, which are prevalent in real-world scenarios [48].

Specifically, the uncertainty-oriented reward-shaping technique is designed to provide dense rewards at each step, guiding the agent to reduce perceptual uncertainty and minimize prediction errors. This approach addresses classical sparse reward strategy [49] in Eq. (9), where the agent can only access the reward at the end of the episode. By incorporating a weighted combination of intermediate rewards, our method allows for more granular and informative feedback to the agent. Formally, we define the dense reward $r_t$ as follows (proof of the equivalence in Appendix A.3):

$$r_t = \mathcal{L}(\hat{y}_{t-1}, y) - \gamma \cdot \mathcal{L}(\hat{y}_t, y), \quad (t > 0) \quad (11)$$

where $\gamma$ is the discount factor. This dense reward structure motivates the policy $\pi$ to seek new observations $o_t$ as feedback from the environment. The continuous feedback loop facilitated by the dense rewards enables the agent to efficiently adapt to new situations and make informed decisions in the face of uncertainty. Moreover, by fairly distributing the reward across each step, we alleviate the challenges of exploration and credit assignment [50], facilitating faster convergence and more efficient learning.

As for the reinforcement learning backbone, we employ Proximal Policy Optimization (PPO) [51] because of its learning efficiency and convergence stability. PPO enables stable policy updates by constraining the size of the policy

---

**Algorithm 1** Training REIN-EAD by Policy Learning

---

**Require:** Training data $\mathcal{D}$, number of iterations $M$, number of epochs $E$, loss function $\mathcal{L}$, perception model $f(\cdot; \boldsymbol{\theta})$, policy model $\pi(\cdot; \boldsymbol{\phi})$.
**Ensure:** The parameters $\boldsymbol{\theta}, \boldsymbol{\phi}$ of the learned EAD model.
1: **for** iteration $\leftarrow 0$ **to** $M - 1$ **do**
2:    ▷ *Roll-out perception $f(\cdot; \boldsymbol{\theta})$ and policy $\pi(\cdot; \boldsymbol{\phi})$ in the environment.* ◁
3:    Collect set of augmented trajectory and label pairs $\mathcal{D}_\tau = \{(\tau, y)\}$ by running policy $\pi(\cdot; \boldsymbol{\phi})$ and perception $f(\cdot; \boldsymbol{\theta})$ on $\mathcal{M}(x)$ with $(x, y) \sim \mathcal{D}$, where

$$\tau = (o'_0, b_0, a_0, r_0, o'_1, b_1, \ldots)$$

4:    Estimate advantages $\hat{A}_t$ using any advantage estimation algorithm
5:    **for** epoch $\leftarrow 0$ **to** $E - 1$ **do**
6:      $\boldsymbol{\phi}_{\text{old}} \leftarrow \boldsymbol{\phi}$
7:      **for all** mini-batch $\mathcal{B}_\tau \in \mathcal{D}_\tau$ **do**
8:        Update $\boldsymbol{\phi}$ by maximizing estimated the PPO-Clip objective $\mathcal{J}_{\text{policy}}(\boldsymbol{\phi})$ defined as:

$$\frac{1}{|\mathcal{B}_\tau|} \sum_{\tau \in \mathcal{B}_\tau} \sum_t \min(R(\boldsymbol{\phi})\hat{A}_t, \text{clip}(R(\boldsymbol{\phi}), 1-\epsilon, 1+\epsilon)\hat{A}_t)$$

       with one-step gradient ascent, where

$$R(\boldsymbol{\phi}) = \frac{\pi(a_t \mid b_t, \boldsymbol{\phi})}{\pi(a_t \mid b_t, \boldsymbol{\phi}_{\text{old}})}$$

9:        Update $\boldsymbol{\theta}$ by minimize the estimated objective in Eq. (9), namely $\mathcal{J}_{\text{percep}}(\boldsymbol{\theta})$:

$$\frac{1}{|\mathcal{B}_\tau|} \sum_{(\tau, y) \in \mathcal{B}_\tau} \sum_t \mathcal{L}(f(o_t, b_t; \boldsymbol{\theta}), y) + \lambda \cdot \mathcal{H}(f(o_t, b_t; \boldsymbol{\theta}))$$

       with one-step gradient descent.
10:      **end for**
11:    **end for**
12: **end for**

---

change at each iteration based on the dense rewards provided by our uncertainty-oriented reward-shaping technique. This incremental learning process ensures that the agent maintains a stable trajectory towards reducing uncertainty and minimizing prediction errors. The detailed training procedure is in Algorithm 1.

### 3.4 Adversary-agnostic Defense against Patch Attacks

The computation of $\mathcal{P}_x$ typically necessitates the resolution of the inner maximization in Eq. (4) by online adversarial training [7]. However, this is not only computationally expensive [52] but also problematic as inadequate assumptions for characterizing adversaries can hinder the model's ability to generalize across diverse, unseen attacks [53].

**Offline adversarial patch approximation (OAPA).** While the USAP approach [18] effectively learns an optimal informative strategy given sufficient training epochs, it is computationally expensive due to the need for extensive sampling to ensure the sampled manifold contains an adequate number of adversarial examples. To improve sampling efficiency while preserving the adversary-agnostic property, we introduce OAPA that employs the projected gradient technique to approximate the adversarial patch manifold $\mathcal{P}_x$ before training the REIN-EAD model. OAPA systematically characterizes

the manifold of adversarial patterns by pre-generating a surrogate set of patches $\tilde{\mathcal{P}} = p_i$, where each patch $p_i$ is derived through projected gradient ascent against the visual backbone. This offline approximation of the adversarial patch manifold allows the REIN-EAD model to learn a compact yet expressive representation of adversarial patterns, enabling it to effectively defend against previously unseen attacks. Our empirical findings suggest that performing this offline approximation maximization is highly effective in developing models robust to a broad spectrum of attacks (refer to Sec. 4). Additionally, because this maximization process occurs offline and only once before training, it substantially boosts training efficiency and renders it competitive with conventional training methods.

## 4 EXPERIMENTS

In this section, we verify the effectiveness of REIN-EAD on different tasks, including face recognition in Sec. 4.1, 3D object classification in Sec. 4.2 and object detection in Sec. 4.3.

### 4.1 Evaluation on Face Recognition

#### 4.1.1 Experimental Settings

**Experimental environment.** To enable the unconstrained navigation and observation collection of REIN-EAD, we construct a manipulable simulation environment for both training and purposes. To align with specific vision tasks, we define the state as a combination of the camera's yaw and pitch, while the action corresponds to the rotation of the camera[2]. This definition establishes the transition function, with the core of the simulation environment revolving around the observation function, which generates a 2D image based on the camera's state. As elaborated in Sec. 3.1, the training of original EAD requires differentiable environmental dynamics. To achieve a fair comparison, we first employ the advanced 3D generative model, EG3D [54], which enables realistic differentiable rendering (refer to Appendix C.1 for simulation fidelity). To further substantiate the superiority of REIN-EAD, we conduct training and testing of REIN-EAD in the same environment to ensure fair comparison.

**Evaluation metrics and protocols.** To rigorously validate the effectiveness of REIN-EAD, we conduct extensive experiments on CelebA-3D, which we meticulously reconstruct from 2D face images in CelebA into 3D representations by utilizing GAN inversion [55] with EG3D [54]. To assess the **standard accuracy**, we sample 2,000 test pairs from the CelebA and follow the well-established evaluation protocol from LFW [56]. To comprehensively evaluate the robustness, we report the **attack success rate (ASR)** on 100 identity pairs, considering both impersonation and dodging attacks [57] under various attack methods in both white-box and black-box settings. The white-box attack methods emcompass MIM [58], EoT [59], GenAP [4] and Face3Dadv (3DAdv) [5]. In the black-box attacks, we employ the transfer-based attack, targeting the surrogate models such as IResNet-18 CosFace [60], [61] and IResNet-50 Softmax with 3DAdv, and the query-based methods including NAttack [62] and

---

2. To prevent the agent from "cheating" by rotating to angles that conceal adversarial patches, specific constraints are imposed on the yaw and pitch.

TABLE 1: The **Standard accuracy** (%) and **attack success rates** (%) on face recognition. [†] denotes the methods using adversarial training. Methods with light blue background do not require a differentiable environment, while the rest do.

| Method | Acc. (%) | White-box | | | | Transfer-based | | Query-based | | Adaptive | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | MIM | EoT | GenAP | 3DAdv | Cos. | Softmax | NAttack | RGF | BPDA | Worst-case |
| *Impersonation Attack* | | | | | | | | | | | |
| Undefended | 88.86 | 100.0 | 100.0 | 99.00 | 89.00 | 28.00 | 23.00 | 100.0 | 100.0 | 100.0 | 100.0 |
| JPEG | 89.98 | 99.00 | 100.0 | 99.00 | 93.00 | 33.00 | 33.00 | 96.00 | 94.00 | 99.00 | 100.0 |
| LGS | 83.50 | 5.10 | 7.21 | 33.67 | 30.61 | 11.63 | 6.98 | 11.63 | 4.65 | 38.37 | 48.83 |
| SAC | 86.83 | 6.06 | 9.09 | 67.68 | 64.64 | 8.70 | 9.78 | 13.04 | 14.13 | 48.00 | 67.68 |
| PZ | 87.58 | 4.17 | 5.21 | 59.38 | 45.83 | 6.45 | 9.68 | 4.30 | 3.26 | 89.76 | 92.86 |
| SAC[†] | 80.55 | 3.16 | 3.16 | 18.94 | 22.11 | 11.11 | 12.36 | 12.22 | 14.44 | 51.73 | 51.73 |
| PZ[†] | 85.85 | 3.13 | 3.16 | 19.14 | 27.34 | 8.24 | 5.00 | 10.58 | 9.41 | 61.01 | 98.99 |
| DOA[†] | 79.55 | 95.50 | 89.89 | 96.63 | 89.89 | 15.73 | 17.97 | 34.83 | 16.86 | 89.89 | 96.63 |
| EAD | **90.45** | 4.12 | 3.09 | 5.15 | **7.21** | 4.17 | 5.20 | 4.12 | 4.12 | 8.33 | 9.38 |
| **REIN-EAD** | 89.03 | **2.10** | **1.06** | **3.15** | 7.37 | **2.10** | **2.08** | **1.05** | **2.12** | **4.21** | **7.37** |
| *Dodging Attack* | | | | | | | | | | | |
| Undefended | 88.86 | 100.0 | 100.0 | 99.00 | 98.00 | 44.00 | 35.00 | 96.00 | 96.00 | 100.0 | 100.0 |
| JPEG | 89.98 | 98.00 | 99.00 | 95.00 | 88.00 | 49.00 | 45.00 | 81.00 | 83.00 | 100.0 | 100.0 |
| LGS | 83.50 | 49.47 | 52.63 | 74.00 | 77.89 | 22.11 | 21.05 | 18.95 | 20.00 | 78.92 | 78.92 |
| SAC | 86.83 | 73.46 | 73.20 | 92.85 | 78.57 | 40.80 | 36.84 | 55.26 | 50.00 | 65.22 | 92.85 |
| PZ | 87.58 | 6.89 | 8.04 | 58.44 | 57.14 | 41.67 | 28.34 | 28.33 | 31.67 | 88.89 | 90.00 |
| SAC[†] | 80.55 | 78.78 | 78.57 | 79.59 | 85.85 | 47.46 | 43.54 | 55.93 | 62.71 | 85.02 | 85.85 |
| PZ[†] | 85.85 | 6.12 | 6.25 | 14.29 | 20.41 | 50.88 | 47.69 | 56.14 | 50.87 | 69.45 | 98.00 |
| DOA[†] | 79.55 | 75.28 | 67.42 | 87.64 | 95.51 | 30.33 | 31.46 | 53.93 | 28.09 | 95.51 | 95.51 |
| EAD | **90.45** | **0.00** | **0.00** | 2.10 | 13.68 | 5.26 | 7.36 | 1.05 | **0.00** | 22.11 | 22.11 |
| **REIN-EAD** | 89.03 | 1.04 | 2.04 | 5.15 | **13.54** | **4.17** | **7.29** | **1.03** | **0.00** | **8.16** | **14.43** |

RGF [63]. Note that 3DAdv utilizes expectation over 3D transformations during the optimization, endowing it with inherent robustness to 3D viewpoint variation within a range of $\pm15°$. More details are described in Appendices C.1 & C.2.

**Implementation details.** To extract discriminative visual features, we employ the pretrained IResNet-50 ArcFace [60] as the visual backbone, leveraging its pretrained weights while keeping them frozen during subsequent training stages. To model the recurrent perception and policy components, we adopt a variant of the Decision Transformer [64], which effectively processes feature sequences extracted by the visual backbone to predict a normalized embedding for FR. We set the maximum horizon length to $H = 4$ in EAD as default. In contrast to EAD, REIN-EAD eliminates the dependency on Back-Propagation Through Time (BPTT), allowing us to extend the horizon length to $H = 16$ without VRAM constraints. Furthermore, we incorporate an additional MLP value head into REIN-EAD to facilitate advantage estimation for PPO. Further details on REIN-EAD can be found in Appendix C.5.

**Defense baselines.** To fully evaluate the effectiveness of REIN-EAD, we benchmark it against a diverse range of state-of-the-art defense methods. These baselines include adversarial training-based Defense against Occlusion Attacks (DOA) [8], and purification-based methods like JPEG compression (JPEG) [21], local gradient smoothing (LGS) [22], segment and complete (SAC) [12], PatchZero (PZ) [13], Patch-Agnostic Defense (PAD) [65] and DIFFender [66]. For DOA, we employ rectangle-shaped PGD patch attacks [7] with 10 iterations and a step size of $2/255$. Note that SAC and PZ require a patch segmenter to locate the area of adversarial patches. Therefore, we train the segmenter using patches of Gaussian noise to ensure the same adversarial-agnostic

setting. Besides, we consider enhanced versions of SAC and PZ, denoted as SAC[†] and PZ[†], which involve training with adversarial patches generated using the EoT technique. More details are presented in Appendix C.4.

### 4.1.2 Experimental Results

**Effectiveness of REIN-EAD.** Table 1 presents a comprehensive evaluation of both the standard accuracy and robust performance against diverse attacks under a white-box setting, with the adversarial patch size set to 8% of the image size. Remarkably, our approach significantly outperforms previous state-of-the-art passive techniques that are agnostic to adversarial examples in both clean accuracy and defense efficacy. For instance, REIN-EAD reduces the attack success rate of 3DAdv by an impressive effect in both scenarios. Furthermore, REIN-EAD also improves the average attack rate reduction when faced with black-box and adaptive attacks compared with EAD by extending the informative greedy policy to an accumulative one. Notely, REIN-EAD even surpasses the performance of the undefended passive model regarding standard accuracy, effectively reconciling the trade-off between robustness and accuracy [67] through embodied perception. Furthermore, our method even outstrips the baselines by incorporating adversarial examples during training. Although SAC[†] and PZ[†] are trained using patches generated via EoT [59], we still obtain superior performance, highlighting the effectiveness of leveraging the environmental feedback in active defense.

Fig. 3 provides a visual illustration of the defense process executed by REIN-EAD. While REIN-EAD may be initially fooled by the adversarial patch, its subsequent active interactions with the environment progressively increase the similarity between the positive pair and decrease the simi-
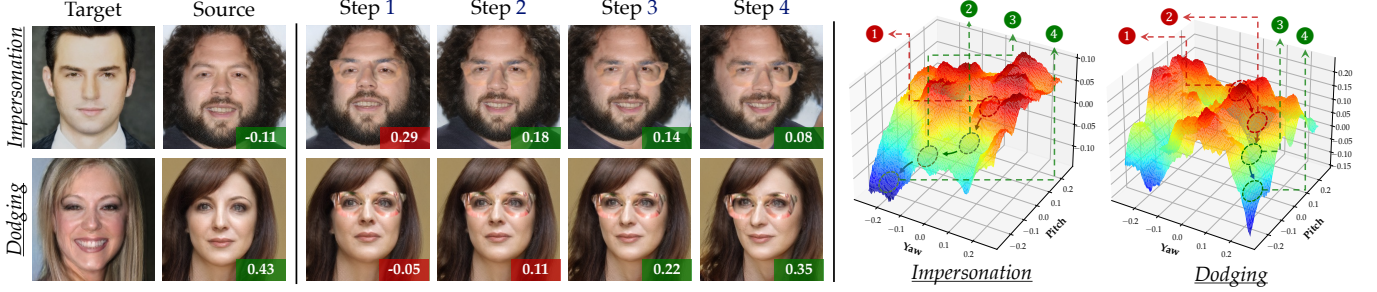
Fig. 3: **Qualitative results** of REIN-EAD. The first two columns present the original image pairs, and the subsequent columns depict the interactive inference steps taken by the model. The defensive trajectory of REIN-EAD is plotted on the loss landscape *w.r.t.* yaw and pitch of the camera, considering the IResNet-50 ArcFace as the target model [60]. The adversarial glasses are generated with 3DAdv, which are robust to 3D viewpoint variation. The computed optimal threshold for distinguishing between positive and negative pairs is set to 0.19 from [-1, 1].
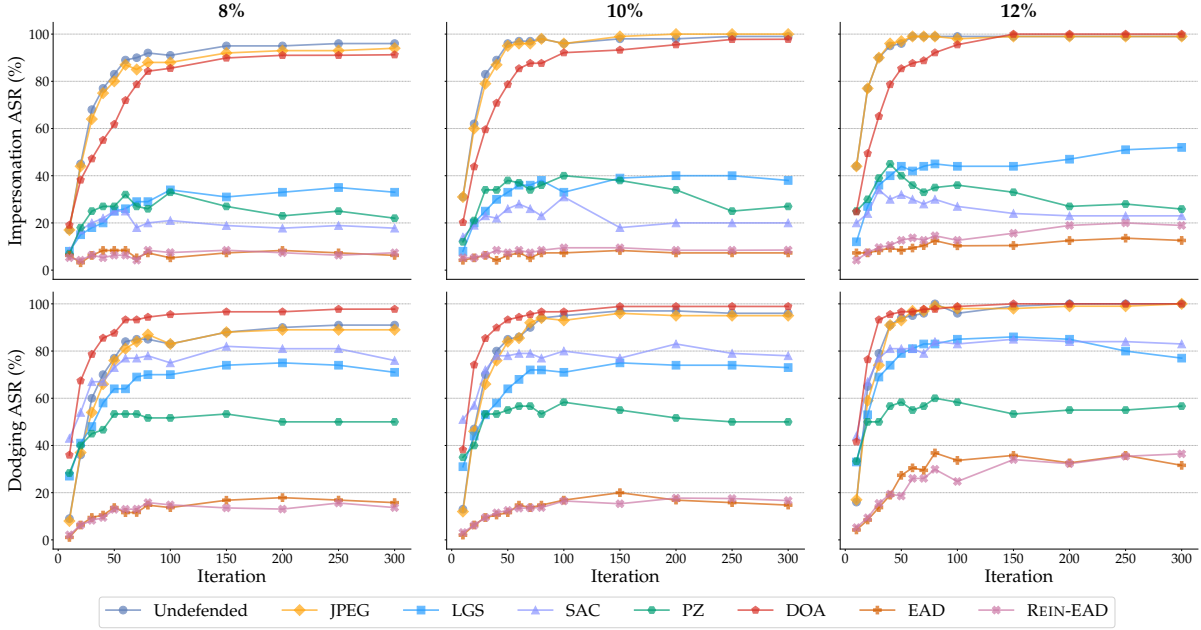


Fig. 4: Comparative evaluation of defense methods across varying attack iterations with different adversarial patch sizes. Note that SAC, PZ and DOA are involved with adversarial training.

larity between the negative pair. Consequently, REIN-EAD effectively mitigates the impact of adversarial hallucination through the proactive acquisition of additional observations, as depicted in the corresponding loss landscapes.

**Effectiveness against adaptive attack.** While the deterministic and differentiable dynamic models could potentially enable backpropagation through the entire inference trajectory of EAD, the computational cost becomes prohibitive due to the rapid consumption of GPU memory as trajectory length $H$ increases. To overcome this, we first adopt an approach similar to the original strategy that approximates the true gradient by computing the expected gradient over a surrogate uniform superset policy distribution. This approximation (BPDA) necessitates an optimized patch to handle diverse action policies. Our adaptive attack implementation builds upon 3DAdv [5] leveraging 3D viewpoint variations. Moreover, we also evaluate more sophisticated adaptive attacks, including 1) attacking the whole pipeline using true gradients obtained by gradient checkpointing, and 2) targeting the perception and policy sub-modules independently. The **worst-case** performance of REIN-EAD

across two adaptive attacks is presented in Table 1 for a reliable evaluation. Moreover, to benchmark other defense methods, we report their worst-case performance under a series of adaptive attacks. More details are presented in Appendix C.3. We can see that REIN-EAD maintains its robustness against the most potent adaptive attacks. This observation shows that the defensive capabilities of REIN-EAD stem from the synergistic integration of its policy and perception models, rather than relying on a short-cut strategy to neutralize adversarial patches from specific viewpoints.

**Generalization of REIN-EAD.** As demonstrated in Table 1, despite no prior knowledge of specific adversaries, ours exhibits remarkable generalization across various **unseen adversarial attack methods**. It is partially attributed to the inherent ability of REIN-EAD to dynamically interact with their environment, enabling them to adapt and respond to novel types of attacks. Additionally, we assess the models' resilience across a wide range of **patch sizes** and **attack iterations**. Despite being trained solely on patches constituting 10% of the image, REIN-EAD consistently maintains a notably low attack success rate, even when subjected to larger

TABLE 2: The **standard accuracy** and **white-box impersonation attack success rates** on ablated models. For the model with stochastic policy, we report the mean and standard deviation across five independent runs to ensure reliability.

| Category | Component | Acc. (%) | Attack Success Rate (%) | | | | |
|---|---|---|---|---|---|---|---|
| | | | MIM | EoT | GenAP | 3DAdv | Adaptive |
| Passive Perception | Undefended | 88.86 | 100.0 | 100.0 | 99.00 | 98.00 | 98.00 |
| | + Random Movement | 90.38 (± 0.12) | 4.17 (± 2.28) | 5.05 (± 1.35) | 8.33 (± 2.21) | 76.77 (± 3.34) | 76.77 (± 3.34) |
| EAD | + Perception Model | 90.22 (± 0.31) | 18.13 (± 4.64) | 18.62 (± 2.24) | 22.19 (± 3.97) | 30.77 (± 1.81) | 31.13 (± 3.01) |
| | + Policy Model | 89.85 | 3.09 | 4.12 | 7.23 | 11.34 | 15.63 |
| REIN-EAD | + Multi-steps Interaction | 89.02 (± 0.18) | 2.12 (± 0.08) | 3.19 (± 0.44) | 5.33 (± 0.87) | 12.77 (± 1.13) | 6.38 (± 0.80) |
| | + OAPA (FGSM [68]) | 88.95 (± 0.19) | 3.22 (± 0.23) | 3.03 (± 0.41) | 5.15 (± 1.06) | 13.18 (± 1.11) | 5.14 (± 0.62) |
| | + OAPA (PGD [7]) | 89.03 (± 0.21) | **2.10** (± 0.42) | **1.06** (± 0.50) | **3.15** (± 1.08) | **7.37** (± 1.32) | **4.21** (± 0.56) |



Fig. 5: Comparative evaluation of computational overhead of defense methods.



Fig. 6: Performance variation along different decision steps.

patch size and increased attack iteration, as shown in Fig. 4. This exceptional resilience can be attributed to REIN-EAD's primary reliance on environmental information, rather than solely depending on patterns of presupposed adversaries. By dynamically updating their perception, REIN-EAD can generalize well to different aspects. The details are available in Appendix C.7.

**Computational overhead.** We further compare the computational overhead of REIN-EAD with passive defense baselines regarding both training and inference times. The evaluation is conducted on an NVIDIA GeForce RTX 3090 Ti with a batch size of 64. SAC [12] and PZ [13] require two-stage segmenter training: initial training with pre-generated adversarial images followed by self-adversarial training. DOA [8] necessitates feature extractor retraining. Our REIN-EAD approach involves offline and online phases without adversarial training. As indicated in Fig. 5, although differential rendering employed by EAD imposes significant computational demands during the online training phase, the total training time of EAD effectively achieves an effective balance between the purely adversarial training method DOA and partially adversarial methods such as SAC and PZ. This efficiency primarily stems from our unique adversary-agnostic approach (OAPA), which eliminate the need to generate adversarial examples *online*, thereby enhancing training efficiency. Despite the larger horizon of REIN-EAD, which inherently increases sampling time, it still demonstrates faster training than EAD. This advantage stems from REIN-EAD employed by the model-free approach avoids the computationally intensive process of backpropagation to obtain policy gradients, thus improving overall learn-
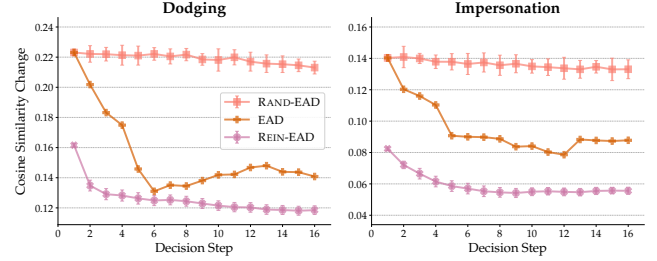
ing efficiency. Regarding model inference, the perception model accounts for 98.4% of this processing time, while the lightweight policy MLP requires only 1.6%. In total, REIN-EAD exhibits superior speed compared to other baselines, such as LGS and SAC. This advantage is attributed to the reliance of LGS and SAC on CPU-intensive, rule-based image preprocessing techniques, which inevitably reduces their inference efficiency. More details about the computational overhead are provided in Appendix C.6.

### 4.1.3 Ablation Study

**Effectiveness of recurrent feedback.** We thoroughly investigate the critical role of recurrent feedback, *i.e.*, reflecting on prior beliefs using a comprehensive fusion model, in achieving robust performance. In Table 2, even when equipped with only the perception model, REIN-EAD significantly surpasses both the undefended baseline and passive FR model that relies on multi-view ensembles (Random Movement). Notably, the multi-view ensemble model fails to counteract the state-of-the-art 3DAdv. This observation corroborates that REIN-EAD's defensive strength is not merely a function of the vulnerability of adversarial examples to viewpoint transformations. Instead, the superior performance of REIN-EAD can be attributed to its ability to actively explore the environment and dynamically adjust its perception.

**Impact of horizon length $H$.** We conduct an analysis of the impact of horizon length $H$ on the performance of REIN-EAD in Fig. 6. By examining the changes in the similarity between face pairs affected by adversarial patches, we consistently demonstrate the effectiveness of REIN-EAD in mitigating the detrimental effects of adversarial attacks. Specifically, for impersonation attacks, the change in similarity implies an
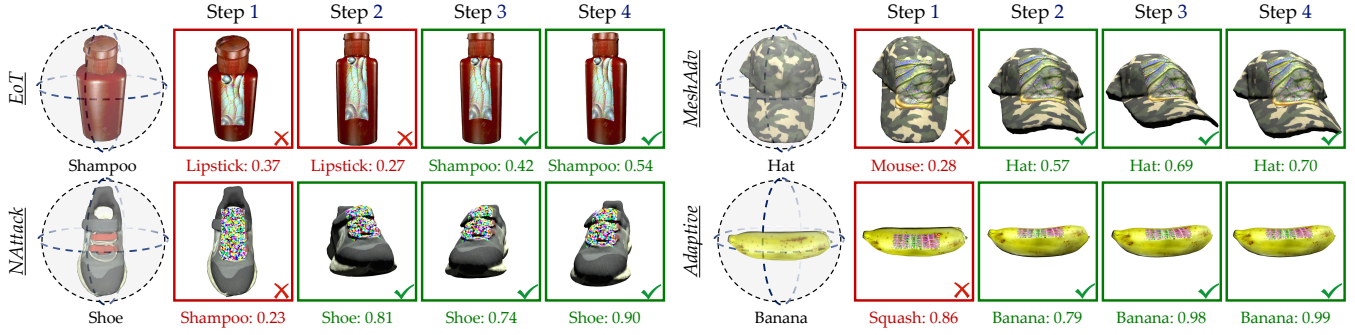
Fig. 7: Qualitative results of REIN-EAD on dynamic OmniObject3D, with the adversarial patch occupying $20\%$ of the object's bounding box in the front view. The state space for object classification is defined as $[-\frac{\pi}{2}, \frac{\pi}{2}] \times [0, \frac{\pi}{2}]$, encompassing a comprehensive range of viewpoints.

TABLE 3: The **Standard accuracy** (%) and **attack success rates** (%) on 3D object classification. $^\dagger$ denotes the methods using adversarial training. Methods with light blue background do not require a differentiable environment, while the rest do.

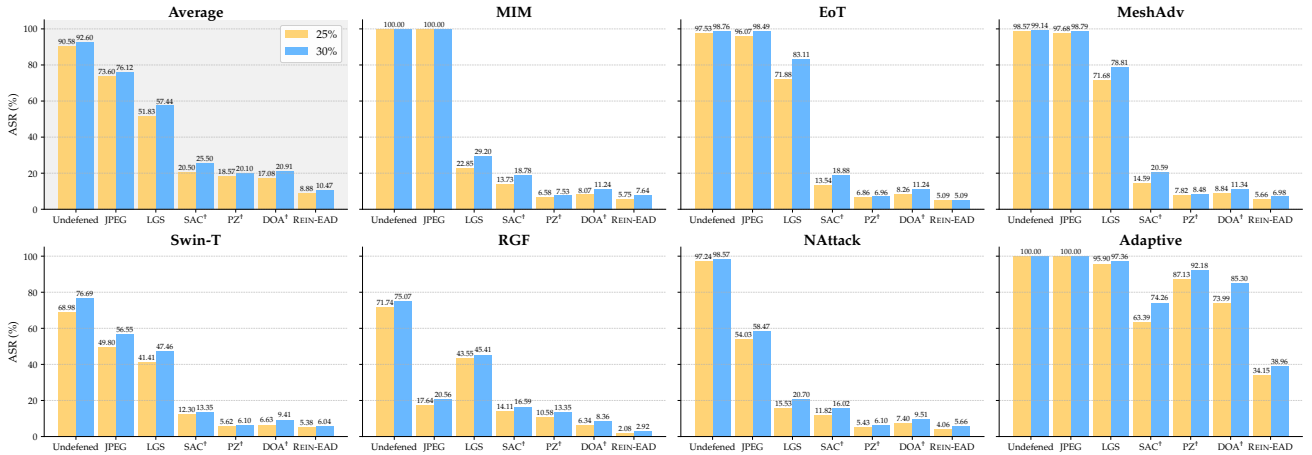| Method | Acc. (%) | White-box | | | Transfer-based | Query-based | | Adaptive |
|---|---|---|---|---|---|---|---|---|
| | | MIM | EoT | MeshAdv | Swin-T | RGF | NAttack | BPDA |
| Undefended | 88.17 | 100.00 | 93.72 | 96.19 | 58.90 | 68.03 | 96.48 | 100.00 |
| JPEG | 83.22 | 99.50 | 92.44 | 94.76 | 43.35 | 15.73 | 49.09 | 100.00 |
| LGS | 85.91 | 18.36 | 61.91 | 64.55 | 33.59 | 38.96 | 13.77 | 93.46 |
| PAD | 87.16 | 25.31 | 27.14 | 29.45 | 18.48 | 17.81 | 27.22 | 90.47 |
| DIFFender | 80.20 | 20.61 | 54.71 | 61.40 | 28.24 | 14.15 | 18.41 | 45.08 |
| SAC$^\dagger$ | 88.00 | 10.30 | 9.53 | 10.49 | 9.72 | 11.15 | 10.20 | 57.01 |
| PZ$^\dagger$ | 88.00 | 5.34 | 5.34 | 7.53 | 4.29 | 9.06 | 4.39 | 80.65 |
| DOA$^\dagger$ | 87.33 | 6.63 | 6.53 | 7.30 | 5.57 | 4.61 | 5.96 | 59.75 |
| **REIN-EAD** | **88.93** | **3.21** | **4.15** | **4.34** | **3.87** | **2.26** | **3.87** | **28.96** |



Fig. 8: Evaluating generalization on 3D object classification models under attacks with different patch sizes.

increase, while a decrease is observed for dodging attacks. This trend suggests that the accumulation of information from additional viewpoints during the decision process effectively attenuates the issues of information loss and model hallucination engendered by adversarial patches.
**Efficiency of learned policy.** We further validate its superiority of the policy empirically by comparing the performance of EAD with two variants: EAD integrated with a random movement policy RAND-EAD and REIN-EAD. All three approaches share identical neural network architecture and parameters. Fig. 6 illustrates that the RAND-EAD cannot mitigate the adversarial effect with random exploration, even when a greater number of actions are employed. Consequently, the exploration efficiency of the random policy is significantly inferior to REIN-EAD. Moreover, REIN-EAD

demonstrates a more stable and continuous deduction of adversarial effect, which can be primarily attributed to its adoption of an accumulative informative policy that consistently reduces perceptual uncertainty, instead of a greedy approach. This strategy effectively avoids the oscillations resulting from temporal inconsistency.
**Influence of patched data.** For REIN-EAD, the inherent limitations of RL in efficiently exploring the vast patch space are addressed by incorporating the OAPA algorithm. We compare OAPA using default PGD with a variant using FGSM adversarial examples, which are less optimized due to taking only a single gradient step. As shown in Table 2, OAPA ensures that REIN-EAD achieves superior robustness after substantial training, with PGD outperforming FGSM and others.
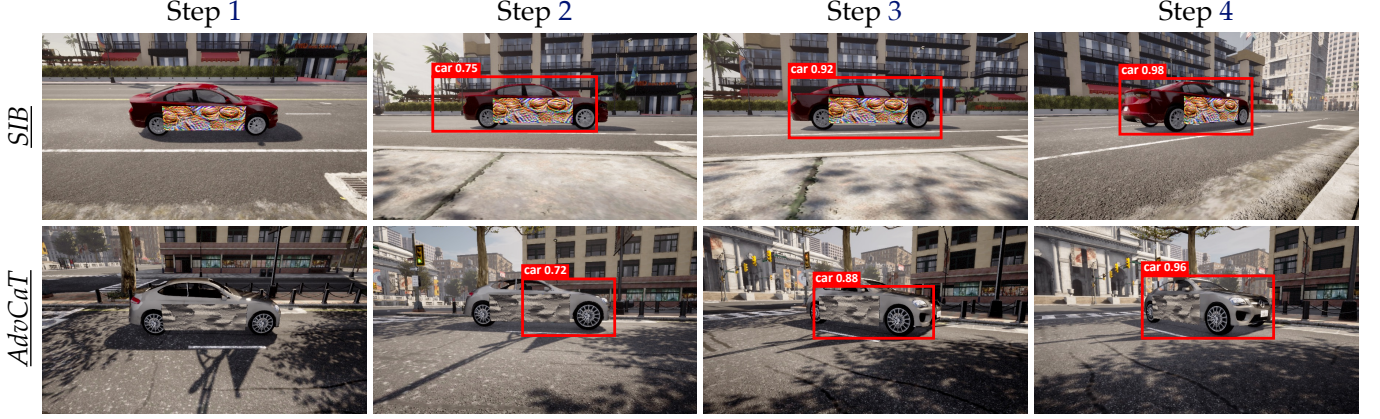
Fig. 9: Qualitative results of REIN-EAD on CARLA. Adversarial patches cover 25% of the object's front-view bounding box.

**Stable convergence of policy training.** To address active defense challenges in dynamic 3D environments, REIN-EAD enhances standard PPO with two key stability features: a two-phase training approach (offline perception pretraining followed by joint online training) that accelerates convergence, and integrated supervised learning that functions as regularization to maintain perception capabilities while preventing erratic policy updates. More experiments are presented Appendix C.9.

**Alternative reward formulations.** We have conducted additional ablation experiments to compare the performance of REIN-EAD with the uncertainty-oriented reward shaping to two alternative reward formulations, named Direct Entropy Deduction and Binary Outcome Reward. More details are presented Appendix C.10. As shown in Table C.7, our proposed reward shaping approach outperforms other methods in terms of both clean accuracy and adversarial robustness against patches. It employs a dense formulation that accelerates convergence and guides the model to learn a policy that maximizes information gain towards accurate perception.

## 4.2 Evaluation on 3D Object Classification

### 4.2.1 Experimental Settings

**Application in non-differential environment.** The classification task, widely employed and inherently vulnerable to patch attacks, often relies on rendering techniques for 3D dataset synthesis due to the scarcity of 3D data [69]. However, the discrete nature of rasterization renders analytical derivation of action transitions non-differentiable [70], limiting the application of EAD that depends on a differential dynamic model. Existing differentiable rendering frameworks [71], [72], [73], [74] either lack precision or compromise rendering quality, failing to meet the training requirements of EAD. To address this, we introduce REIN-EAD, a novel framework for operating effectively in non-differentiable environments.

**Dynamic OmniObject3D.** We leverage the recently proposed OmniObject3D [75], the largest real-scanned 3D dataset. OmniObject3D shares numerous common classes with classic 2D datasets (e.g. ImageNet [76]), making it particularly suitable for evaluating REIN-EAD. The environment is established by Pytorch3D [72] and Gym [77], enabling the rendering of objects at specific viewpoints according to the agent's actions. We refer to this environment as dynamic

TABLE 4: The performance of REIN-EAD with different reward shaping.

| Reward | Acc (%) | Attack Success Rate (%) | | | |
|---|---|---|---|---|---|
| | | MIM | EoT | GenAP | 3DAdv |
| Entropy Deduction | 88.67 | 3.15 | 2.11 | 4.21 | 11.42 |
| Outcome Reward | 88.62 | 3.22 | 3.26 | 5.94 | 10.86 |
| **ours** | **89.03** | **2.10** | **3.15** | **7.37** | **4.21** |

TABLE 5: The performance on object detection in EG3D. † indicates training with adversarial examples.

| Method | Average Precision (%) | | | | |
|---|---|---|---|---|---|
| | Clean | EoT | SIB | UAP | AdvCaT |
| Undefended | 88.55 | 9.06 | 17.30 | 20.29 | 28.38 |
| JPEG | 88.40 | 11.78 | 12.46 | 13.11 | 30.20 |
| LGS | 87.81 | 43.15 | 34.26 | 10.01 | 57.31 |
| SAC | 88.55 | 67.99 | 69.70 | 71.64 | 32.80 |
| PZ | 88.55 | 80.58 | 81.32 | 81.87 | 28.65 |
| SAC† | 88.55 | 70.10 | 71.08 | 74.06 | 40.67 |
| PZ† | 88.55 | 85.31 | 85.43 | 83.53 | 43.36 |
| EAD | 92.50 | 91.61 | **91.47** | 91.02 | 91.34 |
| **REIN-EAD** | **94.26** | **92.09** | 91.45 | **91.14** | **92.13** |

OmniObject3D. Details of the environment establishment can be found in Appendix D.1.

**Adversaries in the texture space.** To address adversarial threats in dynamic environments, we employ Pytorch3D to implement patch attacks on 3D mesh textures through its differential back-propagation pipeline[3]. This approach enables rendering 3D adversarial objects as 2D images from specified viewpoints, ensuring consistent adversarial appearance across multiple views. As the attack modifies the texture space, the patch's shape varies with perspective, challenging defense generalization. MeshAdv [78] accounts for expectations across 3D transformations in its differential rendering, suitable for multi-view 3D mesh attacks. The patch affects 20% of the bounding box area in the object's front view. Further information is detailed in Appendix D.3.

**Implementation details.** For the visual backbone, we employ the pretrained Swin Transformer (Swin-S) [79] on ImageNet and fine-tune it on the dynamic OmniObject3D. The weights of Swin-S are frozen in the subsequent experiments. We

---

3. Pytorch3D facilitates texture differentiation but does not support action transition differentiation for training EAD.

implement REIN-EAD following a paradigm similar to that used in the FR system, leveraging the OAPA algorithm for REIN-EAD with patches that occupy 20% of the object bounding box. More details are in Appendix D.2.

**Defense baseline.** We employ the same defense baselines as in the FR task, making necessary adaptations to the parameters to accommodate the classification settings. Notably, SAC, PZ, and DOA involve adversarial training, which is a widely adopted technique for enhancing robustness. More details are shown in Appendix D.4.

### 4.2.2 Experimental Results

**Effectiveness of REIN-EAD.** We evaluate the robustness of REIN-EAD and five baseline defenses under various white-box, black-box and adaptive attacks on the test set of OmniObject3D in Table 3. In most cases, JPEG and LGS fail to provide an effective defense, while the other baselines provide protection under white-box and black-box settings. Notably, REIN-EAD significantly reduces the attack success rate of various unseen adversaries, without compromising standard accuracy. In contrast, although SAC$^\dagger$ and PZ$^\dagger$ can purify the adversarial patch with the prior knowledge of the EoT adversary, they are inferior to REIN-EAD due to the loss of masked features. Furthermore, the robustness gap is amplified under a stronger adaptive attack targeting the combination of classifier and defense module. We can see that the performances of passive baselines drop significantly under adaptive attack. In contrast, REIN-EAD achieves the strongest robustness among the baselines. These results demonstrate the effectiveness of our model-free learning method and its applicability in general non-differentiable environments. Fig. 7 illustrates that even when initially deceived by adversaries, REIN-EAD can actively explore and observe the environment to correct and refine its predictions.

**Generalization of REIN-EAD.** We further investigate the adaptability of REIN-EAD to various patch sizes. Specifically, the defense methods are trained using patches covering 20% of the bounding box, and their performance is subsequently assessed with patches sized at 25% and 30%. As depicted in Fig. 8, REIN-EAD maintains its robustness when encountering patches larger than those used in the training phase.

## 4.3 Evaluation on Object Detection

### 4.3.1 Experimental Settings

Object detection in autonomous driving is more challenging than face verification and object classification, as the model must distinguish vehicles from intricate backgrounds and accurately regress bounding boxes in non-differentiable real-world environments. To tackle this, we have designed EAD and REIN-EAD for object detection and evaluate them on EG3D [54] that supports differentiable learning, and CARLA [48] belonging to a photorealistic environment used in autonomous driving research.

**Evaluation on differentiable EG3D.** We adopt a differentiable generative framework by EG3D, enabling the generation of diverse vehicle types from controllable perspectives while ensuring 3D consistency. The differentiable environment allows for the verification of both EAD and REIN-EAD paradigms. More details can be found in Appendix E.1.

TABLE 6: The performance on object detection model in CARLA. $^\dagger$ indicates training with adversarial examples.

| Method | Average Precision (%) | | | | |
|---|---|---|---|---|---|
| | Clean | EoT | SIB | UAP | AdvCaT |
| Undefended | 80.97 | 28.47 | 28.61 | 35.85 | 38.87 |
| JPEG | 81.57 | 36.50 | 35.80 | 34.66 | 37.28 |
| LGS | 80.32 | 74.73 | 72.64 | 57.76 | 51.34 |
| SAC | 79.78 | 27.06 | 28.55 | 42.03 | 37.38 |
| PZ | 80.70 | 62.43 | 59.35 | 49.24 | 37.90 |
| SAC$^\dagger$ | 79.28 | 31.68 | 32.95 | 30.70 | 39.09 |
| PZ$^\dagger$ | 80.91 | 76.10 | 75.50 | 75.24 | 40.51 |
| **REIN-EAD** | **83.15** | **82.82** | **81.97** | **82.12** | **82.86** |

**Evaluation on non-differentiable CARLA.** We evaluate the REIN-EAD's defense capabilities in CARLA, a practical and non-differentiable autonomous driving simulator. The model's robustness is assessed on 41 vehicles, encompassing all available asset blueprints in CARLA. The details about the CARLA environment are provided in Appendix E.2.

**Evaluation metric and adversaries.** The **average precision** at intersections over union thresholds from 50% to 95% (AP@50:95) is used as the evaluation metric. The adversarial methods include EoT [59], SIB [80], UAP [81] and AdvCaT [82], which aim to launch a hiding attack that makes the vehicle disappear from the detector. SIB impacts both hidden and final layers, perturbing features before REIN-EAD, while AdvCaT generates inconspicuous environmental mosaic camouflage. All attack methods use the expectation over 3D transformations to enhance patch resilience against viewpoint changes. The patch, attached to the vehicle's side, occupies 25% of the initial perspective's bounding box.

**Implementation details.** The pretrained YOLOv5n [83] serves as the visual backbone, combined with a Decision Transformer to implement REIN-EAD. Additional details on REIN-EAD can be found in Appendix E.3. The defense baselines follow the same approach as in previous tasks. More details can be found in Appendix E.5.

### 4.3.2 Experimental Results

**Effectiveness on EG3D.** We evaluate REIN-EAD and other defense baselines on 200 test data samples from EG3D under four white-box attacks designed to deceive the YOLO detector. Table 5 shows that the undefended model and JPEG defense are breached by all attacks. LGS performs poorly under EoT, SIB, and UAP but moderately resists AdvCaT. SAC and PZ defend well against noisy-patterned attacks, with enhanced versions (SAC$^\dagger$ and PZ$^\dagger$) performing slightly better. PZ$^\dagger$ approximates clean average precision but fails against stealthy AdvCaT. REIN-EAD surpasses others in clean and robust performance, significantly outperforming passive baselines under different attacks.

**Effectiveness on CARLA.** We evaluate the performance of REIN-EAD in the photo-realistic CARLA environment in a non-differential and practical setting. Table 6 reveals a similar trend that aligns with observations across other tasks. Specifically, LGS performs better under EoT and SIB attack, while SAC series degraded significantly due to the severe occlusion of its completion mechanism (visualized in Appendix E.6). REIN-EAD achieves the best results in both clean and adversarial conditions, thanks to recurrent
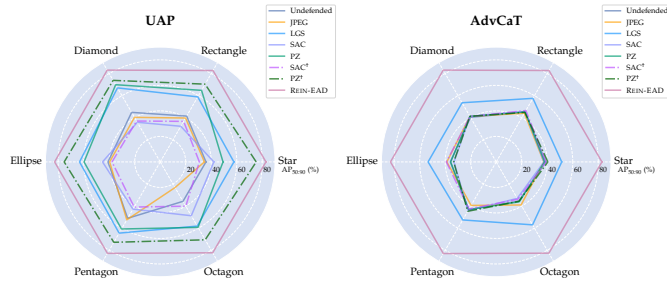
Fig. 10: Comparative evaluation of object detection under different white-box attacks and patch shapes on CARLA.

temporal feedback. Fig. 9 illustrates REIN-EAD's defending process under SIB and AdvCaT attacks, showing improved prediction accuracy and bounding box precision with interaction and environment observation. More qualitative results are presented in Appendix E.6. Furthermore, we discuss some failure cases and analyze their implications in Appendix E.9.

**Generalization on patch shapes.** We evaluate a rectangle-trained REIN-EAD with adversarial patches of 6 varying shapes while maintaining a fixed patch area occupancy. As demonstrated in Fig. 10, our model surpasses other methods in both clean and robust accuracy, even when faced with diverse, unencountered adversarial patch shapes. More results are provided in Appendix E.8. These findings further underscore the exceptional generalization of REIN-EAD in dealing with unknown adversarial attacks.

## 5 CONCLUSION

In this paper, we introduce Reinforced Embodied Active Defense (REIN-EAD), a novel proactive defensive framework that effectively mitigates adversarial patch attacks in real-world 3D environments. REIN-EAD leverages exploration and interaction with the environment to contextualize environmental information and refine its understanding of the target object. It accumulates multi-step interactions for temporal consistency that balances immediate prediction accuracy with long-term entropy minimization. Moreover, REIN-EAD involves an uncertainty-oriented reward-shaping mechanism to improve efficiency without requiring differentiable environments. Extensive experiments demonstrate that REIN-EAD significantly enhances robustness and generalization and obtains strong applicability in complex tasks.

## REFERENCES

[1] T. B. Brown, D. Mané, A. Roy, M. Abadi, and J. Gilmer, "Adversarial patch," *arXiv preprint arXiv:1712.09665*, 2017.
[2] M. Sharif, S. Bhagavatula, L. Bauer, and M. K. Reiter, "Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition," in *CCS*, 2016, pp. 1528–1540.
[3] Z. Zhu, Y. Zhang, H. Chen, Y. Dong, S. Zhao, W. Ding, J. Zhong, and S. Zheng, "Understanding the robustness of 3d object detection with bird's-eye-view representations in autonomous driving," *arXiv preprint arXiv:2303.17297*, 2023.
[4] Z. Xiao, X. Gao, C. Fu, Y. Dong, W. Gao, X. Zhang, J. Zhou, and J. Zhu, "Improving transferability of adversarial patches on face recognition with generative models," in *CVPR*, 2021, pp. 11 845–11 854.
[5] X. Yang, L. Xu, T. Pang, Y. Dong, Y. Wang, H. Su, and J. Zhu, "Face3dadv: Exploiting robust adversarial 3d patches on physical face recognition," *in IJCV*, pp. 1–19, 2024.

[6] D. Song, K. Eykholt, I. Evtimov, E. Fernandes, B. Li, A. Rahmati, F. Tramer, A. Prakash, and T. Kohno, "Physical adversarial examples for object detectors," in *WOOT*, 2018.
[7] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," *arXiv preprint arXiv:1706.06083*, 2017.
[8] T. Wu, L. Tong, and Y. Vorobeychik, "Defending against physically realizable attacks on image classification," *arXiv preprint arXiv:1909.09552*, 2019.
[9] S. Rao, D. Stutz, and B. Schiele, "Adversarial training against location-optimized adversarial patches," in *ECCV*. Springer, 2020, pp. 429–448.
[10] S. Gowal, S.-A. Rebuffi, O. Wiles, F. Stimberg, D. A. Calian, and T. A. Mann, "Improving robustness using generated data," *in NeurIPS*, vol. 34, pp. 4218–4233, 2021.
[11] C. Xiang, A. N. Bhagoji, V. Sehwag, and P. Mittal, "Patchguard: A provably robust defense against adversarial patches via small receptive fields and masking." in *USENIX Security Symposium*, 2021, pp. 2237–2254.
[12] J. Liu, A. Levine, C. P. Lau, R. Chellappa, and S. Feizi, "Segment and complete: Defending object detectors against adversarial patch attacks with robust patch detection," in *CVPR*, 2022, pp. 14 973–14 982.
[13] K. Xu, Y. Xiao, Z. Zheng, K. Cai, and R. Nevatia, "Patchzero: Defending against adversarial patch attacks by detecting and zeroing the patch," in *WACV*, 2023, pp. 4632–4641.
[14] A. Athalye, N. Carlini, and D. Wagner, "Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples," in *ICML*. PMLR, 2018, pp. 274–283.
[15] F. Tramer, N. Carlini, W. Brendel, and A. Madry, "On adaptive attacks to adversarial example defenses," *inNeurIPS*, vol. 33, pp. 1633–1645, 2020.
[16] N. J. Thomas, "Are theories of imagery theories of imagination? an active perception approach to conscious mental content," *Cognitive science*, vol. 23, no. 2, pp. 207–245, 1999.
[17] G. Elsayed, S. Shankar, B. Cheung, N. Papernot, A. Kurakin, I. Goodfellow, and J. Sohl-Dickstein, "Adversarial examples that fool both computer vision and time-limited humans," *in NeurIPS*, vol. 31, 2018.
[18] L. Wu, X. Yang, Y. Dong, X. Liuwei, H. Su, and J. Zhu, "Embodied active defense: Leveraging recurrent feedback to counter adversarial patches," in *ICLR*, 2024.
[19] K. Eykholt, I. Evtimov, E. Fernandes, B. Li, A. Rahmati, C. Xiao, A. Prakash, T. Kohno, and D. Song, "Robust physical-world attacks on deep learning visual classification," in *CVPR*, 2018, pp. 1625–1634.
[20] X. Yang, C. Liu, L. Xu, Y. Wang, Y. Dong, N. Chen, H. Su, and J. Zhu, "Towards effective adversarial textured 3d meshes on physical face recognition," in *CVPR*, 2023, pp. 4119–4128.
[21] G. K. Dziugaite, Z. Ghahramani, and D. M. Roy, "A study of the effect of jpg compression on adversarial images," *arXiv preprint arXiv:1608.00853*, 2016.
[22] M. Naseer, S. Khan, and F. Porikli, "Local gradients smoothing: Defense against localized adversarial attacks," in *WACV*. IEEE, 2019, pp. 1300–1307.
[23] H. Zhang and J. Wang, "Towards adversarially robust object detection," in *ICCV*, 2019, pp. 421–430.
[24] F. Liu, B. Han, T. Liu, C. Gong, G. Niu, M. Zhou, M. Sugiyama *et al.*, "Probabilistic margins for instance reweighting in adversarial training," *in NeurIPS*, vol. 34, pp. 23 258–23 269, 2021.
[25] C. Yu, B. Han, L. Shen, J. Yu, C. Gong, M. Gong, and T. Liu, "Understanding robust overfitting of adversarial training and beyond," in *ICML*. PMLR, 2022, pp. 25 595–25 610.
[26] M. Li, T. Zhou, Z. Huang, J. Yang, J. Yang, and C. Gong, "Dynamic weighted adversarial learning for semi-supervised classification under intersectional class mismatch," *ACM Transactions*, vol. 20, no. 4, pp. 1–24, 2024.
[27] J. Aloimonos, I. Weiss, and A. Bandyopadhyay, "Active vision," *in IJCV*, vol. 1, pp. 333–356, 1988.
[28] R. Bajcsy, "Active perception," *Proceedings of the IEEE*, vol. 76, no. 8, pp. 966–1005, 1988.
[29] K. Kotar and R. Mottaghi, "Interactron: Embodied adaptive object detection," in *CVPR*, 2022, pp. 14 860–14 869.
[30] S. Ruan, Y. Dong, H. Su, J. Peng, N. Chen, and X. Wei, "Improving viewpoint robustness for visual recognition via adversarial training," *arXiv preprint arXiv:2307.11528*, 2023.

[31] H. Ci, M. Liu, X. Pan, F. Zhong, and Y. Wang, "Proactive multi-camera collaboration for 3d human pose estimation," *arXiv preprint arXiv:2303.03767*, 2023.

[32] X. Ma, S. Yong, Z. Zheng, Q. Li, Y. Liang, S.-C. Zhu, and S. Huang, "Sqa3d: Situated question answering in 3d scenes," *arXiv preprint arXiv:2210.07474*, 2022.

[33] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *MICCAI*. Springer, 2015, pp. 234–241.

[34] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *ICCV*, 2017, pp. 2961–2969.

[35] N. Carlini and D. Wagner, "Towards evaluating the robustness of neural networks," in *S&P*. Ieee, 2017, pp. 39–57.

[36] A. R. Cassandra, L. P. Kaelbling, and M. L. Littman, "Acting optimally in partially observable stochastic domains," in *Aaai*, vol. 94, 1994, pp. 1023–1028.

[37] K. Kar, J. Kubilius, K. Schmidt, E. B. Issa, and J. J. DiCarlo, "Evidence that recurrent circuits are critical to the ventral stream's execution of core object recognition behavior," *Nature neuroscience*, vol. 22, no. 6, pp. 974–983, 2019.

[38] A. v. d. Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," *arXiv preprint arXiv:1807.03748*, 2018.

[39] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *ICML*. PMLR, 2021, pp. 8748–8763.

[40] L. Smith and Y. Gal, "Understanding measures of uncertainty for adversarial example detection," *arXiv preprint arXiv:1803.08533*, 2018.

[41] H. Zhang, H. Chen, D. Boning, and C.-J. Hsieh, "Robust reinforcement learning on state observations with learned optimal adversary," *arXiv preprint arXiv:2101.08452*, 2021.

[42] C. Ying, X. Zhou, H. Su, D. Yan, N. Chen, and J. Zhu, "Towards safe reinforcement learning via constraining conditional value-at-risk," *arXiv preprint arXiv:2206.04436*, 2022.

[43] Y. Dong, S. Ruan, H. Su, C. Kang, X. Wei, and J. Zhu, "Viewfool: Evaluating the robustness of visual recognition to adversarial viewpoints," *in NeurIPS*, vol. 35, pp. 36 789–36 803, 2022.

[44] R. S. Sutton, "Dyna, an integrated architecture for learning, planning, and reacting," *ACM Sigart Bulletin*, vol. 2, no. 4, pp. 160–163, 1991.

[45] M. Blondel and V. Roulet, "The elements of differentiable programming," *arXiv preprint arXiv:2403.14606*, 2024.

[46] Y. Bengio, P. Simard, and P. Frasconi, "Learning long-term dependencies with gradient descent is difficult," *IEEE trans. TNN*, vol. 5, no. 2, pp. 157–166, 1994.

[47] R. Antonova, J. Yang, K. M. Jatavallabhula, and J. Bohg, "Rethinking optimization with differentiable simulation from a global perspective," in *Conference on Robot Learning*. PMLR, 2023, pp. 276–286.

[48] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun, "Carla: An open urban driving simulator," in *Conference on robot learning*. PMLR, 2017, pp. 1–16.

[49] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.

[50] A. Y. Ng, D. Harada, and S. Russell, "Policy invariance under reward transformations: Theory and application to reward shaping," in *Icml*, vol. 99, 1999, pp. 278–287.

[51] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.

[52] E. Wong, L. Rice, and J. Z. Kolter, "Fast is better than free: Revisiting adversarial training," *arXiv preprint arXiv:2001.03994*, 2020.

[53] C. Laidlaw, S. Singla, and S. Feizi, "Perceptual adversarial robustness: Defense against unseen threat models," *arXiv preprint arXiv:2006.12655*, 2020.

[54] E. R. Chan, C. Z. Lin, M. A. Chan, K. Nagano, B. Pan, S. De Mello, O. Gallo, L. J. Guibas, J. Tremblay, S. Khamis *et al.*, "Efficient geometry-aware 3d generative adversarial networks," in *CVPR*, 2022, pp. 16 123–16 133.

[55] J.-Y. Zhu, P. Krähenbühl, E. Shechtman, and A. A. Efros, "Generative visual manipulation on the natural image manifold," in *ECCV*. Springer, 2016, pp. 597–613.

[56] G. B. Huang, M. Ramesh, T. Berg, and E. Learned-Miller, "Labeled faces in the wild: A database for studying face recognition in unconstrained environments," University of Massachusetts, Amherst, Tech. Rep. 07-49, October 2007.

[57] X. Yang and J. Zhu, "Adversarial attacks on face recognition," in *Handbook of Face Recognition*. Springer, 2023, pp. 387–404.

[58] Y. Dong, F. Liao, T. Pang, H. Su, J. Zhu, X. Hu, and J. Li, "Boosting adversarial attacks with momentum," in *CVPR*, 2018, pp. 9185–9193.

[59] A. Athalye, L. Engstrom, A. Ilyas, and K. Kwok, "Synthesizing robust adversarial examples," in *ICML*. PMLR, 2018, pp. 284–293.

[60] I. C. Duta, L. Liu, F. Zhu, and L. Shao, "Improved residual networks for image and video recognition," in *ICPR*. IEEE, 2021, pp. 9415–9422.

[61] H. Wang, Y. Wang, Z. Zhou, X. Ji, D. Gong, J. Zhou, Z. Li, and W. Liu, "Cosface: Large margin cosine loss for deep face recognition," in *CVPR*, 2018, pp. 5265–5274.

[62] Y. Li, L. Li, L. Wang, T. Zhang, and B. Gong, "Nattack: Learning the distributions of adversarial examples for an improved black-box attack on deep neural networks," in *ICML*. PMLR, 2019, pp. 3866–3876.

[63] S. Ghadimi and G. Lan, "Stochastic first-and zeroth-order methods for nonconvex stochastic programming," *SIAM journal on optimization*, vol. 23, no. 4, pp. 2341–2368, 2013.

[64] L. Chen, K. Lu, A. Rajeswaran, K. Lee, A. Grover, M. Laskin, P. Abbeel, A. Srinivas, and I. Mordatch, "Decision transformer: Reinforcement learning via sequence modeling," in *NeurIPS*, vol. 34, pp. 15 084–15 097, 2021.

[65] L. Jing, R. Wang, W. Ren, X. Dong, and C. Zou, "Pad: Patch-agnostic defense against adversarial patch attacks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 24 472–24 481.

[66] C. Kang, Y. Dong, Z. Wang, S. Ruan, Y. Chen, H. Su, and X. Wei, "Diffender: Diffusion-based adversarial defense against patch attacks," in *European Conference on Computer Vision*. Springer, 2024, pp. 130–147.

[67] D. Su, H. Zhang, H. Chen, J. Yi, P.-Y. Chen, and Y. Gao, "Is robustness the cost of accuracy?–a comprehensive study on the robustness of 18 deep image classification models," in *ECCV*, 2018, pp. 631–648.

[68] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," *arXiv preprint arXiv:1312.6199*, 2013.

[69] T. Saito and T. Takahashi, "Comprehensible rendering of 3-d shapes," in *SIGGRAPH*, 1990, pp. 197–206.

[70] H. Kato, D. Beker, M. Morariu, T. Ando, T. Matsuoka, W. Kehl, and A. Gaidon, "Differentiable rendering: A survey," *arXiv preprint arXiv:2006.12057*, 2020.

[71] S. Liu, T. Li, W. Chen, and H. Li, "Soft rasterizer: A differentiable renderer for image-based 3d reasoning," in *ICCV*, 2019, pp. 7708–7717.

[72] N. Ravi, J. Reizenstein, D. Novotny, T. Gordon, W.-Y. Lo, J. Johnson, and G. Gkioxari, "Accelerating 3d deep learning with pytorch3d," *arXiv:2007.08501*, 2020.

[73] M. M. Loper and M. J. Black, "Opendr: An approximate differentiable renderer," in *ECCV*. Springer, 2014, pp. 154–169.

[74] H. Rhodin, N. Robertini, C. Richardt, H.-P. Seidel, and C. Theobalt, "A versatile scene model with differentiable visibility applied to generative pose estimation," in *ICCV*, 2015, pp. 765–773.

[75] T. Wu, J. Zhang, X. Fu, Y. Wang, J. Ren, L. Pan, W. Wu, L. Yang, J. Wang, C. Qian, D. Lin, and Z. Liu, "Omniobject3d: Large-vocabulary 3d object dataset for realistic perception, reconstruction and generation," in *CVPR*, 2023.

[76] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *CVPR*, 2009, pp. 248–255.

[77] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, "Openai gym," *arXiv preprint arXiv:1606.01540*, 2016.

[78] C. Xiao, D. Yang, B. Li, J. Deng, and M. Liu, "Meshadv: Adversarial meshes for visual recognition," in *CVPR*, 2019, pp. 6891–6900.

[79] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," *ICCV*, pp. 9992–10 002, 2021.

[80] Y. Zhao, H. Zhu, R. Liang, Q. Shen, S. Zhang, and K. Chen, "Seeing isn't believing: Towards more robust adversarial attack against real world object detectors," in *CCS*, 2019, pp. 1989–2004.

[81] S.-M. Moosavi-Dezfooli, A. Fawzi, O. Fawzi, and P. Frossard, "Universal adversarial perturbations," in *CVPR*, 2017, pp. 86–94.

[82] Z. Hu, W. Chu, X. Zhu, H. Zhang, B. Zhang, and X. Hu, "Physically realizable natural-looking clothing textures evade person detectors via 3d modeling," in *CVPR*, 2023, pp. 16 975–16 984.

[83] G. Jocher, A. Stoken, J. Borovec, A. Chaurasia, L. Changyu, A. Hogan, J. Hajek, L. Diaconu, Y. Kwon, Y. Defretin *et al.*, "ultralytics/yolov5: v5. 0-yolov5-p6 1280 models, aws, supervise. ly and youtube integrations," *Zenodo*, 2021.

[84] D. Barber and F. Agakov, "The im algorithm: a variational approach to information maximization," *in NeurIPS*, vol. 16, no. 320, p. 201, 2004.

[85] Y. Polyanskiy and Y. Wu, "Lecture notes on information theory," *Lecture Notes for ECE563 (UIUC) and*, vol. 6, no. 2012-2016, p. 7, 2014.

[86] Z. Liu, P. Luo, X. Wang, and X. Tang, "Large-scale celebfaces attributes (celeba) dataset," *Retrieved August*, vol. 15, no. 2018, p. 11, 2018.

[87] J. Deng, J. Guo, X. Niannan, and S. Zafeiriou, "Arcface: Additive angular margin loss for deep face recognition," in *CVPR*, 2019.

[88] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[89] Y. Bengio, N. Léonard, and A. Courville, "Estimating or propagating gradients through stochastic neurons for conditional computation," *arXiv preprint arXiv:1308.3432*, 2013.

[90] S. Sabour, Y. Cao, F. Faghri, and D. J. Fleet, "Adversarial manipulation of deep representations," *arXiv preprint arXiv:1511.05122*, 2015.

[91] T. Chen, B. Xu, C. Zhang, and C. Guestrin, "Training deep nets with sublinear memory cost," *arXiv preprint arXiv:1604.06174*, 2016.

[92] Y. Guo, L. Zhang, Y. Hu, X. He, and J. Gao, "Ms-celeb-1m: A dataset and benchmark for large-scale face recognition," in *ECCV*. Springer, 2016, pp. 87–102.

[93] S. Thys, W. Van Ranst, and T. Goedemé, "Fooling automated surveillance cameras: adversarial patches to attack person detection," in *CVPR workshops*, 2019, pp. 0–0.

# APPENDIX A
# PROOFS AND ADDITIONAL THEORY

## A.1 Proof of Theorem 3.1

*Proof.* For a series of observations $\{o_1, \cdots o_t\}$ and a previously maintained belief $b_{t-1}$ determined by the scene $x$, we expand the left-hand side of Eq. (6) as follows:

$$
\begin{aligned}
\mathbb{E}_x I(o_t; y \mid b_{t-1}) &= \mathbb{E}_{x,y} \log \frac{p(b_{t-1})p(b_{t-1}, o_t, y)}{p(b_{t-1}, y)p(b_{t-1}, o_t)} \\
&= \mathbb{E}_{x,y} \log \frac{p(y \mid b_{t-1}, o_t)}{p(y \mid b_{t-1})}.
\end{aligned}
\tag{A.12}
$$

By introducing the variational distribution $q_\theta(y \mid o_1, \cdots o_t)$ as a multiplicative factor in the integrand of Eq. (A.12), we obtain:

$$
\begin{aligned}
\mathbb{E}_x I(o_t; y \mid b_{t-1}) &= \mathbb{E}_{x,y} \log \frac{p(y \mid b_{t-1}, o_t)q_\theta(y \mid b_{t-1}o_t)}{p(y \mid b_{t-1})q_\theta(y \mid b_{t-1}, o_t)} \\
&= \mathbb{E}_{x,y} \log \frac{q_\theta(y \mid b_{t-1}, o_t)}{p(y \mid b_{t-1})} \\
&\quad + \mathbb{E}_x D_{\text{KL}}(p(y \mid b_{t-1}, o_t) \| q_\theta(y \mid b_{t-1}, o_t)).
\end{aligned}
\tag{A.13}
$$

The non-negativity property of the KL-divergence allows us to establish a lower bound for the mutual information:

$$
\begin{aligned}
\mathbb{E}_x I(o_t; y \mid b_{t-1}) &\geq \mathbb{E}_{x,y} \log \frac{q_\theta(y \mid b_{t-1}, o_t)}{p(y \mid b_{t-1})} \\
&= \mathbb{E}_{x,y} \log q_\theta(y \mid b_{t-1}, o_t) + \mathcal{H}(y \mid b_{t-1}),
\end{aligned}
\tag{A.14}
$$

where $\mathcal{H}(y \mid b_{t-1})$ represents the conditional entropy of $y$ given the belief $b_{t-1}$, and Eq. (A.14) is the well-known Barber and Agakov bound [84]. We proceed by selecting an energy-based variational family that incorporates a *critic* $\mathcal{E}_\theta(bt-1, o_t, y)$ and is scaled by the data density $p(b_{t-1}, o_t)$:

$$
q_\theta(y \mid b_{t-1}, o_t) = \frac{p(y \mid b_{t-1})}{Z(b_{t-1}, o_t)} e^{\mathcal{E}_\theta(b_{t-1}, o_t, y)},
\tag{A.15}
$$

where $Z(b_{t-1}, o_t) = \mathbb{E}_y e^{\mathcal{E}_\theta(b_{t-1}, o_t, y)}$. Substituting the distribution defined in Eq. (A.15) into Eq. (A.14) yields:

$$
\begin{aligned}
&\mathbb{E}_x I(o_t; y \mid b_{t-1}) \\
&\geq \mathbb{E}_{x,y} \log q_\theta(y \mid b_{t-1}, o_t) + \mathcal{H}(y \mid b_{t-1}) \\
&= \mathbb{E}_{x,y}[\mathcal{E}_\theta(b_{t-1}, o_t, y)] - \mathbb{E}_x[\log Z(b_{t-1}, o_t)],
\end{aligned}
\tag{A.16}
$$

which represents the unnormalized version of the Barber and Agakov bound. Applying the inequality $\log Z(b_{t-1}, o_t) \leq \frac{Z(b_{t-1}, o_t)}{g(b_{t-1}, o_t)} + \log[g(b_{t-1}, o_t)] - 1$ for any $g(b_{t-1}, o_t) > 0$, with the bound becoming tight when $g(b_{t-1}, o_t) = Z(b_{t-1}, o_t)$, we arrive at a tractable upper bound known as the tractable unnormalized version of the Barber and Agakov lower bound on mutual information:

$$
\begin{aligned}
&\mathbb{E}_{x,y}[\mathcal{E}_\theta(b_{t-1}, o_t, y)] - \mathbb{E}_x[\log Z(b_{t-1}, o_t)] \\
&\geq \mathbb{E}_{x,y}[\mathcal{E}_\theta(b_{t-1}, o_t, y)] - \mathbb{E}_x\left[\frac{\mathbb{E}_y e^{\mathcal{E}_\theta(b_{t-1}, o_t, y)}}{g(b_{t-1}, o_t)}\right] \\
&\quad - \mathbb{E}_x \log[g(b_{t-1}, o_t)] + 1 \\
&= 1 + \mathbb{E}_{x,y}\left[\log \frac{e^{\mathcal{E}_\theta(b_{t-1}, o_t, y)}}{g(b_{t-1}, o_t)}\right] - \mathbb{E}_x\left[\frac{\mathbb{E}_y e^{\mathcal{E}_\theta(b_{t-1}, o_t, y)}}{g(b_{t-1}, o_t)}\right].
\end{aligned}
\tag{A.17}
$$

To mitigate variance, we leverage multiple samples $\{x^{(j)}, y^{(j)}\}j = 1^K$ from $\mathcal{D}$ to implement a low-variance, high-bias estimation of the mutual information. For an observation trajectory originating from a different scene $(x^{(j)}, y^{(j)})$ $(j \neq i)$, with annotations $\{y^{(j)}\}j = 1, j \neq i^K$ independent of $x^{(i)}$ and $y^{(i)}$, we have:

$$
g(b_{t-1}, o_t) = g(b_{t-1}, o_t; y^{(1)}, \cdots, y^{(K)}).
\tag{A.18}
$$

This enables us to utilize the additional samples $\{x^{(j)}, y^{(j)}\}j = 1^K$ to construct a Monte-Carlo estimate of the function $Z(bt - 1, o_t)$:

$$
\begin{aligned}
g(b_{t-1}, o_t; y_1, \cdots y_K) &= m(b_{t-1}, o_t; y^{(1)}, \cdots y^{(K)}) \\
&= \frac{1}{K} \sum_{j=1}^{K} e^{\mathcal{E}_\theta(y^{(j)} \mid b_{t-1}, o_t)}.
\end{aligned}
$$

When estimating the bound over $K$ samples, the last term in Eq. (A.17) reduces to a constant value of 1:

$$
\begin{aligned}
&\mathbb{E}_x\left[\frac{\mathbb{E}_{y^{(1)}, \cdots, y^{(K)}} e^{\mathcal{E}_\theta(b_{t-1}, o_t, y)}}{m(b_{t-1}, o_t; y^{(1)}, \cdots y^{(K)})}\right] \\
&= \mathbb{E}_{x_1}\left[\frac{\frac{1}{K}\sum_{j=1}^K e^{\mathcal{E}_\theta(b_{t-1}, o_t^{(1)}, y^{(j)})}}{m(b_{t-1}, o_t^{(1)}; y^{(1)}, \cdots y^{(K)})}\right] = 1.
\end{aligned}
\tag{A.19}
$$

Applying Eq. (A.19) back to Eq. (A.17) and averaging the bound over $K$ samples (reindexing $x^{(1)}$ as $x^{(j)}$ for each term), we precisely recover the lower bound on mutual information proposed by Oord *et al.* [38]:

$$
\begin{aligned}
&1 + \mathbb{E}_{x^{(j)}, y^{(j)}}\left[\log \frac{e^{\mathcal{E}_\theta(b_{t-1}^{(j)}, y^{(j)})}}{g(b_{t-1}^{(j)}; y^{(1)}, \cdots, y^{(K)})}\right] \\
&\quad - \mathbb{E}_{x^{(j)}}\left[\frac{\mathbb{E}_{y^{(j)}} e^{\mathcal{E}_\theta(b_{t-1}^{(j)}, o_t^{(j)}, \hat{y}^{(j)})}}{g(b_{t-1}^{(j)}, o_t^{(j)}; y^{(1)}, \cdots, (K))}\right] \\
&= \mathbb{E}_{x^{(j)}, y^{(j)}}\left[\log \frac{e^{\mathcal{E}_\theta(b_{t-1}^{(j)}, o_t^{(j)}, y^{(j)})}}{g(b_{t-1}^{(j)}, o_t^{(j)}; y^{(1)}, \cdots, y^{(K)})}\right] \\
&= \mathbb{E}_{x^{(j)}, y^{(j)}}\left[\log \frac{e^{\mathcal{E}_\theta(b_{t-1}^{(j)}, o_t^{(j)}, y^{(j)})}}{\frac{1}{K}\sum_{\hat{y}^{(j)}} e^{\mathcal{E}_\theta(b_{t-1}^{(j)}, o_t^{(j)}, \hat{y}^{(j)})}}\right].
\end{aligned}
\tag{A.20}
$$

Multiplying and dividing the integrand in Eq. (A.12) by $\frac{p(y \mid b_{t-1}^{(j)}, o_{t-1})}{Z(b_{t-1}^{(j)}, o_t)}$ and extracting $\frac{1}{K}$ from the brackets transforms the equation into:

$$
\mathbb{E}_{x^{(j)}, y^{(j)}}\left[\log \frac{q_\theta(b_{t-1}^{(j)}, o_t^{(j)}, y^{(j)})}{\sum_{\hat{y}^{(j)}} q_\theta(b_{t-1}^{(j)}, o_t^{(j)}, \hat{y}^{(j)})}\right] + \log(K).
\tag{A.21}
$$

Therefore, we obtain

$$
\begin{aligned}
&\mathbb{E}_{(x^{(j)}, y^{(j)}) \sim \mathcal{D}}\left[\frac{1}{K}\sum_{j=1}^{K} \log \frac{q_\theta(y^{(j)} \mid b_{t-1}^{(j)}, o_t^{(j)})}{\frac{1}{K}\sum_{k=1}^{K} q_\theta(y^{(k)} \mid b_{t-1}^{(j)}, o_t^{(j)})}\right] \\
&\leq \mathbb{E}_x I(o_t; y \mid b_{t-1}) - \frac{\log(K)}{K}.
\end{aligned}
\tag{A.22}
$$

Given a scene annotation $y$, we quantify the uncertainty of annotation $y$ at time step $t$ using the conditional entropy of $y$ given the series of observations $\{b_{t-1}, o_t\}$, denoted as

$\mathcal{H}(y \mid b_{t-1}, o_t)$. In the following, we demonstrate that the conditional mutual information $I(o_t; y \mid b_{t-1})$ is equivalent to the decrease in conditional entropy:

$$
\begin{aligned}
&I(o_t; y \mid b_{t-1}) \\
=&I(y; b_{t-1}, o_t) - I(y; b_{t-1}) \\
=& \left[\mathcal{H}(y) - I(y; b_{t-1})\right] - \left[\mathcal{H}(y) - I(y; b_{t-1}, o_t)\right] \\
=& \mathcal{H}(y \mid b_{t-1}) - \mathcal{H}(y \mid b_{t-1}, o_t).
\end{aligned}
$$

Here the first equality in the derivation follows from the Kolmogorov identities [85]. $\square$

Theorem 3.1 enables us to forge a connection between the mutual information in Eq. (6) and the greedy informative exploration defined in Definition 3.3. Consequently, we can infer the relationship between the policy model of EAD, which optimizes the InfoNCE objective in Eq. (5), and the principle of greedy informative exploration.

## A.2 Proof of Theorem 3.7

*Proof.* We aim to demonstrate that the accumulative informative policy $\pi^*$ results in a greater or equal reduction in entropy of $y$ compared to the greedy informative policy $\pi^g$, under the condition that the belief update function $f_b$ is bijective.

Recall that the information gain from time 0 to $H$ under a given policy $\pi$ is defined as the reduction in entropy of $y$:

$$\Delta\mathcal{H}_\pi = \mathcal{H}(y) - \mathcal{H}(y \mid b_{H-1}, o_H),$$

where $b_{H-1}$ and $o_H$ are the belief and observation at time $H-1$ and $H$, respectively, under policy $\pi$.

Specifically, for the two policies:

$$\Delta\mathcal{H}_{\pi^*} = \mathcal{H}(y) - \mathcal{H}(y \mid b^*_{H-1}, o^*_H),$$

$$\Delta\mathcal{H}_{\pi^g} = \mathcal{H}(y) - \mathcal{H}(y \mid b^g_{H-1}, o^g_H).$$

We first consider the cumulative information gain over each time step. For each time step $t = 1$ to $H$, define the incremental information gain under policy $\pi$ as:

$$
\begin{aligned}
\Delta\mathcal{H}^\pi_t =& \mathcal{H}(y \mid b^\pi_{t-1}, o^\pi_t) - \mathcal{H}(y \mid b^\pi_t, o^\pi_{t+1}) \\
=& \left[\mathcal{H}(y \mid b^\pi_{t-1}, o^\pi_t) - \mathcal{H}(y \mid b^\pi_t)\right] \\
&+ \left[\mathcal{H}(y \mid b^\pi_t) - \mathcal{H}(y \mid b^\pi_t, o^\pi_{t+1})\right].
\end{aligned}
\tag{A.23}
$$

Since $f_b$ is bijective, each observation $o_t$ uniquely determines the subsequent belief $b_t$, and vice versa. This bijectivity ensures that the mapping between beliefs and observations preserves information about $y$ without loss:

$$\mathcal{H}(y \mid b^\pi_{t-1}, o^\pi_t) - \mathcal{H}(y \mid b^\pi_t) = 0. \tag{A.24}$$

Thereby, we have

$$\Delta\mathcal{H}_\pi = \sum_{t=1}^H \Delta\mathcal{H}^\pi_t = \sum_{t=1}^H \left[\mathcal{H}(y \mid b^\pi_t) - \mathcal{H}(y \mid b^\pi_t, o^\pi_{t+1})\right]. \tag{A.25}$$

To this end, we compare the trajectory information gain between these two policies:

$$
\begin{aligned}
&\Delta\mathcal{H}_{\pi^*} - \Delta\mathcal{H}_{\pi^g} \\
=& \left(\sum_{t=1}^H \Delta\left[\mathcal{H}(y \mid b^*_t) - \mathcal{H}(y \mid b^*_t, o^*_{t+1})\right]\right) \\
&- \left(\sum_{t=1}^H \Delta\left[\mathcal{H}(y \mid b^g_t) - \mathcal{H}(y \mid b^g_t, o^g_{t+1})\right]\right).
\end{aligned}
$$

Since $\pi^*$ maximizes $\Delta\mathcal{H}$ considering the $H$-steps trajectory, we have:

$$
\begin{aligned}
&\left(\sum_{t=1}^H \Delta\left[\mathcal{H}(y \mid b^*_t) - \mathcal{H}(y \mid b^*_t, o^*_{t+1})\right]\right) \\
&- \left(\sum_{t=1}^H \Delta\left[\mathcal{H}(y \mid b^g_t) - \mathcal{H}(y \mid b^g_t, o^g_{t+1})\right]\right) \geq 0,
\end{aligned}
\tag{A.26}
$$

with equality if and only if for all $t$, $\Delta\mathcal{H}^*_t = \Delta\mathcal{H}^g_t$, which occurs precisely when the problem exhibits optimal substructure. In such cases, the greedy policy $\pi^g$ inherently accumulates information gains as effectively as the accumulative policy $\pi^*$. $\square$

## A.3 Derivation of Reward

Given the definition of the dense reward $r_t$ from Eq. (11), we substitute into the expression for trajectory reward $\mathcal{R}(\tau)$:

$$
\begin{aligned}
\mathcal{R}(\tau) &= \sum_{t=1}^H \gamma^{t-1}\left[\mathcal{L}(\hat{y}_{t-1}, y) - \gamma \cdot \mathcal{L}(\hat{y}_t, y)\right] \\
&= \sum_{t=1}^H \gamma^{t-1}\mathcal{L}(\hat{y}_{t-1}, y) - \sum_{t=1}^H \gamma^t\mathcal{L}(\hat{y}_t, y).
\end{aligned}
\tag{A.27}
$$

Observe that the second summation now spans from $t = 1$ to $t = H$, which matches the first summation's indexing shifted by one. Therefore, the two summations can be aligned as follows:

$$
\begin{aligned}
\mathcal{R}(\tau) =& \gamma^0\mathcal{L}(\hat{y}_0, y) + \gamma^1\mathcal{L}(\hat{y}_1, y) + \cdots + \gamma^{H-1}\mathcal{L}(\hat{y}_{H-1}, y) \\
&- \left(\gamma^1\mathcal{L}(\hat{y}_1, y) + \gamma^2\mathcal{L}(\hat{y}_2, y) + \cdots + \gamma^H\mathcal{L}(\hat{y}_H, y)\right).
\end{aligned}
$$

All intermediate terms cancel out due to the telescoping nature of the series:

$$\mathcal{R}(\tau) = \mathcal{L}(\hat{y}_0, y) - \gamma^H\mathcal{L}(\hat{y}_H, y). \tag{A.28}$$

Thus, the cumulative discounted reward $\mathcal{R}(\tau)$ is equivalent to the initial loss minus the discounted final loss. Since the initial loss $\mathcal{L}(\hat{y}_0, y)$ is solely determined by the initial state distribution $\rho_0$ and is independent of the policy, it can be considered a constant during policy optimization. This completes the proof that this form of reward aligns with the objective in Eq.( 9).

# APPENDIX B
## EXPERIMENT DETAILS FOR SIMULATION ENVIRONMENT

**Environmental dynamics.** We formally define the state $s_t = (h_t, v_t)$ as a combination of the camera's yaw $h_t \in \mathbb{R}$ and pitch $v_t \in \mathbb{R}$ at moment $t$, while the action is defined as a continuous rotation denoted by $a_t = (\Delta h, \Delta v)$. Consequently, the transition function is expressed as $T(s_t, a_t, x) = s_t + a_t$, indicating that the next state is obtained by adding the action's rotation to the current state. The observation function is reformulated using a 3D generative model, denoted as $O(s_t, x) = \mathcal{R}(s_t, x)$, where $\mathcal{R}(\cdot)$ represents a renderer (*e.g.*, a 3D generative model or graphic engine) that generates a 2D image observation $o_t$ given the camera parameters determined by the state $s_t$.

In computational graphics, the detailed formulation for the renderer is presented as:

$$o_t = \mathcal{R}'(\boldsymbol{E}_t, \boldsymbol{I}, x), \tag{B.29}$$

where $\boldsymbol{E}_t \in \mathbb{R}^{4 \times 4}$ represents the camera's extrinsic matrix determined by the state $s_t$, and $\boldsymbol{I} \in \mathbb{R}^{3 \times 3}$ is the pre-defined camera intrinsic matrix. To utilize the renderer for generating 2D images, we need to calculate the camera's extrinsic matrix $\boldsymbol{E}_t$ based on the state $s_t$. Assuming a right-handed coordinate system and column vectors, we have:

$$\boldsymbol{E}_t = \begin{bmatrix} \boldsymbol{R}_t & \boldsymbol{T} \\ \boldsymbol{0} & 1 \end{bmatrix}, \tag{B.30}$$

where $\boldsymbol{R}_t \in \mathbb{R}^{3 \times 3}$ is the rotation matrix determined by $s_t$, and $\boldsymbol{T} \in \mathbb{R}^{3 \times 1}$ is the invariant translation vector. The rotation matrices for yaw $h_t$ and pitch $v_t$ are given by:

$$\boldsymbol{R}^y(h_t) = \begin{bmatrix} \cos(h_t) & 0 & \sin(h_t) \\ 0 & 1 & 0 \\ -\sin(h_t) & 0 & \cos(h_t) \end{bmatrix}, \tag{B.31}$$

$$\boldsymbol{R}^x(v_t) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos(v_t) & -\sin(v_t) \\ 0 & \sin(v_t) & \cos(v_t) \end{bmatrix}. \tag{B.32}$$

The combined rotation $\boldsymbol{R}_t$ is obtained by multiplying the individual rotation matrices: $\boldsymbol{R}_t = \boldsymbol{R}^y(h_t) \times \boldsymbol{R}^x(v_t)$. The complete extrinsic matrix is then constructed as:

$$\boldsymbol{E}_t = \begin{bmatrix} \boldsymbol{R}^y(h_t) \times \boldsymbol{R}^x(v_t) & \boldsymbol{T} \\ \boldsymbol{0} & 1 \end{bmatrix}. \tag{B.33}$$

**Applying function.** In the experiments, the adversarial patch is attached to a flat surface, such as eyeglasses for face recognition and billboards for object detection. By utilizing the known corner coordinates of the adversarial patch in the world coordinate system, both the extrinsic matrix $\boldsymbol{E}_t$ and the intrinsic matrix $\boldsymbol{K}$ are employed to render image observations containing the adversarial patch. The projection process of the 3D patch, as described by Zhu *et al.* [3], is followed to construct the applying function. The projection matrix $\boldsymbol{M}_{\text{3d-2d}} \in \mathbb{R}^{4 \times 4}$ is specified as:

$$\boldsymbol{M}_{\text{3d-2d}} = \begin{bmatrix} \boldsymbol{K} & \boldsymbol{0} \\ \boldsymbol{0} & 1 \end{bmatrix} \times \boldsymbol{E}_t. \tag{B.34}$$

This process is differentiable, enabling the optimization of the adversarial patches.

TABLE C.1: quantitative evaluation for CelebA-3D. The image size is $112 \times 112$.

| | PSNR↑ | SSIM↑ | LPIPS↓ | ID↑ |
|---|---|---|---|---|
| CelebA-3D | 21.28 | .7601 | .1314 | .5771 |

In summary, a deterministic environmental model is proposed, which is applicable to all the experimental environments (*e.g.*, EG3D, CARLA):

$$
\begin{aligned}
\text{State} \quad & s_t = (h_t, v_t) \in \mathbb{R}^2, \\
\text{Action} \quad & a_t = (\Delta h, \Delta v) \in \mathbb{R}^2, \\
\text{Transition Function} \quad & T(s_t, a_t, x) = s_t + a_t, \\
\text{Observation Function} \quad & Z(s_t, x) = \mathcal{R}(s_t, x).
\end{aligned}
$$

The primary distinction between simulations for different tasks lies in the feasible viewpoint regions. These regions are detailed in the implementation sections for each task, specifically in Appendices C, D & E.

# APPENDIX C
## EXPERIMENT DETAILS FOR FACE RECOGNITION
### C.1 CelebA-3D

We employ an unofficial implementation of GAN Inversion with EG3D (https://github.com/oneThousand1000/EG3D-projector), utilizing default parameters to transform 2D images from the CelebA dataset [86] into 3D latent representation $w^+$. As the 3D generative model prior, we leverage the EG3D models pre-trained on the FFHQ dataset, which are officially released at https://catalog.ngc.nvidia.com/orgs/nvidia/teams/research/models/eg3d. To reduce computational overhead, we forgo the super-resolution module of EG3D and directly render RGB images with a resolution of $112 \times 112$ pixels using its neural renderer.

We conduct a comprehensive evaluation of the reconstructed CelebA-3D dataset to assess its quality. Image quality is measured using PSNR, SSIM and LPIPS between the original images and the EG3D-rendered images from the same viewpoint. These metrics provide a multi-faceted assessment of the fidelity of the reconstructed images compared to their 2D counterparts. Furthermore, we introduce a modified identity consistency (ID) metric, slightly deviating from the one presented in [54], to evaluate the identity consistency between the reconstructed 3D faces and their original 2D faces. Our ID metric calculates the mean ArcFace [87] cosine similarity score between pairs of views of the face rendered from random camera poses and its original image from CelebA.

The results presented in Table C.1 demonstrate that the learned 3D prior over FFHQ enables remarkably high-quality single-view geometry recovery. Consequently, our reconstructed CelebA-3D dataset exhibits high image quality and sufficient identity consistency with its original 2D form, making it suitable for subsequent experiments. A selection of reconstructed multi-view faces is shown in Fig. C.1.

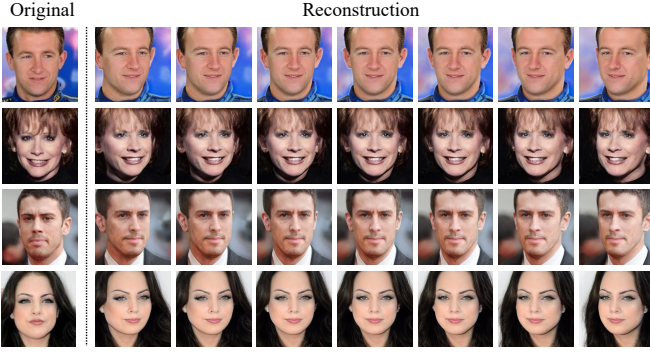The CelebA-3D dataset incorporates the rich annotations provided by the widely-used CelebA dataset, which can be accessed at https://mmlab.ie.cuhk.edu.hk/projects/CelebA.html

Original    Reconstruction



Fig. C.1: The first column presents the source face from CelebA, and the succeeding columns demonstrate the rendered multiview faces from inverted $w^+$ with EG3D. The dimensions of each rendered facial image are $112 \times 112$.

## C.2 Details for Attacks

**Impersonation and dodging in adversarial attacks against FR system.** In adversarial attacks targeting FR systems, two primary objectives emerge impersonation and dodging. Impersonation entails the manipulation of an image to deceive FR algorithms into erroneously identifying an individual as a different person. Dodging focuses on preventing accurate identification by the FR system through strategic modifications to the image. These alterations are designed to either prevent the FR system from detecting a face altogether or to obscure the association between the detected face and its corresponding authentic identity within the system's database. The nature of these subtasks underscores the profound challenges they pose to the security and integrity of FR systems.

**Attacks in pixel space.** The Momentum Iterative Method (MIM) [58] and Expectation over Transformation (EoT) [59] represent state-of-the-art techniques for refining adversarial patches within the RGB pixel space. MIM enhances the generation of adversarial examples by incorporating momentum into the optimization process, thereby facilitating the transferability of these examples across different models and settings. In contrast, EoT employs a diverse array of transformations, such as rotations and variations in illumination, to boost the robustness of attacks under physical conditions. To ensure optimal performance, we adhered to the recommended parameters as detailed in [4], setting the number of iterations to $N = 150$, the learning rate to $\alpha = 1.5/255$, and the decay factor to $\mu = 1$. These parameters are consistently maintained across all experimental conditions. Furthermore, the sampling frequency for EoT is established at $M = 10$ to balance the computational efficiency and attack effectiveness.

**Attacks in latent space of the generative model.** GenAP [4] and Face3DAdv [5] represent pioneering approaches that shift the focus of adversarial patches from the pixel space to the latent space of a Generative Adversarial Network (GAN). By utilizing the generative capabilities, GANs develop adversarial patches that can deceive facial recognition systems, accounting for 3D variations. Additionally, we employ the Adam optimizer [88] for the latent space of EG3D for patch optimization, with a learning rate of $\eta = 0.01$ and an iteration count of $N = 150$. The sampling frequency is set at $M = 10$.

## C.3 Details for Adaptive Attacks

**Adaptive attack for defense baselines.** To launch adaptive attacks against parameter-free, purification-based defenses such as JPEG and LGS, we employ Backward Pass Differentiable Approximation (BPDA) as proposed by Athalye *et al.* [14]. This method assumes that the output from each defense mechanism closely approximates the original input. For adaptive attacks on SAC and PZ, we use their official implementations [12], incorporating Straight-Through Estimators (STE) [89] for backpropagation through thresholding operations.

**Adaptive attack with uniform superset policy.** In adaptive attacks for REIN-EAD, we leverage uniform superset approximation for the policy model. Thus, we have the surrogate policy

$$\tilde{\pi} := \mathcal{U}(h_{\min}, h_{\max}) \times \mathcal{U}(v_{\min}, v_{\max}), \qquad \text{(C.35)}$$

where $a_t \in [h_{\min}, h_{\max}] \times [v_{\min}, v_{\max}]$, and $[h_{\min}, h_{\max}]$, $[v_{\min}, v_{\max}]$ separately denotes the pre-defined feasible region for horizontal rotation (yaw) and vertical rotation (pitch). The optimization objective is outlined as follows, with a simplified sequential representation for clarity[4]:

$$\max_{p} \quad \mathbb{E}_{s_0 \sim \rho_0, a_i \sim \tilde{\pi}} \mathcal{L}(\overline{y}_\tau, y),$$

$$\text{with} \quad \{\overline{y}_\tau, b_\tau\} = f(\{A(o_i, p; s_0 + \sum_{j=0}^{i-1} a_j)\}_{i=0}^{\tau}; \boldsymbol{\theta}) \quad \text{(C.36)}$$

$$\text{s.t.} \quad p \in [0, 1]^{H_p \times W_p \times C},$$

where $\rho_0$ denotes the distribution of initial state $s_0$, and $\mathcal{L}$ is the task-specific loss function.

**Adaptive Attack Against Sub-Modules.** An end-to-end attack may not always be the most effective strategy, particularly against defenses with complex forward passes. Targeting the weakest component is often sufficient. Therefore, we propose two separate adaptive attacks: one against the perception model and another against the policy model. The attack on the perception model aims to generate an adversarial patch that corrupts the internal belief $b_t$ [90]. The optimization objective for this attack is to maximize the Euclidean distance between the corrupted belief $b\tau$ and the benign belief $b\tau^+$, formulated as follows:

$$\max_{p} \quad \mathbb{E}_{s_0 \sim \rho_0, a_i \sim \tilde{\pi}} \|b_\tau - b_\tau^+\|_2^2,$$

$$\text{with} \quad \{\overline{y}_\tau, b_\tau\} = f(\{A(o_i, p; s_0 + \sum_{j=0}^{i-1} a_j)\}_{i=1}^{\tau}; \boldsymbol{\theta}), \quad \text{(C.37)}$$

$$\{\overline{y}_\tau^+, b_\tau^+\} = f(\{o_i\}_{i=0}^{\tau}; \boldsymbol{\theta}),$$

$$\text{s.t.} \quad p \in [0, 1]^{H_p \times W_p \times C}.$$

For the attack against the policy model, the goal is to create an adversarial patch that induces the policy model to output a zero action $a_i = \pi(b_i; \boldsymbol{\phi}) = 0$, thereby keeping the model stationary with an invariant state $s_i = s_1$ and generating erroneous predictions $\overline{y}_\tau$. While the original problem can be challenging with policy output as a constraint, we employ Lagrangian relaxation to incorporate the constraint into the objective and address the following problem:

---

4. The recurrent inference procedure is presented sequentially in this section for simplicity.

$$\max_p \quad \mathbb{E}_{s_0 \sim \rho_0} \mathcal{L}(\overline{y}_\tau, y) + c \cdot \|\pi(\{A(o_i, p; s_i)\}_{i=0}^\tau; \boldsymbol{\phi})\|_2^2,$$

$$\text{with} \quad \{\overline{y}_\tau, b_\tau\} = f(\{A(o_i, p; s_i)\}_{i=0}^\tau; \boldsymbol{\theta})$$

$$\text{s.t.} \quad p \in [0, 1]^{H_p \times W_p \times C},$$

$$(C.38)$$

where $c > 0$ is a constant that yields an adversarial example ensuring the model outputs zero actions.

**Adaptive Attack for the Entire Pipeline.** Attacking the model through backpropagation is infeasible due to the rapid consumption of GPU memory as trajectory length increases (*e.g.*, 4 steps require nearly 90 GB of video memory). To mitigate this, we use gradient checkpointing [91] to reduce memory consumption. By selectively recomputing parts of the computation graph defined by the $H$-step inference procedure, instead of storing them, this technique effectively reduces memory costs at the expense of additional computation. Using this method, we successfully attack the entire pipeline along a 4-step trajectory using an NVIDIA RTX 3090 Ti, but still fail to extend it to attack 16-step REIN-EAD.

Regarding implementation, we adopt the same hyper-parameters as Face3DAdv and consider the action bounds defined in Appendix C.5. For the constant $c$ in the adaptive attack against the policy model, we employ a bisection search to identify the optimal value as per Carlini *et al.* [35], finding $c = 100$ to be the most effective. Additionally, the evaluation results from these adaptive attacks and analysis are detailed in Appendix C.7 with main results in Table C.6.

## C.4 Details for Defenses

**JPEG compression.** we set the quality parameter to 75.
**Local gradients smoothing.** We adopt the implementation at https://github.com/fabiobrau/local_gradients_smoothing, and maintain the default hyper-parameters claimed in [22].
**Segment and complete.** We use official implementation at https://github.com/joellliu/SegmentAndComplete, and retrain the patch segmenter with adversarial patches optimized by EoT [14] and USAP technique separately. We adopt the same hyper-parameters and training process claimed in [12], except for the prior patch sizes, which we resize them proportionally to the input image size.
**Patchzero.** For Patchzero [13], we directly utilize the trained patch segmenter of SAC, for they share almost the same training pipeline.
**Defense against Occlusion Attacks.** DOA is an adversarial training-based method. We adopt the DOA training paradigm to fine-tune the same pre-trained IResNet-50 and training data as REIN-EAD with the code at https://github.com/P2333/Bag-of-Tricks-for-AT. As for hyper-parameters, we adopt the default ones in [8] with the training patch size scaled proportionally.

## C.5 Details for Implementations

**Model details.** In the context of face recognition, we implement REIN-EAD model with a composition of a pre-trained face recognition feature extractor and a Decision transformer, specifically, we select IResNet-50 as the the visual backbone to extract feature. For each time step $t > 0$, we use the IResNet-50 to map the current observation input of dimensions $112 \times 112 \times 3$ into an embedding with a dimensionality of $512$. This embedding is then concatenated with the previously extracted embedding sequence of dimensions $(t - 1) \times 512$, thus forming the temporal sequence of observation embeddings ($t \times 512$) for Decision Transformer input. The Decision Transformer subsequently outputs the temporal-fused face embedding for inference as well as the predicted action. For the training process, we further map the fused face embedding into logits using a linear projection layer. For the sake of simplicity, we directly employ the Softmax loss function for model training.

In experiments, we use pre-trained IResNet-50 with ArcFace margin on MS1MV3 [92] from `InsightFace` [87], which is available at https://github.com/deepinsight/insightface/tree/master/model_zoo.
**Curriculum training.** It's observed that the training suffers from considerable instability when simultaneously training perception and policy models from scratch. A primary concern is that the perception model, in its early training stages, cannot provide accurate supervision signals, leading the policy network to generate irrational actions and hindering the overall learning process. To mitigate this issue, we initially train the perception model independently using frames obtained from a random action policy, namely *offline phase*. Once achieving a stable performance from the perception model, we proceed to the *online phase* and jointly train both the perception and policy networks, employing Algorithm 1,thereby ensuring their effective coordination and learning.

Meanwhile, learning offline with pre-collected data in the first phase proves to be significantly more efficient than online learning through interactive data collection from the environment. By dividing the training process into two distinct phases *offline* and *online*, we substantially enhance training efficiency and reduce computational costs.
**Training details.** To train REIN-EAD for face recognition, we randomly sample images from $2,500$ distinct identities from the training set of CelebA-3D. we adopt the previously demonstrated phased training paradigm with hyper-parameters listed in Table C.2.

## C.6 Computational Overhead

This section evaluates our method's computational overhead compared to other passive defense baselines in facial recognition systems. Performance evaluation is performed on an NVIDIA GeForce RTX 3090 Ti and an AMD EPYC 7302 16-Core Processor, using a training batch size of 64. SAC and PZ necessitate training a segmenter to identify the patch area, entailing two stages: initial training with pre-generated adversarial images and subsequent self-adversarial training [12], [13]. DOA, an adversarial training-based approach, requires retraining the feature extractor [8]. Additionally, REIN-EAD's training involves offline and online phases, without involving adversarial training.

As indicated in Table C.3, although differential rendering imposes significant computational demands during the online training phase, the total training time of our REIN-EAD model is effectively balanced between the pure adversarial

TABLE C.2: Hyper-parameters of REIN-EAD for face recognition.

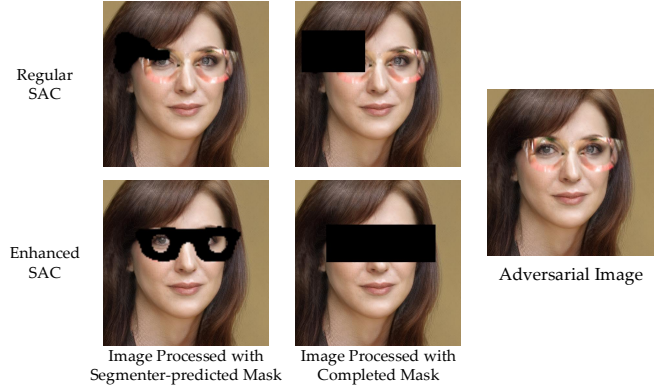| Hyper-parameter | Value |
|---|---|
| Lower bound for horizontal rotation ($h_{min}$) | $-0.35$ |
| Upper bound for horizontal rotation ($h_{max}$) | $0.35$ |
| Lower bound for vertical rotation ($v_{min}$) | $-0.25$ |
| Upper bound for vertical rotation ($v_{max}$) | $0.25$ |
| Ratio of patched data ($r_{patch}$) | $0.4$ |
| Training epochs for offline phase ($lr_{offline}$) | $50$ |
| learning rate for offline phase ($lr_{offline}$) | $10^{-3}$ |
| batch size for offline phase ($b_{offline}$) | $64$ |
| Learning rate for online phase ($lr_{online}$) | $2.5 \times 10^{-4}$ |
| Batch size for online phase ($b_{online}$) | $256$ |
| Return attenuation factor ($\gamma$) | $0.95$ |
| Updates per iteration ($n$) | $2$ |



Fig. C.2: Qualitative results of SAC trained with different data. The first column presents the adversarial image processed by regular SAC which is trained with a patch filled with Gaussian noise, while the subsequent column demonstrates the one processed by enhanced SAC. The adversarial patches are generated with 3DAdv and occupy $8\%$ of the image.

training method DOA and the partially adversarial methods like SAC and PZ. This efficiency stems mainly from our unique USAP approach, which bypasses the need to generate adversarial examples, thereby boosting training efficiency. In terms of model inference, our REIN-EAD, along with PZ and DOA, demonstrates superior speed compared to LGS and SAC. This is attributed to the latter methods requiring CPU-intensive, rule-based image preprocessing, which diminishes their inference efficiency.

Regarding detailed training, the REIN-EAD model was trained following the configuration in Appendix C.5. The offline training utilized 4 NVIDIA GeForce RTX 3090 Ti for approximately 2.5 hours (150 minutes). Due to the increased sampling requirements of reinforcement learning, the online training phase was extended to about 12 hours (728 minutes) and utilized eight NVIDIA GeForce RTX 3090 Ti GPUs.

Regarding the computational costs of its key components, the end-to-end inference pipeline achieves 11.5 ms per instance on the face recognition task. The perception model accounts for 98.4% of this processing time (IResNet50 backbone: 8.13 ms; causal transformer: 3.19 ms), while the lightweight policy MLP requires only 0.18 ms. Besides, the Offline Patch Approximation (OAPA) stage occurs before training, and does not add computational burden to either the training or inference phases. The entire training set of approximately 50,000 instances was processed in an acceptable timeframe of around 20 minutes using eight NVIDIA RTX 3090 GPUs via batch processing. The results show that the perception model accounts for the majority of the computational costs in REIN-EAD during inference. These findings suggest that future work on optimizing the computational efficiency of REIN-EAD should focus primarily on the perception model.

## C.7 More Evaluation Results

**Evaluation with different patch sizes.** To further assess the generalizability of the REIN-EAD model across varying patch sizes and attack methods, we conduct experiments featuring both impersonation and dodging attacks. These attacks share similarities with the setup illustrated in Table 1. Although with different patch sizes, the results in Table C.4

and Table C.5 bear a considerable resemblance to those displayed in Table 1. This congruence further supports the adaptability of the REIN-EAD model in tackling unseen attack methods and accommodating diverse patch sizes.

**Evaluation with different adaptive attacks.** As Table C.6 demonstrates, our original adaptive attack using USP was more effective than tracing the authentic policy of EAD (overall). This may be attributed to vanishing or exploding gradients [14] that impede optimization. This problem is potentially mitigated by our approach to computing expectations over a uniform policy distribution. In the meantime, The results reaffirm the robustness of REIN-EAD against a spectrum of adaptive attacks. It further shows that REIN-EAD's defensive capabilities arise from the synergistic integration of its policy and perception models, facilitating strategic observation collection rather than learning a shortcut strategy to neutralize adversarial patches from specific viewpoints.

## C.8 More qualitative results

**Qualitative comparison of different versions of SAC.** SAC is a preprocessing-based method that adopts a segmentation model to detect patch areas, followed by a "shape completion" technique to extend the predicted area into a larger square, and remove the suspicious area [12]. As shown in Fig. C.2, the enhanced SAC, while exhibiting superior segmentation performance in scenarios like face recognition, inadvertently increases the likelihood of masking critical facial features such as eyes and noses. This leads to a reduced ability of the face recognition model to correctly identify individuals, thus impacting its performance in dodging attacks.

## C.9 Stable Convergence of Policy Training

REIN-EAD has incorporated several key design choices into the standard PPO algorithm for promoting stable convergence during policy training. First, to enhance training stability and efficiency, we avoid training perception and

TABLE C.3: Computational overhead comparison of different defense methods in face recognition. We report the training and inference time of defense on a NVIDIA GeForce RTX 3090 Ti and an AMD EPYC 7302 16-Core Processor with the training batch size as $64$.

| Method | # Params (M) | Parametric Model | Training Epochs | Training Time per Batch (s) | Total Training Time (GPU hours) | Inference Time per Instance (ms) |
|---|---|---|---|---|---|---|
| JPEG | - | non-parametric | - | - | - | 9.65 |
| LGS | - | non-parametric | - | - | - | 26.22 |
| SAC | 44.71 | segmenter | 50 + 10 | 0.152 / 4.018 | 104 | 26.43 |
| PZ | | | | | | 11.88 |
| DOA | 43.63 | feature extractor | 100 | 1.732 | 376 | 8.10 |
| EAD | 57.30 | Perception & | 50 + 50 | 0.595 / 1.021 | 175 | 11.51 |
| REIN-EAD | 57.52 | Policy Model | 50 + 40 | 0.595 / 1.096 | 188 | |

TABLE C.4: The **white-box attack success rates** (%) on face recognition models with different patch sizes. $\dagger$ denotes methods are trained with adversarial examples.

| Method | 8 % | | | | 10 % | | | | 12 % | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | MIM | EoT | GenAP | 3DAdv | MIM | EoT | GenAP | 3DAdv | MIM | EoT | GenAP | 3DAdv |
| | | | | | | *Impersonation Attack* | | | | | | |
| Undefended | 100.0 | 100.0 | 99.00 | 98.00 | 100.0 | 100.0 | 100.0 | 99.00 | 100.0 | 100.0 | 100.0 | 99.00 |
| JPEG | 99.00 | 100.0 | 99.00 | 93.00 | 100.0 | 100.0 | 99.00 | 99.00 | 100.0 | 100.0 | 99.00 | 99.00 |
| LGS | 5.10 | 7.21 | 33.67 | 30.61 | 6.19 | 7.29 | 41.23 | 36.08 | 7.21 | 12.37 | 61.85 | 49.48 |
| SAC | 6.06 | 9.09 | 67.68 | 64.64 | **1.01** | 3.03 | 67.34 | 63.26 | 5.05 | 4.08 | 69.70 | 66.32 |
| PZ | 4.17 | 5.21 | 59.38 | 45.83 | 2.08 | 3.13 | 60.63 | 58.51 | 4.17 | **3.13** | 60.63 | 58.33 |
| SAC$^\dagger$ | 3.16 | 3.16 | 18.94 | 22.11 | 2.10 | 3.16 | 21.05 | 16.84 | 3.16 | 4.21 | 15.78 | 18.95 |
| PZ$^\dagger$ | 3.13 | 3.16 | 19.14 | 27.37 | 2.11 | 3.13 | 20.00 | 30.53 | 5.26 | 5.26 | 18.95 | 28.42 |
| DOA$^\dagger$ | 95.50 | 89.89 | 96.63 | 89.89 | 95.50 | 93.26 | 100.0 | 96.63 | 94.38 | 93.26 | 100.0 | 100.0 |
| **EAD** | 4.12 | 3.09 | 5.15 | **7.21** | 3.09 | 2.06 | **4.17** | **8.33** | 3.09 | 5.15 | **8.33** | **10.42** |
| **REIN-EAD** | **2.10** | **1.06** | **3.15** | 7.37 | 1.04 | **1.06** | 6.32 | 9.47 | 3.15 | 4.21 | 7.37 | 15.79 |
| | | | | | | *Dodging Attack* | | | | | | |
| Undefended | 100.0 | 100.0 | 99.00 | 89.00 | 100.0 | 100.0 | 100.0 | 95.00 | 100.0 | 100.0 | 100.0 | 99.00 |
| JPEG | 98.00 | 99.00 | 95.00 | 88.00 | 100.0 | 100.0 | 99.00 | 95.00 | 100.0 | 100.0 | 100.0 | 98.00 |
| LGS | 49.47 | 52.63 | 74.00 | 77.89 | 48.93 | 52.63 | 89.47 | 75.78 | 55.78 | 54.73 | 100.0 | 89.47 |
| SAC | 73.46 | 73.20 | 92.85 | 78.57 | 80.06 | 78.57 | 92.85 | 91.83 | 76.53 | 77.55 | 92.85 | 92.92 |
| PZ | 6.89 | 8.04 | 58.44 | 57.14 | 8.04 | 8.04 | 60.52 | 65.78 | 13.79 | 12.64 | 68.49 | 75.71 |
| SAC$^\dagger$ | 78.78 | 78.57 | 79.59 | 85.85 | 81.65 | 80.80 | 82.82 | 86.73 | 80.61 | 84.69 | 87.87 | 87.75 |
| PZ$^\dagger$ | 6.12 | 6.25 | 14.29 | 20.41 | 7.14 | 6.12 | 21.43 | 25.51 | 11.22 | 10.20 | 24.49 | **30.61** |
| DOA$^\dagger$ | 75.28 | 67.42 | 87.64 | 95.51 | 78.65 | 75.28 | 97.75 | 98.88 | 80.90 | 82.02 | 94.38 | 100.0 |
| **EAD** | **0.00** | **0.00** | **2.10** | **13.68** | 2.11 | 1.05 | **6.32** | **16.84** | 2.10 | 3.16 | **12.64** | 34.84 |
| **REIN-EAD** | 1.04 | 2.04 | 5.15 | 13.54 | 1.03 | 2.02 | 8.16 | 14.58 | 6.18 | 5.15 | 13.54 | 34.02 |

policy models simultaneously from the outset. Early perception models provide unreliable guidance, leading to erratic policy actions and hampering progress. Our solution involves two stages: an initial offline phase where the perception model is trained independently on data from a random policy until stable, followed by an online phase where both networks are trained jointly. This separation enables smoother learning and is significantly more resource-efficient, as the offline pre-training is faster and less computationally demanding than immediate online learning through interaction. As illustrated in Fig. C.3, the pre-trained EAD model ($Pretrained$) demonstrates a better starting point than the model without pre-training ($FromScratch$) and reaches a $10^{-2}$ loss level when trained with the same iterations. Second, alongside the reinforcement learning objective, we incorporate supervised learning signals derived from ground-truth perceptual annotations. Given that active defense naturally extends passive perception, it is intuitive to retain the original supervised learning objective when training the perception model alongside the policy model. Furthermore, this supervised signal acts as a form of regularization, sustaining the EAD model with a considerable perception capability and preventing excessively large or detrimental changes driven solely by the potentially noisy RL reward signal. It guides the policy towards known effective behaviors, further enhancing stability. As shown by Fig. C.3, the un-regularized curve ($w/o\ J_{percep}$) rapidly increases, fluctuates wildly, and fails to converge.

## C.10 Alternative Reward Shaping

**Direct Entropy Deduction.** To explicitly encourage directly reducing the uncertainty in predicting the target, we define the reward as the weighted reduction of the entropy term at every step:

$$\hat{r}_t = \mathcal{H}(\hat{y}_{t-1} \mid b_{t-1}) - \gamma \cdot \mathcal{H}(\hat{y}_t \mid b_{t-1}, o_t). \qquad (C.39)$$

TABLE C.5: The **black-box attack success rates** (%) on face recognition models with different patch sizes. [†] denotes methods are trained with adversarial examples.

| Method | 8 % | | | | 10 % | | | | 12 % | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Cos. | Softmax | NAttack | RGF | Cos. | Softmax | NAttack | RGF | Cos. | Softmax | NAttack | RGF |
| *Impersonation Attack* | | | | | | | | | | | | |
| Undefended | 28.00 | 23.00 | 100.00 | 100.00 | 41.00 | 36.00 | 100.00 | 100.00 | 55.00 | 45.00 | 100.00 | 100.00 |
| JPEG | 33.00 | 33.00 | 96.00 | 94.00 | 42.00 | 40.00 | 98.00 | 98.00 | 62.00 | 48.00 | 99.00 | 99.00 |
| LGS | 11.63 | 6.98 | 11.63 | 4.65 | 12.79 | 11.63 | 9.30 | 3.49 | 15.12 | 17.44 | 9.30 | 4.65 |
| SAC | 8.70 | 9.78 | 13.04 | 14.13 | 13.04 | 11.97 | 14.13 | 15.22 | 17.39 | 13.04 | 13.04 | 17.39 |
| PZ | 6.45 | 9.68 | 4.30 | 3.26 | 5.38 | 7.63 | 4.30 | 4.31 | 8.60 | 5.52 | 4.30 | 4.15 |
| SAC[†] | 11.11 | 12.36 | 12.22 | 14.44 | 10.00 | 19.10 | 13.33 | 15.56 | 12.22 | 17.98 | 15.56 | 14.44 |
| PZ[†] | 8.24 | 5.00 | 10.58 | 9.41 | 10.59 | 3.75 | 9.41 | 8.25 | 11.77 | **5.00** | 9.41 | 8.25 |
| DOA[†] | 15.73 | 17.97 | 34.83 | 16.86 | 15.73 | 14.61 | 32.58 | 15.72 | 16.85 | 17.98 | 31.46 | 10.11 |
| **EAD** | 4.17 | 5.20 | 4.12 | 4.12 | 5.21 | 5.21 | **3.09** | 3.08 | 7.29 | 5.21 | **2.06** | **4.12** |
| **REIN-EAD** | **2.10** | **2.08** | **1.05** | **2.10** | 4.21 | 3.16 | 3.12 | **2.10** | **5.26** | 8.33 | 4.21 | 4.21 |
| *Dodging Attack* | | | | | | | | | | | | |
| Undefended | 44.00 | 35.00 | 96.00 | 96.00 | 53.00 | 43.00 | 100.00 | 99.00 | 68.00 | 61.00 | 100.00 | 100.00 |
| JPEG | 49.00 | 45.00 | 81.00 | 83.00 | 58.00 | 51.00 | 94.00 | 96.00 | 71.00 | 72.00 | 99.00 | 99.00 |
| LGS | 22.11 | 21.05 | 18.95 | 20.00 | 22.11 | 21.05 | 20.00 | 16.84 | 29.48 | 31.58 | 20.00 | 24.21 |
| SAC | 40.80 | 36.84 | 55.26 | 50.00 | 47.37 | 44.74 | 53.95 | 52.63 | 50.00 | 50.00 | 59.21 | 47.37 |
| PZ | 41.67 | 28.34 | 28.33 | 31.67 | 38.33 | 36.67 | 28.33 | 26.67 | 41.67 | 28.33 | 30.00 | 30.00 |
| SAC[†] | 47.46 | 43.54 | 55.93 | 62.71 | 44.07 | 50.00 | 47.46 | 66.10 | 44.07 | 46.77 | 44.07 | 57.63 |
| PZ[†] | 50.88 | 47.69 | 56.14 | 50.87 | 47.37 | 52.30 | 50.88 | 56.14 | 43.86 | 53.84 | 52.63 | 49.12 |
| DOA[†] | 30.33 | 31.46 | 53.93 | 28.09 | 30.33 | 31.46 | 61.80 | 25.85 | 31.46 | 34.83 | 62.93 | 23.60 |
| **EAD** | 5.26 | 7.36 | 1.05 | **0.00** | **4.21** | 6.31 | **0.00** | **0.00** | **10.52** | 16.84 | **0.00** | 3.16 |
| **REIN-EAD** | **4.17** | **7.29** | **1.03** | **0.00** | 6.25 | 8.42 | 1.03 | 2.06 | 12.50 | **13.54** | 5.15 | **3.06** |

TABLE C.6: Evaluation of adaptive attacks on EAD and REIN-EAD. Columns with *USP* represent results obtained by optimizing the patch with expected gradients over the Uniform Superset Policy (USP). And *perception* and *policy* separately represent adaptive attacks against a single submodule. And *overall* denotes attacking the model by following the gradients for along the overall (4 steps) trajectory with gradient-checkpointing.

| Method | Attack Success Rate (%) | | | |
|---|---|---|---|---|
| | USP | Perception | Policy | Overall |
| *Impersonation Attack* | | | | |
| EAD | 8.33 | 1.04 | **9.38** | 7.29 |
| REIN-EAD | **4.21** | 2.11 | 2.17 | - |
| *Dodging Attack* | | | | |
| EAD | **22.11** | 10.11 | 16.84 | 15.79 |
| REIN-EAD | **8.16** | 3.06 | 2.06 | - |



Fig. C.3: Episodic Final Loss across training iterations for three model configurations.

**Binary Outcome Reward.** As a sparse and unbiased baseline, the binary outcome reward indicates whether the agent finishes the task successfully. Under the case of active perception, we define the binary outcome reward as the predicted probability of the target exceeding the pre-defined threshold. Formally, it's defined as

$$\hat{r}_t = \mathbb{I}(\hat{y}_t = y), \tag{C.40}$$

where $\mathbb{I}(\cdot)$ is the indicator function. When the perception model uses probabilistic modeling, which maps from previous belief $b_{t-1}$ and observation $o_t$ to a distribution over target $y$, we regard the episode as successful once the posterior probability for the target exceeds a confidence threshold $\kappa$.
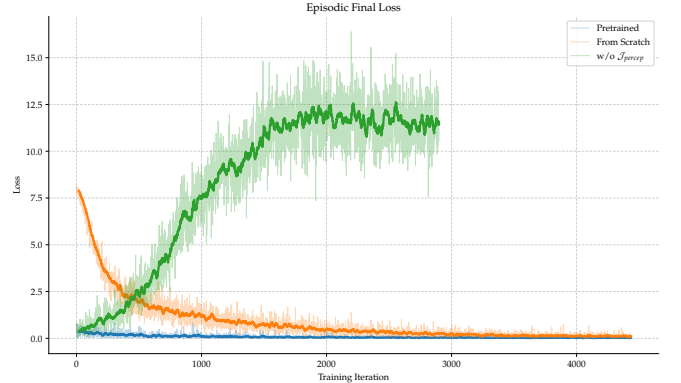
The reward is therefore

$$\hat{r}_t = \mathbb{I}(f(y \mid b_{t-1}, o_t; \boldsymbol{\theta}) > \kappa). \tag{C.41}$$

In implementation, we select $\kappa = 0.95$ for face recognition during training.

**Evaluation Results.** As shown in Table C.7, our proposed reward shaping approach outperforms other methods in terms of both clean accuracy and adversarial robustness against patches. The Direct Entropy Deduction encourages the policy to select actions that reduce uncertainty; however, as it does not leverage ground-truth labels, it may inadvertently promote confident yet incorrect predictions. Consequently, this method is vulnerable to stronger adversarial attacks, such as 3DAdv, resulting in a higher attack success rate. Meanwhile, the Binary Outcome Reward, due to its sparsity, requires

TABLE C.7: The performance of REIN-EAD with different reward shaping.

| Reward | Acc (%) | Attack Success Rate (%) | | | |
|---|---|---|---|---|---|
| | | MIM | EoT | GenAP | 3DAdv |
| Entropy Deduction | 88.67 | 3.15 | 2.11 | 4.21 | 11.42 |
| Outcome Reward | 88.62 | 3.22 | 3.26 | 5.94 | 10.86 |
| **ours** | **89.03** | **2.10** | **3.15** | **7.37** | **4.21** |



Fig. C.4: The loss convergence of REIN-EAD with different rewards.

more iterations to converge, leading to comparatively weaker performance given the same number of training iterations. Furthermore, because it lacks uncertainty-informed guidance, its capability for proactive information-seeking to counter adversarial patches is limited. In contrast, our reward shaping employs a dense formulation that accelerates convergence and guides the model to unbiasedly learn a policy that maximizes information gain towards accurate perception. Furthermore, we present the loss convergence behavior of REIN-EAD under different reward designs as shown in Fig. C.4. The Binary Outcome Reward method demonstrates significant instability and slow, suboptimal convergence, ultimately plateauing at a high loss value. In contrast, both the Direct Entropy Deduction method and our proposed reward shaping technique achieve substantially faster and more stable convergence, with the loss for both methods approaching near-zero values within approximately 4,000 training iterations.

# APPENDIX D
## EXPERIMENT DETAILS FOR OBJECT CLASSIFICATION

### D.1  Details on Dynamic OmniObject3D

The original version of OmniObject3D is accessible at https://omniobject3d.github.io. We preprocess the dataset since it is raw and does not differentiate any splits for classification tasks. For the dataset following a long tail distribution, we drop categories with less than 10 instances, which do not contribute much to our experiments, and split the training and test data in each class at a ratio of 4 to 1. The final dataset has 176 classes, with 4409 objects for training and 1192 objects for testing. We use Pytorch3D (https://github.com/facebookresearch/pytorch3d) as the

simulation engine since it provides efficient API for batch rendering and a differential pipeline for implementing adversarial attacks on the texture of objects. The meshes of the original OmniObject3D dataset have millions of faces, and the scale of them varies significantly. To lower computational overhead and facilitate batch rendering, we use Pymeshlab (https://pymeshlab.readthedocs.io) to process the meshes, simplifying the face number to $10,000$ and normalizing the scale without compromising any rendering quality. The image is rendered with a perspective camera of FoV $60$ and resolution $256 \times 256$ by a hard Phong shader. We use Gym [77] to warp Pytorch3D as the environment for REIN-EAD.

### D.2  Details for implementation

**Model details.** For the experiment conducted on dynamic OmniObject3D, we implement REIN-EAD for classification with a combination of Swin Transformer and Decision Transformer. We use the pretained Swin-Small Transformer from PyTorch Image Models (https://github.com/huggingface/pytorch-image-models) and finetune it with a head of 176 classes on the training set of dynamic OmniObject3D. To implement REIN-EAD, We replace the head of the Swin Transformer with a temporal-fusing module and a Decision Transformer. At each step $t > 0$, the feature embedding of length 768 extracted by the Swin-Small backbone is fed to the temporal-fusing module where it concatenates with the previous extracted observation sequence of dimension $(t-1) \times 768$ to form a temporal sequence of visual features. This sequence is then temporal-fused by the Decision Transformer, which outputs a refined feature embedding. Finally, a feature decoder, an action decoder and a value decoder implemented by shallow MLP are employed to decode the embedding into refined label, predicted action and value respectively. The output of the action decoder is a predicted view angle. In the training stage, we use it as the mean to sample from a multivariate Gaussian distribution with fixed variance for RL exploration, while in testing, we use it directly as the actual action value.

**Training details.** We adopt a similar two-phase training paradigm as REIN-EAD for face recognition. In the offline phase, we freeze the Swin-Small backbone and train the Decision Transformer and feature decoder with a random action policy to collect observations. In the online phase, we employ Algorithm 1 to incorporate the training of the policy network. We approximate the surrogate set of patches with PGD, named OAPA, which attacks the Swin-Small backbone in single view. The learning rate and a number of iterations are $\alpha = 8/255$, $N = 30$, consistent with DOA†. Note that the training of OAPA for REIN-EAD is offline and only done once, while DOA† generates a surrogate patch in every iteration. The hyper-parameters of REIN-EAD for object classification are shown in Table D.8.

### D.3  Details for attack

**Attack in texture space for 3D environment.** We implement adversarial attacks to mislead the classifier to output a wrong label. The attack methods on the classification task are similar to those in FR system, but are implemented in texture space. We leverage the traceable rendering computing

TABLE D.8: Hyper-parameters of REIN-EAD for object classification.

| Hyper-parameter | Value |
| --- | --- |
| Lower bound for horizontal rotation ($h_{\min}$) | $-90°$ |
| Lower bound for horizontal rotation ($h_{\max}$) | $90°$ |
| Lower bound for vertical rotation ($v_{\min}$) | $0°$ |
| Upper bound for vertical rotation ($v_{\max}$) | $90°$ |
| Ratio of patched data ($r_{\text{patch}}$) | 0.8 |
| Training epochs for offline phase ($N_{\text{offline}}$) | 100 |
| Learning rate for offline phase ($\text{lr}_{\text{offline}}$) | $2 \times 10^{-4}$ |
| Batch size for offline phase ($b_{\text{offline}}$) | 128 |
| Total Episodes for online phase | $150,000$ |
| Learning rate for online phase ($\text{lr}_{\text{online}}$) | $1 \times 10^{-4}$ |
| Batch size for online phase ($b_{\text{online}}$) | 128 |
| Return attenuation factor ($\gamma$) | 0.95 |
| Updates per iteration ($n$) | 2 |

graph of Pytorch3D to generate the adversarial patch directly on the masked texture of the object, which fits the non-planar surface of the object of diverse categories and ensures multi-view consistency. Since the shape of objects varies significantly in the context of the classification task, to ensure that the patch can be fully attached to the object, we set the patch size to $20\%$ of the bounding box of the object and locate it on the center of the bounding box. The adversarial sample for the robustness test is generated on test datasets that are unseen for both REIN-EAD and other defense baselines.

**White-box Attack.** We use MIM [58] as a single-view white box adversary, which is enhanced by momentum item. The decay factor of MIM is set at $\mu = 1.0$, consistent with FR attack. The multi-view threats are established by EoT [59] and MeshAdv [78]. EoT adopts a 2D batch data augmentation that includes shifting, rotating and flipping the rendered image and averages the gradient in image space. In contrast, MeshAdv employs a batch of 3D transformations across the action range of REIN-EAD and backpropagates the gradient through the differential rendering pipeline to texture space, making it more robust to view change. For the white-box attacks, we set the learning rate and the number of iterations at $N = 100$ and $\alpha = 8/255$. The sampling frequency for EoT and MeshAdv is established at $M = 128$.

**Black-box Attack.** For query-based attack, we use RGF [63] and N attack [62]. We set the maximum number of queries at $N = 10000$ and sampling frequency at $M = 100$. The learning rates for RGF and N attack are $\alpha = 0.05$ and $\alpha = 0.1$ respectively. For transfer-based adversaries, We utilize MeshAdv with parameters the same as the white box version and fine-tune a Swin-Tiny Transformer as the surrogate model to launch the transfer attack.

**Adaptive attack.** Following the face recognition setting, we adaptively attack JPEG and LGS with BPDA [14] technique, SAC[†] and PZ[†] with STE [89] technique. For REIN-EAD, we launch the adaptive attack with a uniform superset policy which is provided to be most effective in the FR task.

## D.4 Details for defense

The implementation of JPEG [21] and LGS [22] are consistent with the FR task. For SAC[†] [12] and PZ[†] [13], we retrain the patch segmenter on the training set of dynamic OmniObject3D with adversarial patches optimized by EoT [59] using the same code with FR task. For DOA[†] [8], we follow its training paradigm to fine-tune the same Swin-Small Transformer backbone used by REIN-EAD on the training set of dynamic OmniObject3D. Specifically, we utilize PGD with learning rate $\alpha = 8/255$ and number of iterations $N = 30$ for adversarial training and search the patch location using the gradient-based method with the top candidate number $C = 10$ as described in their paper. The patches used for training SAC[†], PZ[†] and DOA[†] occupy $20\%$ of the bounding box of the object, which is the same with REIN-EAD.

## APPENDIX E
## EXPERIMENT DETAILS FOR OBJECT DETECTION
### E.1 Details on EG3D

For object detection, we use a pre-trained EG3D model on ShapeNet Cars at https://catalog.ngc.nvidia.com/orgs/nvidia/teams/research/models/eg3d to generate multi-view car images with the resolution of $256 \times 256$. We generate 1000 cars with different appearances and split the data into 800 training data and 200 test data. The latent seeds are recorded to ensure that the identity of each car remains consistent in all the experiments. Since the background of the generated image is blank, we are able to annotate the bounding box automatically. The online environment for REIN-EAD is wrapped by Gym [77].

### E.2 Details on CARLA

We utilize CARLA 0.9.14 [48] for a more complex object detection experiment. The training data covers all 41 different vehicle blueprints provided by CARLA. For each blueprint we generate vehicles of different color versions and collect multi-view samples in the different backgrounds on Town10 as the offline datasets for training model and adversarial samples. The training and testing procedures for REIN-EAD are conducted online in CARLA, which allows the agent to explore every possible view within the action range. During the test experiment, we use all 41 vehicles and place them in locations different from the training set. We use the Python API provided by CARLA to automatically collect the data at a resolution of $256 \times 444$ and annotate the label. We warp CARLA with Gym [77] as an online environment for REIN-EAD.

### E.3 Details of implementations

**Model details.** For the object detection task on EG3D, we implement REIN-EAD with a combination of YOLOv5n and Decision Transformer. We use the pretained YOLOv5n from the official implementation (https://github.com/ultralytics/yolov5) and fine-tune it as a single class detection model on the training set of each environment respectively. For each time step $t > 0$, given the current observation of dimensions $256 \times 256 \times 3$ as input, the feature maps of dimension $32 \times 32 \times 64$ output by the second Cross Stage Partial

TABLE E.9: Hyper-parameters of REIN-EAD for object detection on EG3D.

| Hyper-parameter | Value |
|---|---|
| Lower bound for horizontal rotation ($h_{min}$) | $-60°$ |
| Upper bound for horizontal rotation ($h_{max}$) | $60°$ |
| Lower bound for vertical rotation ($v_{min}$) | $0°$ |
| Upper bound for vertical rotation ($v_{max}$) | $30°$ |
| Ratio of patched data ($r_{patch}$) | 0.4 |
| Training epochs for offline phase ($N_{offline}$) | 50 |
| Learning rate for offline phase ($lr_{offline}$) | $2 \times 10^{-4}$ |
| Batch size for offline phase ($b_{offline}$) | 128 |
| Total episodes for online phase ($N_{online}$) | 10,000 |
| Learning rate for online phase ($lr_{online}$) | $1 \times 10^{-4}$ |
| Batch size for online phase ($b_{online}$) | 64 |
| Return attenuation factor ($\gamma$) | 0.95 |
| Updates per iteration ($n$) | 2 |

TABLE E.10: Hyper-parameters of REIN-EAD for object detection on CARLA.

| Hyper-parameter | Value |
|---|---|
| Lower bound for horizontal rotation ($h_{min}$) | $-60°$ |
| Upper bound for horizontal rotation ($h_{max}$) | $60°$ |
| Lower bound for vertical rotation ($v_{min}$) | $5°$ |
| Upper bound for vertical rotation ($v_{max}$) | $35°$ |
| Ratio of patched data ($r_{patch}$) | 0.4 |
| Training epochs for offline phase ($lr_{offline}$) | 50 |
| Learning rate for offline phase ($lr_{offline}$) | $2 \times 10^{-4}$ |
| Batch size for offline phase ($b_{offline}$) | 128 |
| Total Episodes for online phase | 10,000 |
| Learning rate for online phase ($lr_{online}$) | $1 \times 10^{-4}$ |
| Batch size for online phase ($b_{online}$) | 64 |
| Return attenuation factor ($\gamma$) | 0.95 |
| Updates per iteration ($n$) | 2 |

Networks are utilized. We find feature maps at this level to be computationally efficient for REIN-EAD as the later stages of YOLOv5n form a large-scale concatenated feature pyramid. These maps are then reshaped into a sequence with dimensions $64 \times 1024$. To concatenate it with the previous extracted observation sequence $(t - 1) \times 64 \times 1024$, we have a temporal sequence of visual features as the input of Decision Transformer, and it outputs the temporal-fused visual feature sequence $64 \times 1024$ and predicted action. For REIN-EAD, an extra value decoder is utilized to estimate the advantage value. To predict the bounding boxes and objectness score which is required in object detection, we reshape the temporal-fused visual sequence back to its original shape and utilize it as input for the later stage in YOLOv5n.

**Training details.** We adopt a similar training paradigm as REIN-EAD for previous tasks and set the hyper-parameters of REIN-EAD for object detection as Table E.9 and Table E.10

## E.4 Details for attack

**Multi-view hiding attack for vehicle object detection.** In the context of a single-class vehicle object detection task, the goal of the adversary is to place a patch on a vehicle and make it disappear from the object detector. The patch is placed on the side of the vehicle and occupies $25\%$ of the bounding box. We follow the adversarial loss designed for YOLO [93] to minimize the objectness score to achieve a hiding attack. To enhance the multi-view robustness, we train adversarial patches with a batch of images with different views in each iteration. The sampling frequency is set at 100.

**Attack in pixel space.** EoT is used as a baseline multi-view adversary, which incorporates the expectation of view transformation described above. To achieve a more generalized attack, we utilize Universal Adversarial Perturbations (UAP) [81] to generate a single adversarial patch for all the vehicles in the dataset, which is able to hide any vehicles from the detector. The learning rate and the number of iterations for both methods are set at $N = 500$ and $\alpha = 8/255$ .

**Attack in the hidden layer.** As the EAD module is plugged at the middle stage of YOLOv5n, we implement SIB [80] , which attacks the feature in the hidden layer. We set the feature to be perturbed as the input of the EAD module and maximize the difference between the adversarial feature and the original feature with their feature-interference reinforcement loss. The coefficients for objectness loss and feature loss are $\mu_1 = 0.5$ and $\mu_2 = 0.5$ respectively. The learning rate and the number of iterations are consistent with EoT and UAP.

**Attack in latent space.** To achieve an inconspicuous patch attack for both human and detector, we utilized adversarial camouflage textures (AdvCaT) [82], which adopts Voronoi diagram and Gumbel-softmax trick to generate the semantic polygon camouflage with control points and latent seeds. We use K-means with 4 clusters to extract the base colors for camouflage from the environment. We set the learning rate for control points and latent seeds at $\alpha_1 = 0.0005$ and $\alpha_2 = 0.005$ respectively. The number of iterations is $N = 200$.

## E.5 Details for defense

The implementation of JPEG [21] and LGS [22] are consistent with the FR task. For SAC[†] [12] and PZ[†] [13], we retrain the patch segmenter on the training set of the corresponding environment with adversarial patches optimized by EoT using the same code with FR task. The patches used for training SAC[†] and PZ[†] occupy $25\%$ of the bounding box of the vehicle, which is the same as those used for training REIN-EAD.

## E.6 Qualitative comparison of defense baselines in CARLA.

We present the qualitative defense result for baseline methods in Fig. E.5. The smoothing mechanism of LGS is relatively weak and unable to completely eliminate adversarial noise. In contrast, the segmentation backbone of SAC[†] and PZ[†] can perfectly distinguish and remove the patch attacks with noisy patterns, but fail to detect the environmental mosaic AdvCaT. Additionally, the completion mechanism of SAC[†] leads to excessive occlusion, which complicates subsequent detection tasks.
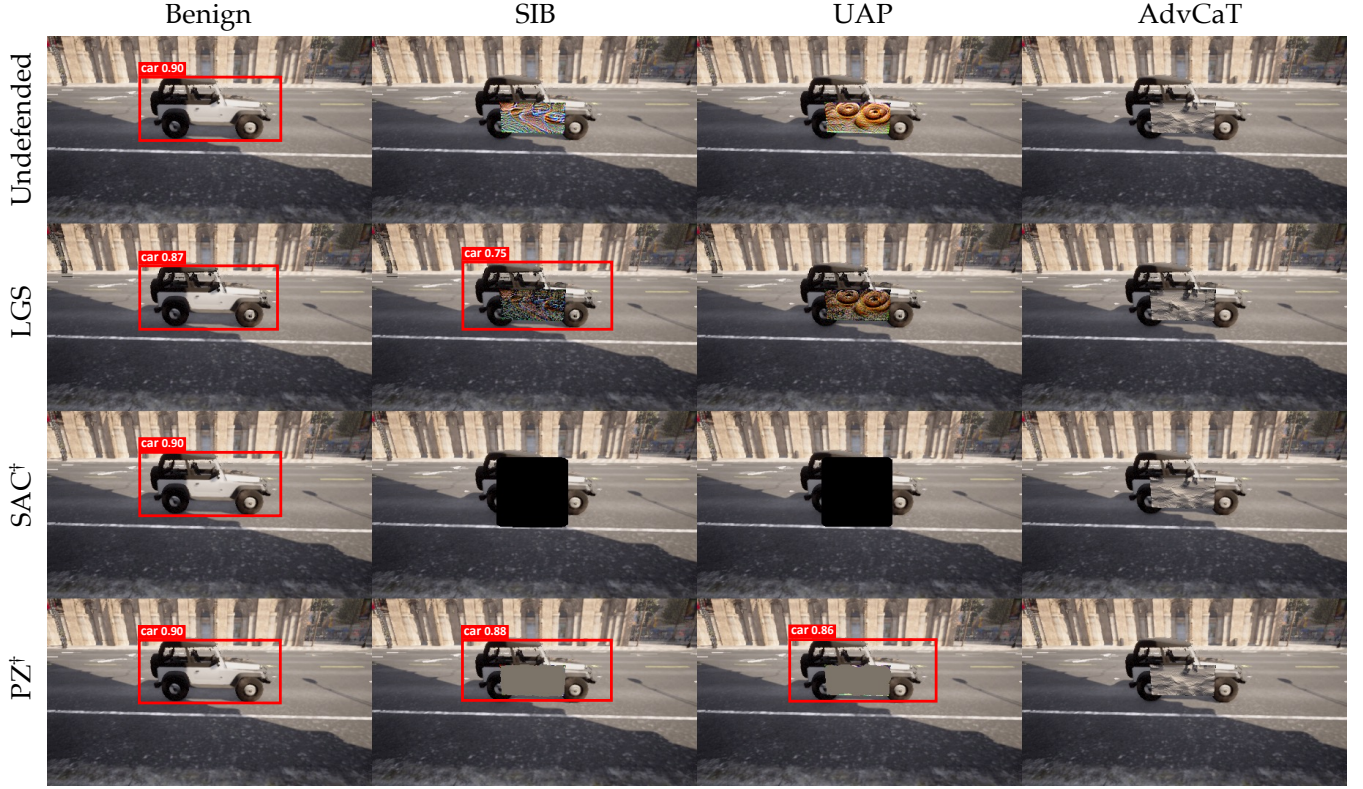
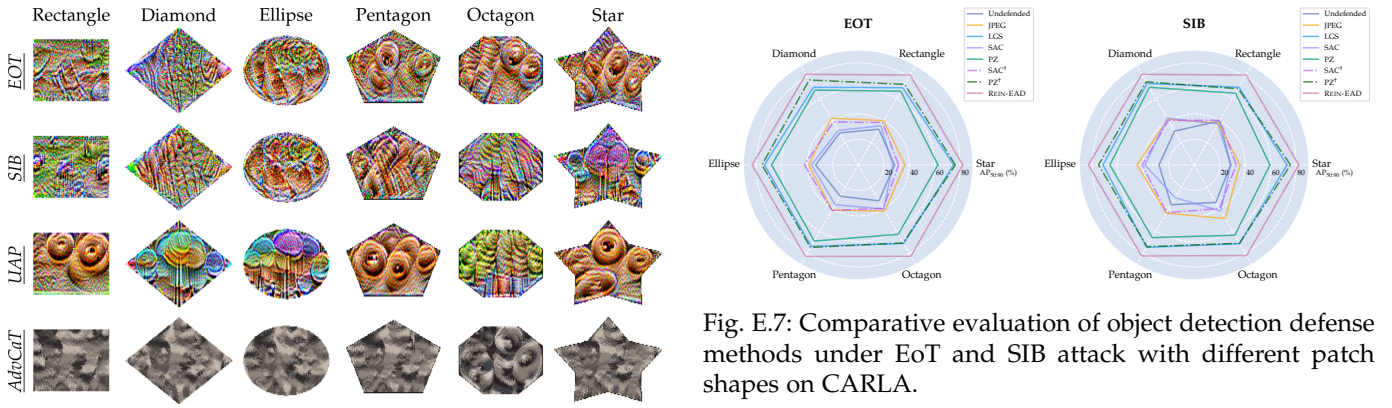Fig. E.5: Visualization of defense baseline on CARLA. $^\dagger$ denotes methods are trained with adversarial examples.



Fig. E.6: Visualization of the generated patch under different attacks and shapes on CARLA.



Fig. E.7: Comparative evaluation of object detection defense methods under EoT and SIB attack with different patch shapes on CARLA.

## E.7 Visualization of patches on CARLA

We generate patches with the shapes of diamond, ellipse, pentagon, octagon and star. They are unseen for both REIN-EAD and other defense baselines. A visualization of these patches are shown in Fig. E.6.

## E.8 More results on patch shapes

More comparative evaluations for EoT and SIB attacks with different patch shapes are illustrated in Fig. E.7.

## E.9 Failure cases

The failure cases reveal several notable limitations of REIN-EAD. First, the system is vulnerable to strategically positioned adversarial patches that occlude critical object features. As illustrated in Fig. E.8, when an adversarial patch obscures the facial region of a green doll, the model consistently misclassifies the object as broccoli. This suggests that the occlusion of key discriminative features compromises the model's ability to acquire adequate information for accurate recognition. Second, the framework exhibits degraded performance when simultaneously subjected to adversarial attacks and natural out-of-distribution interference. The multi-step exploration mechanism, while generally effective, fails to resolve conflicting signals in such compound uncertainty scenarios. A representative example, shown in Fig. E.8, demonstrates persistent prediction uncertainty when
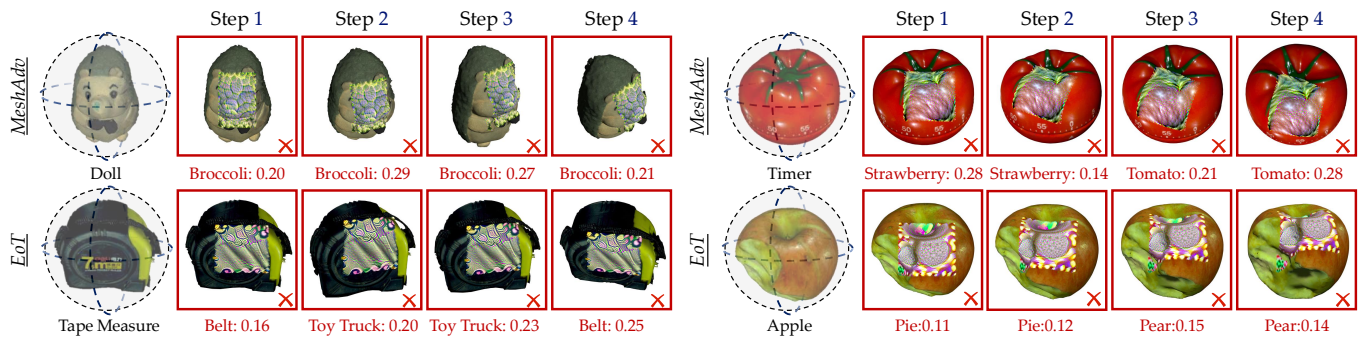
Fig. E.8: Visualization of failure cases by adopting REIN-EAD on dynamic OmniObject3D, with the adversarial patch occupying 20% of the object's bounding box in the front view.

processing a partially consumed apple, where bite-induced shape deformation interacts with adversarial perturbations. These limitations highlight potential areas for future research. Potential avenues include developing patch placement strategies that prioritize important visual areas and incorporating more diverse 3D object datasets for training to enhance the model's robustness to occlusion and out-of-distribution samples.