# Semantic Adversarial Deep Learning

A paper of Tommaso Dreossi, Somesh Jha, and Sanjit A. Seshia
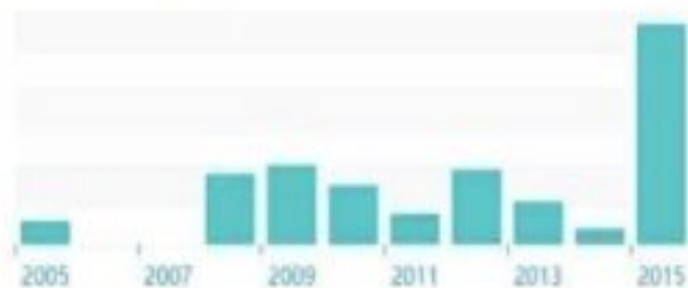
Thanks to Somesh Jha and Sanjit A. Seshia for some of the slides
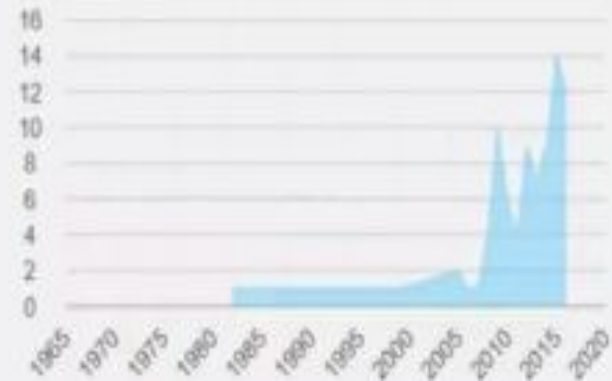
Presented by Yilin Han
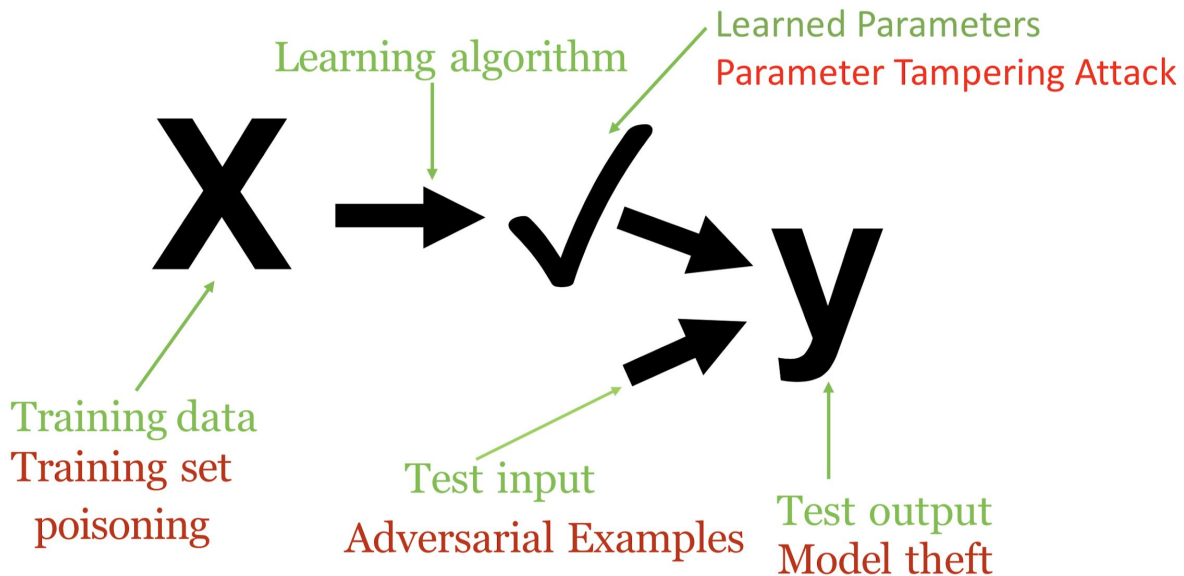
# Background



Adversarial attack

Popularity Over Time:



Publish Year

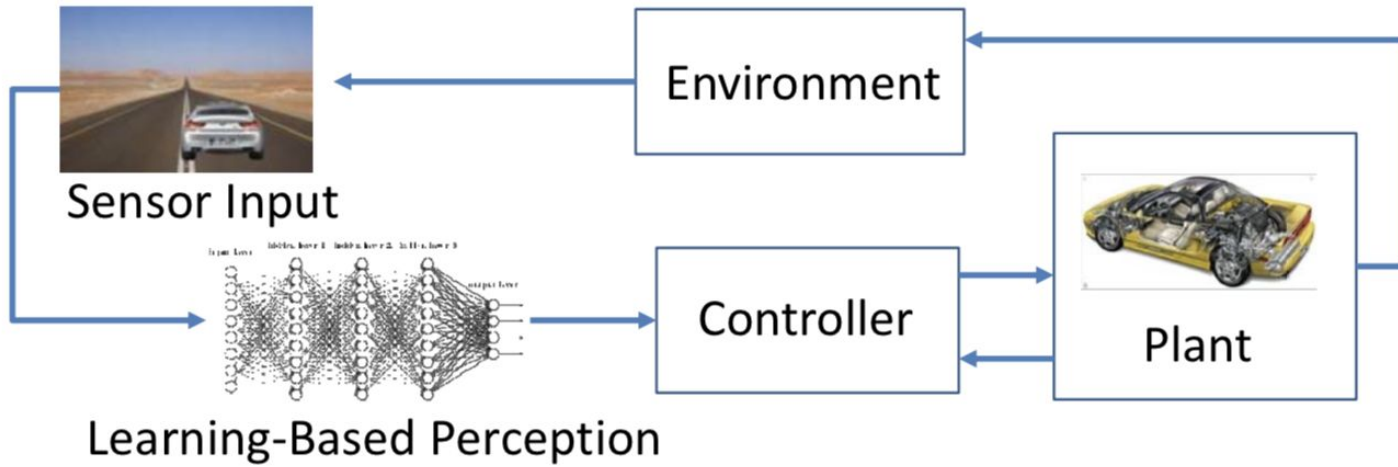# Adversarial Attacks on ML

# Definition

"Adversarial examples are inputs to machine learning models that an attacker has intentionally designed to cause the <span style="color:red">model</span> to make a mistake" (Goodfellow et al 2017)

"Adversarial examples are inputs to machine learning models that an attacker has intentionally designed to cause the <span style="color:red">entire system</span> to make a mistake"

# Automatic Emergency Braking System(AEBS)

Goal: Brake whenever an obstacle is detected

# Semantic Adversarial Analysis and Training
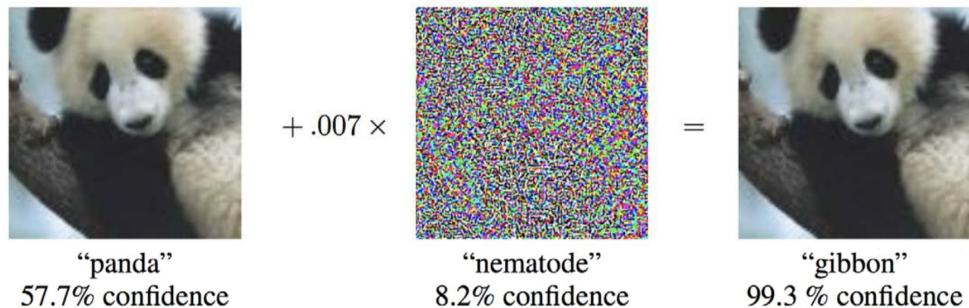
Goal: DNN analysis must be more semantic

- Semantic modification
- System-level specification
- Semantic (re-)training
- Confidence-based analysis

# Semantic Modification

- Allow "Noise": Add a vector δ
- Allow richer transformations:

  Translation, Cloudy
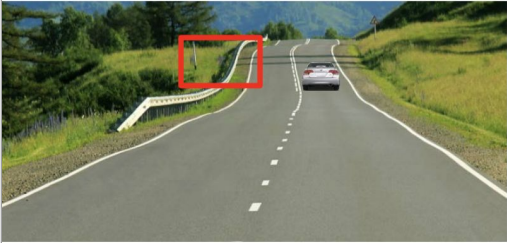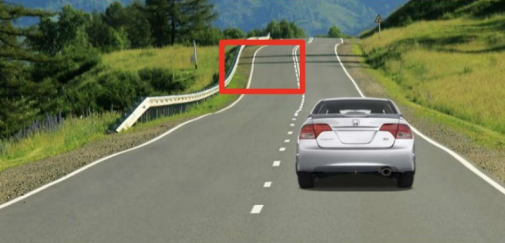  background…..

- DSL that relates to applications



$+ .007 \times$

$=$

"panda"
57.7% confidence

"nematode"
8.2% confidence

"gibbon"
99.3 % confidence

Non-semantic perturbation (i.e., noise)



Semantic perturbation (i.e., translation) 61

# System-level Specification

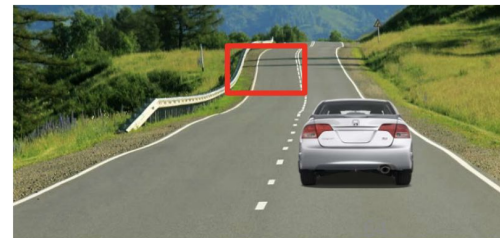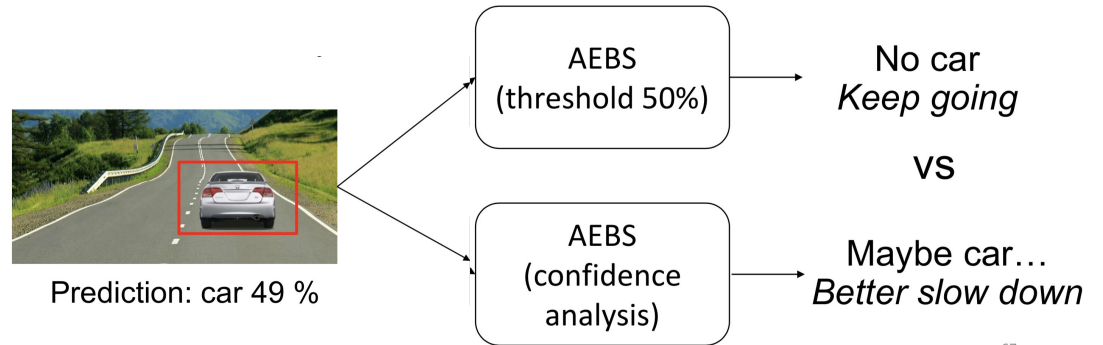|  |  |  |
|---|---|---|
| Perception-level-spec: "Detect cars" | NO | NO |
| Perception-level-spec: "Do not crash" | YES | NO |

# Semantic Training

Semantic augmentation

# Confidence-based Analysis

They argue that confidence levels must be used within the design of ML-based system.



Prediction: car 49 %

AEBS
(threshold 50%)

No car
*Keep going*

VS

AEBS
(confidence analysis)

Maybe car…
*Better slow down*

# Problem

Can we generate adversarial examples that cause system-level failure?


Can we use formal method to verify CPS?

# Compositional Falsification

Statement:

given a formal specification φ (say in a formalism such as signal temporal logic) and a CPS+ML model M, find an input for which M does not satisfy φ.
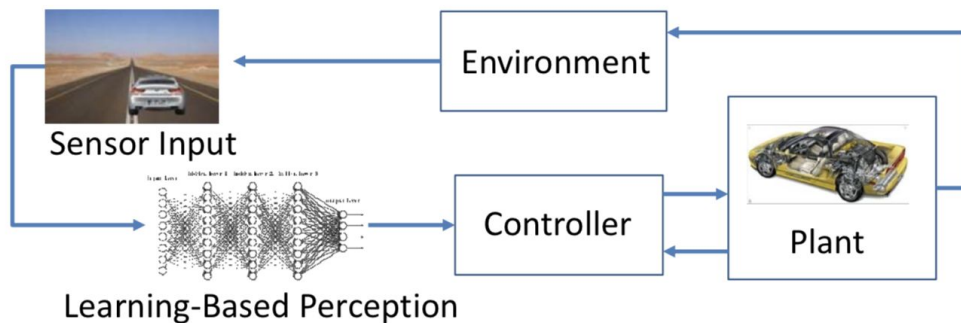
Example:  $\mathbf{G}_{0,T}(\|\mathbf{x}_{\text{ego}} - \mathbf{x}_{\text{obs}}\|_2 \geq 2).$

Problem:

How to deal with ML component?

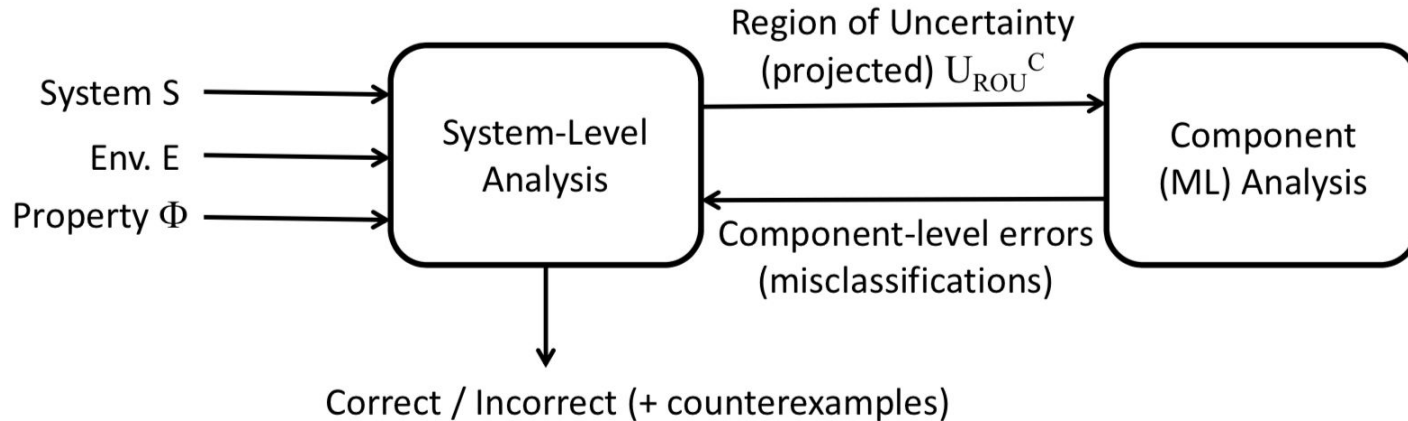# Compositional Falsification

- Standard solution: use compositional verification (Modular)
- However, no formal specification for neural network component.

# Approach: Use a System-Level Specification and Combine CPS Falsifier with ML Analyzer
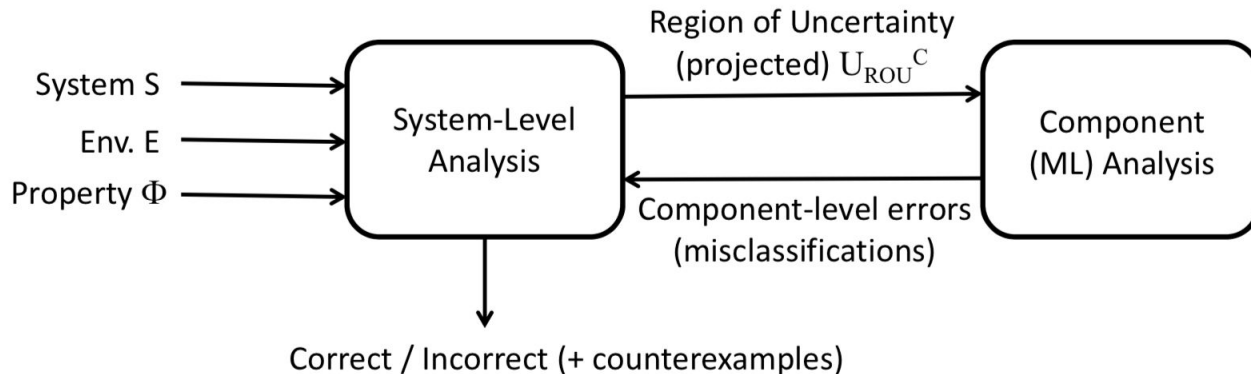
Do not verify the DNN Object Detector.

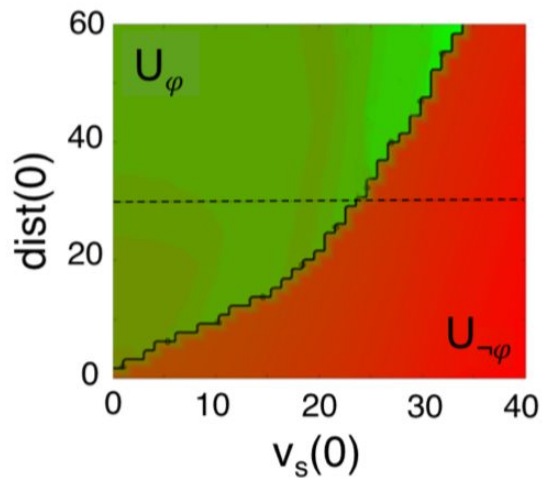Verify the system containing the DNN model

# System-Level Analysis

CPS Falsifier uses abstraction of ML component:

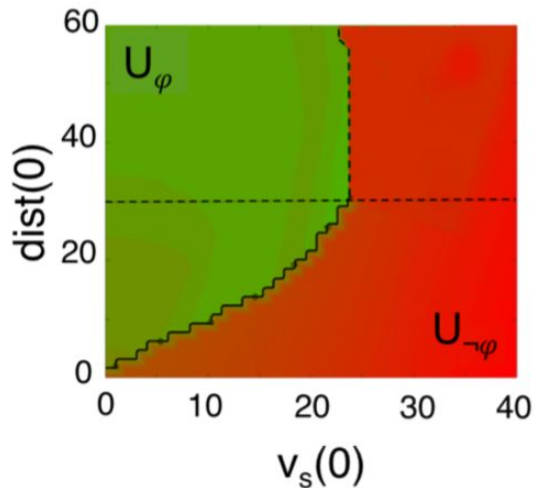- Optimistic analysis: assume ML classifier is always correct
- Pessimistic analysis: assume it is always wrong
- Difference is the region of uncertainty where output of the ML component "matters"



Region of Uncertainty (projected) $U_{ROU}{}^C$

System S → System-Level Analysis

Env. E →

Property $\Phi$ →

Component-level errors (misclassifications)

Component (ML) Analysis

Correct / Incorrect (+ counterexamples)

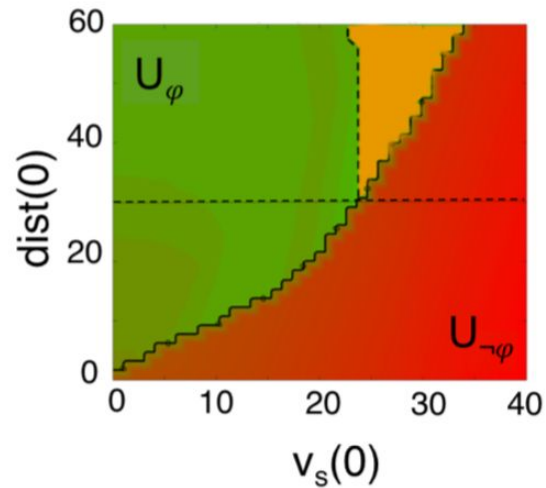# Identify Region of Uncertainty



ML always correct

ML always wrong

Region of uncertainty
(yellow)

# Machine Learning Analyzer



brightness · car z-pos · (1,1,1) · (0,1,0) · car x-pos · (0,0,0) · (1,0,0)

*Semantic modification space $\tilde{X}$*

*Abstraction map*

brightness · car z-pos · car x-pos

Systematic
Sampling (low-discrepancy sampling)

$\gamma(a)$

$a$

$x$

*Abstract space A*

Neural network
$y \in \{car, \neg car\}$

# Sample Results



Misclassifications

Misclassification not of concern

Potential hazard (system-level safety violation)

Corner case

# System Training

- Train the model with both original and counterexample test sets.
- That is they trained a model with hinge loss with a bias factor <span style="color:red">k</span>, which means no penalty for a misclassification if confidence level + k < truth label

$$l(\hat{y}) = \max(0, k + \max_{i \neq l}(\hat{y}_i - \hat{y}_l))$$

Result:

- Accuracy increases on counterexamples testing set.
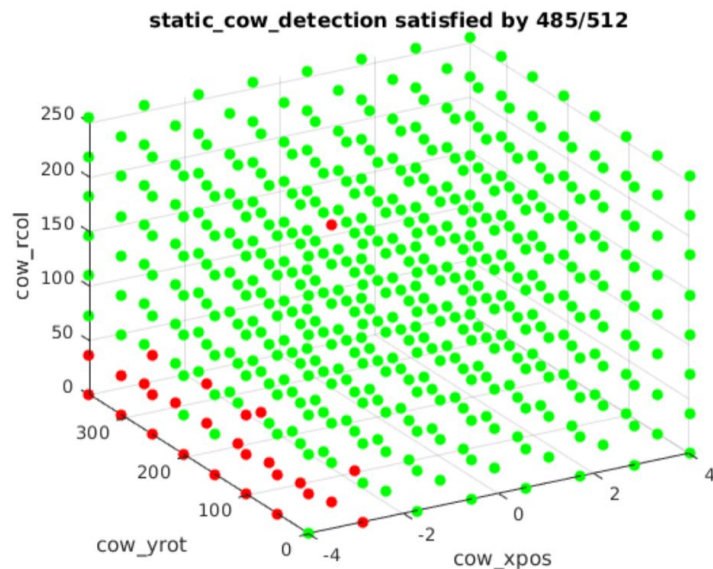- Accuracy decreases on original testing set.

# System Training

- Train the model with both original and counterexample test sets. Where the counterexamples are generated from composition falsification framework.
- Red dots are Semantic counterexamples

Result:

Still specification violation

However, obstacles get detected earlier.



static_cow_detection satisfied by 485/512

# Conclusion

- Formal Methods can apply to cyber-physical system with high assurance.
- Compositional falsification framework can also be used for verification