# That Person Moves Like A Car: Misclassification Attack Detection for Autonomous Systems Using Spatiotemporal Consistency
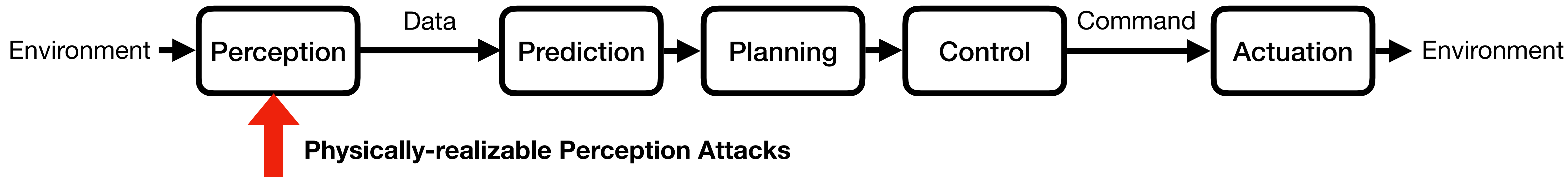
Yanmao Man#, Raymond Muller§, Ming Li#, Z. Berkey Celik§, Ryan Gerdes‡

#University of Arizona   §Purdue University   ‡Virginia Tech

# Autonomous Systems

Environment → **Perception** — Data → **Prediction** → **Planning** → **Control** — Command → **Actuation** → Environment

# Perception Security

Environment → **Perception** → *Data* → **Prediction** → **Planning** → **Control** → *Command* → **Actuation** → Environment

**Physically-realizable Perception Attacks**

These perception attacks alter the data at the source, hence bypassing traditional digital defenses
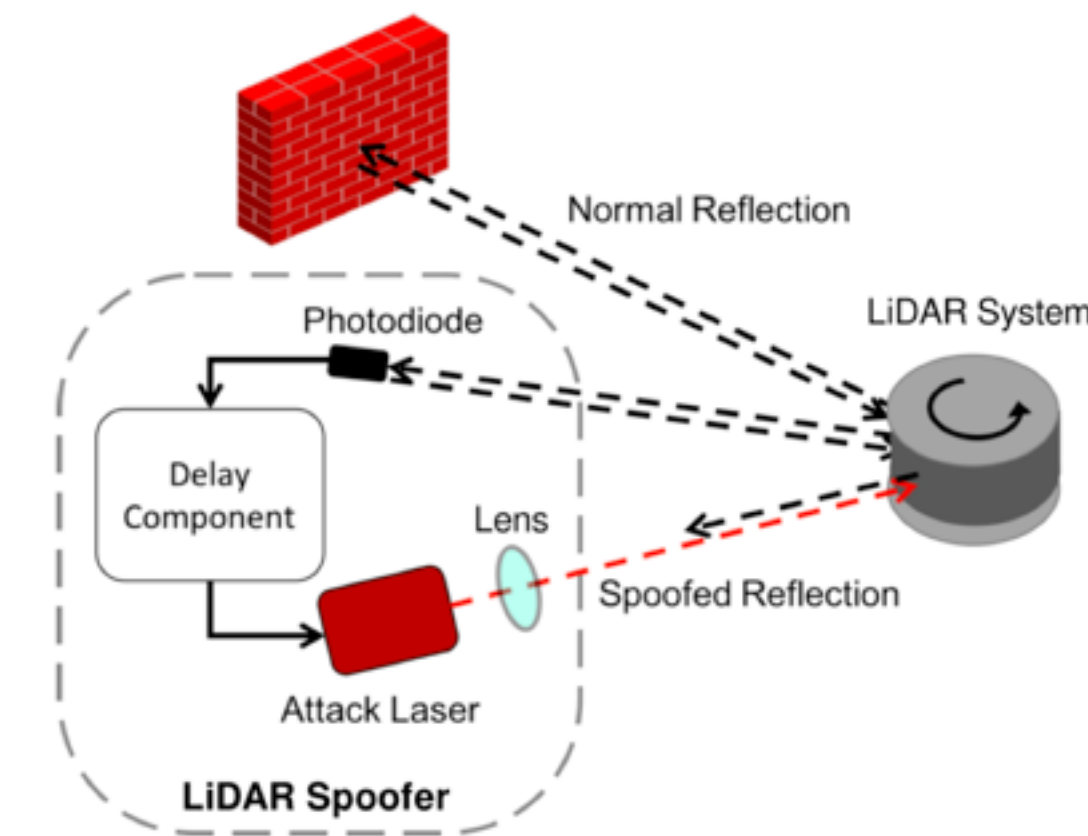
Stop Sign Sticker

Phantom Attack

stop sign
STOP

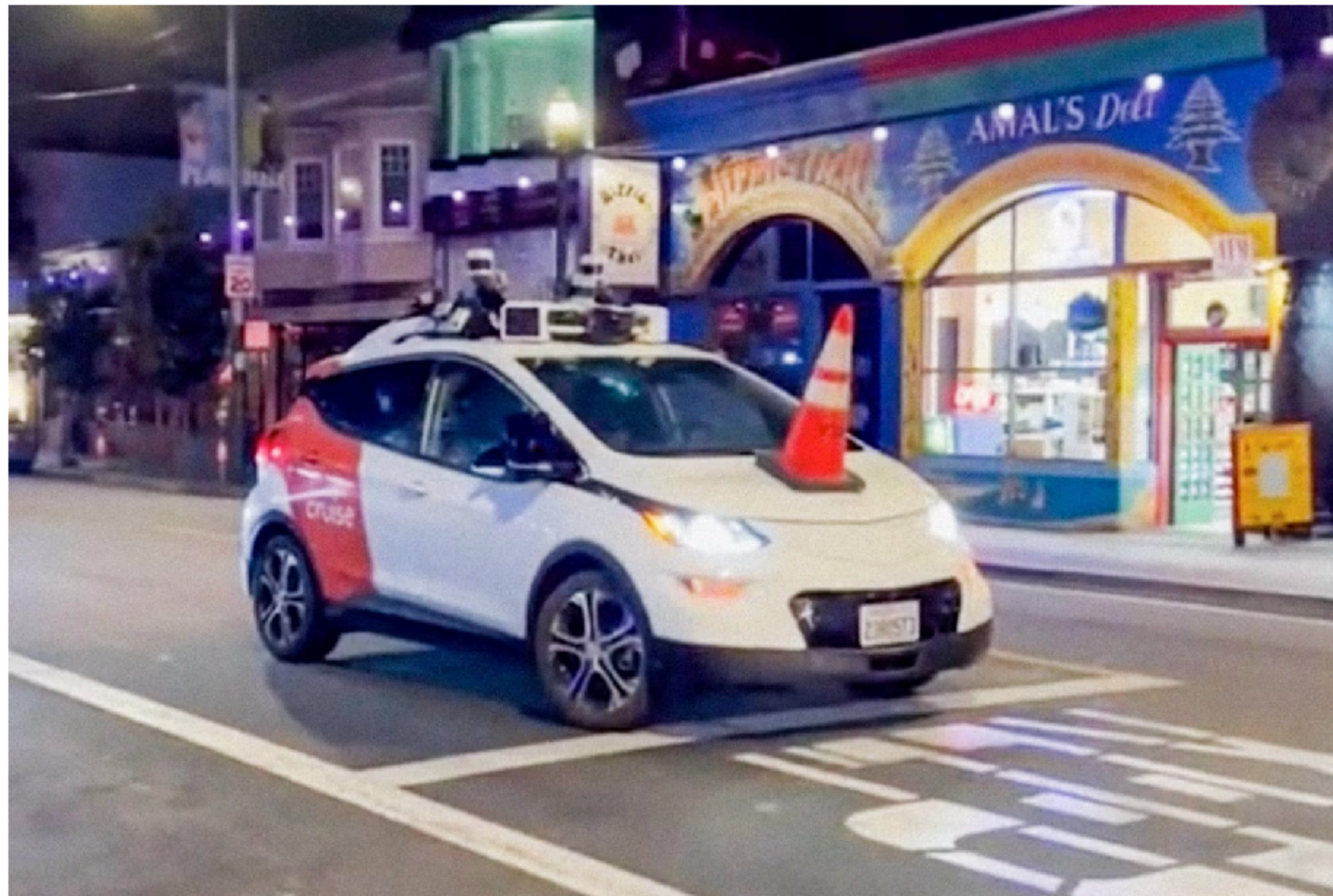GhostImage Camera Attack

LiDAR Spoofing

1. Eykholt, Kevin, et al. "Robust physical-world attacks on deep learning visual classification." *CVPR 2018*.
2. Nassi, Ben, et al. "Phantom of the adas: Securing advanced driver-assistance systems from split-second phantom attacks." CCS 2020
3. Man, Yanmao et al."GhostImage: Remote perception attacks against camera-based image classification systems." RAID 2020
4. Cao, Yulong, et al. "Adversarial sensor attack on lidar-based perception in autonomous driving." CCS 2019.

# The Self-Driving Cars Wearing a Cone of Shame

There's a brilliant activist campaign to stop San Francisco's autonomous vehicles in their tracks.

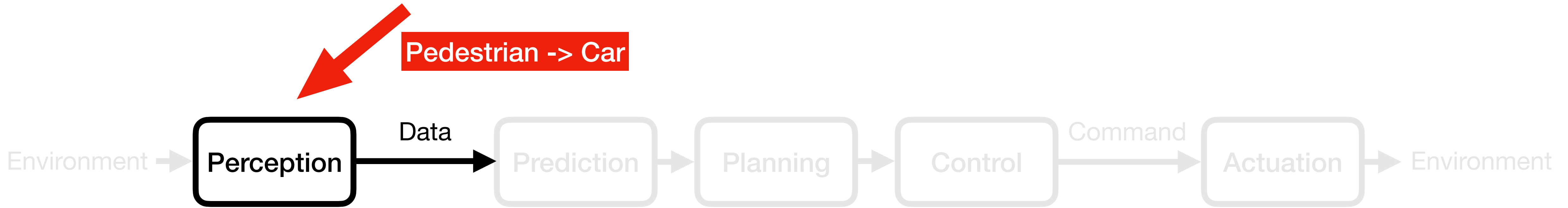BY ALISON GRISWOLD                          JULY 11, 2023 • 10:45 AM



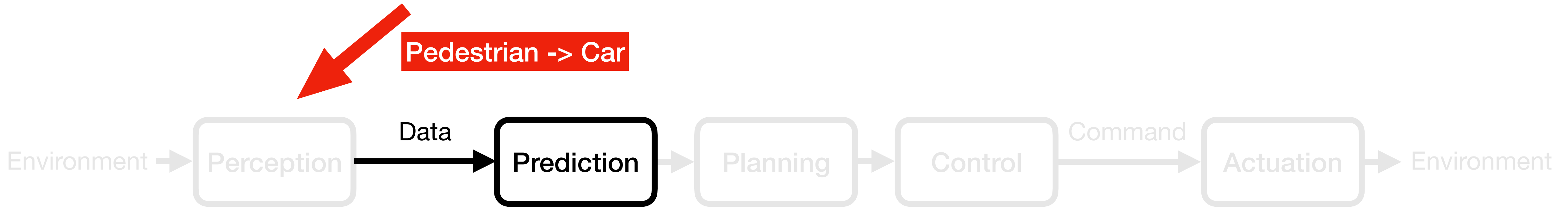It looks like a sad unicorn (which, in a way, it is).   Screengrab from TikTok/Safe Street Rebel
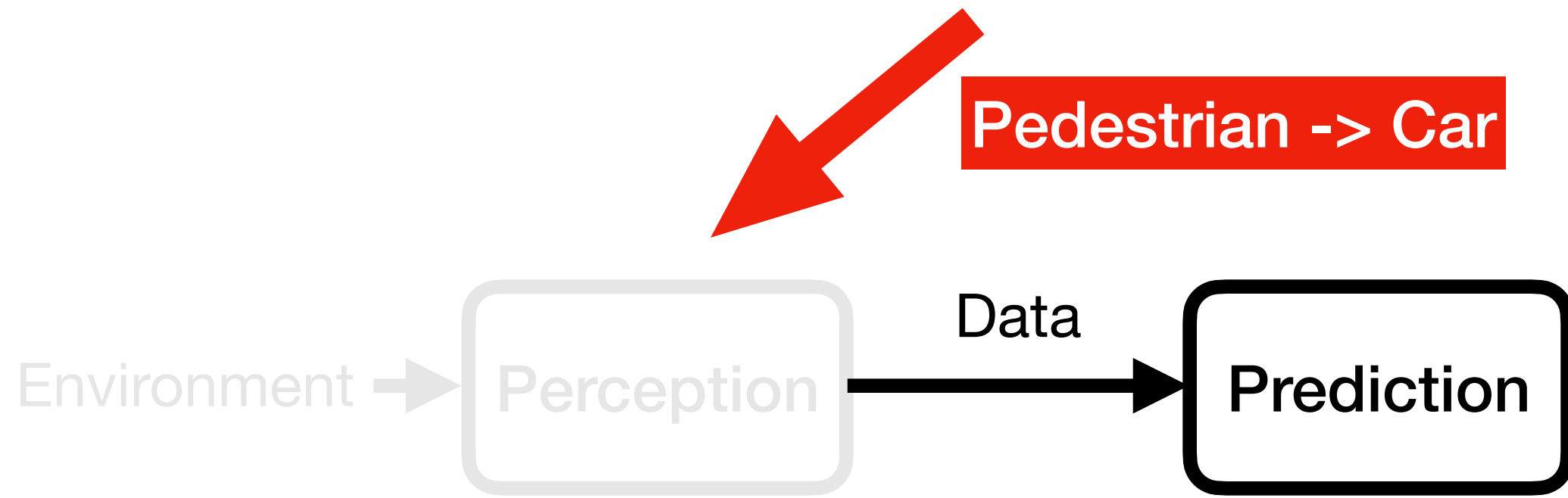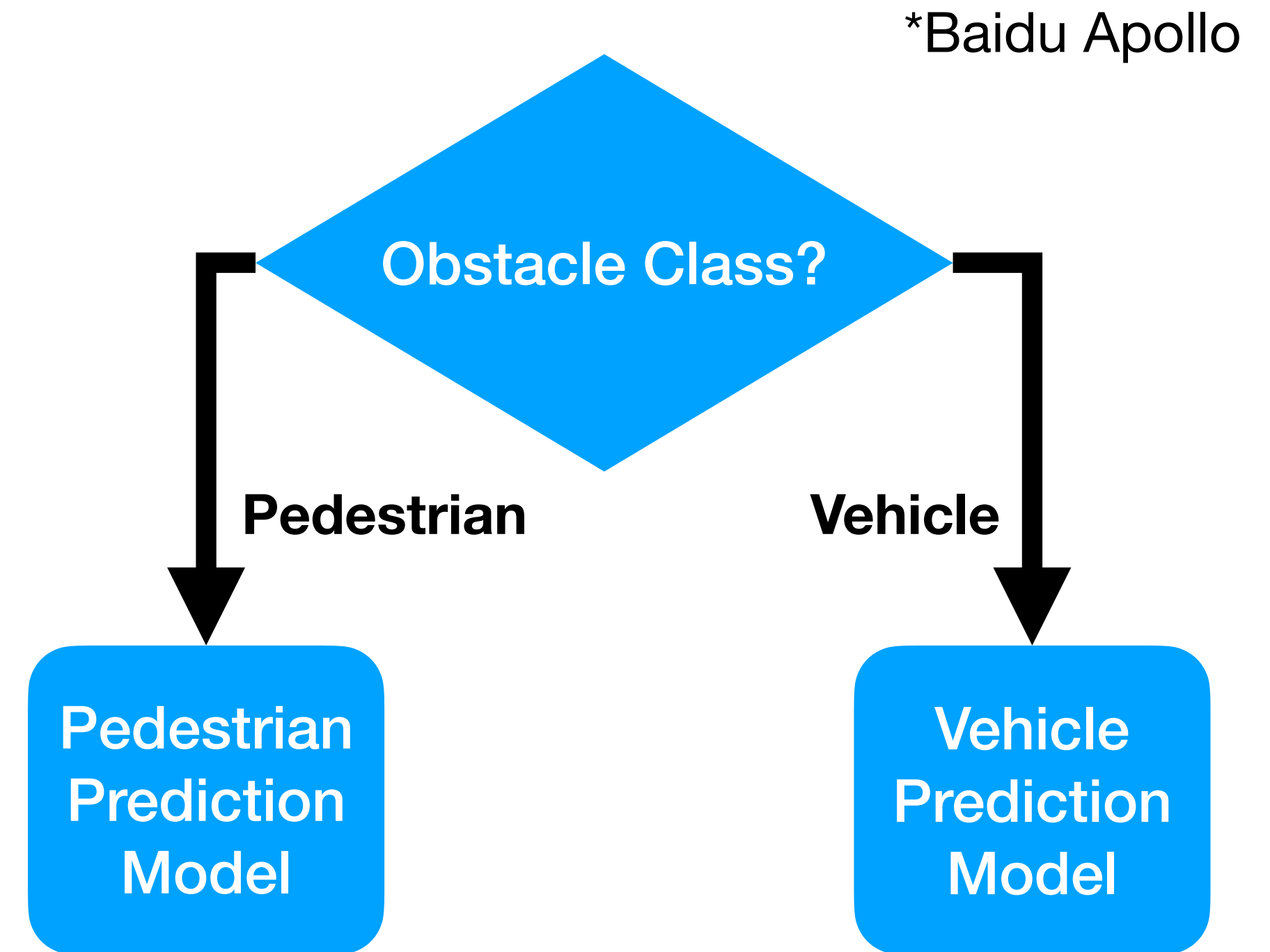
# Perception Security



Environment → Perception → Data → Prediction → Planning → Control → Command → Actuation → Environment

# Misclassification Attacks

Pedestrian -> Car

Environment → | Perception | → Data → | Prediction | → | Planning | → | Control | → Command → | Actuation | → Environment

# Misclassification Attacks

Pedestrian -> Car

Environment → Perception → Data → Prediction → Planning → Control → Command → Actuation → Environment

# Misclassification Attacks

Pedestrian -> Car

*Baidu Apollo

Environment ➡ Perception —Data→ Prediction

Obstacle Class?

Pedestrian — Pedestrian Prediction Model

Vehicle — Vehicle Prediction Model

# Misclassification Attacks



https://toocooltrafficschool.com/following-distance/

Obstacle Class?

Pedestrian

Vehicle

Pedestrian Prediction Model

Vehicle Prediction Model

Emergency Brake

Keep Following

Person
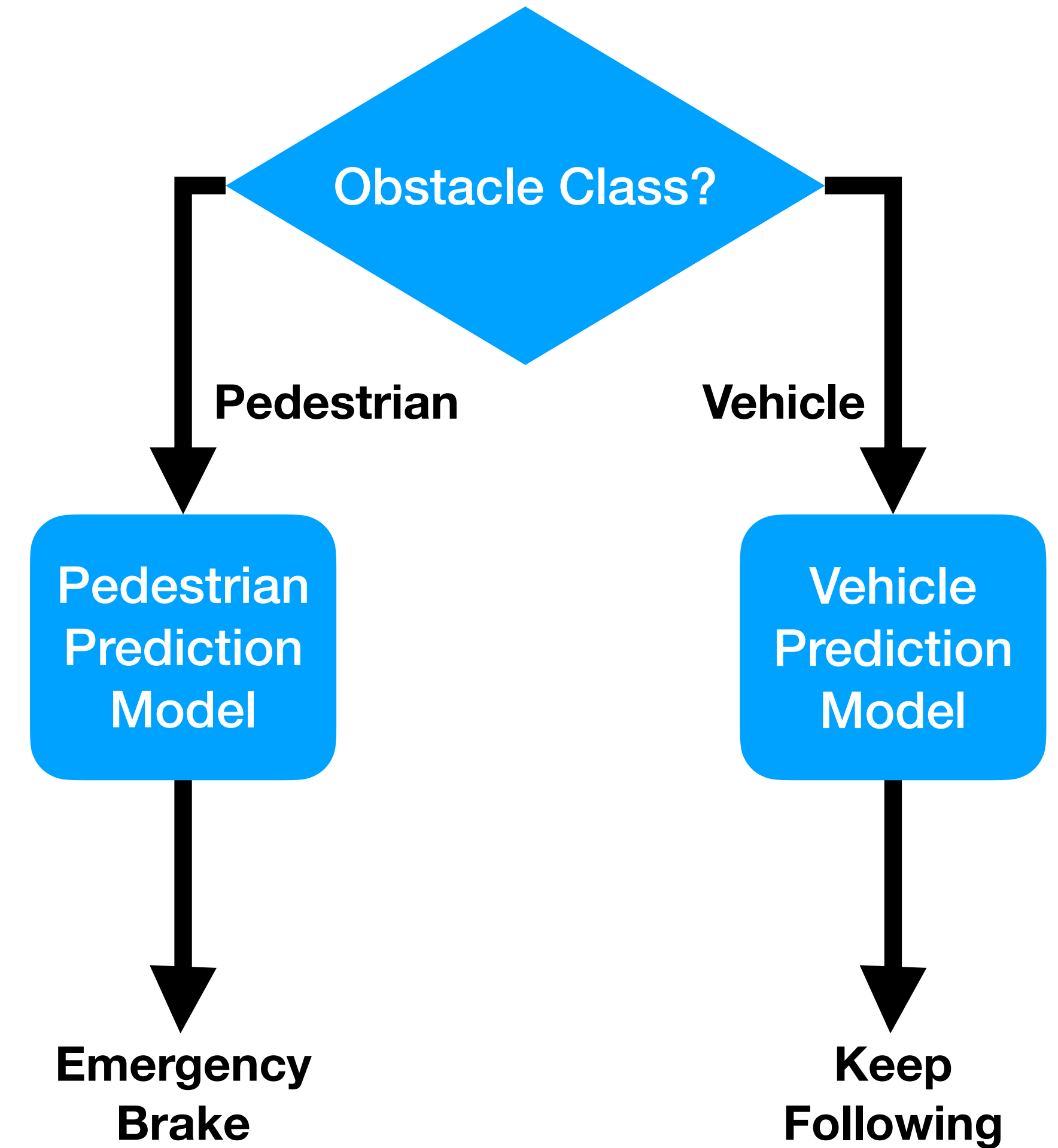
[5] Man, Yanmao, et al. "Evaluating perception attacks on prediction and planning of autonomous vehicles." *USENIX Security Symposium Poster Session*. 2022.

# Misclassification Attacks

**Exorcising "Wraith": Protecting LiDAR-based Object Detector in Automated Driving System from Appearing Attacks**

USENIX Security 2023

**Towards Robust LiDAR-based Perception in Autonomous Driving: General Black-box Adversarial Sensor Attack and Countermeasures**

USENIX Security 2020

**Drift with Devil: Security of Multi-Sensor Fusion based Localization in High-Level Autonomous Driving under GPS Spoofing**

USENIX Security 2020

**Anomaly Detection Against GPS Spoofing Attacks on Connected and Autonomous Vehicles Using Learning From Demonstration**

IEEE T-ITS 2023

**SAVIOR: Securing Autonomous Vehicles with Robust Physical Invariants**

USENIX Security 2020

**ObjectSeeker: Certifiably Robust Object Detection against Patch Hiding Attacks via Patch-agnostic Masking**

IEEE S&P 2023

**AdvIT: Adversarial Frames Identifier Based on Temporal Consistency In Videos**

IEEE ICCV 2019

Chaowei Xiao [1] *    Ruizhi Deng [2]    Bo Li [3]    Taesung Lee [4]
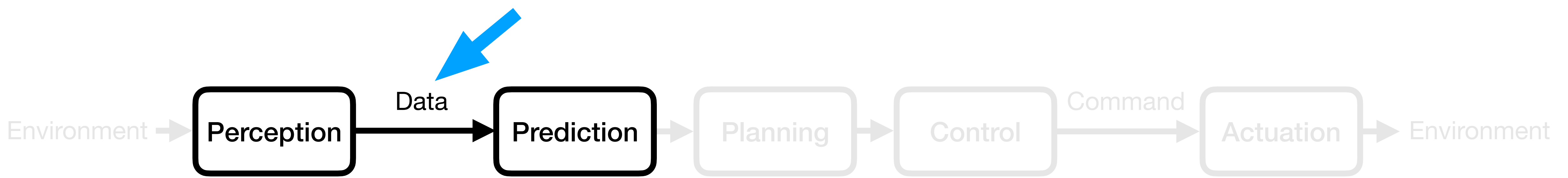Benjamin Edwards[4]    Jinfeng Yi [5]    Dawn Song [6]    Mingyan Liu [1]    Ian Molloy[4]
[1] University of Michigan, Ann Arbor [2] Simon Fraser University [3] UIUC
[4] IBM Research AI [5] JD.com [6] UC Berkeley

PercepGuard aims to detect misclassification attacks

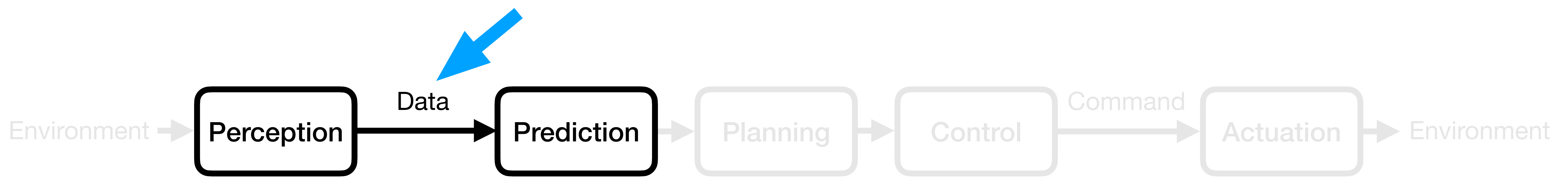Existing defenses against perception attacks are either

- Specific to some sensing modality
  - LiDARs
  - GPS
  - IMU
- Specific to some attack methodology
  - Adversarial Patch
  - Norm-bounded

Environment → **Perception** → Data → **Prediction** → Planning → Control → Command → Actuation → Environment

Agnostic to
- Attack methodologies
- Object detection and tracking algorithms

PercepGuard aims to detect misclassification attacks

Environment → **Perception** → Data → **Prediction** → Planning → Control → Command → Actuation → Environment

PercepGuard aims to detect misclassification attacks by verifying the spatiotemporal consistency of the perception result

# Spatio-temporal Consistency

t=1

Traffic Sign

Pedestrian

Incoming Car

Preceding Car

# Spatio-temporal Consistency

# Spatio-temporal Consistency

t=3

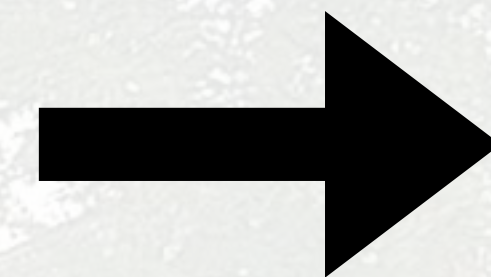# Spatio-temporal Consistency

# Spatio-temporal Consistency

**Research Questions:**
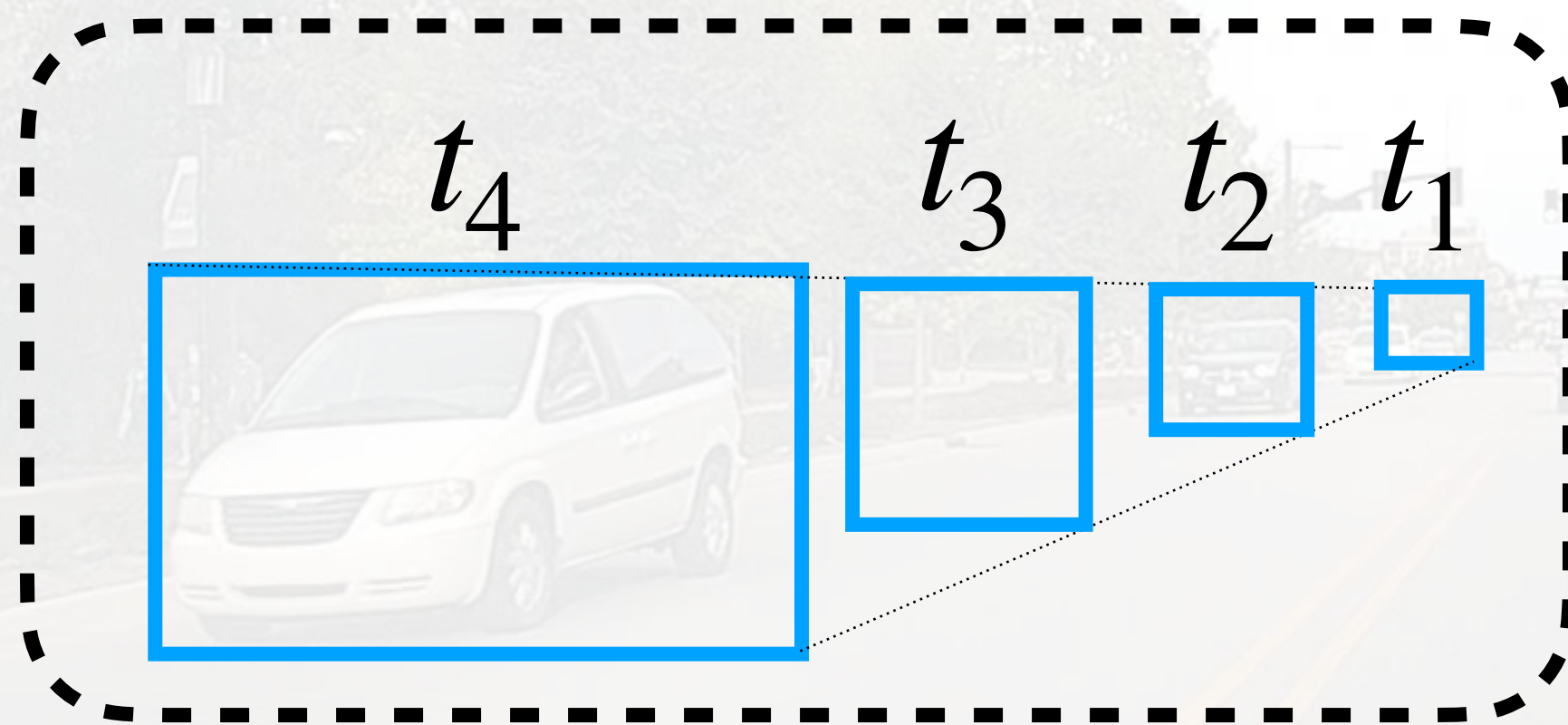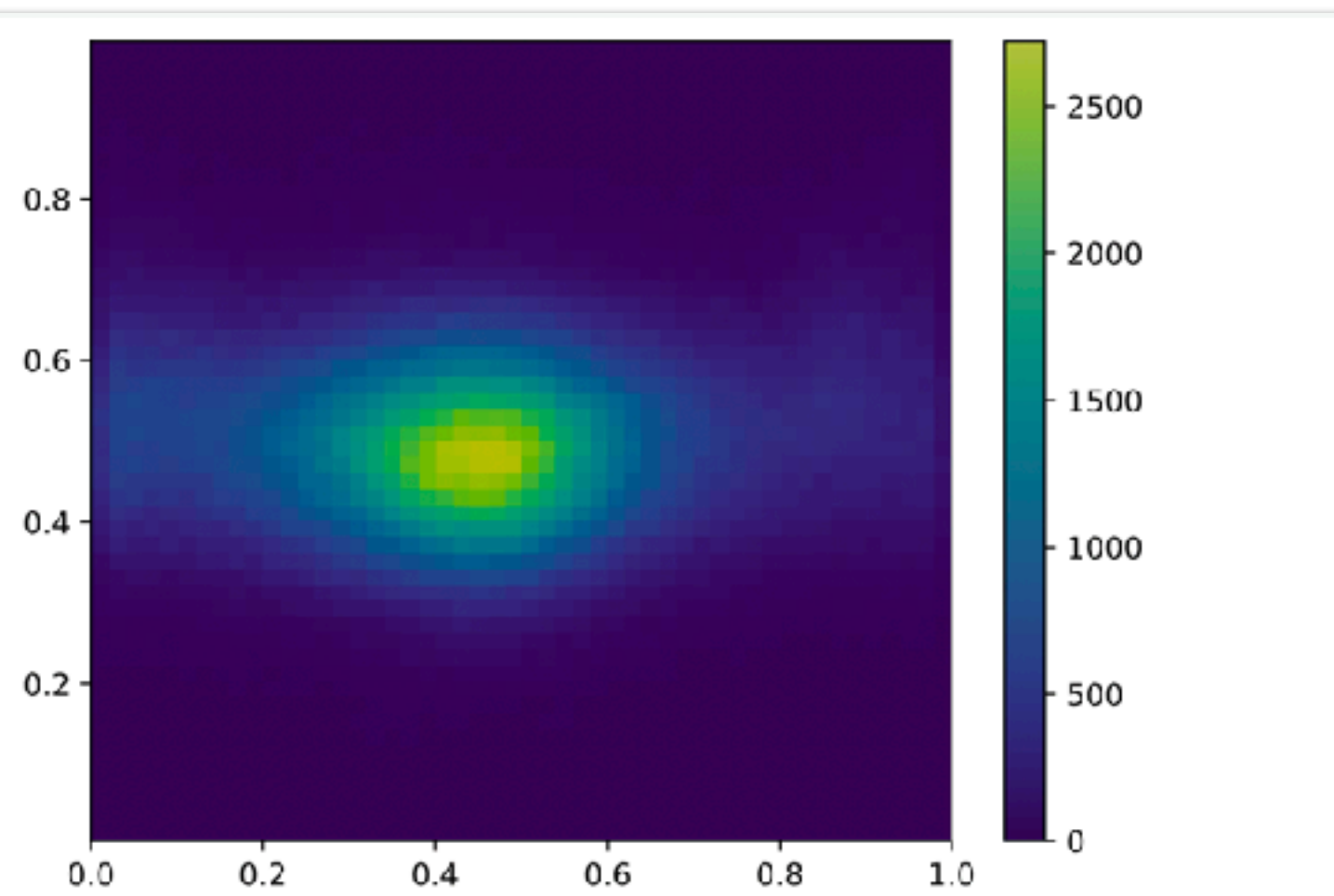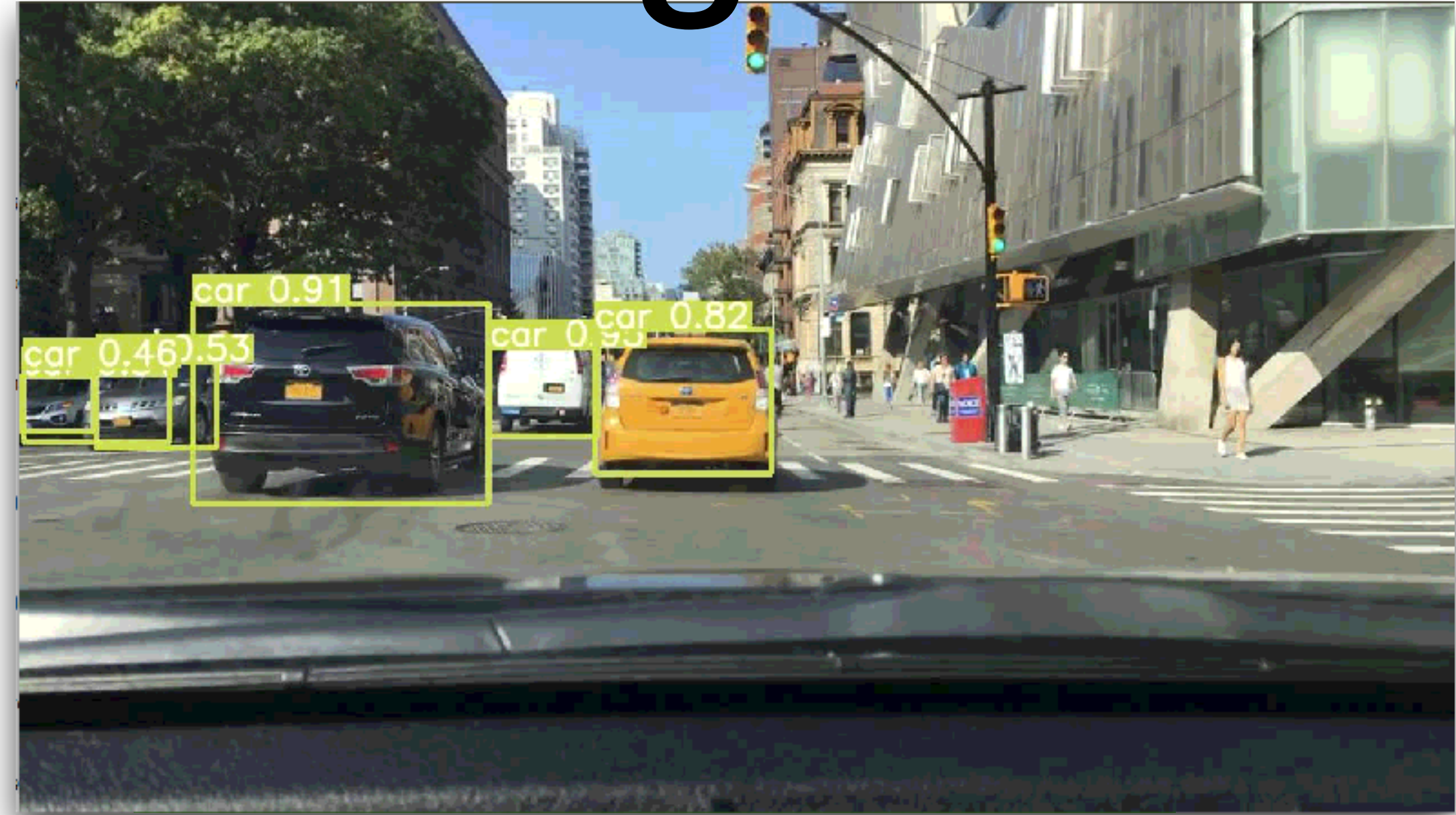
- Do bounding boxes provide statistically-sufficient information?

- Does the detection algorithm produce low false positive and negative rates?

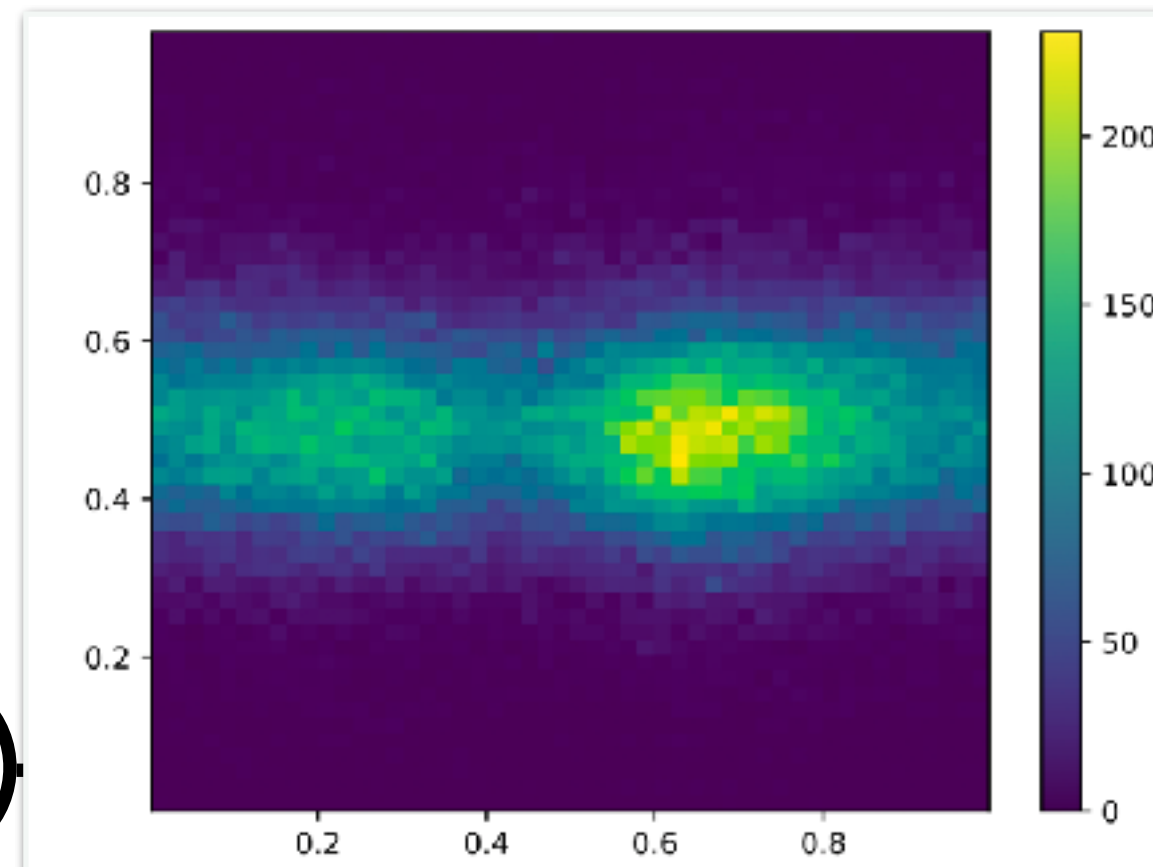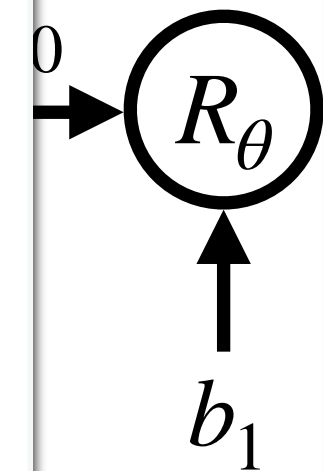- Is it robust against adaptive attackers?

# PercepGuard Design



**Research Questions:**

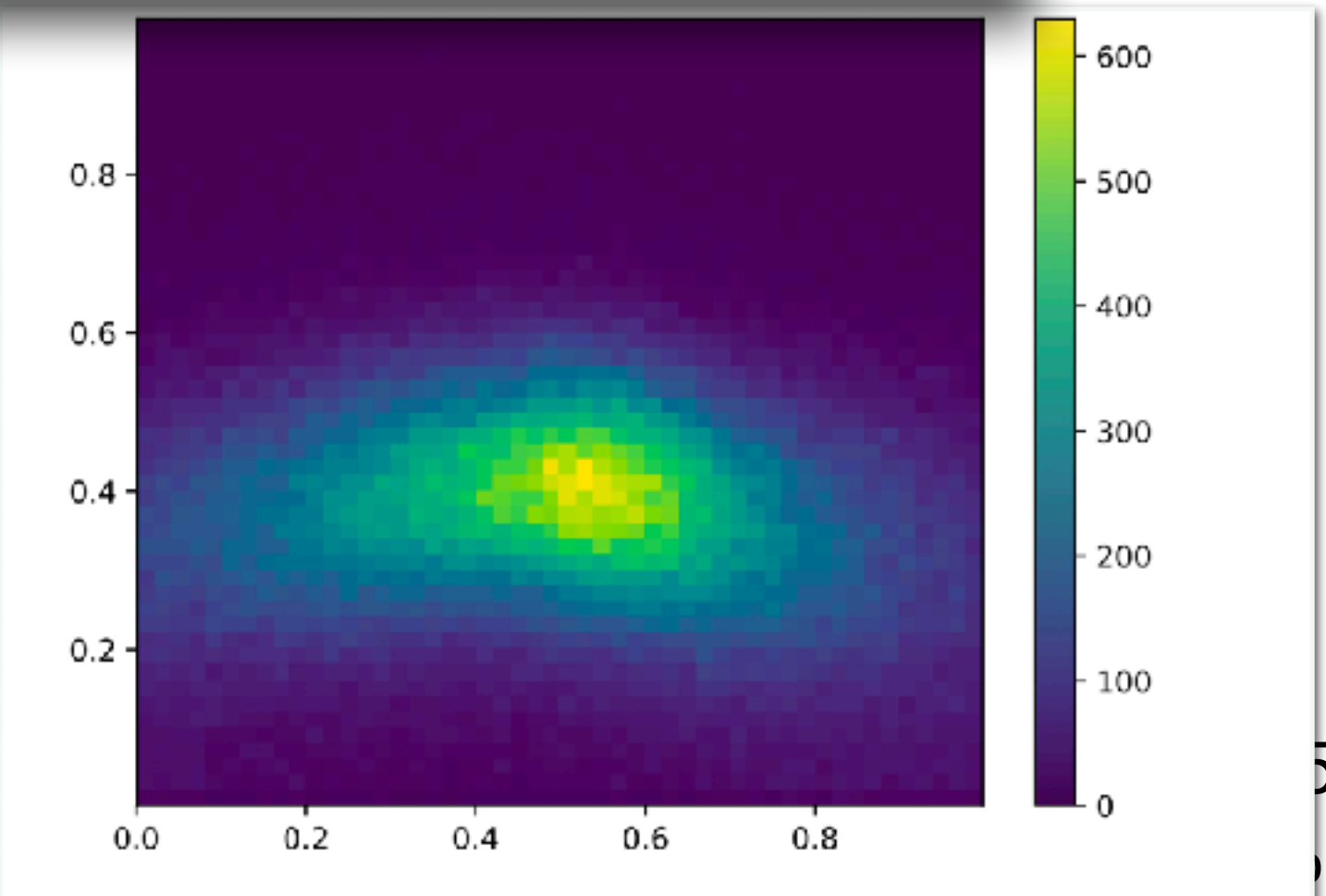- Do bounding boxes provide statistically-sufficient information?



(a) Ground truth: Histogram of cars' locations.

$g_0$ → $R_\theta$ ← $b_1$



(a) Ground truth: Histogram of locatic

Figure 2: Persons' locations.

$g_{N-1}$ → $R$ ← $b_1$



(a) Ground truth: Histogram of locations.

Figure 3: Traffic Signs' locations.

5%

**Research Questions:**

- Do bounding boxes provide statistically-sufficient information?

- Does the detection algorithm produce low false positive and negative rates?

- Is it robust against adaptive attackers?

# Evaluation

**Research Questions:**

- Do bounding boxes provide statistically-sufficient information?

- Does the detection algorithm produce low false positive and negative rates?
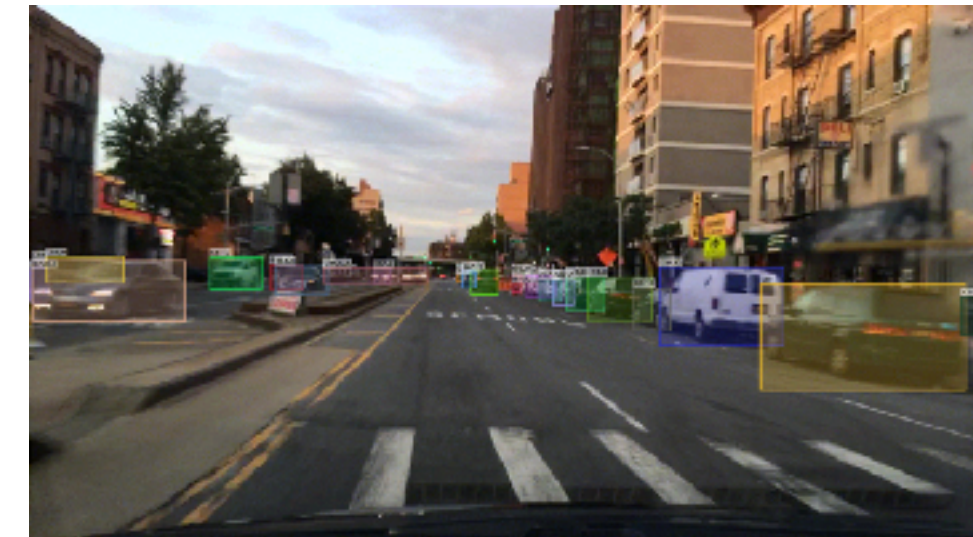
- Is it robust against adaptive attackers?



**Dataset:**

- Berkeley Driving Dataset (BDD)

- Five object classes:

  - bike, bus, car, pedestrian, truck

Classification Accuracy: 95%

False Negative Rate: 5%

# Evaluation



$$\underset{\Delta}{\text{minimize}} \quad \|\Delta\| \qquad \text{such that} \quad \bar{c} = c''$$
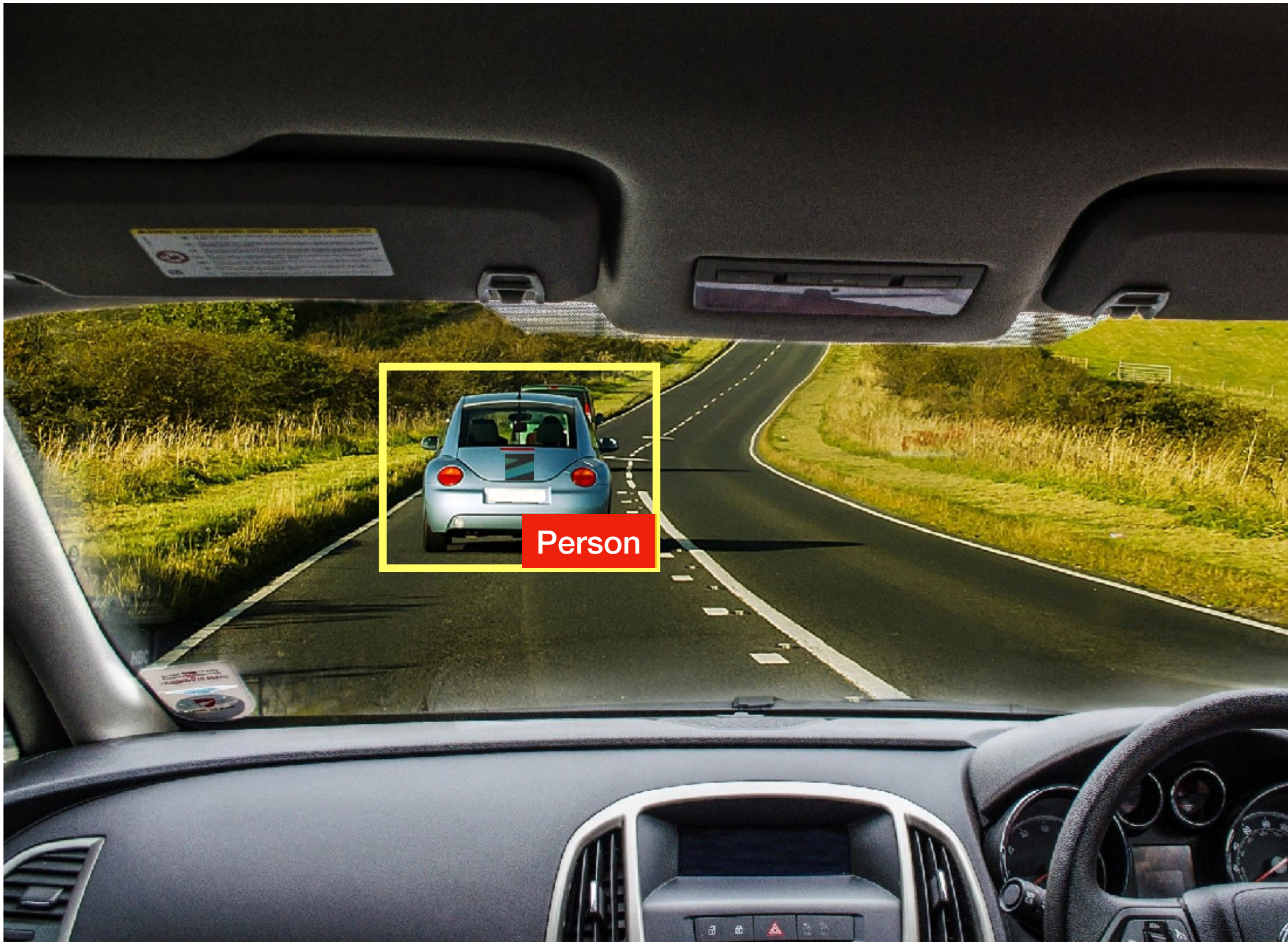
"person"

YOLO's classification result

**Attack Model:**

- *Attack Goal*: Causing the rear car to recognize the front car as a person thus decides to stop (e.g., on a highway)

- *Attacker's Capability*: They utilize the adversarial machine learning to generate adversarial patches with white-box knowledge of the object detection algorithm (e.g., YOLO)

True Positive Rate: ??

False Negative Rate: 5%

# Evaluation

$$\underset{\Delta}{\text{minimize}} \quad \|\Delta\| \qquad \text{such that} \quad \bar{c} = c''$$

$$\text{"person"} \rightarrow \bar{c} = c' \leftarrow \text{Our classification result}$$

Table I: Adversarial patch attacks with BDD100K

| Attack Type | Patch Size | A.M.R. | T.P.R. | A.S.R. |
|---|---|---|---|---|
| Defense-unaware | 20 × 20 | 83.47% | 99.63% | 0.3% |
| | 40 × 40 | 89.41% | 100% | 0% |
| | 60 × 60 | 92.94% | 100% | 0% |

But, what about adaptive attackers, who are aware of our defense and try to evade it?

True Positive Rate: Above 99%

False Negative Rate: 5%

# Evaluation

$$\underset{\Delta}{\text{minimize}} \quad \|\Delta\| \qquad \text{such that} \quad \bar{c} = c''$$
$$\bar{c} = c'$$

Table I: Adversarial patch attacks with BDD100K

| Attack Type | Patch Size | A.M.R. | T.P.R. | A.S.R. |
|---|---|---|---|---|
| Defense-unaware | 20 × 20 | 83.47% | 99.63% | 0.3% |
| | 40 × 40 | 89.41% | 100% | 0% |
| | 60 × 60 | 92.94% | 100% | 0% |

**Research Questions:**

- Do bounding boxes provide statistically-sufficient information?

- Does the detection algorithm produce low false positive and negative rates?

- Is it robust against adaptive attackers?

But, what about adaptive attackers, who are aware of our defense and try to evade it?

True Positive Rate: Above 99%

False Negative Rate: 5%

# Evaluation

$$\underset{\Delta}{\text{minimize}} \quad \|\Delta\| \qquad \text{such that} \quad \bar{c} = c''$$
$$\bar{c} = c'$$

Table I: Adversarial patch attacks with BDD100K

**Research Questions:**

- Do bounding boxes provide statistically-sufficient information?

- Does the detection algorithm produce low false positive and negative rates?

- Is it robust against adaptive attackers?

| Attack Type | Patch Size | A.M.R. | T.P.R. | A.S.R. |
|---|---|---|---|---|
| Defense-unaware | $20 \times 20$ | 83.47% | 99.63% | 0.3% |
| | $40 \times 40$ | 89.41% | 100% | 0% |
| | $60 \times 60$ | 92.94% | 100% | 0% |

True Positive Rate: Above 99%
False Negative Rate: 5%

# Evaluation

$$\underset{\Delta}{\text{minimize}} \quad \|\Delta\| \qquad \text{such that} \quad \bar{c} = c''$$
$$\bar{c} = c'$$

**Research Questions:**

- Do bounding boxes provide statistically-sufficient information?

- Does the detection algorithm produce low false positive and negative rates?

- Is it robust against adaptive attackers?

Table I: Adversarial patch attacks with BDD100K

| Attack Type | Patch Size | A.M.R. | T.P.R. | A.S.R. |
|---|---|---|---|---|
| Defense-unaware | $20 \times 20$ | 83.47% | 99.63% | 0.3% |
| | $40 \times 40$ | 89.41% | 100% | 0% |
| | $60 \times 60$ | 92.94% | 100% | 0% |
| Defense-aware | $20 \times 20$ | 73.25% | 98.74% | 0.92% |
| | $40 \times 40$ | 80.49% | 90.33% | 7.78% |
| | $60 \times 60$ | 87.6% | 85.67% | 12.55% |

True Positive Rate: Above 99%

False Negative Rate: 5%

# Evaluation

$$\underset{\Delta}{\text{minimize}} \quad \|\Delta\| \qquad \text{such that} \quad \bar{c} = c''$$
$$\bar{c} = c'$$

**Research Questions:**

- Do bounding boxes provide statistically-sufficient information?

- Does the detection algorithm produce low false positive and negative rates?

- Is it robust against adaptive attackers?

Table I: Adversarial patch attacks with BDD100K

| Attack Type | Patch Size | A.M.R. | T.P.R. | A.S.R. |
|---|---|---|---|---|
| Defense-unaware | 20 × 20 | 83.47% | 99.63% | 0.3% |
| | 40 × 40 | 89.41% | 100% | 0% |
| | 60 × 60 | 92.94% | 100% | 0% |
| Defense-aware | 20 × 20 | 73.25% | 98.74% | 0.92% |
| | 40 × 40 | 80.49% | 90.33% | 7.78% |
| | 60 × 60 | 87.6% | 85.67% | 12.55% |

True Positive Rate: Above ~~99%~~ 85%

False Negative Rate: 5%

# Evaluation

**Research Questions:**

- Do bounding boxes provide statistically-sufficient information?

- Does the detection algorithm produce low false positive and negative rates?
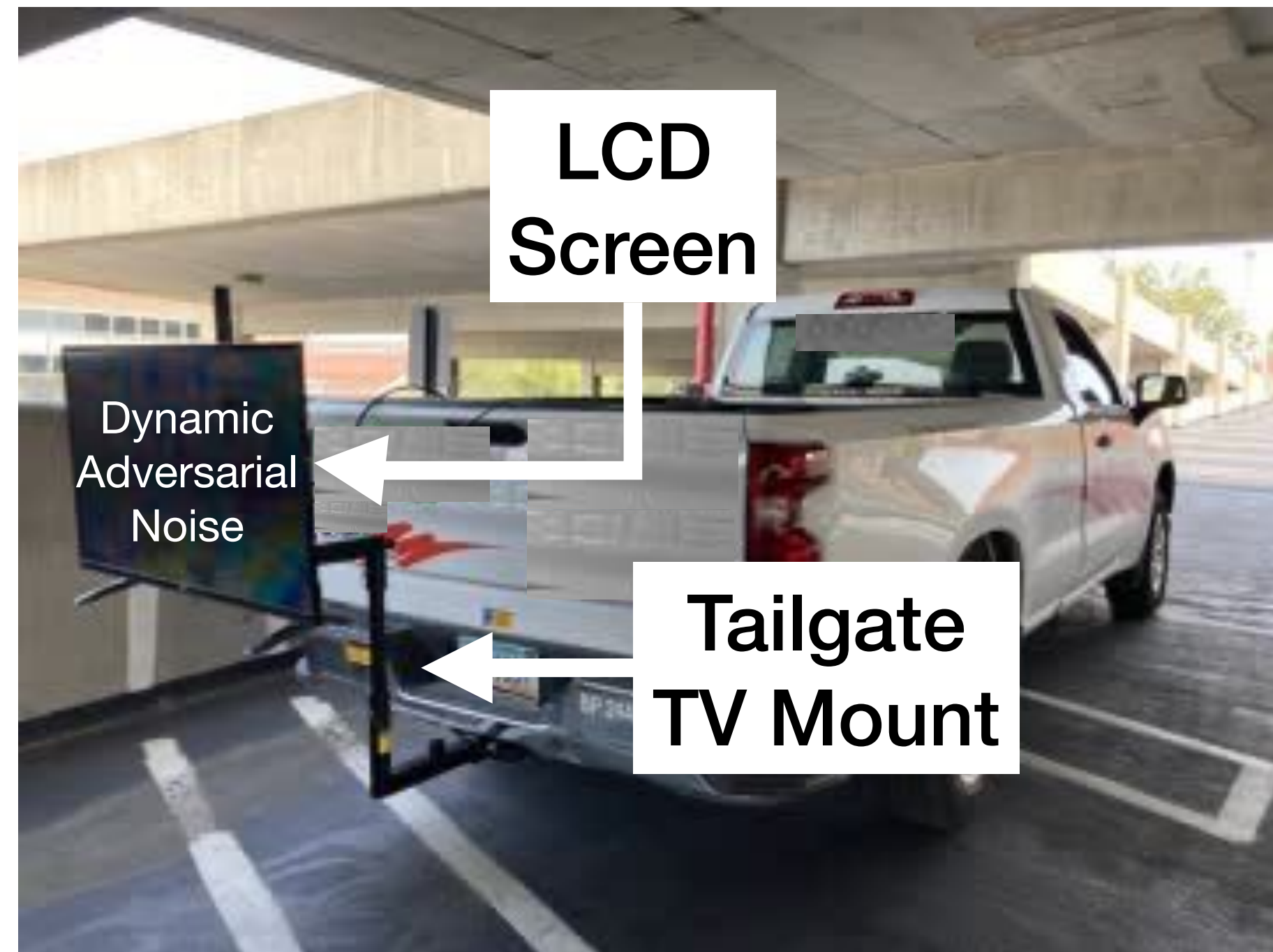
- Is it robust against adaptive attackers?

**A sequence of bounding boxes** → **Object Class**

RNN

+

**Contextual Information**

- Ego-vehicle velocity
- Relative velocity to the object

| | | | | |
|---|---|---|---|---|
| Defense-aware | $20 \times 20$ | 73.25% | 98.74% | 0.92% |
| | $40 \times 40$ | 80.49% | 90.33% | 7.78% |
| | $60 \times 60$ | 87.6% | 85.67% | 12.55% |
| with contexts | $60 \times 60$ | 88.76% | 99.35% | 0.6% |

True Positive Rate: Above ~~99%~~ 85%

False Negative Rate: 5%

# Evaluation

**Research Questions:**

- Do bounding boxes provide statistically-sufficient information?

- Does the detection algorithm produce low false positive and negative rates?

- Is it robust against adaptive attackers?

**A sequence of bounding boxes** → **RNN** → **Object Class**

**+**

**Contextual Information**

- Ego-vehicle velocity
- Relative velocity to the object

| | | | | |
|---|---|---|---|---|
| Defense-aware | $20 \times 20$ | 73.25% | 98.74% | 0.92% |
| | $40 \times 40$ | 80.49% | 90.33% | 7.78% |
| | $60 \times 60$ | 87.6% | 85.67% | 12.55% |
| with contexts | $60 \times 60$ | 88.76% | 99.35% | 0.6% |

True Positive Rate: Above ~~99%~~ ~~85%~~

False Negative Rate: 5%        99%

# Real-world Experiments

# Real-world Experiments

# Real-world Experiments



(a) Person  (b) Stop sign on monitor  (c) Projected stop sign  (d) Adversarial patch

Table 4: Real image attacks in the real-world

| Real Images of | Device | ARR | TPR | ASR |
|---|---|---|---|---|
| People | Monitor | 63.2% | 83.3% | 10.6% |
|  | Projector | 58.8% | 100% | 0% |
| Stop Signs | Monitor | 40.0% | 100% | 0% |
|  | Projector | 20.0% | 100% | 0% |

Table 5: Adversarial patch attacks in the real-world

| Attack Type | Device | AMR | TPR | ASR |
|---|---|---|---|---|
| Defense-unaware | Monitor | 52.2% | 100% | 0% |
|  | Projector | 27.3% | 100% | 0% |
| Defense-aware | Monitor | 45.5% | 100% | 0% |
|  | Projector | 20.0% | 100% | 0% |

# More Evaluation

- Baseline comparison

- Sensitivity analysis

- Alternative operating points
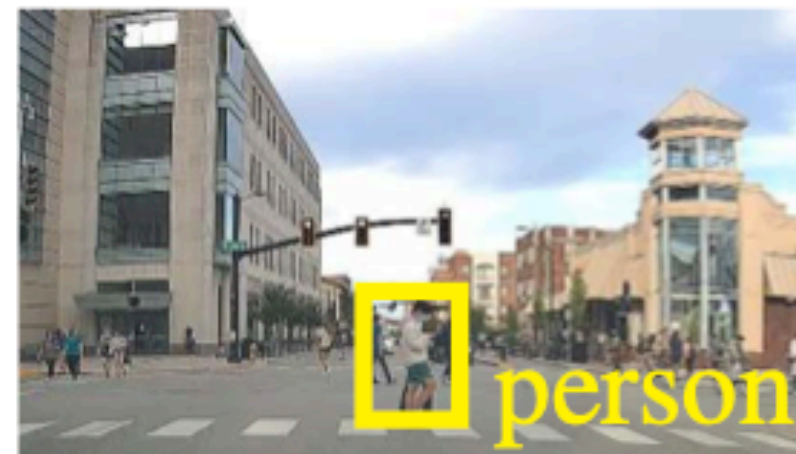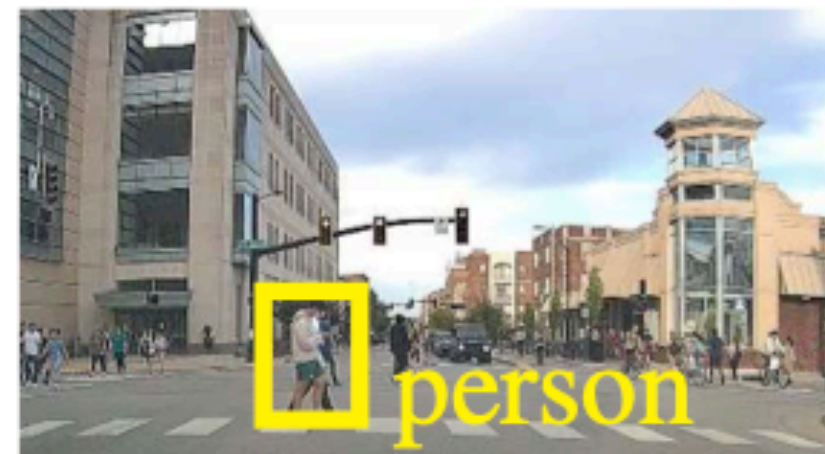
- Additional features

# That Person Moves Like A Car:
# Misclassification Attack Detection
# for Autonomous Systems
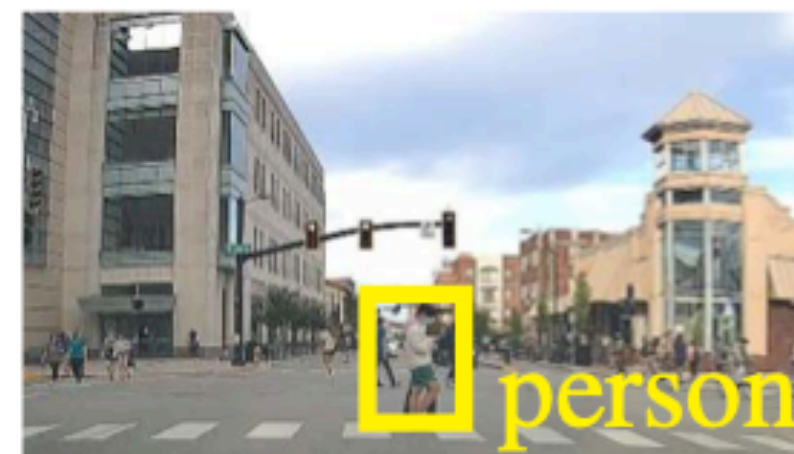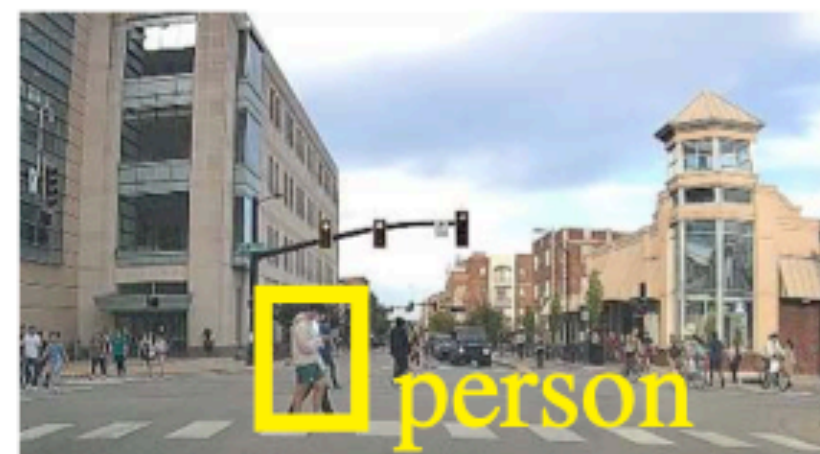# Using Spatiotemporal Consistency

# That Person Moves Like A Car:
## Misclassification Attack Detection
## for Autonomous Systems
## Using Spatiotemporal Consistency

# That Person Moves Like A Car:
# Misclassification Attack Detection
# for Autonomous Systems
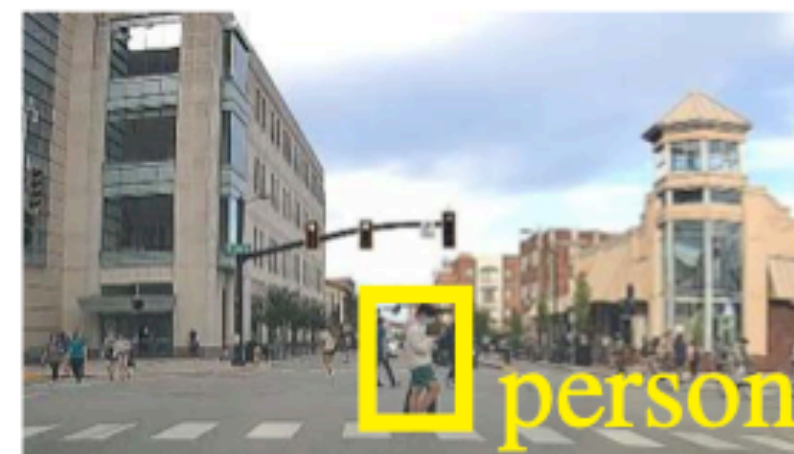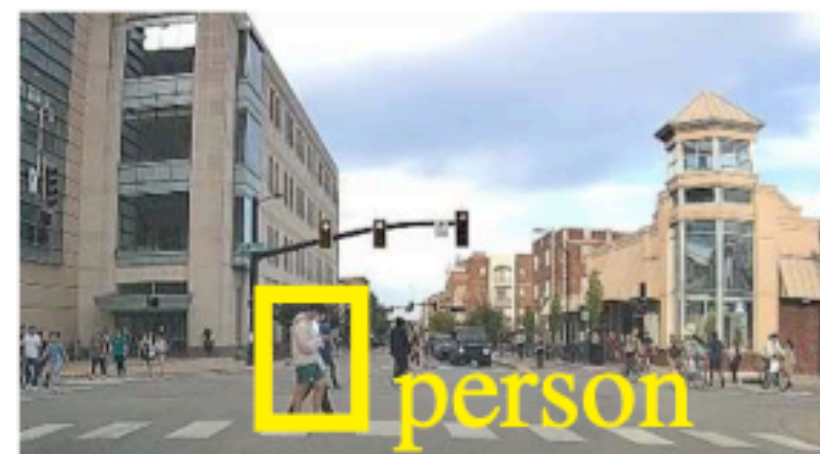# Using Spatiotemporal Consistency

# That Person Moves Like A Car: Misclassification Attack Detection for Autonomous Systems Using **Spatiotemporal Consistency**



- Adaptive Attacks

- Contextual information

- True positive rate: above 99%

- False negative rate: 5%
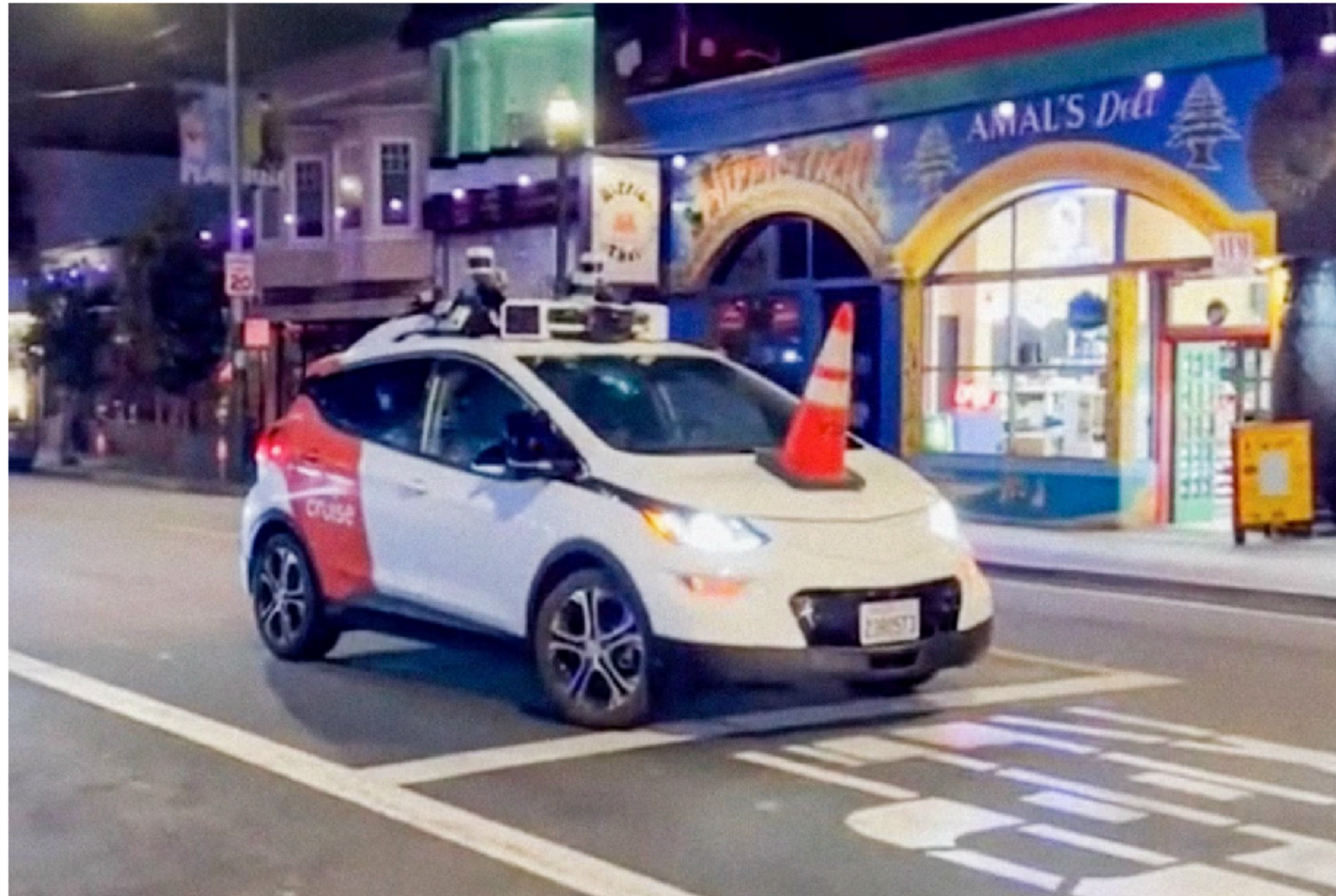
# Future Work



- More spatiotemporal features
  - Different Sensors
  - Semantic Segmentation
- Detecting object creation attacks
- Attention-based
- Sensor Configuration Randomization

# The Self-Driving Cars Wearing a Cone of Shame

There's a brilliant activist campaign to stop San Francisco's autonomous vehicles in their tracks.

BY ALISON GRISWOLD

JULY 11, 2023 • 10:45 AM



It looks like a sad unicorn (which, in a way, it is).   Screengrab from TikTok/Safe Street Rebel

slate.com

# That Person Moves Like A Car: Misclassification Attack Detection for Autonomous Systems Using Spatiotemporal Consistency

Yanmao Man[#], Raymond Muller[§], Ming Li[#], Z. Berkey Celik[§], Ryan Gerdes[‡]

[#]University of Arizona    [§]Purdue University    [‡]Virginia Tech