

General Black-box Adversarial Sensor Attack and Countermeasures

Robust LiDAR-based Perception in Autonomous Driving

(Charles) Chengzeng You
Connected and autonomous vehicles
Department of computing
Email: chengzeng.you19@imperial.ac.uk

Content



6. Attack Conclusion



7. Physics-Informed Anomaly Detection



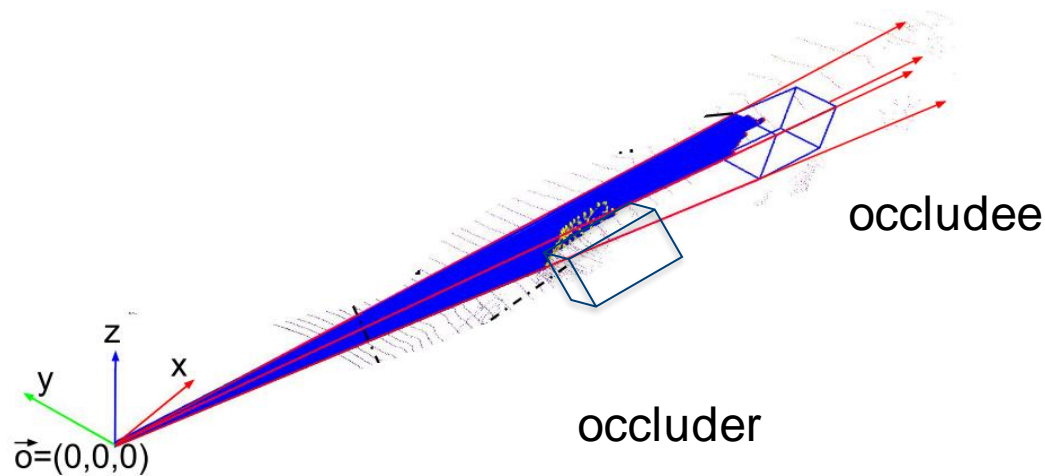
8. Physics-Embedded Perception Architecture



9. Limitations

6.1 Occlusion patterns

- state-of-the-art 3D object detection model designs generally ignore the occlusion patterns in LiDAR point clouds.



Content



6. Attack Conclusion



7. Physics-Informed Anomaly Detection



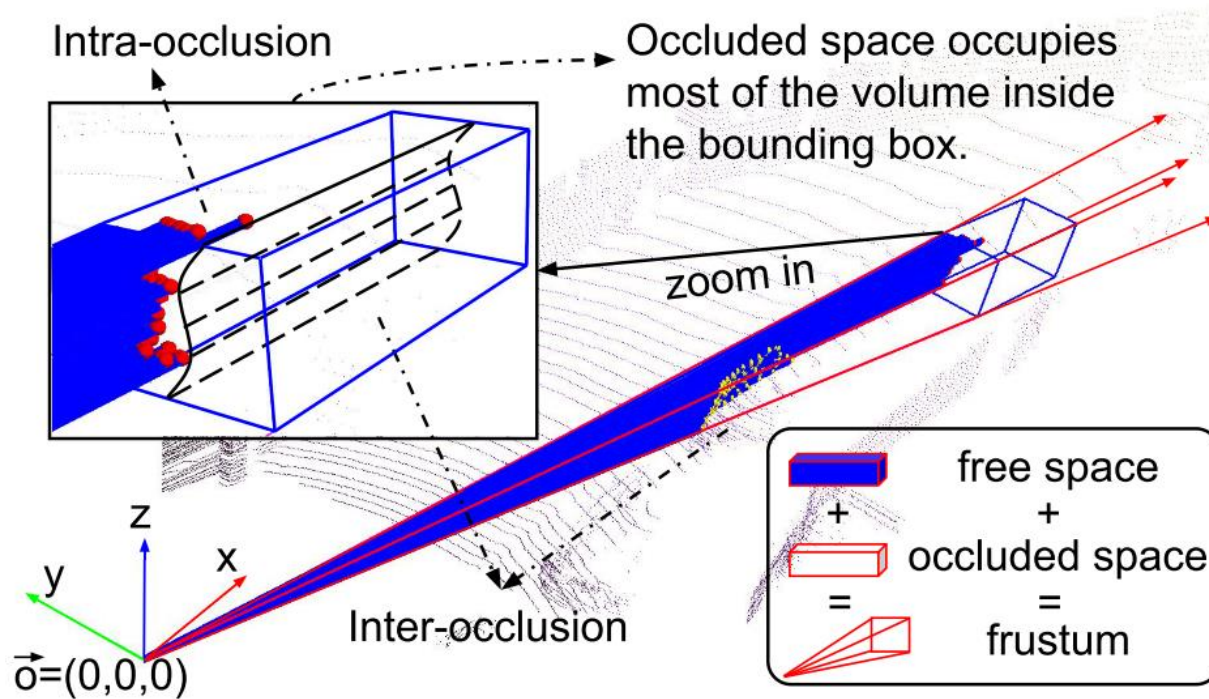
8. Physics-Embedded Perception Architecture



9. Limitations

7.1 CARLO (occlusion-Aware hierarchy anomaly detection)

7.1.1 Free Space Detection (FSD)



7.1 CARLO (occlusion-Aware hierarchy anomaly detection)

7.1.1 Free Space Detection (FSD)

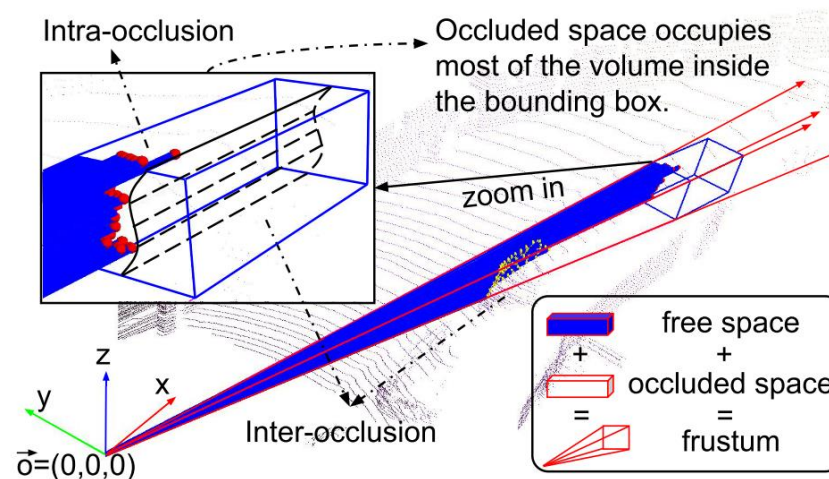
Observation:

f = the volume of FS / the volume of a detected bounding box

$\exists b \in (0,1)$, implying $f \in (0,b]$

For fake vehicles:

$\exists a \in (0,1)$ such that $f' \in [a,1)$.



7.1 CARLO (Occlusion-Aware hierarchy anomaly detection)

7.1.1 Free Space Detection (FSD)

$$f_B = \frac{\sum_{c \in B} \mathbb{1} \cdot FS(c)}{|B|}$$

$FS(c)$ indicates whether the cell c is free or not

$|B|$ denotes the total number of cells in the bounding box B

7.1 CARLO (Occlusion-Aware hierarchy anomaly detection)

7.1.1 Free Space Detection (FSD)

100 ms/bounding box

Setup

cell size: 0.25

Dataset: KITTI training set

600 generated attack traces

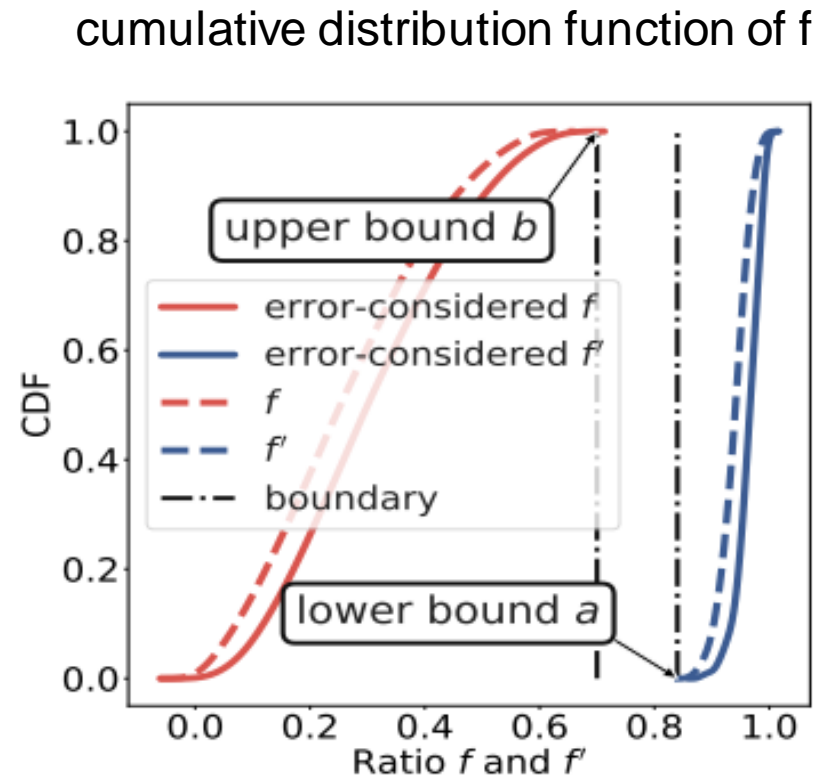
cumulative distribution function:

$$F_X(x) = P(X \leq x)$$

$F_X(x)$ = function of X

X = real value variable

P = probability that X will have a value less than or equal to x



7.1 CARLO (Occlusion-Aware hierarchy anomaly detection)

7.1.2 Laser Penetration Detection (LPD)

three spaces of one frustum:

- 1) the space between the LiDAR sensor and the bounding box B ↑
- 2) the space inside the bounding box B
- 3) **the space behind the bounding box B** ↓

7.1 CARLO (Occlusion-Aware hierarchy anomaly detection)

7.1.2 Laser Penetration Detection (LPD)

5 ms/bounding box

$$g_B = \frac{\sum_{\vec{p} \in B_{\downarrow}} \mathbb{1}(\vec{p})}{\sum_{\vec{p} \in B \cup B_{\downarrow} \cup B_{\uparrow}} \mathbb{1}(\vec{p})}$$

For benign vehicles:

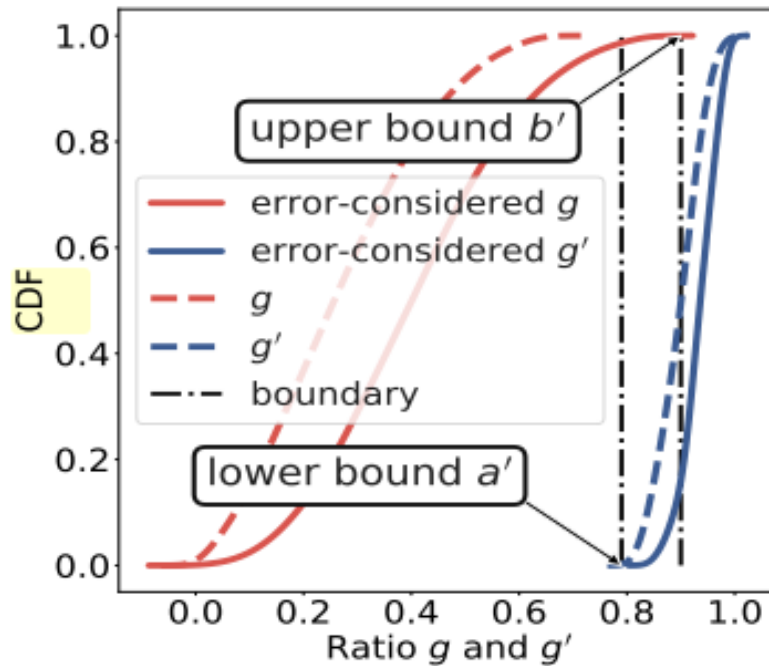
g should be upper bounded by $\exists b' \in (0,1)$.

For spoofed vehicles:

g' should be lower bounded by $\exists a' \in (0,1)$.

7.1 CARLO (Occlusion-Aware hierarchy anomaly detection)

7.1.2 Laser Penetration Detection (LPD)



$b' > a'$
Might cause erroneous
detection of potential
anomalies

7.1 CARLO (Occlusion-Aware hierarchy anomaly detection)

7.1.3 Hierarchy Design

around 8.5 ms / each vehicle

Algorithm 1: CARLO

input: Detected bounding boxes $\mathbf{B} = \{\mathbf{B}\}$;
LiDAR laser ray directions $\mathbf{L} = \{\mathbf{L}\}$;
3D point cloud $\mathbf{X} = \{\vec{p}\}$;
1 Threshold of FSD $\frac{a+b}{2}$;
Thresholds of LPD $b' + \epsilon, a' - \epsilon$;
output: Valid bounding boxes $\mathbf{B}_{\text{valid}} = \{\mathbf{B}\}$;
Adversarial bounding boxes $\mathbf{B}_{\text{adv}} = \{\mathbf{B}\}$;
2 **Initialization** : $\mathbf{B}_{\text{valid}} \leftarrow \emptyset, \mathbf{B}_{\text{adv}} \leftarrow \emptyset, g \leftarrow 0, f \leftarrow 0$;
/* Initiate parameters. */

7.1 CARLO (Occlusion-Aware hierarchy anomaly detection)

7.1.3 Hierarchy Design

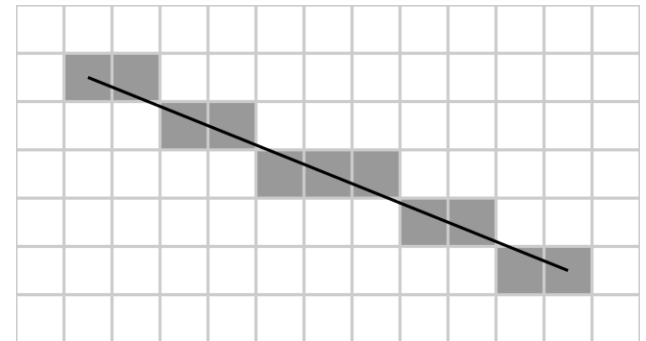
```
3 for  $B \in \mathbf{B}$  do
    /* Initiate parameters, where  $FS(\cdot)$  is the
       free space and  $F_B$  is the frustum of  $B$ . */
4    $F_B \leftarrow \emptyset, FS(\cdot) \leftarrow \emptyset;$ 
5   for  $L \in \mathbf{L}$  do
       /* Predict whether  $L$  will intersect with
           $B$ . */
6       if  $L \cap B$  then
7            $\vec{p}_L \leftarrow L;$ 
8            $F_B.append([L, \vec{p}_L]);$ 
       /* Extract the frustum  $F_B$  of  $B$ . */
9   end
10   $g \leftarrow \text{Equation 6};$ 
    /* Calculate  $g$  by  $F_B$  for  $B$  (LPD). */
```

7.1 CARLO (occlusion-Aware hierarchy anomaly detection)

7.1.3 Hierarchy Design

```
11  if  $g < a' - \varepsilon$  then
12    |  $B_{\text{valid}}.\text{append}(B)$ ;
13    | /* Certainly valid vehicles. */
14  else if  $g > b' + \varepsilon$  then
15    |  $B_{\text{adv}}.\text{append}(B)$ ;
16    | /* Certainly spoofed vehicles. */
17  else
18    | /* Calculate  $f$  by  $F_B$  for  $B$  (FSD). */
19    | for  $[L, \vec{p}_L] \in F_B$  do
20    |   |  $FS(L) \leftarrow \text{Bresenham}([L, \vec{p}_L])[16]$ ;
21    |   |  $FS(B) \leftarrow FS(B) \cup FS(L)$ ;
22    | end
23    |  $f \leftarrow \text{Equation 5}$ ;
24    | if  $f < \frac{a+b}{2}$  then
25    |   |  $B_{\text{valid}}.\text{append}(B)$ ;
26    |   |  $B_{\text{adv}}.\text{append}(B)$ ;
27  end
28  Return :  $B_{\text{valid}}, B_{\text{adv}}$ 
```

Illustration of the result of
Bresenham's line



7.2 CARLO Evaluation

7.2.1 Experimental setup

Dateset K: existing real-world traces

Dateset R: simulated attack traces

Adversarial model: Adv-LiDAR [17] on Apollo 5.0

Target location: 5-8 meters in front of the victim

CARLO-guarded models: CARLO(M(.))

Defense goal: successfully detect the spoofed fake vehicles from the output bounding boxes without hurting the original performance of the target models

7.2 CARLO Evaluation

7.2.2 Evaluation metrics

$$ASR = \frac{\text{\# of successful attacks}}{\text{\# of total point cloud samples}}$$

$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

$$AP = \int_0^1 p(r)dr$$

TP = True positive

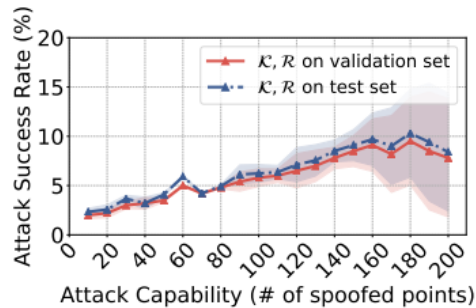
TN = True negative

FP = False positive

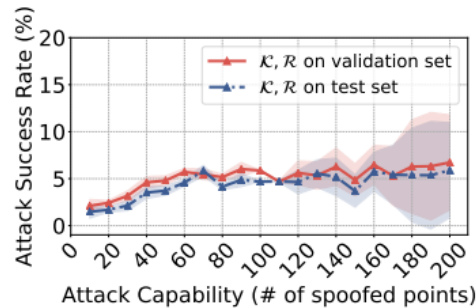
FN = False negative

7.2 CARLO Evaluation

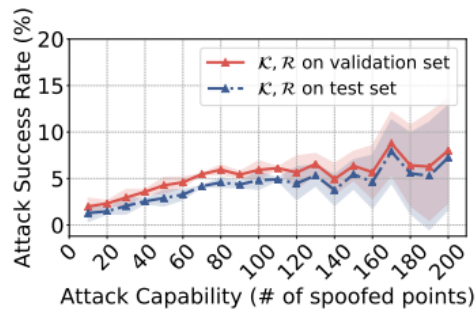
7.2.2 Evaluation metrics



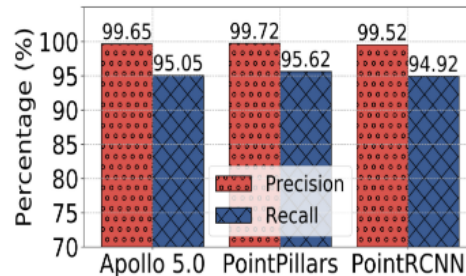
(a) CARLO-guarded Apollo 5.0.



(b) CARLO-guarded PointPillars.



(c) CARLO-guarded PointRCNN.



(d) Precision and recall of CARLO.

reduce the mean ASR to around 5.5%.

the remaining 5.5% comes from the detection errors .

Figure 12: Attack success rates (ASRs) of proposed black-box spoofing attacks on three CARLO-guarded models.

7.2 CARLO Evaluation

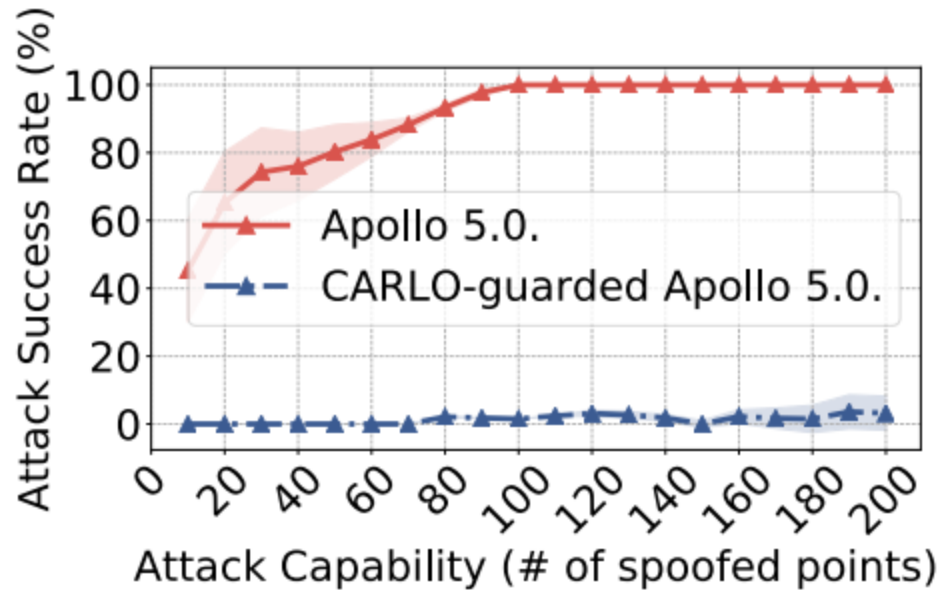
7.2.2 Evaluation metrics

Table 1: PointPillars' and PointRCNN's APs (%) of 3D car detection on the KITTI validation set.

Model	PointPillars			PointRCNN		
	Easy	Mod.	Hard	Easy	Mod.	Hard
Original	86.56	76.87	72.09	88.80	78.58	77.64
Attack	74.06	56.69	53.98	84.51	71.17	68.06
CARLO	86.57	78.60	73.55	88.91	78.61	77.63

7.2 CARLO Evaluation

7.2.3 Defense against White-box Attacks



7.2 CARLO Evaluation

7.2.4 Defense against Adaptive Attacks

$$f_B = \frac{\sum_{c \in B} \mathbb{1} \cdot FS(c)}{|B|}$$

absolute free space

7.2 CARLO Evaluation

7.2.4 Defense against Adaptive Attacks

Attack goal:

- 1.to spoof a vehicle at target locations
- 2.minimize the size of the bounding box

$$\min_{\theta, \tau} \quad \mathcal{L}(x \oplus V \cdot H(\theta, \tau)^T) + \lambda \cdot \mathcal{V}_B(V \cdot H(\theta, \tau)^T)$$

$\mathcal{L}(\cdot)$: the adversarial loss, given a detector

x : the corresponding input feature matrix

V : adversarial spoofed input feature matrix

$\oplus(\cdot)$: merging function

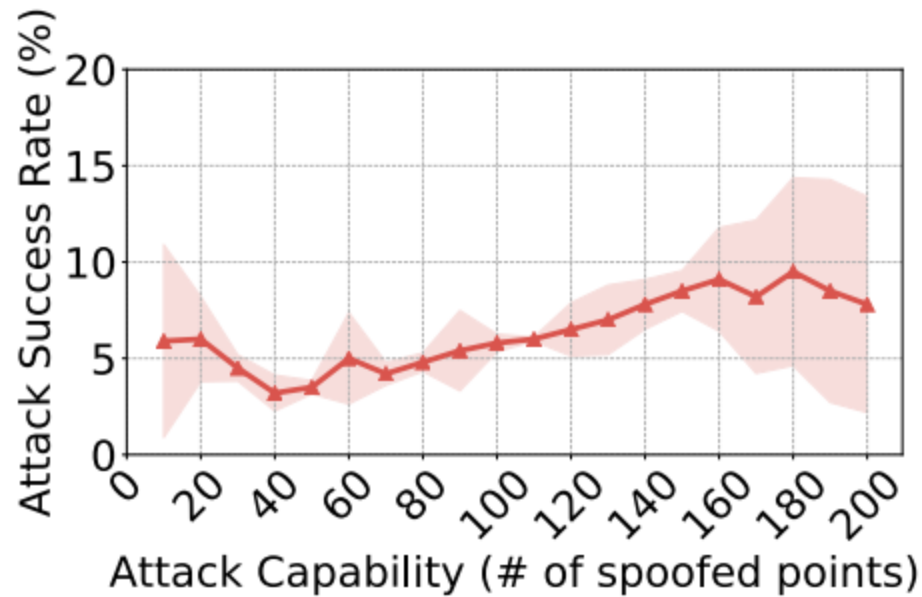
$H(\theta, \tau)$: transformation matrix

$\mathcal{V}_B(\cdot)$: the volume of the target bounding box B ,

λ : a hyper-parameter

7.2 CARLO Evaluation

7.2.4 Defense against Adaptive Attacks



Content



6. Attack Conclusion



7. Physics-Informed Anomaly Detection



8. Physics-Embedded Perception Architecture



9. Limitations

8.1 Sequential View Fusions(SVF)

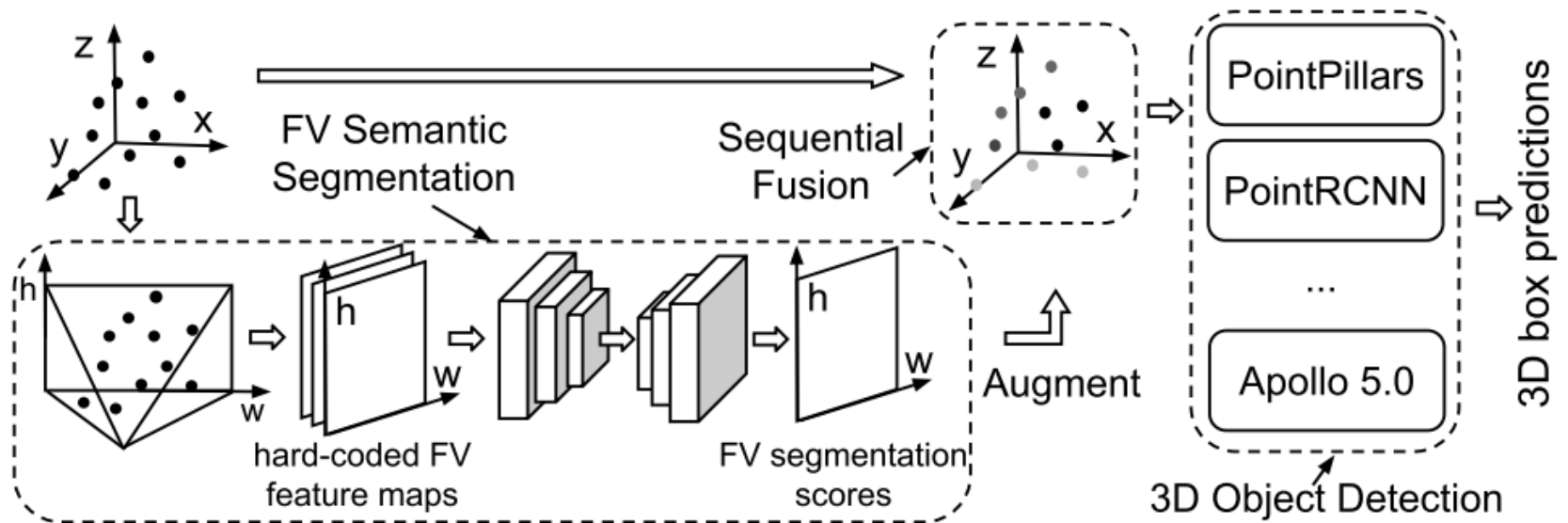


Figure 17: Sequential view fusion (SVF) architecture.

8.2 SVF Evaluation

8.2.1 Experimental setup

Target model: trained SVF-PointPillars and SVF- PointRCNN

Dataset : KITTI training set.

Evaluate SVF against Adv-LiDAR [17] on Apollo 5.0 and the adaptive attacks.

8.2 SVF Evaluation

8.2.2 Evaluation metrics

leverage ASR to test their robustness against LiDAR spoofing attacks

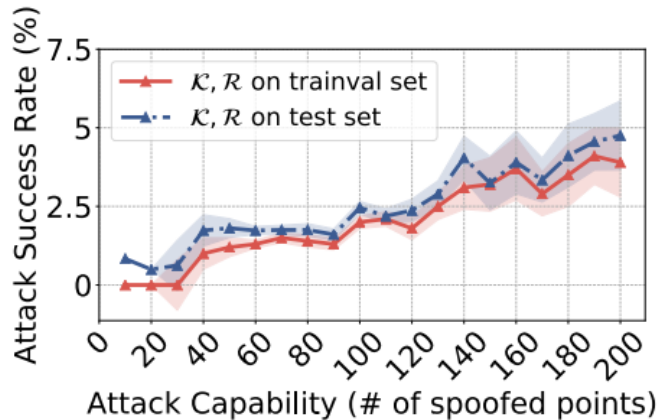
Table 3: SVF-PointPillars' and SVF-PointRCNN's APs (%) of 3D car detection on the KITTI validation set.

Model	Car Detection		
	Easy	Moderate	Hard
SVF-PointPillars	85.93	74.12	70.19
SVF-PointRCNN	88.12	76.56	74.81

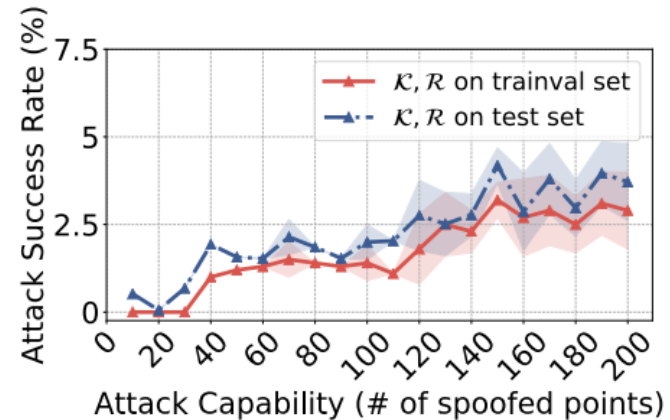
8.2 SVF Evaluation

8.2.2 Evaluation metrics

leverage ASR to test their robustness against LiDAR spoofing attacks



(a) ASR of SVF-PointPillars.



(b) ASR of SVF-PointRCNN.

Figure 18: Attack success rates (ASRs) of proposed black-box spoofing attack on SVF models.

8.2 SVF Evaluation

8.2.3 Defense against White-box Attacks

Since SVF requires re-training for the model, cannot directly evaluate Adv-LiDAR on SVF-Apollo

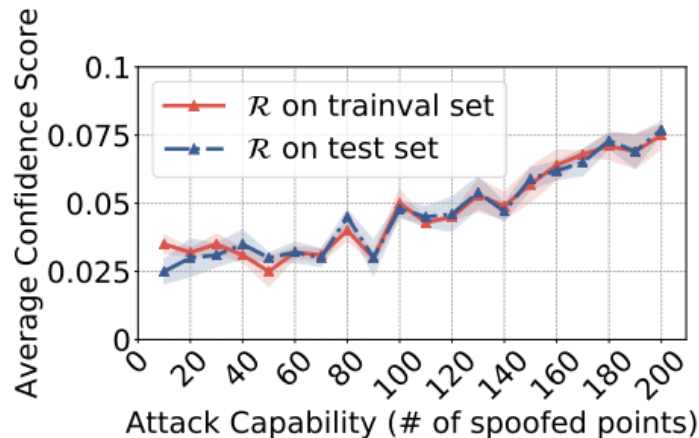


Figure 19: Average confidence score of Adv-LiDAR on the segmentation network.

8.2 SVF Evaluation

8.2.4 Defense against Adaptive Attacks

Assume the adversaries are aware of the SVF architecture.

Attack goal: both fool the semantic segmentation and 3D object detection modules.

$$\min_{\theta, \tau} \quad -\mathcal{L}_{seg}(x \odot V \cdot H(\theta, \tau)^T)$$

8.2 SVF Evaluation

8.2.4 Defense against Adaptive Attacks

Assume the adversaries are aware of the SVF architecture.

Attack goal: both fool the semantic segmentation and 3D object detection modules.

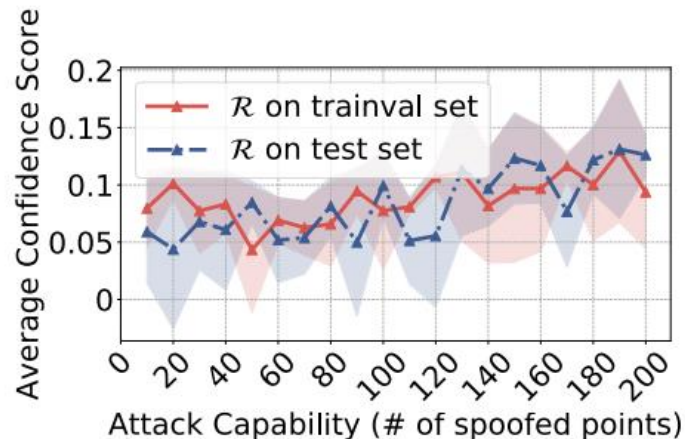


Figure 20: Average confidence score of the adaptive attack on the segmentation network.

Content



6. Attack Conclusion



7. Physics-Informed Anomaly Detection



8. Physics-Embedded Perception Architecture



9. Limitations

9 Limitations

1. **same pattern** at the **target location**, whether this defense strategy is transferable to other LiDAR spoofing attacks (e.g. other patterns) still remains unexplored.
2. the identified vulnerability does not provide completeness, there may exist other potential vulnerabilities hidden in the AD systems to be discovered and exploited
3. Lack of attacks on other types of objects, such as pedestrians
4. Lack of real-world validations

Thank you
