

---

# Fooling Detection Alone is Not Enough: First Adversarial Attack against Multiple Object Tracking

---

Yunhan Jia\*, Yantao Lu\*, Junjie Shen<sup>†</sup>, Qi Alfred Chen<sup>†</sup>, Zhenyu Zhong, Tao Wei  
Baidu X-Lab, <sup>†</sup>UC Irvine  
{jiayunhan, yantaolu, edwardzhong, lenx}@baidu.com, junjies1, alfchen@uci.edu

## Abstract

Recent work in adversarial machine learning started to focus on the visual perception in autonomous driving and studied Adversarial Examples (AEs) for object detection models. However, in such visual perception pipeline the detected objects must also be tracked, in a process called Multiple Object Tracking (MOT), to build the moving trajectories of surrounding obstacles. Since MOT is designed to be robust against errors in object detection, it poses a general challenge to existing attack techniques that blindly target objection detection: we find that a success rate of over 98% is needed for them to actually affect the tracking results, a requirement that no existing attack technique can satisfy. In this paper, we are the first to study adversarial machine learning attacks against the complete visual perception pipeline in autonomous driving, and discover a novel attack technique, tracker hijacking, that can effectively fool MOT using AEs on object detection. Using our technique, successful AEs on as few as one single frame can move an existing object in to or out of the headway of an autonomous vehicle to cause potential safety hazards. We perform evaluation using the Berkeley Deep Drive dataset and find that on average when 3 frames are attacked, our attack can have a nearly 100% success rate while attacks that blindly target object detection only have up to 25%.

## 1 Introduction

Since the first Adversarial Example (AE) against traffic sign image classification discovered by Eykholt *et al.* [10], several research work in adversarial machine learning [9, 30, 15, 16, 35, 6] started to focus on the context of visual perception in autonomous driving, and studied AEs on object detection models. For example, Eykholt *et al.* [9] and Zhong *et al.* [36] studied AEs in the form of adversarial stickers on stop signs or the back of front cars against YOLO object detectors [23], and performed indoor experiments to demonstrate the attack feasibility in the real world. Building upon these work, most recently Zhao *et al.* [35] leveraged image transformation techniques to improve the robustness of such adversarial sticker attacks in outdoor settings, and were able to achieve a 72% attack success rate with a car running at a constant speed of 30 km/h on real roads.

While these results from prior work are alarming, object detection is in fact only the first half of the visual perception pipeline in autonomous driving, or in robotic systems in general — in the second half, the detected objects must also be tracked, in a process called *Multiple Object Tracking (MOT)*, to build the moving trajectories, called *trackers*, of surrounding obstacles. This is *required* for the subsequent driving decision making process, which needs the built trajectories to predict future moving trajectories for these obstacles and then plan a driving path accordingly to avoid collisions with them. To ensure high tracking accuracy and robustness against errors in object detection, in MOT only the detection results with sufficient consistency and stability across multiple frames can be included in the tracking results and actually influence the driving decisions. Thus, MOT in the visual perception of autonomous driving poses a general challenge to existing attack techniques that blindly

---

\*Equal contribution

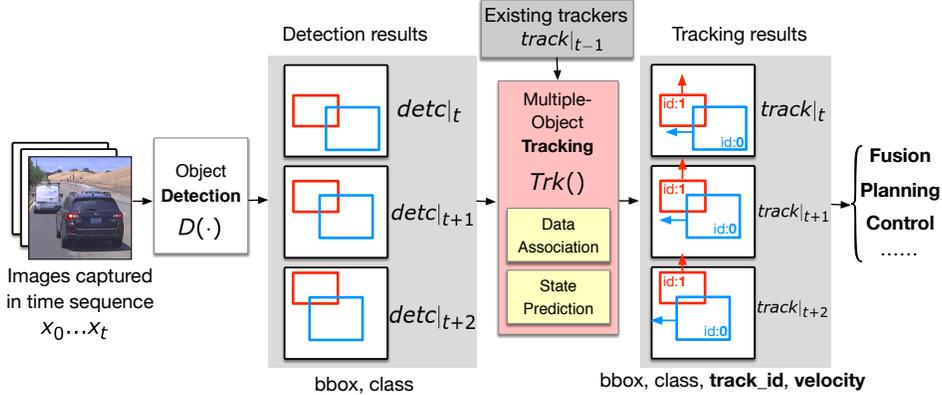


Figure 1: The complete visual perception pipeline in autonomous driving, i.e., both object detection and Multiple Object Tracking (MOT) [3, 13, 12, 34, 8, 19, 28].

target objection detection. For example, as shown by our analysis later in §4, an attack on objection detection needs to succeed consecutively for at least 60 frames to fool a representative MOT process, which requires an at least 98% attack success rate (§4). To the best of our knowledge, no existing attacks on objection detection can achieve such a high success rate [9, 30, 15, 16, 35, 6].

In this paper, we are the first to study adversarial machine learning attacks considering the *complete* visual perception pipeline in autonomous driving, i.e., both object detection and object tracking, and discover a novel attack technique, called *tracker hijacking*, that can effectively fool the MOT process using AEs on object detection. Our key insight is that although it is highly difficult to directly create a tracker for fake objects or delete a tracker for existing objects, we can carefully design AEs to attack the tracking error reduction process in MOT to deviate the tracking results of existing objects towards an attacker-desired moving direction. Such process is designed for increasing the robustness and accuracy of the tracking results, but ironically, we find that it can be exploited by attackers to substantially alter the tracking results. Leveraging such attack technique, successful AEs on as few as *one single frame* is enough to move an existing object in to or out of the headway of an autonomous vehicle and thus may cause potential safety hazards.

We select 20 out of 100 randomly sampled video clips from the Berkeley Deep Drive dataset to evaluate our attack technique. Under recommended MOT algorithm configurations in practice [37] and normal measurement noise levels, we find that our attack can succeed with successful AEs on as few as *one frame*, and 2 to 3 consecutive frames on average. We also reproduce and compare with previous attacks that blindly target object detection, and find that when attacking 3 consecutive frames, our attack has a nearly 100% success rate while attacks that blindly target object detection only have up to 25%.

**Contributions.** In summary, this paper makes the following contributions:

- We are the first to study adversarial machine learning attacks considering the complete visual perception pipeline in autonomous driving, i.e., both object detection and MOT. We find that without considering MOT, an attack blindly targeting object detection needs at least a success rate of 98% to actually affect the complete visual perception pipeline in autonomous driving, which is a requirement that no existing attack technique can satisfy.
- We discover a novel attack technique, tracker hijacking, that can effectively fool MOT using AEs on object detection. This technique exploits the tracking error reduction process in MOT, and can enable successful AEs on as few as one single frame to move an existing object in to or out of the headway of an autonomous vehicle to cause potential safety hazards.
- The attack evaluation using the Berkeley Deep Drive dataset shows that our attack can succeed with successful AEs on as few as one frame, and only 2 to 3 consecutive frames on average, and when 3 consecutive frames are attacked, our attack has a nearly 100% success rate while attacks that blindly target object detection only have up to 25%.
- Code and evaluation data are all available at an anonymized GitHub repository [1].

## 2 Background and Related Work

**Adversarial examples for object detection.** Since the first physical adversarial examples against traffic sign classifier demonstrated by Eykholt *et al.* [10], several work in adversarial machine learning [9, 30, 15, 16, 35, 6] have been focused on the visual perception task in autonomous driving, and more specifically, the object detection models. To achieve high attack effectiveness in practice, the key challenge is how to design robust attacks that can survive distortions in real-world driving scenarios such as different viewing angles, distances, lighting conditions, and camera limitations. For example, Lu *et al.* [15] shows that AEs against Faster-RCNN [25] generalize well across a sequence of images in digital space, but fail in most of the sequence in physical world; Eykholt *et al.* [9] generates adversarial stickers that, when attached to stop sign, can fool YOLOv2 [23] object detector, while it is only demonstrated in indoor experiment within short distance; Chen *et al.* [6] generates AEs based on expectation over transformation techniques, while their evaluation shows that the AEs are not robust to multiple angles, probably due to not considering perspective transformations [35]. It was not until recently that physical adversarial attacks against object detectors achieve a decent success rate (70%) in fixed-speed (6 km/h and 30 km/h) road test [35].

While the current progress in attacking object detection is indeed impressive, in this paper we argue that in the actual visual perception pipeline of autonomous driving, object tracking, or more specifically MOT, is an integral step, and without considering it, existing adversarial attacks against object detection still cannot affect the visual perception results even with high attack success rate. As shown in our evaluation in §4, with a common setup of MOT, an attack on object detection needs to reliably fool at least 60 consecutive frames to erase one object (e.g., stop sign) from the tracking results, in which case even a 98% attack success rate on object detectors is not enough (§4).

**MOT background.** MOT aims to identify objects and their trajectories in video frame sequence. With the recent advances in object detection, *tracking-by-detection* [18] has become the dominant MOT paradigm, where the detection step identifies the objects in the images and the tracking step links the objects to the trajectories (*i.e.*, trackers). Such paradigm is widely adopted in autonomous driving systems today [3, 13, 12, 34, 8, 19, 28], and a more detailed illustration is in Fig. 1. As shown, each detected objects at time  $t$  will be associated with a dynamic state model (e.g., position, velocity), which represents the past trajectory of the object ( $track|_{t-1}$ ). A per-track Kalman filter [3, 13, 11, 20, 32] is used to maintain the state model, which operates in a recursive *predict-update* loop: the predict step estimates current object state according to a motion model, and the update step takes the detection results  $detect|_t$  as *measurement* to update its state estimation result  $track|_t$ .

The association between detected objects with existing trackers is formulated as a bipartite matching problem [27, 11, 20] based on the pairwise similarity costs between the trackers and detected objects, and the most commonly used similarity metric is the spatial-based cost, which measures the overlapping between bounding boxes, or bboxes [3, 14, 29, 27, 11, 20, 37, 32, 4, 5]. To reduce errors in this association, an accurate velocity estimation is necessary in the Kalman filter prediction [7, 31]. Due to the discreteness of camera frames, Kalman filter uses the velocity model to estimate the location of the tracked object in the next frame in order to compensate the object motion between frames. However, as described later in §3, such error reduction process unexpectedly makes it possible to perform tracker hijacking.

MOT manages tracker creation and deletion with two thresholds. Specifically a new tracker will be created only when the object has been constantly detected for a certain number of frames, this threshold will be referred to as the *hit count*, or  $H$  in the rest of the paper. This helps to filter out occasional false positives produced by object detectors. On the other hand, a tracker will be deleted if no objects is associated with for a duration of  $R$  frames, or called a *reserved age*. It prevents the tracks from being accidentally deleted due to infrequent false negatives of object detectors. The configuration of  $R$  and  $H$  usually depends on both the accuracy of detection models, and the frame rate (fps). Previous work suggest a configuration of  $R = 2 \cdot \text{fps}$ , and  $H = 0.2 \cdot \text{fps}$  [37], which gives a  $R = 60$  frames and  $H = 6$  frames for a common 30 fps visual perception system. We will show in §4 that an attack that blindly targeting object detection needs to constantly fool at least 60 frames ( $R$ ) to erase an object, while our proposed tracker hijacking attack can fabricate object that last for  $R$  frames and vanish target object for  $H$  frames in the tracking result by attacking as few as one frame, and only 2~3 frames on average ( $S4$ ).

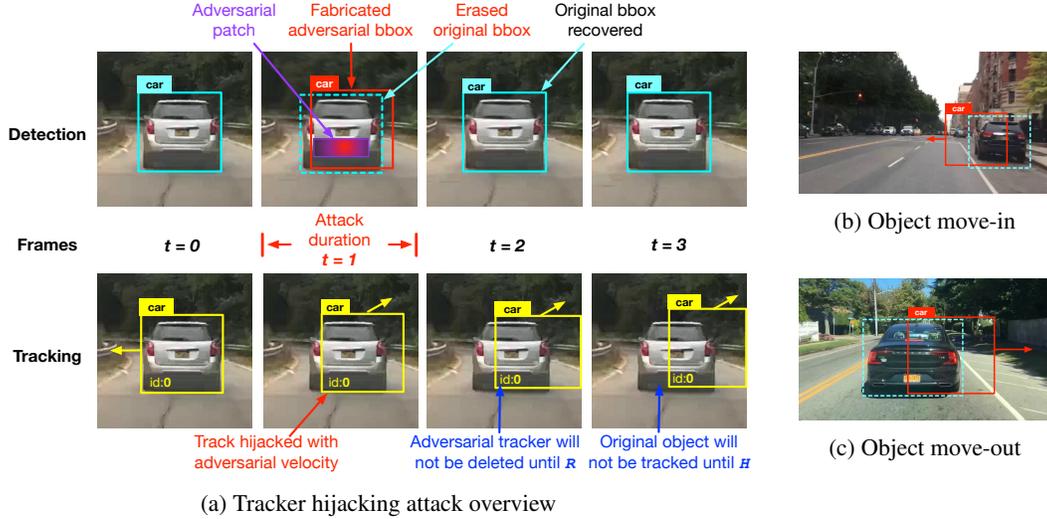


Figure 2: Description of the tracker hijacking attack flow (a), and two different attack scenarios: object move-in (b) and move-out (c), where tracker hijacking may lead to severe safety consequences including emergency stop and rear-end crashes.

### 3 Tracker Hijacking Attack

**Overview.** Fig. 2a illustrates the tracker hijacking attack discovered in this paper, in which an AE for object detection (e.g., in the form of adversarial patches on the front car) that can fool the detection result for as few as one frame can largely deviate the tracker of a target object (e.g., a front car) in MOT. As shown, the target car is originally tracked with a predicted velocity to the left at  $t_0$ . The attack starts at time  $t_1$  by applying an adversarial patch onto the back of the car. The patch is carefully generated to fool the object detector with two adversarial goals: (1) erase the bounding box of target object from detection result, and (2) fabricate a bounding box with similar shape that is shifted a little bit towards an attacker-specified direction. The fabricated bounding box (red one in detection result at  $t_1$ ) will be associated with the original tracker of target object in the tracking result, which we call a *hijacking* of the tracker, and thus would give a fake velocity towards the attacker-desired direction to the tracker. The tracker hijacking shown in Fig. 2a lasts for only one frame, but its adversarial effects could last tens of frames, depending on the MOT parameter  $R$  and  $H$  (introduced in §2). For example, at time  $t_2$  after the attack, all detection bounding boxes are back to normal, however, two adversarial effects persist: (1) the tracker that has been hijacked with attacker-induced velocity *will not be deleted until a reserved age ( $R$ ) has passed*, and (2) the target object, though is recovered in the detection result, *will not be tracked until a hit count ( $H$ ) has reached*, and before that the object remains missing in the tracking result. However, it’s important to note that our attack may not always succeed with one frame in practice, as the recovered object may still be associated with its original tracker, if the tracker is not deviated far enough from the object’s true position during a short attack duration. Our empirical results show that our attack usually achieves a nearly 100% success rate when 3 consecutive frames are successfully attacked using AE (§4).

Such persistent adversarial effects may cause severe safety consequences in self-driving scenarios. We highlight two attack scenarios that can cause emergency stop or even a rear-end crashes:

**Attack scenario 1: Target object move-in.** Shown in Fig. 2b, an adversarial patch can be placed on roadside objects, e.g., a parked vehicle to deceive visual perception of autonomous vehicles passing by. The adversarial patch is generated to cause a translation of the target bounding box towards the center of the road in the detection result, and the hijacked tracker will appear as a moving vehicle cutting in front in the perception of the victim vehicle. This tracker would last for 2 seconds if  $R$  is configured as 2-fps as suggested in [37], and tracker hijacking in this scenario could cause an emergency stop and potentially a rear-end crash.

**Attack scenario 2: Target object move-out.** Similarly, tracker hijacking attack can also deviate objects in front of the victim autonomous vehicle away from the road to cause a crash as shown in Fig. 2c. Adversarial patch applied on the back of front car could deceive MOT of autonomous vehicle behind into believing that the object is moving out of its way, and the front car will be missing from

the tracking result for a duration of  $200ms$ , if  $H$  uses the recommended configuration of  $0.2\cdot$  fps [37]. This may cause the victim autonomous vehicle to crash into the front car.

### 3.1 Attack Methodology

---

#### Algorithm 1 Tracker Hijacking Attack

---

**Input:** Video image sequence  $X = [x_0, x_1, \dots, x_n]$ ; object detector  $D(\cdot)$ ; MOT algorithm  $Trk(\cdot)$ ;

**Input:** Index of target object to be hijacked  $K$ , attacker-desired directional velocity  $\vec{v}$ , adversarial patch area as a mask matrix  $patch$ .

**Output:** Sequence of adversarial examples  $X' = [x'_1, \dots, x'_r]$  required for a successful attack.

**Initialization**  $X' \leftarrow \{\}$ ,  $detc|_0 \leftarrow D(x_0)$ ,  $track|_0 \leftarrow \{current\_tracks\}$

```

1: for  $t = 1$  to  $n$  do
2:    $detc|_t \leftarrow D(x_t)$ 
3:   if  $detc|_t[K]$  matches  $track|_{t-1}[K]$  then  $\triangleright$  target object matches with an existing tracker
4:     find position  $pos$  to place fabricated bbox with Eq. 1
            $pos \leftarrow \text{FINDPOS}(Trk(\cdot), track|_{t-1}, K, \vec{v}, patch)$  see SuppAlg.1.1
5:     generate adversarial frame  $x'$  with Eq. 2  $\triangleright$  attack object detector with specialized loss
            $x'_t \leftarrow \text{GENERATEADV}(x, D(\cdot), pos, K, patch)$  see SuppAlg.1.2
6:      $X' \leftarrow^+ x'_t$ 
7:   else
8:     return  $X'$   $\triangleright$  attack succeeds when target object is not associated with original tracker
9:   end if
10:   $track|_t \leftarrow Trk(track|_{t-1}, D(x'_t))$   $\triangleright$  update current tracker with adversarial frame
11: end for

```

---

**Targeted MOT design.** Our attack targets the most common MOT pipeline described in §2. Specifically, we target first-order Kalman filter which predicts a state vector containing position and velocity of detected objects over time. For the data association, we adopt the mostly widely used Intersection over Union (IoU) as the similarity metric, and the IoU between bounding boxes are calculated by Hungarian matching algorithm [17] to solve the bipartite matching problem that associates bounding boxes detected in consecutive frames with existing trackers. Such combination of algorithms in the MOT is the most common in previous work [14, 29, 27] and real-world systems [3].

We now describe our methodology of generating an adversarial patch that manipulates detection results to hijack a tracker. As detailed in Alg. 1, given a targeted video image sequence, the attack iteratively finds the minimum required frames to perturb for a successful track hijack, and generates the adversarial patches for these frames. In each attack iteration, an image frame in the original video clip is processed, and given the index of target objects  $K$ , the algorithm finds an optimal position to place the adversarial bounding box  $pos$  in order to hijack the tracker of target object by solving Eq. 1. The attack then constructs adversarial frame against object detection model with an adversarial patch, using Eq. 2 as the loss function to erase the original bounding box of target object and fabricate the adversarial bounding box at the given location. The tracker is then updated with the adversarial frame that deviates the tracker from its original direction. If the target object in the next frame is not associate with its original tracker by the MOT algorithm, attack has succeeded; otherwise, this process is repeated for the next frame. We discuss two critical steps in this algorithm below, and please refer to our supplementary material for the complete implementation of the algorithm.

**Finding optimal position for adversarial bounding box.** To deviate the tracker of a target object  $K$ , besides removing its original bounding box  $detc|_t[K]$ , the attack also needs to fabricate an adversarial box with a shift  $\delta$  towards a specified direction. This turns into an optimization problem (Eq. 1) of finding the translation vector  $\delta$  that maximizes the cost of Hungarian matching ( $\mathcal{M}(\cdot)$ ) between the detection box and the existing tracker so that the bounding box is still associated with its original tracker ( $\mathcal{M} \leq \lambda$ ), but the shift is large enough to give an adversarial velocity to the tracker. Note that we also limit the shifted bounding box to be overlapped with the  $patch$  to facilitate

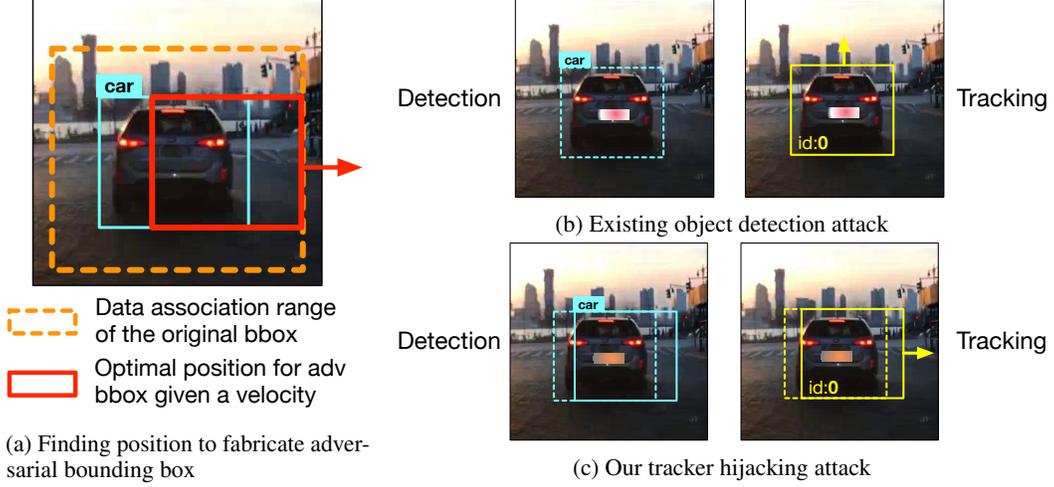


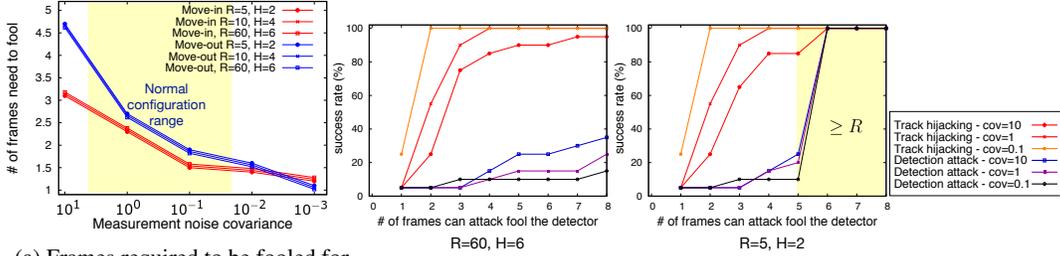
Figure 3: Comparison between previous object detection attack and our tracker hijacking attack. Previous attack that simply erase the bbox has no impact on the tracking output (b), while tracker hijacking attack that fabricates bbox with carefully chosen position successfully redirects the tracker towards attacker-specified direction (c).

adversarial example generation, as it's often easier for adversarial perturbations to affect prediction results in its proximity, especially in physical settings [6].

$$\begin{aligned} & \max_{\delta} \mathcal{M}(detc|_t[K] + \delta, track|_{t-1}[K]) \\ & s.t. \mathcal{M} \leq \lambda, IoU(detc|_t[K] + \delta, patch) \geq \gamma \end{aligned} \quad (1)$$

**Generating adversarial patch against object detection.** Similar to the existing adversarial attacks against object detection models [6, 10, 35], we also formulate the adversarial patch generation as an optimization problem shown in Eq. 2. Existing attacks without considering MOT directly minimize the probability of target class (e.g., a stop sign) to erase the target from detection result. However, as shown in Fig. 3b, such AEs are highly ineffective in fooling MOT as the tracker will still track for  $R$  frames even after the detection bounding box is erased. Instead, the loss function of our tracker hijacking attack incorporates two loss terms:  $\mathcal{L}_1$  minimizes the target class probability at given location to erase the target bounding box, where  $\sum_{i=0}^B \mathbb{1}_i^{obj}$  identifies all bounding boxes ( $B$ ) before non-max suppression [21], who contain the center location  $(cx_t, cy_t)$  of  $pos$ , while  $C_i$  is the confidence score of bounding boxes;  $\mathcal{L}_2$  controls the fabrication of adversarial bounding box at given center location  $(cx_t, cy_t)$  with given shape  $(w_t, h_t)$  to hijack the tracker. In the implementation, we use Adam optimizer to minimize the loss by iteratively perturbing the pixels along the gradient directions within the patch area, and the generation process stops when an adversarial patch that satisfies the requirements is generated. Note that the fabrication loss  $\mathcal{L}_2$  needs only to be used when generating the first adversarial frame in a sequence to give the tracker an attacker-desired velocity  $\vec{v}$ , and then  $\lambda$  can be set to 0 to only focus on erasing target bounding box similar to previous work. Thus, our attack wouldn't add much difficulty to the optimization. Details of our algorithm can be found in the supplementary material, and the implementation can be found at [1].

$$\begin{aligned} & \min_{\Delta \in patch} \mathcal{L}_1(x_t + \Delta) + \lambda \cdot \mathcal{L}_2(x_t + \Delta) \\ \mathcal{L}_1 &= \sum_{i=0}^B \mathbb{1}_i^{obj} \cdot [C_i^2 - CrossEntropy(p_i, class_t)] \\ \mathcal{L}_2 &= \sum_{i=0}^B \mathbb{1}_i^{obj} \cdot \{ [(cx_i - cx_t)^2 + (cy_i - cy_t)^2] + [(\sqrt{w_i} - \sqrt{w_t})^2 + (\sqrt{h_i} - \sqrt{h_t})^2] \\ & \quad + (1 - C_i)^2 + CrossEntropy(p_i, class_t) \} \end{aligned} \quad (2)$$



(a) Frames required to be fooled for a successful tracker hijack

(b) Attack success rate at  $R = 60, H = 6$ , and  $R = 5, H = 2$

Figure 4: In normal measurement noise covariance range (a), our tracker hijacking attack would require the adversarial example to fool only 2~3 consecutive frames on average to successfully deviate the target tracker despite the  $(R, H)$  settings. Moreover we compare the success rate of tracker hijacking with previous adversarial attack against object detectors only under different attacker capabilities, *i.e.*, the number of consecutive frames the adversarial example can reliably fool the object detector (b). Tracker hijacking achieves superior attack success rate (100%) even by fooling as few as 3 frames, while previous attack is only effective when the adversarial example can reliably fools at least  $R$  consecutive frames.

## 4 Attack Evaluation

In this section, we describe our experiment settings for evaluating the effectiveness of our tracker hijacking attack, and comparing it with previous attacks that blindly attack object detection in detail.

### 4.1 Experiment Methodology

**Evaluation metrics.** We define a successful attack as that *the detected bounding box of target object can no longer be associated with any of the existing trackers when attack has stopped*. We measure the effectiveness of our track hijacking attack using the minimum number of frames that the AEs on object detection need to succeed. The attack effectiveness highly depends on the difference between the direction vector of the original tracker and adversary’s objective. For example, attacker can cause a large shift on tracker with only one frame if choosing the adversarial direction to be opposite to its original direction, while it would be much harder to deviate the tracker from its established track, if the adversarial direction happens to be the same as the target’s original direction. To control the variable, we measure the number of frames required for our attack in two previous defined attack scenarios: target object move-in and move-out. Specifically, in all move-in scenarios, we choose the vehicle parked along the road as target, and the attack objective is to move the tracker to the center, while in all move-out scenarios, we choose vehicles that are moving forward, and the attack objective is to move the target tracker off the road.

**Dataset selection.** We randomly sampled 100 video clips from Berkeley Deep Drive dataset [33], and then manually selected 10 suitable for the object move-in scenario, and another 10 for the object move-out scenario. For each clip, we manually label a target vehicle and annotate the patch region to be a small area at the back of it as shown in Fig. 3c. All videos have the same frame rate of 30 fps.

**Implementation details.** We implement our targeted visual perception pipeline using Python, with YOLOv3 [24] as the object detection model since it is among the most popular detectors used by real-time systems. For the MOT implementation, we use the Hungarian matching implementation called `linear_assignment` in the `sklearn` package for the data association, and we provide a reference implementation of Kalman filter based on the one used in OpenCV [22].

The effectiveness of attack depends on a configuration parameter of Kalman filter, called *measurement noise covariance* (*cov*). *cov* is an estimation about how much noise is in the system, a low *cov* value would give Kalman filter more confidence on the detection result at time  $t$  when updating the tracker, while a high *cov* value would make Kalman filter to place trust more on its own previous prediction at time  $t - 1$  than that at time  $t$ . We give a detailed introduction of configurable parameters in Kalman filter in §2 of our supplementary material. This measurement noise covariance is often tuned based on the performance of detection models in practice. We evaluate our approach under different *cov* configurations ranging from very small ( $10^{-3}$ ) to very large (10) as shown in Fig. 4a, while *cov* is usually set between 0.01 and 10 in practice [3, 13].

## 4.2 Evaluation Results

**Effectiveness of tracker hijacking attack.** Fig. 4a shows the average number of frames that the AEs on object detection need to fool for a successful track hijacking over the 20 video clips in the evaluation. Although a configuration with  $R = 60$  and  $H = 6$  is recommended when fps is 30 [37], we still test different reserved age ( $R$ ) and hit count ( $H$ ) combinations as real-world deployment are usually more conservative and use smaller  $R$  and  $H$  [3, 13]. The results show that tracker hijacking attack only requires successful AEs on object detection in 2 to 3 consecutive frames on average to succeed despite the  $(R, H)$  configurations. We also find that even with a successful AE on only one frame, our attack still has 50% and 30% success rates when  $cov$  is 0.1 and 0.01 respectively.

Interestingly, we find that object move-in generally requires less frames compared with object move-out. The reason is that the parked vehicles in move-in scenarios (Fig. 2b) naturally have a moving-away velocity relative to the autonomous vehicle. Thus, compared to move-out attack, move-in attack triggers a larger difference between the attacker-desired velocity and the original velocity. This makes the original object, once recovered, harder to associate correctly, making hijacking easier.

**Comparison with attacks that blindly target object detection.** Fig. 4b shows the success rate of our attack and previous attacks that blindly target object detection, which we denote as *detection attack*. We reproduced the recent adversarial patch attack on object detection from Jia *et al.* [36] in 2018, which targets the autonomous driving context and has validated attack effectiveness using real-world car testing. In this attack, the objective is to erase the target class from the detection result of each frame. Evaluated under two  $(R, H)$  settings, we find that tracker hijacking attack achieves superior attack success rate (100%) even by attacking as few as 3 frames, while the detection attack needs to reliably fool at least  $R$  consecutive frames to guarantee success. When  $R$  is set to 60 according to the frame rate of 30 fps, the detection attack needs to have an adversarial patch that can constantly succeed at least 60 frames while the victim autonomous vehicle is driving. It translates to an over 98.3% ( $\frac{59}{60}$ ) AE success rate, which has never been achieved or even got close to in previous work [35, 9, 6, 15]. Note that the detection attack still can have up to  $\sim 25\%$  success rate before  $R$ . This is because the detection attack causes the object to disappear for some frames, and when the vehicle heading changes during such disappearing period, it is still possible to cause the original object, when recovered, to misalign with the tracker predication in the original tracker. However, since our attack is designed to intentionally mislead the tracker predication in MOT, our success rate is substantially higher (3-4 $\times$ ) and can reach 100% with as few as 3 frames attacked.

## 5 Discussion and Future Work

**Implications for future research in this area.** Today, adversarial machine learning research targeting the visual perception in autonomous driving, no matter on attack or defense, uses the accuracy of objection detection as the *de facto* evaluation metric [18]. However, as concretely shown in our work, without considering MOT, successful attacks on the detection results alone do not have direct implication on equally or even closely successful attacks on the MOT results, the final output of the visual perception task in real-world autonomous driving [3, 13]. Thus, we argue that future research in this area should consider: (1) using the MOT accuracy as the evaluation metric, and (2) instead of solely focusing on object detection, also studying weaknesses specific to MOT or interactions between MOT and object detection, which is a highly under-explored research space today. This paper marks the first research effort towards both directions.

**Practicality improvement.** Our evaluation currently are all conducted digitally with captured video frames, while our method should still be effective when applied to generate physical patches. For example, our proposed adversarial patch generation method can be naturally combined with different techniques proposed by previous work to enhance reliability of AEs in the physical world (*e.g.*, non-printable loss [26] and expectation-over-transformation [2]). We leave this as future work.

**Generality improvement.** Though in this work we focused on MOT algorithm that uses IoU based data association, our approach of finding location to place adversarial bounding box is generally applicable to other association mechanisms (*e.g.*, appearance-based matching). Our AE generation algorithm against YOLOv3 should also be applicable to other object detection models with modest adaptations. We plan to provide reference implementations of more real-world end-to-end visual perception pipelines to pave the way for future adversarial learning research in self-driving scenarios.

## 6 Conclusion

In this work, We are the first to study adversarial machine learning attacks against the complete visual perception pipeline in autonomous driving, i.e., both object detection and MOT. We discover a novel attack technique, tracker hijacking, that exploits the tracking error reduction process in MOT and can enable successful AEs on as few as one frame to move an existing object in to or out of the headway of an autonomous vehicle to cause potential safety hazards. The evaluation results show that on average when 3 frames are attacked, our attack can have a nearly 100% success rate while attacks that blindly target object detection only have up to 25%. The source code and data is all available at [1].

Our discovery and results strongly suggest that MOT should be systematically considered and incorporated into future adversarial machine learning research targeting the visual perception in autonomous driving. Our work initiates the first research effort along this direction, and we hope that it can inspire more future research into this largely overlooked research perspective.

## References

- [1] Anonymized. Anonymized github repository for the source code of our attack and evaluation data. <https://github.com/anonymousjack/hijacking>.
- [2] A. Athalye, L. Engstrom, A. Ilyas, and K. Kwok. Synthesizing robust adversarial examples. *arXiv preprint arXiv:1707.07397*, 2017.
- [3] Baidu. Baidu Apollo. <https://github.com/ApolloAuto/apollo>.
- [4] P. Bergmann, T. Meinhardt, and L. Leal-Taixé. Tracking without bells and whistles. *CoRR*, abs/1903.05625, 2019.
- [5] A. Bewley, Z. Ge, L. Ott, F. Ramos, and B. Upcroft. Simple online and realtime tracking. In *2016 IEEE International Conference on Image Processing (ICIP)*, pages 3464–3468. IEEE, 2016.
- [6] S.-T. Chen, C. Cornelius, J. Martin, and D. H. P. Chau. Shapeshifter: Robust physical adversarial attack on faster r-cnn object detector. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 52–68. Springer, 2018.
- [7] W. Choi. Near-online multi-target tracking with aggregated local flow descriptor. In *Proceedings of the IEEE international conference on computer vision*, pages 3029–3037, 2015.
- [8] A. Ess, K. Schindler, B. Leibe, and L. Van Gool. Object detection and tracking for autonomous navigation in dynamic environments. *The International Journal of Robotics Research*, 29(14):1707–1725, 2010.
- [9] 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 models. *arXiv preprint arXiv:1707.08945*, 2017.
- [10] 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 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1625–1634, 2018.
- [11] W. Feng, Z. Hu, W. Wu, J. Yan, and W. Ouyang. Multi-object tracking with multiple cues and switcher-aware classification. *arXiv preprint arXiv:1901.06129*, 2019.
- [12] S. Kato, E. Takeuchi, Y. Ishiguro, Y. Ninomiya, K. Takeda, and T. Hamada. An open approach to autonomous vehicles. *IEEE Micro*, 35(6):60–68, 2015.
- [13] S. Kato, S. Tokunaga, Y. Maruyama, S. Maeda, M. Hirabayashi, Y. Kitsukawa, A. Monrroy, T. Ando, Y. Fujii, and T. Azumi. Autoware on board: enabling autonomous vehicles with embedded systems. In *ICCPs’18*, pages 287–296. IEEE Press, 2018.
- [14] C. Long, A. Haizhou, Z. Zijie, and S. Chong. Real-time multiple people tracking with deeply learned candidate selection and person re-identification. In *ICME*, 2018.
- [15] J. Lu, H. Sibai, and E. Fabry. Adversarial examples that fool detectors. *arXiv preprint arXiv:1712.02494*, 2017.
- [16] J. Lu, H. Sibai, E. Fabry, and D. Forsyth. Standard detectors aren’t (currently) fooled by physical adversarial stop signs. *arXiv preprint arXiv:1710.03337*, 2017.
- [17] F. Luetteke, X. Zhang, and J. Franke. Implementation of the hungarian method for object tracking on a camera monitored transportation system. In *ROBOTIK 2012; 7th German Conference on Robotics*, pages 1–6. VDE, 2012.
- [18] W. Luo, J. Xing, A. Milan, X. Zhang, W. Liu, X. Zhao, and T.-K. Kim. Multiple object tracking: A literature review. *arXiv preprint arXiv:1409.7618*, 2014.
- [19] MathWorks. Automated driving toolbox. <https://www.mathworks.com/products/automated-driving.html>.
- [20] S. Murray. Real-time multiple object tracking—a study on the importance of speed. *arXiv preprint arXiv:1709.03572*, 2017.
- [21] A. Neubeck and L. Van Gool. Efficient non-maximum suppression. In *18th International Conference on Pattern Recognition (ICPR’06)*, volume 3, pages 850–855. IEEE, 2006.
- [22] OpenCV. Kalman Filter Class Reference. [https://docs.opencv.org/3.4.1/dd/d6a/classcv\\_1\\_1KalmanFilter.html](https://docs.opencv.org/3.4.1/dd/d6a/classcv_1_1KalmanFilter.html).

- [23] J. Redmon and A. Farhadi. Yolo9000: better, faster, stronger. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7263–7271, 2017.
- [24] J. Redmon and A. Farhadi. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018.
- [25] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015.
- [26] 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 *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, pages 1528–1540. ACM, 2016.
- [27] S. Sharma, J. A. Ansari, J. K. Murthy, and K. M. Krishna. Beyond pixels: Leveraging geometry and shape cues for online multi-object tracking. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3508–3515. IEEE, 2018.
- [28] Udacity. Self-driving car engineer nanodegree program. <https://www.udacity.com/course/self-driving-car-engineer-nanodegree--nd013>.
- [29] Y. Xiang, A. Alahi, and S. Savarese. Learning to track: Online multi-object tracking by decision making. In *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 4705–4713, Dec 2015.
- [30] C. Xie, J. Wang, Z. Zhang, Y. Zhou, L. Xie, and A. Yuille. Adversarial examples for semantic segmentation and object detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1369–1378, 2017.
- [31] A. Yilmaz, O. Javed, and M. Shah. Object tracking: A survey. *Acm computing surveys (CSUR)*, 38(4):13, 2006.
- [32] J. H. Yoon, C. Lee, M. Yang, and K. Yoon. Online multi-object tracking via structural constraint event aggregation. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [33] F. Yu, W. Xian, Y. Chen, F. Liu, M. Liao, V. Madhavan, and T. Darrell. Bdd100k: A diverse driving video database with scalable annotation tooling. *arXiv preprint arXiv:1805.04687*, 2018.
- [34] D. Zhao, H. Fu, L. Xiao, T. Wu, and B. Dai. Multi-object tracking with correlation filter for autonomous vehicle. *Sensors*, 18(7):2004, 2018.
- [35] Y. Zhao, H. Zhu, Q. Shen, R. Liang, K. Chen, and S. Zhang. Practical adversarial attack against object detector. *arXiv preprint arXiv:1812.10217*, 2018.
- [36] Z. Zhong, W. Xu, Y. Jia, and T. Wei. Perception Deception: Physical Adversarial Attack Challenges and Tactics for DNN-Based Object Detection. In *Black Hat Europe*, 2018.
- [37] J. Zhu, H. Yang, N. Liu, M. Kim, W. Zhang, and M.-H. Yang. Online multi-object tracking with dual matching attention networks. In *Computer Vision – ECCV 2018*, pages 379–396, Cham, 2018. Springer International Publishing.