# Learning When to Use Adaptive Adversarial Image Perturbations Against Autonomous Vehicles

Hyung-Jin Yoon , Hamidreza Jafarnejadsani , and Petros Voulgaris , *Fellow, IEEE*

*Abstract*—Deep neural network (DNN) models are widely used in autonomous vehicles for object detection using camera images. However, these models are vulnerable to adversarial image perturbations. Existing methods for generating these perturbations use the image frame as the decision variable, resulting in a computationally expensive optimization process that starts over for each new image. Few approaches have been developed for attacking online image streams while considering the physical dynamics of autonomous vehicles, their mission, and the environment. To address these challenges, we propose a multi-level stochastic optimization framework that monitors the attacker's capability to generate adversarial perturbations. Our framework introduces a binary decision attack/not attack based on the attacker's capability level to enhance its effectiveness. We evaluate our proposed framework using simulations for vision-guided autonomous vehicles and actual tests with a small indoor drone in an office environment. Our results demonstrate that our method is capable of generating real-time image attacks while monitoring the attacker's proficiency given state estimates.

*Index Terms*—Adversarial machine learning, reinforcement learning, autonomous vehicle.

## I. INTRODUCTION

**M**ACHINE learning (ML) tools that detect objects using high-dimensional sensors, such as camera images [1] or point clouds measured by LiDAR [2], are extensively used in autonomous vehicles [3], [4]. As vision-based autonomous vehicles become more integrated into society, it is crucial to ensure the robustness of these systems, which rely on various sensor signals in uncertain environments. Analyzing worst-case scenarios within uncertainties has been a useful approach to robustify control systems [5] and reinforcement learning [6]. To follow this approach, researchers have revealed the vulnerability of machine learning methods, especially deep learning tools developed for computer vision tasks such as object detection and classification, to data perturbed by adversaries. For instance, small perturbations can be added to images that are unnoticeable to human eyes but result in incorrect image classifications [7],

[8], [9]. Moreover, recent works have demonstrated adversarial image perturbations against autonomous vehicles, including (1) modifying physical objects, such as putting stickers on a road [10] or a road sign [11], to fool an ML image classifier or end-to-end vision-based autonomous car; and (2) fooling object tracking algorithms in autonomous driving systems [12]. Adversarial machine learning commonly focuses on creating stealthy and natural-looking perturbations to evade human detection. Such attacks are designed to resemble out-of-distribution samples that may occur in real-world environments. As a consequence, ensuring the robustness of ML-based autonomous vehicle systems against adversarial attacks has become increasingly critical.

While the aforementioned adversarial image perturbations against autonomous cars [10], [11], [12] successfully reveal weaknesses in vision-guided navigation in autonomous vehicles, these perturbed images are generated offline. However, offline methods [11], [12] do not consider the effect of real-time attacks on dynamically changing environments during driving or flight of the vehicles. To prevent accidents [13] caused by vision-guided autonomous vehicles due to defective perception systems and their vulnerabilities, we need to study attack and defense techniques that go beyond offline methods for deep neural networks.

There are two approaches to generating adversarial image perturbations, depending on the attacker's access to the victim perception model. In the *white-box* attack approach, the attacker has full access to the victim ML classifier (or object detector) and generates adversarial image perturbations through iterative optimization [8], [12]. In this method, images are the decision variables of optimization, and the training loss function is reused with incorrect labels set by the attacker. The optimization takes iterative gradient descent optimization steps with respect to the image variables calculated using back-propagation through the known victim ML classifier [8] (or object detector [12]). On the other hand, in the *black-box* approaches [14], [15], the attacker only has access to input and output pairs of the victim model and must estimate the gradient. However, estimating the gradients in *black-box* attacks requires a large number of samples, which may not be available from autonomous systems operating in dynamic environments.

*Statement of Contribution:* To our knowledge, this letter is the first to propose a stealthy attack scheme on image streams used for object detection/tracking in autonomous vehicles (e.g., self-driving cars and drones) that can be deployed *online*, and the *physical dynamics of the system* and the *varying surrounding environment* are taken into account in the optimization phase of the attack scheme. In this letter, we present a framework that utilizes generative adversarial networks (GANs) to generate adversarial images in real-time scenarios without the need for

(a) Adversarial patch example.

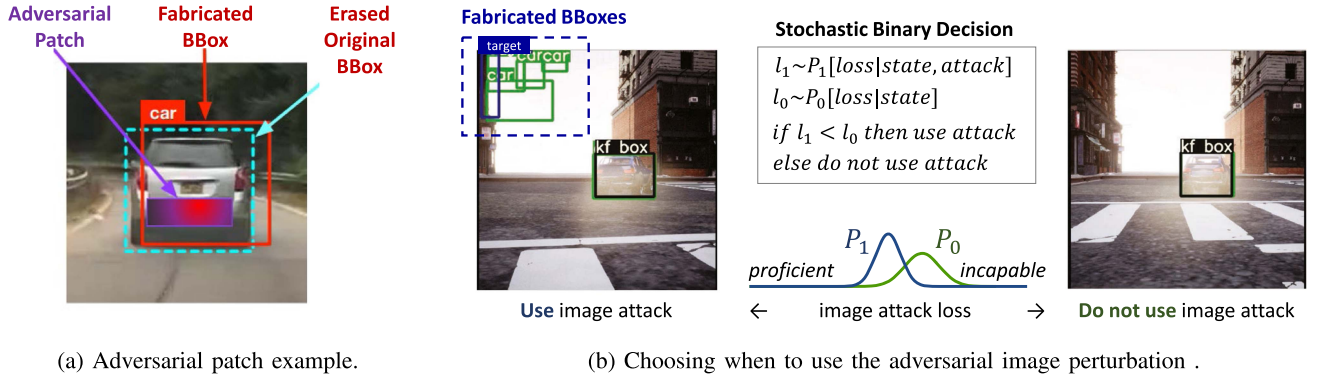(b) Choosing when to use the adversarial image perturbation .

Fig. 1. Image attacks: (a) Adversarial patch in [12], (b) Adversarial perturbation with binary decision in this letter. *kf box* denotes Kalman filtered bounding box (BBox). An illustrative video is available at https://youtu.be/rWq9JR4PY5o.

iterative steps. Building on the approach outlined in [16], our proposed multi-level framework consists of several components. First, the GAN functions as an online image generator. Second, a reinforcement learning agent is trained to misguide the vehicle according to the adversary's objective. Lastly, a binary decision-maker determines when to use image attacks based on the proficiency of the image attack generator, given the current state estimate. Our framework provides a more efficient and practical alternative to iterative *white-box* methods for generating adversarial images.

Our contributions can be summarized as follows:

- We propose a *real-time* adversarial image perturbation framework that allows for implementation on real-world robots, in contrast to existing offline methods.
- We introduce a *state estimation-based reinforcement learning* approach that learns to decide on the image frame area to fabricate bounding boxes. This approach eliminates the need for manual annotation of patch areas.
- We incorporate a constraint on the strength of the image perturbation, making the attacked image frame *less noticeable and more stealthy* compared to existing methods. This is demonstrated in Fig. 1.

## II. RELATED WORKS

Adversarial image perturbations have been extensively studied to attack autonomous vehicles that rely on camera images for navigation [12], [17], [18]. For instance, in [17], an optimization problem was formulated to place black marks on the road, which caused an end-to-end autonomous driving car to veer off the road in a virtual reality environment. This approach was inspired by the demonstration of attacking Tesla's autonomous driving systems with just three small stickers [10]. In another work [12], the authors demonstrated the effectiveness of a *white-box* adversarial image perturbation method on object tracking of an autonomous system that uses *Kalman* filter (KF) to disrupt the object tracking. This method was also shown to be effective in attacking an industry-level perception module that uses vision-based object detection fused with LIDAR, GPS, and IMU [18].

The aforementioned *white-box* methods [12], [18] require full sets of iterative optimization computations for every new image, rendering them unsuitable for dynamic environments with evolving situations and control loops of autonomous vehicles. These approaches do not consider the varying computation
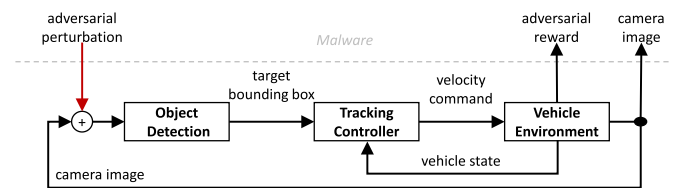


Fig. 2. Attacker (malware) and victim system (guidance).

time of the iterative optimizations, which might have different termination steps for online applications. Additionally, the attack methods in [12], [18] often require additional state information that is not always readily available, unlike the image stream. For instance, generating the adversarial patch in Fig. 1(a) in [12] requires the attacker to know the exact anchor index associated with the target bounding box (BBox) and the location to place the BBox. As mentioned in the *open review* [19] by the authors in [12], the adversarial patch area was manually annotated in each video frame.

Although there are various other adversarial image attack methods available, many of them are offline methods that require additional information such as labeled training data to generate adversarial images.

## III. REAL-TIME ADVERSARIAL IMAGE ATTACK

Our goal is to develop a real-time solution that can learn to generate adversarial image perturbations and decide when to use the attack based on the proficiency of the attack generator, as illustrated in Fig. 1(b). The adversarial image perturbations are designed to manipulate the perception of autonomous vehicles to misguide them according to the adversary's objectives, such as causing collisions or making the vehicle deviate from its original path. To formally formulate the problem, we consider the following assumptions and settings.

### A. Problem Description and Proposed Framework

We focus on an autonomous vehicle that utilizes an object detection ML method to track a target object using camera images, as shown in Fig. 2. We used a recent version of the *YOLO* object detection model [1], which was downloaded from [20],
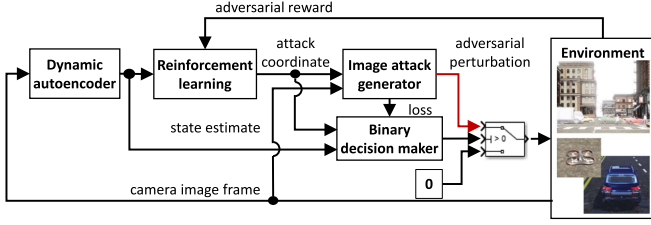
Fig. 3.    Image attack framework with binary decision maker.

for our experiments.[1] The output of the object detector network is a multi-dimensional tensor that is processed using non-max suppression [1] to obtain a list of bounding box coordinates. The box with the highest confidence score for the target class is then selected from the list of detected bounding boxes to generate tracking control commands. The autonomous guidance system uses the vehicle's actuators, including the acceleration pedal, brake, and steering wheel, to keep the target's bounding box centered in the camera view. Consequently, the vehicle moves toward the target object.

We assume the adversary's objective is to disrupt the target tracking control in Fig. 2. The attacker is assumed to be embedded as *Malware* and has access to the image stream, enabling them to perturb the input to the object detection module of the victim system, as illustrated in Fig. 2. Given the image streams denoted as $\mathbf{x}_0, \mathbf{x}_1, \ldots, \mathbf{x}_t$, the attacker's goal is to generate adversarial image perturbations $\mathbf{w}_0, \mathbf{w}_1, \ldots, \mathbf{w}_t$ that mislead the victim vehicle according to adversarial objectives expressed in terms of adversarial rewards $r_1, r_2, \ldots, r_t$. The reward function is based on the vehicle's state, such as position, velocity, or collision states, and actions that involve the coordinates used to fabricate the bounding boxes through the image attack generator, as shown in Fig. 3. These rewards are crucial for applying *reinforcement learning* (RL), which learns the correlation between actions and rewards for different states of the system. In this framework, a binary decision-maker determines when to attack based on the attack proficiency (represented as loss in Fig. 3). The problem addressed in this framework can be summarized as follows:

*Problem:* Develop machine learning methods that learn to increase the sum of rewards $\{r_t\}$ for the adversary by generating adversarial perturbation $\{\mathbf{w}_t\}$ while selecting when to use the attack at the time step $t$, as shown in Fig. 3. The ML method assumes to use only the image stream $\{\mathbf{x}_t\}$ from the autonomous vehicle that has a guidance system and malware shown in Fig. 2.

### B. Online Image Attack With Binary Decision Making

Our framework involves binary decision-making that depends on the proficiency of the image attack generator. This type of decision-making belongs to the multi-armed bandit class of problems [22], where the decision-maker selects the most profitable action. However, unlike the classical multi-armed bandit, where the rewards are generated from independent-stationary distributions, our decision-maker must consider non-stationary system state distributions. Specifically, given the attack coordinate (the position and size of the fabricated bounding box)

chosen by RL and the state estimate from the dynamic autoencoder, the decision-maker must determine whether using the attack is profitable or not. To tackle this challenge, the authors in [23] used a deep neural network (DNN) to learn the correlation between the state, decision, and profit. The profit in this context is aligned with the attacker's capability to generate the perturbation given the attack coordinate. They also employed random dropout [24] with the DNN to estimate the profit distributions for each decision. This multi-armed bandit algorithm, which uses DNN with random dropout, is known as *Neural Thompson Sampling* (NTS).

We sought to implement NTS for binary decision-making, using the proficiency of the image attack generator as the profits in the multi-armed bandit. While the direct application of NTS to our framework is appealing, there is a causality issue to consider. Specifically, the proficiency of the image attack generator depends on how many similar pairs of the attack coordinates and the image frames are experienced by the generator throughout the training. In other words, proficiency would not be closely correlated to the state estimate that is input to the NTS implementations. In our experiments, we tested NTS, but it did not demonstrate the desired behavior of selecting to attack when the attack generator is proficient.

Therefore, we propose an alternative method to NTS that involves comparing two conditional expectations. Specifically, our method compares $E[l_t | \mathbf{h}_t, \mathbf{a}_t]$ with $E[l_t | \mathbf{h}_t]$. Here, $l_t$ represents the loss function used to measure the proficiency of the image attack GAN. The state estimates that are low-dimensional representations of all previous observations (camera image stream), denoted by $\mathbf{h}_t$, are obtained using the dynamic autoencoder shown in Fig. 3. A single camera frame $\mathbf{x}_t$ does not capture the entire scene. Furthermore, a single camera frame does not provide the velocity information of the vehicles. The state estimate $\mathbf{h}_t$ provides better information on the state of the environment $\mathbf{s}_t$ than a single camera frame $\mathbf{x}_t$. Additionally, $\mathbf{a}_t$ represents the action determined by reinforcement learning agent in Fig. 3. This action is the attack coordinate, which is the position and size of the fabricated bounding box. The goal of this approach is to compare the expected loss given the attack coordinate $\mathbf{a}_t$ suggested by RL with the expected loss averaged over all other possible attack coordinates. If the expected loss given $\mathbf{a}_t$ is lower than the average loss, then $\mathbf{a}_t$ is considered a promising attack coordinate to be used at this point. We refer to this decision-making method as *Conditional Sampling* (CS). The loss estimation and decision-making procedure of CS are as follows:

*Estimation:* The estimation for CS involves the following optimizations:

$$\arg\min_{\theta^{\text{dec}}} \|l_t - \hat{l}_0(\mathbf{h}_t; \theta^{\text{dec}})\|^2,$$
$$\arg\min_{\theta^{\text{dec}}} \|l_t - \hat{l}_1(\mathbf{h}_t, \mathbf{a}_t; \theta^{\text{dec}})\|^2, \tag{1}$$

where $\hat{l}_0$ and $\hat{l}_1$ are DNNs trained to predict the loss functions values $l_t$ given the state estimate $\mathbf{h}_t$ and the attack coordinate $\mathbf{a}_t$ respectively. The DNNs have parameters denoted as $\theta^{\text{dec}}$ that need to be optimized.

*Decision Making:* The decision to launch an attack is determined by random sampling. To obtain sample image attack losses $\tilde{l}_0$ and $\tilde{l}_1$, we follow the same approach as in NTS, using the current state estimate $\mathbf{h}_t$ and the attack coordinate $\mathbf{a}_t$. Specifically, we generate output samples by applying random dropout in

---

[1]Another popular object detection model, *Faster R-CNN*, can be attacked using similar White box attack method as in [21]. Hence, our proposed method can be implemented with *Faster R-CNN*.
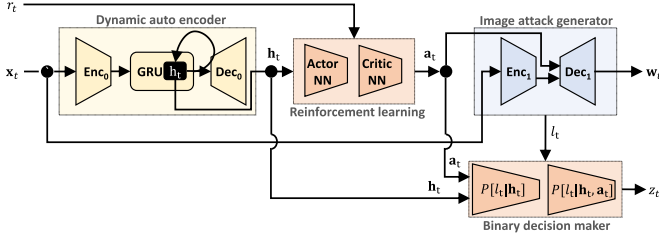
Fig. 4. Multi-level image attack computation network.

DNNs, i.e., $\hat{l}_0$ and $\hat{l}_1$, and estimate Gaussian distributions based on these samples. We then sample from the estimated Gaussian distributions to obtain $\tilde{l}_0$ and $\tilde{l}_1$, which can be expressed as:

$$\tilde{l}_0 \sim P_0[l_t|\mathbf{h}_t] \quad \text{and} \quad \tilde{l}_1 \sim P_1[l_t|\mathbf{h}_t, \mathbf{a}_t], \quad (2)$$

where $P_0[l_t|\mathbf{h}_t]$ and $P_1[l_t|\mathbf{h}_t, \mathbf{a}_t]$ represent the conditional probability distributions associated with sampling from $\hat{l}_0(\mathbf{h}_t; \theta^{\text{dec}})$ and $\hat{l}_1(\mathbf{h}_t, \mathbf{a}_t; \theta^{\text{dec}})$ in Fig. 4. Further details of the sampling procedure can be found in [23]. To make a decision, we select the option with the lower loss value, i.e., if $\tilde{l}_0 < \tilde{l}_1$, then the attack will not be launched in the time step; otherwise, the image attack will be performed until the next time step $t + 1$.

The proposed framework integrates estimation models for binary decision-making within computation networks consisting of DNNs, as depicted in Fig. 4. A major advantage of this framework is the ability to generate real-time adversarial image perturbations through recursive computations. The process involves feeding the image $\mathbf{x}_t$ at time $t$ into encoder networks, $\mathbf{Enc}_0$ and $\mathbf{Enc}_1$, for dimension reduction. The dynamic autoencoder comprises $\mathbf{Enc}_0$, $\mathbf{GRU}$ (gated recurrent unit), and $\mathbf{Dec}_0$ (decoder). The $\mathbf{GRU}$ recursively updates the hidden state $\mathbf{h}_t$ using the encoded image $\mathbf{Enc}_0(\mathbf{x}_t)$ and the high-level attack action $\mathbf{a}_t$, as shown in Fig. 4. The estimated state information in $\mathbf{h}_t$ is then used by the actor (policy) to generate the high-level action, i.e., $\mathbf{a}_t = \mathbf{Actor}(\mathbf{h}_t)$.

As shown in Fig. 4, the adversarial image perturbation $\mathbf{w}_t$ is generated by $\mathbf{Dec}_1$ using the high-level attack action $\mathbf{a}_t$ and another encoded image from $\mathbf{Enc}_1$, i.e., $\mathbf{w}_t = \mathbf{Dec}_1(\mathbf{Enc}_1(\mathbf{x}_t), \mathbf{a}_t)$. The perturbed image frame is obtained by applying the perturbation to the original image with a scale factor $\alpha$, i.e., $\tilde{\mathbf{x}}_t = \max(\min(\mathbf{x}_t + \alpha\mathbf{w}_t, 1), 0)$. The scale factor $\alpha$ can be set to a low value (e.g., 5% or 10%) to improve the stealthiness of the image attack. The binary decision maker selects the decision variable $z_t$ using the conditional sampling (CS) described in (2), where $z_t = 1$ indicates that the attack is used and $z_t = 0$ indicates that the attack is not used (when $\tilde{l}_0 < \tilde{l}_1$).

The recursive process of generating adversarial image perturbation using only camera image is summarized in Algorithm 1. The entire computation at each time-step uses only the current observation and the state values in the previous time-step without iterative optimization, enabling real-time generation of image attacks.

### C. Multi-Time Scale Optimization to Train the Attacker

We employ a multi-level stochastic optimization approach that separates the time scales of the updates for the various components depicted in Fig. 3. Our stochastic optimization method

---

**Algorithm 1:** Recursive Image Attack.

  **Initialize:** $t \leftarrow 0$; load the pre-trained parameters of the recursive attack networks.
  **repeat**
    Generate attack command using RL policy (Actor)
      $\mathbf{a}_t \leftarrow \mathbf{Actor}(\mathbf{h}_t)$
    Encode the observed image $\mathbf{x}_t$ from the environment
      $\zeta_t \leftarrow \mathbf{Enc}_1(\mathbf{x}_t)$
    Generate adversarial image perturbation
      $\mathbf{w}_t \leftarrow \mathbf{Dec}_1(\zeta_t, \mathbf{a}_t)$
    Sample from the conditional distribution as in (2), i.e.,
      $\tilde{l}_0 \sim P_0[l_t|\mathbf{h}_t] \quad \text{and} \quad \tilde{l}_1 \sim P_1[l_t|\mathbf{h}_t, \mathbf{a}_t]$.
    **if** $\tilde{l}_1 < \tilde{l}_0$ **then**
      $\mathbf{x}_{t+1}, \mathbf{s}_{t+1}, r_t, \text{done} \leftarrow \mathbf{Environment}(\mathbf{s}_t, \mathbf{w}_t)$
    **else**
      $\mathbf{x}_{t+1}, \mathbf{s}_{t+1}, r_t, \text{done} \leftarrow \mathbf{Environment}(\mathbf{s}_t, 0)$
    **end if**
    Recursively update the state predictor $\mathbf{h}_{t+1}$ with $\mathbf{x}_{t+1}$
      $\mathbf{h}_{t+1} \leftarrow \mathbf{GRU}(\mathbf{h}_t, \mathbf{Enc}_0(\mathbf{x}_{t+1}), \mathbf{a}_t)$
  **until** done is True

---

trains the multi-level image attack computational networks illustrated in Fig. 4. During training, the learning components and the environment are coupled and update their parameters simultaneously. The choice of time scales for the updates can have a significant impact on the behavior of the multi-time scale optimization process. For instance, in actor-critic [25], the critic has a faster update rate than the actor. In the generative adversarial network [26], the discriminator network has a faster update rate than the generative network. Following the heuristics and theories described in [25], [26], we set slower parameter update rates for the lower-level components.

Let us denote the parameters of the various components as follows: $\theta_n^{\text{img}}$ represents the parameters of $\mathbf{Enc}_1(\cdot)$ and $\mathbf{Dec}_1(\cdot)$, $\theta_n^{\text{sys}}$ represents the parameters of the dynamic autoencoder comprising $\mathbf{Enc}_0(\cdot)$, $\mathbf{GRU}(\cdot)$, and $\mathbf{Dec}_0(\cdot)$, $\theta_n^{\text{actor}}$ represents the parameters of the actor denoted by $\mathbf{Actor}(\cdot)$, and $\theta_n^{\text{critic}}$ represents the parameters of the critic denoted by $Q(\cdot,\cdot)$, which is the action-value function for policy evaluation. We update these parameters using different step sizes, based on the pace of the update rates, as follows:

$$\theta_{n+1}^{\text{img}} = \theta_n^{\text{img}} + \varepsilon_n^{\text{img}} S_n^{\text{img}}(\mathcal{M}_{\text{trajectory}})$$

$$\theta_{n+1}^{\text{dec}} = \theta_n^{\text{dec}} + \varepsilon_n^{\text{dec}} S_n^{\text{dec}}(\mathcal{M}_{\text{decision}})$$

$$\theta_{n+1}^{\text{actor}} = \theta_n^{\text{actor}} + \varepsilon_n^{\text{actor}} S_n^{\text{actor}}(\mathcal{M}_{\text{transition}})$$

$$\theta_{n+1}^{\text{critic}} = \theta_n^{\text{critic}} + \varepsilon_n^{\text{critic}} S_n^{\text{critic}}(\mathcal{M}_{\text{transition}})$$

$$\theta_{n+1}^{\text{sys}} = \theta_n^{\text{sys}} + \varepsilon_n^{\text{sys}} S_n^{\text{sys}}(\mathcal{M}_{\text{trajectory}}) \quad (3)$$

where the update functions $S_n^{\text{img}}$, $S_n^{\text{sys}}$, $S_n^{\text{actor}}$, $S_n^{\text{critic}}$ and $S_n^{\text{dec}}$ are stochastic gradients with loss functions (to be described in following sections) calculated with data samples from the replay buffers, i.e., $\mathcal{M}_{\text{trajectory}}$, $\mathcal{M}_{\text{transition}}$, and $\mathcal{M}_{\text{decision}}$. The replay buffers store a finite number of recently observed tuples of $(\mathbf{x}_t, \mathbf{a}_t)$, $(\mathbf{h}_{t-1}, \mathbf{a}_t, r_t, \mathbf{h}_t)$, and $(\mathbf{h}_t, \mathbf{a}_t, l_t)$ in $\mathcal{M}_{\text{trajectory}}$, $\mathcal{M}_{\text{transition}}$, and $\mathcal{M}_{\text{decision}}$ respectively.

The step sizes for the various components of our multi-time scale optimization, namely $\varepsilon_n^{\text{img}}$, $\varepsilon_n^{\text{dec}}$, $\varepsilon_n^{\text{actor}}$, $\varepsilon_n^{\text{critic}}$, and $\varepsilon_n^{\text{sys}}$, are

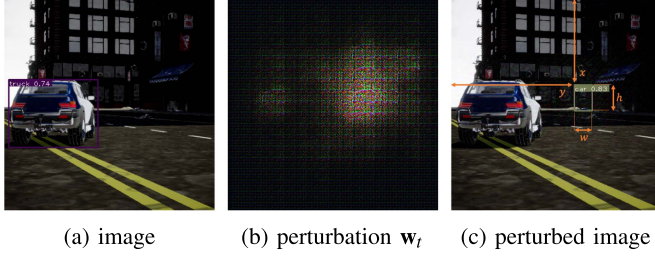(a) image  (b) perturbation $\mathbf{w}_t$  (c) perturbed image

Fig. 5. Fabrication of the bounding box at $(x, y, w, h)$ with $\mathbf{w}_t$.

determined as follows. The generation of an adversarial image perturbation depends on the generator with parameter $\theta_n^{\mathrm{img}}$, the actor that determines the attack coordinates with parameter $\theta_n^{\mathrm{actor}}$, and the binary decision maker that chooses whether to use the adversarial perturbation or not with parameter $\theta_n^{\mathrm{dec}}$. As the generation of the adversarial image perturbation and its use are governed by a policy with parameters $\theta_n^{\mathrm{img}}$, $\theta_n^{\mathrm{actor}}$, and $\theta_n^{\mathrm{dec}}$, we set faster update rates for the parameters that are relevant to policy evaluation, i.e., $\theta_n^{\mathrm{critic}}$ and $\theta_n^{\mathrm{sys}}$. Hence, the step size follows the diminishing rules as $n \to \infty$

$$\varepsilon_n^{\mathrm{img}}/\varepsilon_n^{\mathrm{dec}} \to 0, \ \varepsilon_n^{\mathrm{dec}}/\varepsilon_n^{\mathrm{actor}} \to 0, \ \varepsilon_n^{\mathrm{actor}}/\varepsilon_n^{\mathrm{critic}} \to 0, \ \varepsilon_n^{\mathrm{critic}}/\varepsilon_n^{\mathrm{sys}} \to 0.$$

This is because we intend to set slower update rates for the lower-level components of the policy that generate data for the upper-level components.

We describe the loss functions of the stochastic gradients for the multi-level stochastic optimization as follows:

*1) Image Attack Generator:* We utilize a *white box* model as a proxy object detector to train the attack generator. Specifically, we use the recently released version of *YOLO*, called *YOLOv5* [20], for this purpose.

Our image attack generator aims to fabricate bounding boxes at the target coordinates by injecting adversarial perturbations as shown in Fig. 1(b). Using reinforcement learning, the high-level attacker (reinforcement learning agent) selects the target coordinates to place the fabricated bounding boxes accordingly. Given the high-level attack $\mathbf{a}_t \in [0,1] \times [0,1] \times [0,1] \times [0,1]$ representing the coordinates $x$ and $y$ of a bounding box, its width and height, the image attack network aims to delete all other bounding boxes but keep the one corresponding to the high-level attack as illustrated in Fig. 5. By performing optimization iterations (500 iterations in Fig. 5), we can delete the existing bounding box and place a bounding box according to the target coordinates.

The iterative optimization approach that creates a bounding box for online image attacks, as shown in Fig. 5, is not suitable as it must be performed within a fixed time step of the control loop in the autonomous system. In our framework, we instead train an attack generator that minimizes the loss function through the image attack generator:

$$\arg\min_{\theta^{\mathrm{img}}} l^{img}\left(\mathbf{w}(\mathbf{x}; \theta^{\mathrm{img}}), \mathbf{x}, \mathbf{a}\right), \tag{4}$$

where $\mathbf{w}(\mathbf{x}; \theta^{\mathrm{img}}) = \mathbf{Dec}_1(\mathbf{Enc}_1(\mathbf{x}; \theta^{\mathrm{img}}), \mathbf{a}; \theta^{\mathrm{img}})$, and $\mathbf{x}$ and $\mathbf{a}$ are sampled from $\mathscr{M}_{\mathrm{trajectory}}$. This approach differs from the optimization over image space that is suitable for one-time use, i.e., $\arg\min_{\mathbf{w}} l^{img}(\mathbf{w}, \mathbf{x}, \mathbf{a})$. The stochastic gradient $S_n^{\mathrm{img}}(\mathscr{M}_{\mathrm{trajectory}})$ in (3) is associated with the loss function in (4). To fabricate the detected bounding box, i.e., inverted mapping

from the fabricated detection to the input image, we use the loss function employed from [1], where the same loss function is used to train the YOLO detector. The specific loss function in (4) from [1] is described in the appendix of the extended version [27].

*2) System Identification for State Estimation:* Due to incomplete observations of the state $\mathbf{s}_t$ through image stream $\mathbf{x}_t$, we need to identify the system to construct the estimator. The system identification determines the parameter that maximizes the state estimate's likelihood. We maximize the likelihood of state predictor by minimizing the cross-entropy error between true image streams and the predicted image streams by a stochastic optimization which samples trajectories saved in the memory buffer denoted by $\mathscr{M}_{\mathrm{trajectory}}$ with a loss function to minimize.

The loss function $l^{\mathrm{sys}}(\cdot)$ is calculated using the sampled trajectories from $\mathscr{M}_{\mathrm{trajectory}}$. We calculate the loss function as

$$l^{\mathrm{sys}}(\mathscr{M}_{\mathrm{trajectory}}; \theta_{\mathrm{sys}}) = \frac{1}{M}\sum_{m=1}^{M} H(\mathbf{X}_m, \hat{\mathbf{X}}_m) \tag{5}$$

where $\mathbf{X}_m = (\mathbf{x}_0, \ldots, \mathbf{x}_T)_m$ is the $m^{\mathrm{th}}$ sample image stream with time length $T$. Here, $H(\cdot, \cdot)$ is average of the binary cross-entropy $h(\cdot, \cdot)$ between the original image stream $\mathbf{X}_m$ and the predicted image stream $\hat{\mathbf{X}}_m$, which is computed over the RGB pixel values of the image streams.

We generate the predicted trajectory, $\hat{\mathbf{X}}_m = \{\hat{\mathbf{x}}_1, \ldots, \hat{\mathbf{x}}_T\}$, given the original trajectory with image stream $\mathbf{X}_m = (\mathbf{x}_0, \ldots, \mathbf{x}_T)_m$ and action stream $(\mathbf{a}_0, \ldots, \mathbf{a}_T)_m$ by processing them through the encoder, GRU, and the decoder as

$$\mathbf{h}_{t+1} = \mathrm{GRU}(\mathbf{h}_t, \mathrm{Enc}_1(\mathbf{x}_t), \mathbf{a}_t), \quad \mathbf{h}_0 \sim \mathcal{N}(0, \mathbf{I}),$$

$$\hat{\mathbf{x}}_{t+1} = \mathrm{Dec}_1(\mathbf{h}_{t+1}).$$

With the loss function in (5), the stochastic gradient for the optimization is defined as $S_n^{\mathrm{sys}} = -\nabla_{\theta_{\mathrm{sys}}} l^{\mathrm{sys}}(\mathscr{M}_{\mathrm{trajectory}}; \theta_{\mathrm{sys}})$.

*3) Actor-Critic Policy Improvement:* The attack coordinate $\mathbf{a}_t$ is determined by the policy, i.e., $\mathbf{a}_t = \mu(\mathbf{h}_t; \theta_{\mathrm{actor}})$ that maps the state estimate $\mathbf{h}_t$ into an action $\mathbf{a}_t$. To improve the policy, the critic evaluates the policy relying on the principle of optimality [28]. We employed an actor-critic method [29] for the reinforcement learning agent in the proposed framework. The critic network is updated using the state estimate $\mathbf{h}_t$ to apply the optimality principle with the following stochastic gradient as $S_n^{\mathrm{critic}} = -\nabla_{\theta_{\mathrm{critic}}} l^{\mathrm{critic}}(\mathscr{M}_{\mathrm{transition}}; \theta_{\mathrm{critic}})$ with the following loss function

$$l^{\mathrm{critic}}(\mathscr{M}_{\mathrm{transition}}; \theta^{\mathrm{critic}}) = \frac{1}{M}\sum_{m=1}^{M}(Q(\mathbf{h}_m, \mathbf{a}_m; \theta^{\mathrm{critic}}) - Q_m^{\mathrm{target}})^2,$$

where $Q_m^{\mathrm{target}} = r_m + \gamma Q(\mathbf{h}'_m, \mu(\mathbf{h}'_m; \theta_{\mathrm{actor}}); \theta_{\mathrm{critic}})$ and the state transition samples, i.e., $((\mathbf{h}, \mathbf{a}, \mathbf{h}', r)_0, \ldots, (\mathbf{h}, \mathbf{a}, \mathbf{h}', r)_M)$, are sampled from the replay buffer $\mathscr{M}_{\mathrm{transition}}$. With the same state transition data samples, we calculate the stochastic gradient for the policy update as $S_n^{\mathrm{actor}} = \nabla_{\theta_{\mathrm{actor}}} J(\mathscr{M}_{\mathrm{transition}}; \theta_{\mathrm{actor}})$ with the following estimated value function as

$$J = \frac{1}{M}\sum_{m=1}^{M} Q(\mathbf{h}_m, \mu(\mathbf{h}_m; \theta_{\mathrm{actor}}); \theta_{\mathrm{critic}}),$$

where $Q(\mathbf{h}, \mathbf{a}; \theta_{\mathrm{critic}})$ indicates the value of taking action $\mathbf{a}$ at the state estimated as $\mathbf{h}$.

(a) Drone to a car

(b) Cars and trucks



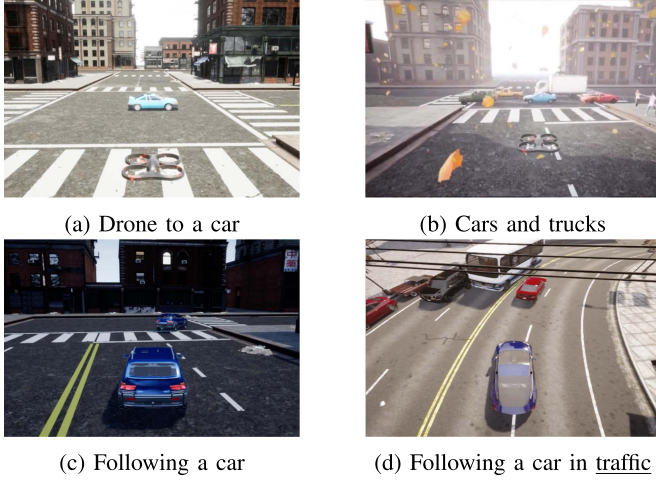(c) Following a car

(d) Following a car in traffic

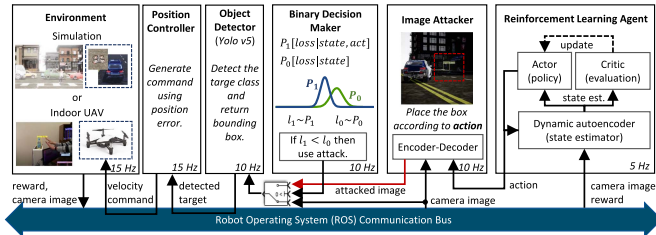Fig. 6.    Simulation environments.



Fig. 7: The framework implemented using *ROS*.

Fig. 7.    The framework implemented using *ROS*.

TABLE I
LIST OF COMPONENTS FOR ABLATION STUDY

| Methods | Generative network (Y/N) | State estimator (Y/N) | Attack switch (Y/N) |
|---|---|---|---|
| Iterative optimization | N | N | N |
| Generative attack | Y | N | N |
| Recursive attack | Y | Y | N |
| Neural Thompson sampling | Y | Y | Y |
| Conditional sampling | Y | Y | Y |

*4) Loss Estimators for the Binary Decision Making:* The stochastic gradient $S_n^{\text{dec}}$ in (3) is associated with the two optimizations described in (1). The entire stochastic optimization with the aforementioned stochastic gradients is summarized as Algorithm 2 in the appendix of the extended version [27].

## IV. EXPERIMENTS

We tested the proposed attack method to determine its ability to mislead autonomous vehicles in line with adversarial objectives. The adversary relied solely on image frames as sensing input and an uncertain actuator, in the form of an adversarial perturbation, to manipulate the paths of autonomous vehicles. Despite the adversary's limited sensor and uncertain actuator, our proposed algorithm successfully misled the autonomous
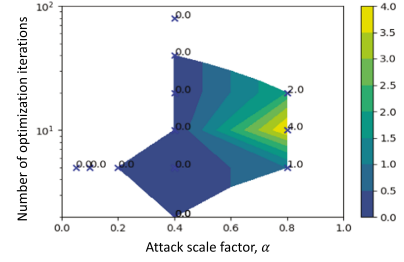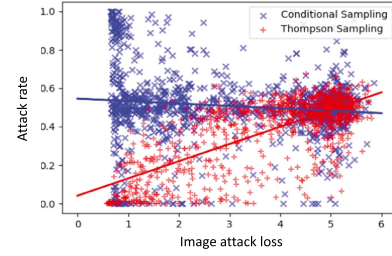


Fig. 8.    Terminal rewards in the third environment (Fig. 6(c)) of the baseline method with different of $\alpha$ and the iterations.



Fig. 9.    Attack rate vs. image attack loss of the 5 training experiments with the second environment (Fig. 6(b)).
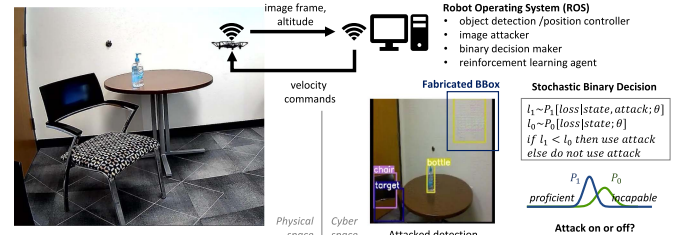


Fig. 10.    Attacking UAV's visual-tracking of a bottle.

vehicles in various simulation environments shown in Fig. 6 (and in illustrative videos[2]).

All our experiments consider attacking a vision-based guidance system depicted in Fig. 2 that uses *YOLOv5* object detector [20]. To simulate the vehicle environment, we employed a game development editor (*Unreal*) that is capable of building *photo-realistic* 3D environments, along with plug-in tools such as *AirSim* [30] and *CARLA* [31]. For the attack algorithm implementation, we used the robot operating system (*ROS*) to simultaneously implement the attack model learning and executing the attack using multiple modules (nodes) as illustrated in Fig. 7. All experiments were conducted on a desktop computer equipped with a GPU capable of rendering the 3D environments and performing DNN training.

### A. Baseline Method

Our proposed framework was compared to an image attack method presented in [12], [18], which uses iterative optimization

---

[2](a) Drone to a car at https://youtu.be/sjgQGgyLR8Y; (b) Cars and trucks at https://youtu.be/Xx9hH6mP0PE; (c) Following a car at https://youtu.be/mOPfPDEXkdM; (d) Following a car in traffic at https://youtu.be/G5gvGbqOPQM

TABLE II
ABLATION STUDY WITH THE LAST 10 EPISODES IN 5 TRAINING EXPERIMENTS, I.E., $N = 50$

| Environments | Object Tracking | Methods | Attack rate (%) | SSIM loss ($10^{-2}$, avg) | L2 loss ($10^{-4}$, avg) | Collision rate (%) | Terminal reward (avg±stdev) | Time AVG reward (avg±stdev) |
|---|---|---|---|---|---|---|---|---|
| Drone to a car (Figure 6a) | YOLO detection only | Normal | - | - | - | 0 | 15.7±0.0 | 1.29±0.07 |
| | | Iterative optimization | 100 | 13.6 | 12.1 | 0 | 15.7±0.1 | 1.29±0.07 |
| | | Generative attack | 100 | 8.74 | 3.94 | 74 | **30.1±8.3** | **2.58±0.80** |
| | | Recursive attack | 100 | 17.4 | 7.61 | **76** | 29.7±8.1 | 2.16±0.45 |
| | | Neural Thompson | **31.7** | **6.81** | **2.58** | 0 | 21.5±10.5 | 1.64±0.49 |
| | | Conditional sampling | 55.3 | 9.15 | 4.11 | 52 | 27.8±9.2 | 2.16±0.57 |
| Cars and trucks (Figure 6b) | YOLO detection only | Normal | - | - | - | 0 | 1.8±1.8 | 0.10±0.15 |
| | | Iterative optimization | 100 | 16.7 | 30.3 | 0 | −0.4±2.3 | −0.03±0.13 |
| | | Generative attack | 100 | 7.52 | 2.68 | 18 | 5.4±13.3 | 0.29±0.75 |
| | | Recursive attack | 100 | 7.26 | 2.74 | **36** | **9.4±6.3** | **0.87±0.43** |
| | | Neural Thompson | 49.3 | 4.94 | 1.77 | 8 | 6.0±9.2 | 0.25±0.46 |
| | | Conditional sampling | **42.7** | **1.97** | **0.81** | 20 | 2.6±10.6 | 0.10±0.53 |
| Following a car (Figure 6c) | YOLO detection and Kalman filter | Normal | - | - | - | 0 | 0±0 | 4.25±0.04 |
| | | Iterative optimization | 100 | 9.17 | 7.58 | 30 | 3.0±4.6 | 4.17±0.23 |
| | | Generative attack | 100 | 5.96 | 3.16 | 66 | 7.0±4.4 | 2.92±1.10 |
| | | Recursive attack | 100 | 4.88 | 2.61 | 74 | 7.4±4.4 | 2.92±1.33 |
| | | Neural Thompson | **33.5** | **1.66** | **0.96** | 52 | 5.2±5.0 | 4.01±0.39 |
| | | Conditional sampling | 72.8 | 3.10 | 2.06 | **76** | **7.6±4.3** | **4.27±3.54** |
| Following a car in traffic (Figure 6d) | YOLO detection and multi-object tracking (SORT) | Normal | - | - | - | 6 | 3.7±4.2 | **2.28±0.30** |
| | | Iterative optimization | 100 | 0.0 | 0.0 | 4 | 3.2±3.5 | 2.24±0.34 |
| | | Generative attack | 100 | 20.1 | 9.1 | **30** | **10.4±3.7** | 0.97±0.72 |
| | | Recursive attack | 100 | 23.5 | 10.6 | 2 | 9.9±2.2 | 0.63±0.60 |
| | | Neural Thompson | **34.6** | **8.1** | **3.7** | 10 | 8.1±5.3 | 1.39±0.83 |
| | | Conditional Sampling | 49.6 | 12.4 | 5.56 | 10 | 9.3±4.8 | 1.09±0.79 |

The stealthiness is evaluated by attack rate, ssim loss, and l2 loss. the attacker's performance is evaluated by collision rate, terminal reward, and time averaged reward.

with the image tensor as the decision variable. The effectiveness of these methods depends on the number of iterations and the scale factor $\alpha$, which limits the size of the adversarial perturbation as $\tilde{\mathbf{x}}_t = \max(\min(\mathbf{x}_t + \alpha \mathbf{w}_t, 1), 0)$. Previous works [12], [18] developed such methods as offline approaches. When an infinite number of iterations are allowed, the offline method can arbitrarily fabricate the bounding box, as illustrated in Fig. 5. To the best of our knowledge, no online image attack methods have been applied to autonomous vehicles. Therefore, we set a baseline method by applying the iterative optimization method to online applications with autonomous vehicles. In addition to the number of iterations, the scale factor $\alpha$ is a critical hyperparameter, as a higher value can increase the attacker's ability to fabricate bounding boxes, but it also reduces the stealthiness of the attack. For example, we collect the performance of the base line method with varying number of iteration and the scale factor as shown in Fig. 8. For the first three experiments (the first three rows in Table II), we set the hyperparameter $\alpha = 0.8$ and 20 iterations for the baseline method (Iterative optimization). The optimization method generated image perturbations approximately every 1 s. However, for our proposed methods, including Generative Attack, Recursive Attack, Neural Thompson, and Conditional Sampling, we set $\alpha = 0.05$ for stealthiness of the image attack. Additionally, the baseline method required manual annotation of the bounding box to be fabricated, as described in [19]. Therefore, we manually placed the target area to fabricate the bounding box according to the adversarial objectives, such as placing the bounding box to the left when the adversary needs to move the vehicle to the left. In the last experiment, we set the scale factor of the baseline method equal to that of our proposed methods, i.e., $\alpha = 0.05$.

## B. Ablation Study

We evaluated the proposed framework by conducting ablations, as presented in Table I. The proposed methods that learn when to use the image attack are denoted as *Neural Thompson*

*sampling* and *Conditional sampling* in Table I. The proposed methods have a generative network for the real-time adversarial image perturbation, a state estimator due to incomplete state observation through image streams, and an attack switch (the binary decision maker) that decides whether to use the attack or not. From top to bottom, the methods gradually add more components to the baseline method and become the proposed ones for the ablation study.

To test the effectiveness of each method, we conducted experiments in four different environments illustrated in Fig. 6. In the first environment (Fig. 6(a)), we set the reward as the distance between the ego-vehicle and the target object being tracked by the ego-vehicle. Thus, the adversary aims to increase the distance to move the vehicle away from the target while learning to increase rewards. In the second environment (Fig. 6(b)), the reward is set as the horizontal coordinate of ego-vehicle with respect to the target object. In this scenario, learning to increase the horizontal coordinate would lead to a crash. In the third environment (Fig. 6(c)), the adversary learns to increase the distance from the front car. In the fourth environment, we rewarded the learning agents when the distance between the ego-vehicle and the front car was greater than 50 meters. As shown in Table II, the four environments have different object-tracking methods. The first two environments use only **YOLO** detection. In the third environment, a Kalman filter smooths out the changes in the bounding boxes that are the outcomes of the detection using the double integrator dynamics. In the last environment, an multi-object tracking (SORT [32]) is implemented to deal with multiple cars in the traffic. We report the performance of the trained attackers (listed in Table I) with the last ten episodes of the entire 200 training episodes in Table II.

Compared to the baseline iterative optimization method, all the other methods, including the recursive attack method and the proposed conditional sampling method in Table II, demonstrated higher terminal rewards, collision rates, and terminal distances. Incorporating the state estimator improved the overall performance regarding terminal rewards, as shown in Table II. Since

we desire infrequent use of the image attack for stealthiness, we measure how frequently the attacks are used, i.e., attack rate. The utilization of binary decision-makers such as Neural Thompson Sampling (NTS) or the proposed conditional sampling resulted in decreased attack rates and the difference between the attacked images and the original images in terms of L2 norm and SSIM loss. Moreover, our proposed conditional sampling method showed greater terminal rewards compared to NTS.

The conditional sampling shows the greater use of the attack when the image attack loss is lower as we intended as shown in Fig. 9. In contrast, the Neural Thompson sampling (NTS) shows the opposite correlation, i.e., using the image attack when the loss values are higher. Further information regarding the simulations can be found in the appendix of the extended version [27].

### C. Real Robot Experiment

We implemented the proposed framework with a miniature drone (*DJI Tello*) to validate its efficacy in real-world scenarios, as depicted in Fig. 7. The drone employs IMU, optical-flow, and barometer for velocity estimation and to follow velocity commands. The drone was connected to a desktop computer via wifi-networks, as shown in Fig. 10. The objective of the online image training was to crash the UAV by teaching it to fabricate the bounding box. In the linked video at https://youtu.be/4w0pvQRCVHc, the online training lasted for ten minutes, and the UAV crashed successfully.

## V. CONCLUSION

This work showed a new online image attack framework that improves the iterative optimization-based methods that are more suitable for offline attack generation. In our proposed framework, the image attacks can be generated in real-time using only the image stream collected from the autonomous vehicle. Furthermore, the proposed conditional sampling for the binary decision making whether to use the attack (or not) improves the stealthiness by waiting until the proficiency increases. This work will serve as a stepping stone towards strengthening the perception in autonomous vehicles by learning worst-case attack scenarios.

## REFERENCES

[1] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 779–788.

[2] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "PointNet: Deep learning on point sets for 3D classification and segmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 652–660.

[3] "DJI mavic active tracking," 2016. Accessed: Dec. 01, 2021. [Online]. Available: https://youtu.be/ss0J0dAI1DM

[4] "What's next," 2022. [Online]. Available: https://waymo.com/intl/en_us/dataset-whats-next/

[5] T. Başar and P. Bernhard, *H-Infinity Optimal Control and Related Minimax Design Problems: A Dynamic Game Approach.* Berlin, Germany: Springer, 2008.

[6] L. Pinto, J. Davidson, R. Sukthankar, and A. Gupta, "Robust adversarial reinforcement learning," in *Proc. 34th Int. Conf. Mach. Learn.*, 2017, pp. 2817–2826.

[7] A. Kurakin, I. J. Goodfellow, and S. Bengio, "Adversarial machine learning at scale," in *Proc. 5th Int. Conf. Learn. Representations*, 2017.

[8] A. Kurakin, I. J. Goodfellow, and S. Bengio, "Adversarial examples in the physical world," in *Proc. 5th Int. Conf. Learn. Representations*, 2017.

[9] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," in *Proc. 3rd Int. Conf. Learn. Representations*, 2015, *arXiv:1412.6572*.

[10] E. Ackerman, "Three small stickers in intersection can cause Tesla autopilot to swerve into wrong lane," *IEEE Spectr.*, 2019. [Online]. Available: https://spectrum.ieee.org/three-small-stickers-on-road-can-steer-tesla-autopilot-into-oncoming-lane

[11] K. Eykholt et al., "Robust physical-world attacks on deep learning visual classification," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 1625–1634.

[12] Y. Jia, Y. Lu, J. Shen, Q. A. Chen, Z. Zhong, and T. Wei, "Fooling detection alone is not enough: First adversarial attack against multiple object tracking," in *Proc. Int. Conf. Learn. Representations*, 2020.

[13] "Tesla crash driver posted videos of himself riding without hands on wheel," *Guardian*, 2021. Accessed: Mar. 15, 2021. [Online]. Available: https://www.theguardian.com/us-news/2021/may/15/tesla-fatal-california-crash-autopilot

[14] A. Ilyas, L. Engstrom, A. Athalye, and J. Lin, "Black-box adversarial attacks with limited queries and information," in *Proc. 34th Int. Conf. Mach. Learn.*, 2018, pp. 2137–2146.

[15] Z. Wei et al., "Heuristic black-box adversarial attacks on video recognition models," in *Proc. AAAI Conf. Artif. Intell.*, 2020, pp. 12338–12345.

[16] C. Xiao, B. Li, J.-Y. Zhu, W. He, M. Liu, and D. Song, "Generating adversarial examples with adversarial networks," in *Proc. 27th Int. Joint Conf. Artif. Intell.*, 2018, pp. 3905–3911.

[17] A. Boloor, K. Garimella, X. He, C. Gill, Y. Vorobeychik, and X. Zhang, "Attacking vision-based perception in end-to-end autonomous driving models," *J. Syst. Architecture*, vol. 110, 2020, Art. no. 101766.

[18] S. Jha et al., "ML-driven malware that targets AV safety," in *Proc. Int. Conf. Dependable Syst. Netw.*, 2020, pp. 113–124.

[19] "Fooling detection alone is not enough: Adversarial attack against multiple object tracking," 2021. Accessed: Dec. 01, 2021. [Online]. Available: https://openreview.net/forum?id=rJl31TNYPr

[20] "YOLOv5│Pytorch," 2021. Accessed: Dec. 01, 2021. [Online]. Available: https://pytorch.org/hub/ultralytics_yolov5

[21] Y. Wang, K. Wang, Z. Zhu, and F.-Y. Wang, "Adversarial attacks on faster R-CNN object detector," *Neurocomputing*, vol. 382, pp. 87–95, 2020.

[22] J. Vermorel and M. Mohri, "Multi-armed bandit algorithms and empirical evaluation," in *Proc. Eur. Conf. Mach. Learn.*, 2005, pp. 437–448.

[23] W. Zhang, D. Zhou, L. Li, and Q. Gu, "Neural Thompson sampling," in *Proc. Int. Conf. Learn. Representations*, 2021.

[24] Y. Gal and Z. Ghahramani, "Dropout as a Bayesian approximation: Representing model uncertainty in deep learning," in *Proc. 34th Int. Conf. Mach. Learn.*, 2016, pp. 1050–1059.

[25] V. R. Konda and J. N. Tsitsiklis, "Actor-critic algorithms," in *Proc. Adv. Neural Inf. Process. Syst. 12*. 2000, pp. 1008–1014.

[26] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter, "GANs trained by a two time-scale update rule converge to a local Nash equilibrium," in *Proc. Annu. Conf. Neural Inf. Process. Syst.*, 2017, pp. 6626–6637.

[27] H.-J. Yoon, H. Jafarnejadsani, and P. Voulgaris, "Learning when to use adaptive adversarial image perturbations against autonomous vehicles," 2022, *arXiv:2212.13667*.

[28] R. Bellman, "Dynamic programming," *Science*, vol. 153, no. 3731, pp. 34–37, 1966.

[29] T. P. Lillicrap et al., "Continuous control with deep reinforcement learning," in *Proc. Int. Conf. Learn. Representations*, 2016.

[30] S. Shah, D. Dey, C. Lovett, and A. Kapoor, "Airsim: High-fidelity visual and physical simulation for autonomous vehicles," in *Field and Service Robotics*. Berlin, Germany: Springer, 2017, pp. 621–635 .

[31] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun, "CARLA: An open urban driving simulator," in *Proc. 1st Annu. Conf. Robot Learn.*, 2017, pp. 1–16.

[32] A. Bewley, Z. Ge, L. Ott, F. Ramos, and B. Upcroft, "Simple online and realtime tracking," in *Proc. IEEE Int. Conf. Image Process.*, 2016, pp. 3464–3468.