

Towards Verified Artificial Intelligence

Sanjit A. Seshia

Professor

EECS, UC Berkeley

Joint work with

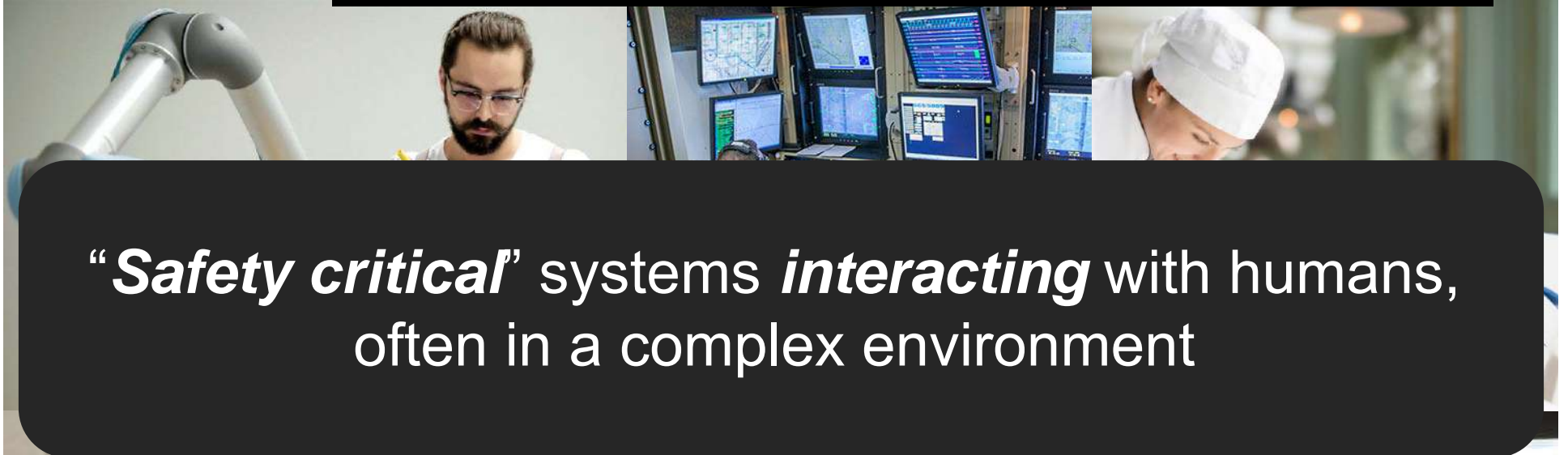
Dorsa Sadigh, Tommaso Dreossi, Alexander Donze,
Anca Dragan, S. Shankar Sastry



SETTA 2017
October 24, 2017



Human Cyber-Physical Systems (h-CPS)

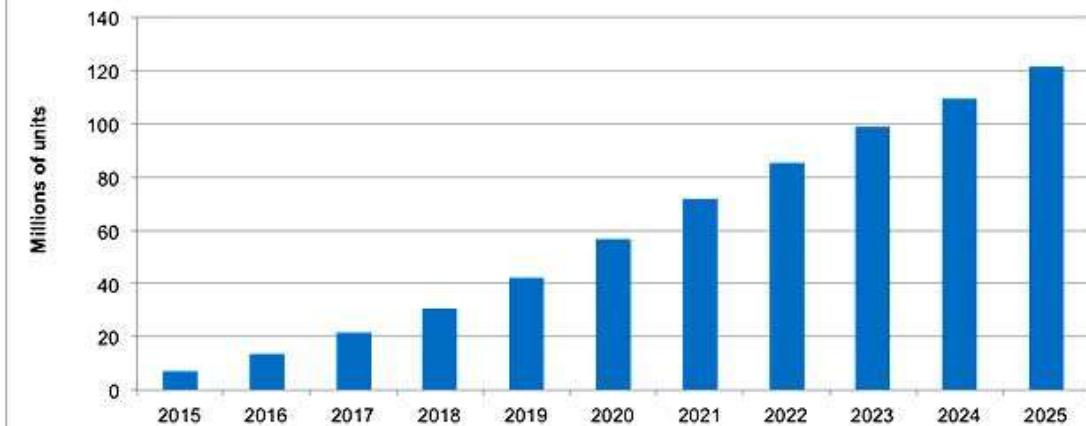


“Safety critical” systems *interacting* with humans,
often in a complex environment



Growing Use of Machine Learning/AI in Cyber-Physical Systems

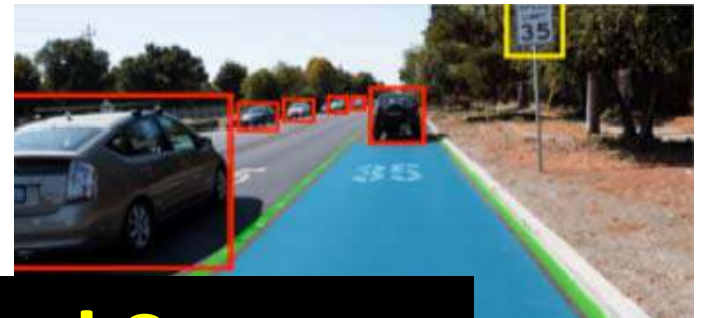
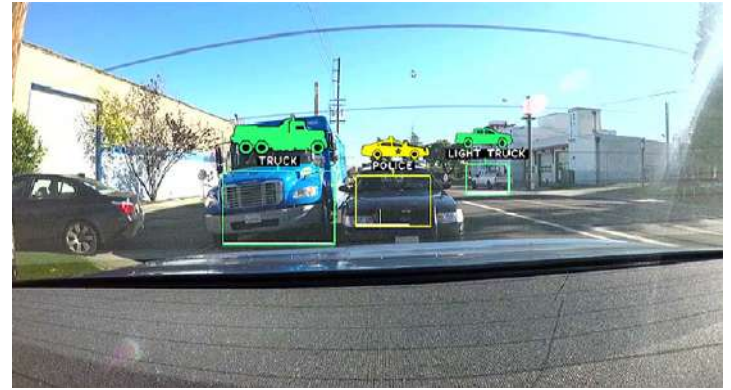
Artificial Intelligence based systems for automotive



Notes: Includes: infotainment (virtual assistance, gesture and speech recognition) and autonomous driving applications (object detection and freespace detection)

Source: IHS Technology - Automotive Electronics Roadmap Report, H1 2016

© 2016 IHS



Many Safety-Critical Systems



Artificial Intelligence (AI)

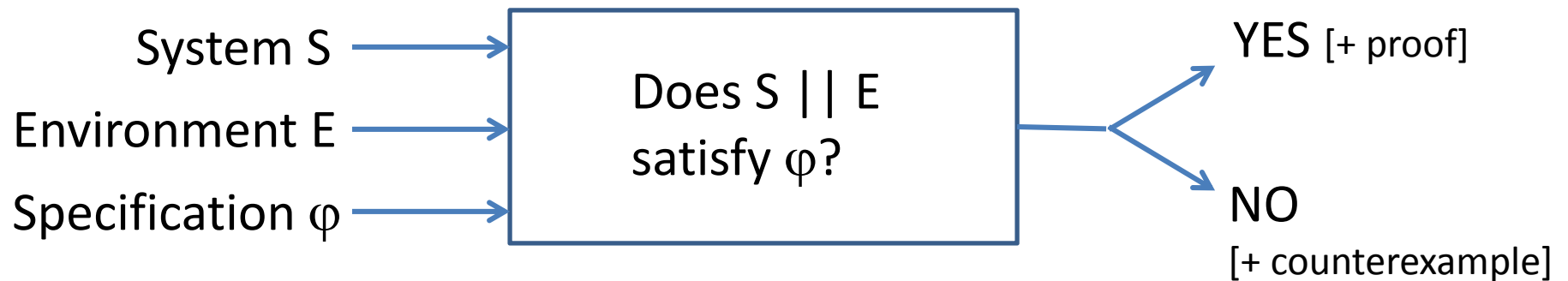
Computational Systems that attempt to **mimic aspects of human intelligence**, including especially the ability to **learn from experience**.

How do we ensure that AI-based systems are Dependable?

The Formal Methods Lens

- Formal Methods \approx Computational Proof methods
 - Specification/Modeling \approx Statement of Conjecture/Theorem
 - Verification \approx Proving/Disproving the Conjecture
 - Synthesis \approx Generating (parts of) Conjecture/Proof
 - Tools/techniques: SAT / SMT solvers, model checkers, theorem provers, simulation-based falsification, ...

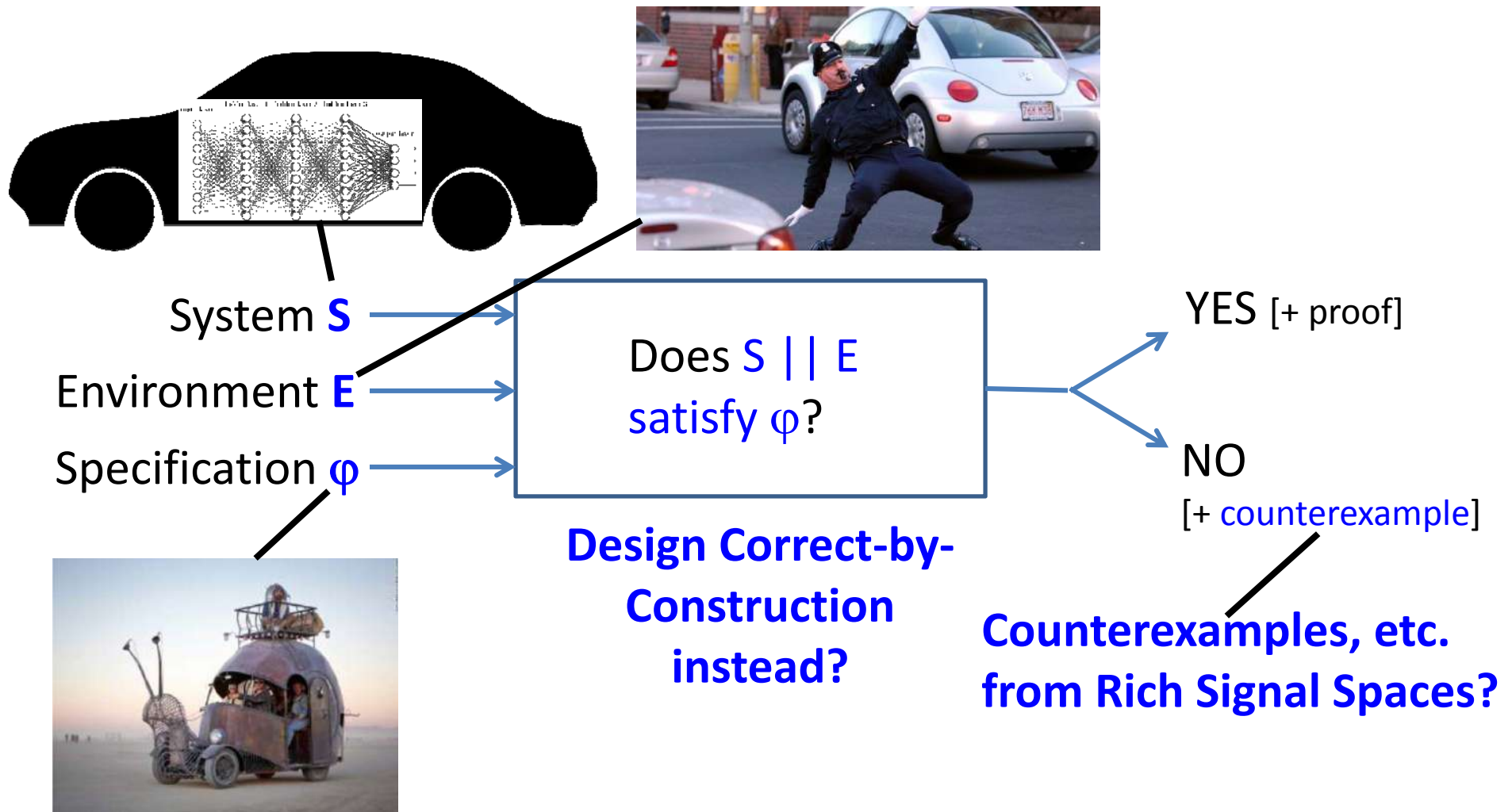
Verification:



Challenges for Verified AI

S. A. Seshia, D. Sadigh, S. S. Sastry.

Towards Verified Artificial Intelligence. July 2016. <https://arxiv.org/abs/1606.08514>.



Talk Outline

- Environment Modeling Challenge
 - Interaction-Aware Control for Human-CPS
- Specification (& Verification) Challenge
 - Verifying Robustness (of Interaction-Aware Controller)
 - Falsification for Deep Learning based CPS
- Conclusions and Future Directions
 - Towards a New Design Methodology for AI-based Systems

Environment Modeling Challenge – Uncertainty and Unknowns

Self-Driving Vehicles: Interact with Humans in Complex Environments;
Significant use of machine learning!



Known Unknowns and
Unknown Unknowns!!

Cannot represent all possible
environment scenarios

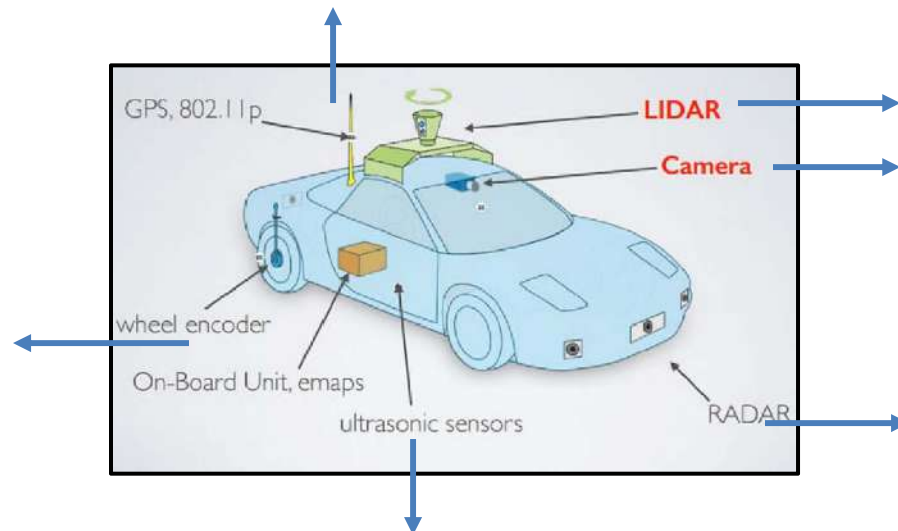
Idea 1: **Introspective** Environment Modeling



Impossible to model
all possible scenarios

Approach: ***Introspect on System to Model the Environment***

Identify: (i) **Interface** between System & Environment,
(ii) (Weakest) **Assumptions** needed to Guarantee Safety/Correctness



Algorithmic techniques to
*generate weakest interface
assumptions and monitor them
at run-time* for potential
violation/mitigation

[Li, Sadigh, Sastry, Seshia; TACAS'14]

Idea 2: **Active** Data Gathering and Learning

***Monitor and Interact with the Environment,
Offline and Online, to Model It.***

Google's Driverless Cars Run Into Problem: Cars With Drivers

By MATT RICHTEL and CONOR DOUGHERTY SEPT. 1, 2015



MOUNTAIN VIEW, Calif. — [Google](#), a leader in

“One of the biggest challenges facing automated cars is *blending them into a world in which humans don't behave by the book.*”

it can be tough to get around if you are a stickler for the rules. One Google car, in a test in 2009, couldn't get through a four-way stop because its sensors kept waiting for other (human) drivers

The Google self-driving car, with Eric Schmidt, left, the company's executive chairman, and Transportation Secretary Anthony Foxx. Justin Sullivan/Getty Images

Challenge: Environment (Human) Modeling

Interaction-Aware Control

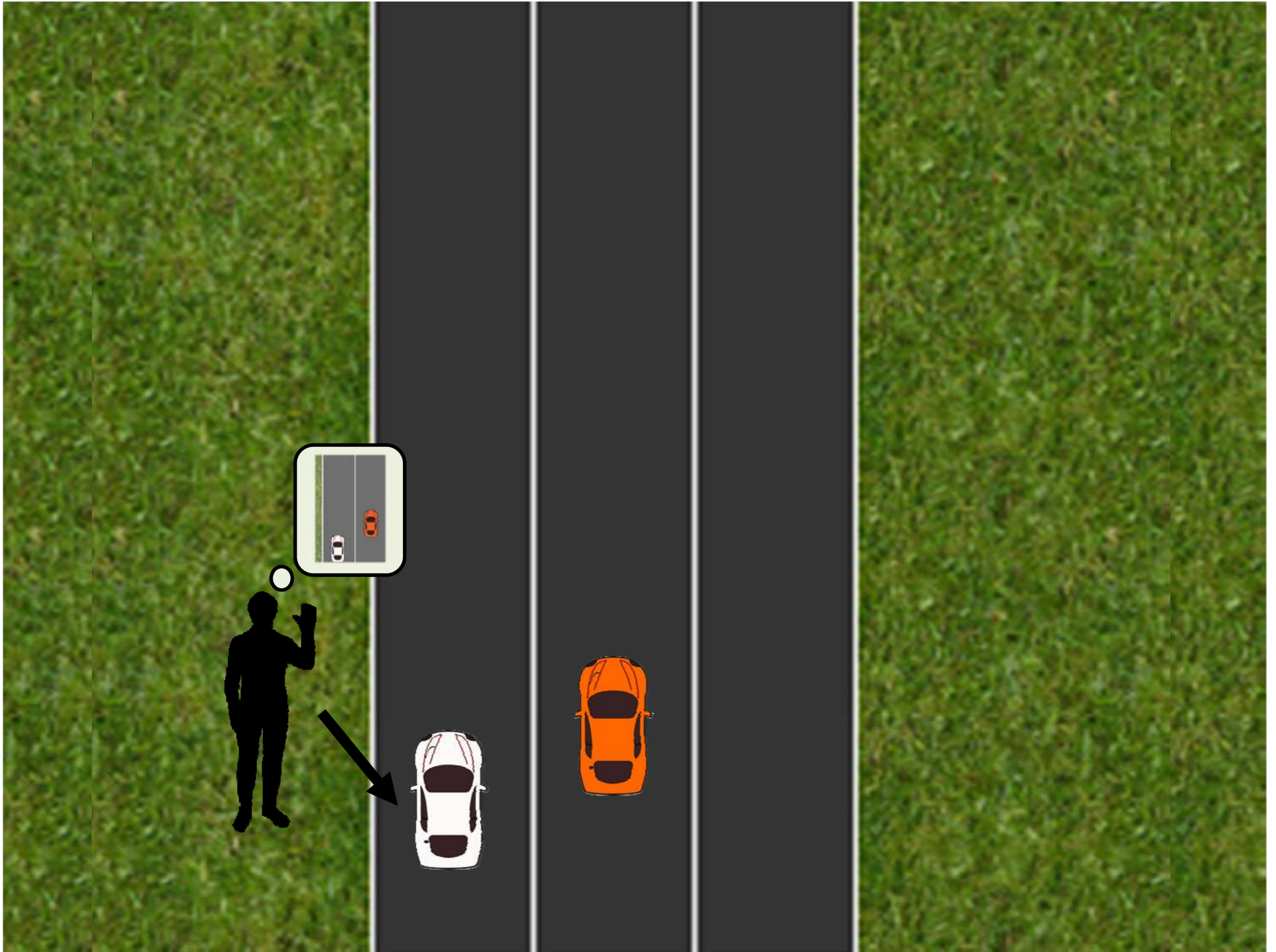
- D. Sadigh, S. Sastry, S. A. Seshia, A. Dragan. *Planning for Autonomous Cars that Leverages Effects on Human Actions*. In RSS, 2016.
- D. Sadigh, S. Sastry, S. A. Seshia, A. Dragan. *Information Gathering Actions over Internal Human State*. In IROS, 2016.
- D. Sadigh, A. Dragan, S. Sastry, S. A. Seshia. *Active Preference-Based Learning of Reward Functions*. In RSS, 2017.

Lane Change on a Highway

Human



Robot



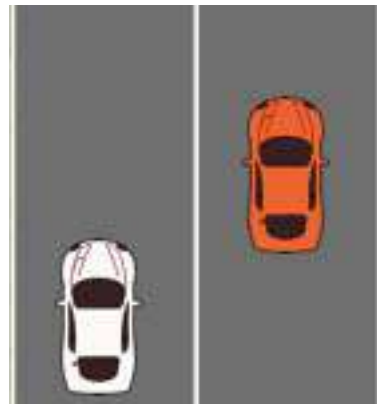


Interaction as a Dynamical System

$$x^{t+1} = f_{\mathcal{H}}(f_{\mathcal{R}}(x^t, u_{\mathcal{R}}^t), u_{\mathcal{H}}^t)$$

Robot actions

u_R



Human
actions u_H

Model the problem as a *Stackelberg (turn-based) Game*.
Robot moves first.

Assumptions/Simplifications

Model Predictive (Receding Horizon) Control:

Optimize over short time horizon N , replan at every step t .

$$R_{\mathcal{R}}(x, \mathbf{u}_{\mathcal{R}}, \mathbf{u}_{\mathcal{H}}) = \sum_{t=1}^N r_{\mathcal{R}}(x^t, \mathbf{u}_{\mathcal{R}}^t, \mathbf{u}_{\mathcal{H}}^t) \quad R_{\mathcal{H}}(x, \mathbf{u}_{\mathcal{R}}, \mathbf{u}_{\mathcal{H}}) = \sum_{t=1}^N r_{\mathcal{H}}(x^t, \mathbf{u}_{\mathcal{R}}^t, \mathbf{u}_{\mathcal{H}}^t)$$

Assume *deterministic* “*rational*” human model,
human optimizes reward function which is a linear
combination of “features”.

Human has full access to $\mathbf{u}_{\mathcal{R}}$ for the short time horizon.

$$\mathbf{u}_{\mathcal{H}}^*(x_0, \mathbf{u}_{\mathcal{R}}) = \operatorname{argmax}_{\mathbf{u}_{\mathcal{H}}} R_{\mathcal{H}}(x_0, \mathbf{u}_{\mathcal{R}}, \mathbf{u}_{\mathcal{H}})$$

Learning (Human) Driver Models

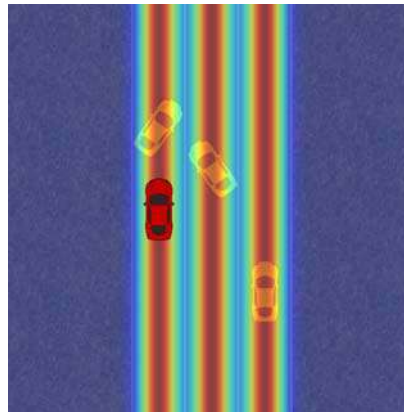
Learn Human's reward function based on **Inverse Reinforcement Learning** [Ziebart et al, AAAI'08; Levine & Koltun, 2012].

Assume structure of human reward function:

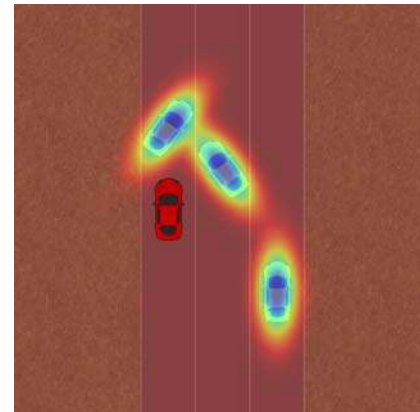
$$r_H(x^t, u_R^t, u_H^t) = w^\top \phi(x^t, u_R^t, u_H^t)$$



(a) Features for the boundaries of the road



(b) Feature for staying inside the lanes.

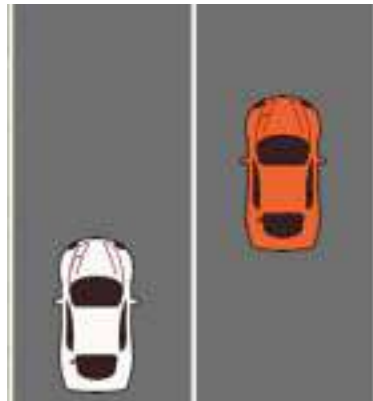


(c) Features for avoiding other vehicles.

Interaction as a Dynamical System

$$\mathbf{u}_R^* = \operatorname{argmax}_{\mathbf{u}_R} R_R(x_0, \mathbf{u}_R, \mathbf{u}_H^*(x_0, \mathbf{u}_R))$$

Model \mathbf{u}_H^* as optimizing the human reward function R_H .



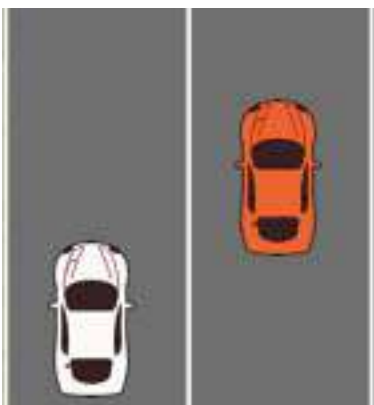
Find optimal actions for the autonomous vehicle while accounting for the human response \mathbf{u}_H^* .

$$\mathbf{u}_H^*(x_0, \mathbf{u}_R) = \operatorname{argmax}_{\mathbf{u}_H} R_H(x_0, \mathbf{u}_R, \mathbf{u}_H)$$

Solution of Nested Optimization

$$u_{\mathcal{R}}^* = \operatorname{argmax}_{u_{\mathcal{R}}} R_{\mathcal{R}}(x, u_{\mathcal{R}}, u_{\mathcal{H}}^*(x, u_{\mathcal{R}}))$$

$$R_{\mathcal{R}}(x, u_{\mathcal{R}}, u_{\mathcal{H}}) = \sum_{t=1}^N r_{\mathcal{R}}(x^t, u_{\mathcal{R}}^t, u_{\mathcal{H}}^t)$$



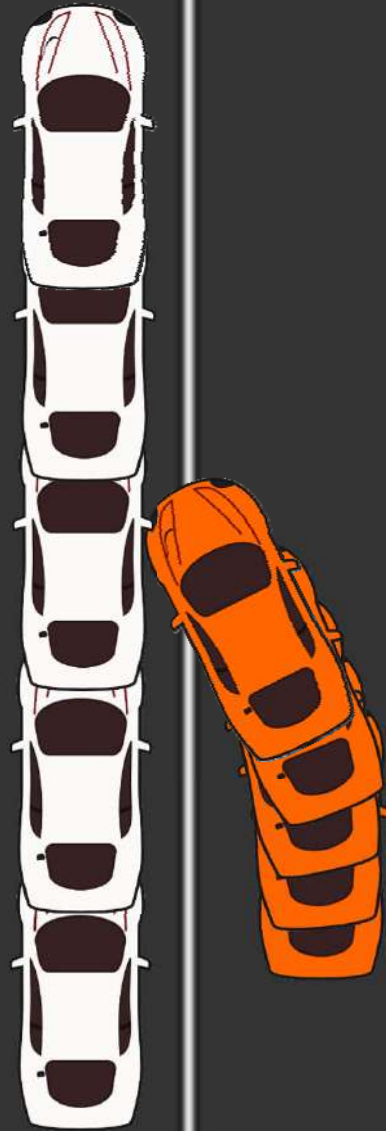
Gradient-Based Method (Quasi-Newton): (solve using L-BFGS technique)

$$\left\{ \begin{array}{l} R_{\mathcal{R}}(x, u_{\mathcal{R}}, u_{\mathcal{H}}^*) \\ \frac{\partial R_{\mathcal{R}}}{\partial u_{\mathcal{R}}} = \frac{\partial R_{\mathcal{R}}}{\partial u_{\mathcal{H}}} \frac{\partial u_{\mathcal{H}}^*}{\partial u_{\mathcal{R}}} + \frac{\partial R_{\mathcal{R}}}{\partial u_{\mathcal{R}}} \end{array} \right.$$

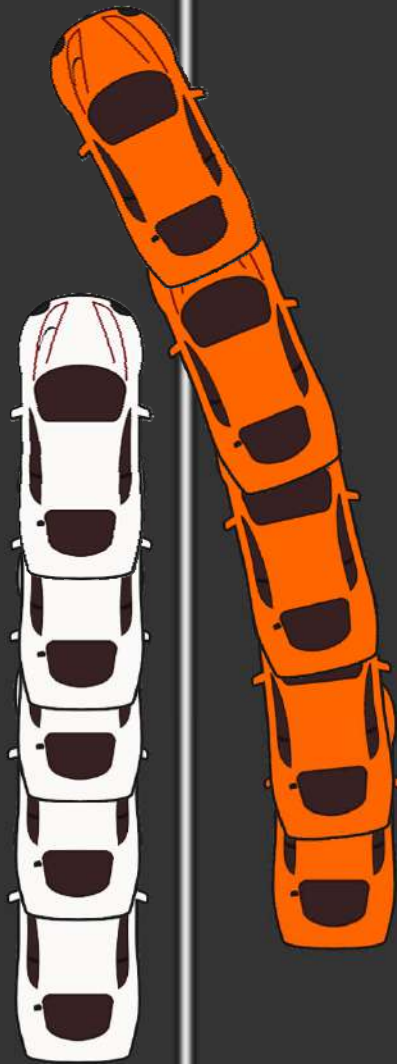
$$u_{\mathcal{H}}^*(x, u_{\mathcal{R}}) \approx \operatorname{argmax}_{u_{\mathcal{H}}} R_{\mathcal{H}}(x, u_{\mathcal{R}}, u_{\mathcal{H}})$$

$$R_{\mathcal{H}}(x, u_{\mathcal{R}}, u_{\mathcal{H}}) = \sum_{t=1}^N r_{\mathcal{H}}(x^t, u_{\mathcal{R}}^t, u_{\mathcal{H}}^t)$$

Cautious Lane Change



Interaction-Aware Lane Change





Aggressive Driver



Distracted Driver



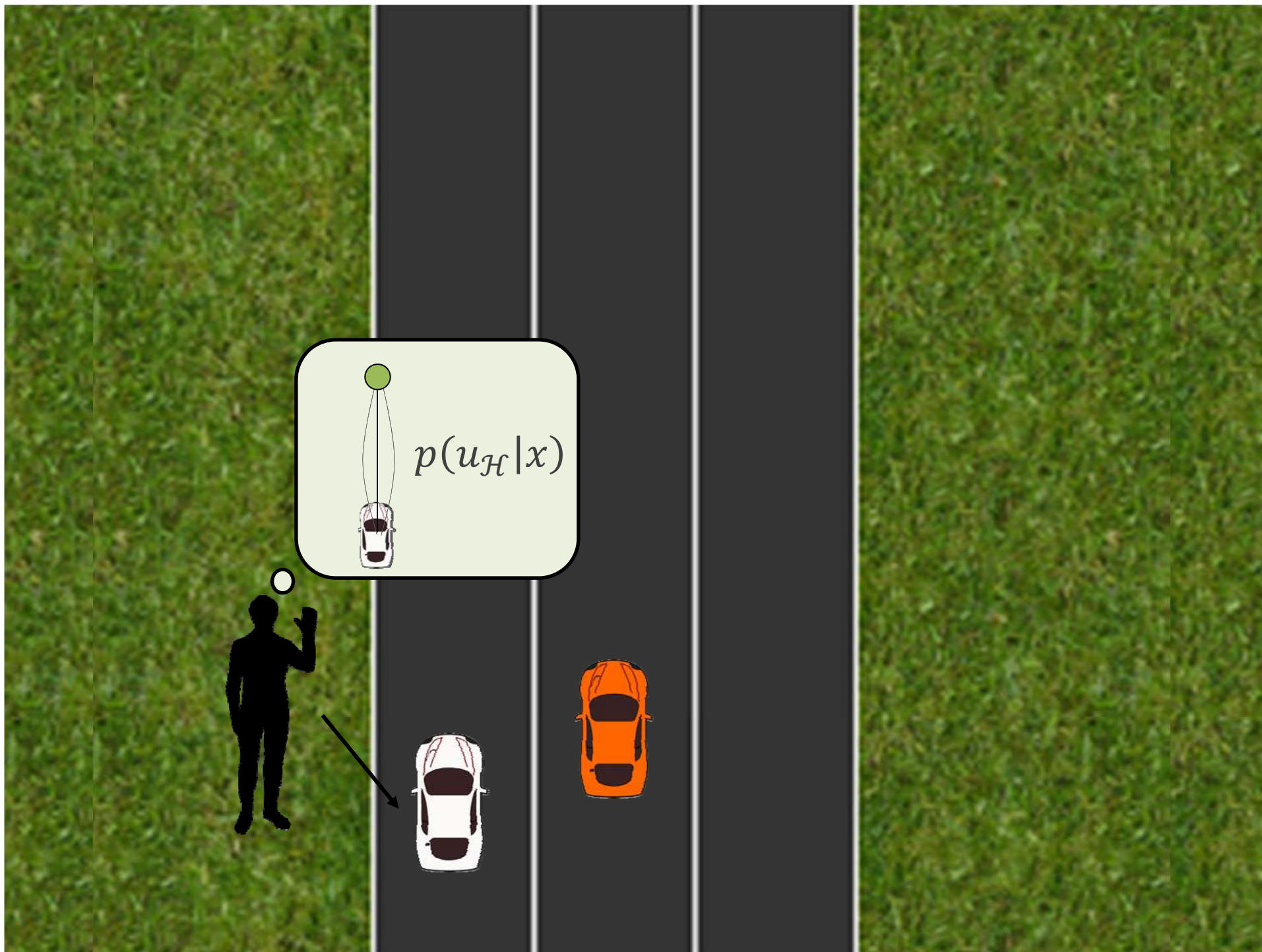
Cautious Driver



Attentive Driver

We can't rely on a
single driver model.

We need to **differentiate**
between different drivers.





$$p(u_{\mathcal{H}}|x) \propto \exp(R_{\mathcal{H}}(x, u_{\mathcal{H}}))$$





$$p(u_{\mathcal{H}}|x, \theta) \propto \exp(R_{\mathcal{H}}(x, u_{\mathcal{H}}, \theta))$$



$$b_{t+1}(\theta) \propto b_t(\theta) \cdot p(u_{\mathcal{H}}|x_t, \theta)$$

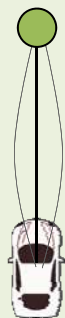




$$p(u_{\mathcal{H}}|x, \theta) \propto \exp(R_{\mathcal{H}}(x, u_{\mathcal{H}}, \theta))$$



$$b_{t+1}(\theta) \propto b_t(\theta) \cdot p(u_{\mathcal{H}}|x_t, \theta)$$



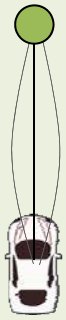
$$p(u_{\mathcal{H}}|x, \theta, u_{\mathcal{R}}) \propto \exp(R_{\mathcal{H}}(x, u_{\mathcal{H}}, \theta, u_{\mathcal{R}}))$$



$$b_{t+1}(\theta) \propto b_t(\theta) \cdot p(u_{\mathcal{H}}|x_t, \theta, u_{\mathcal{R}})$$

$$u_{\mathcal{R}} = \operatorname{argmax}_{u_{\mathcal{R}}} R_{\mathcal{R}}$$





$$p(u_{\mathcal{H}}|x, \theta, u_{\mathcal{R}}) \propto \exp(R_{\mathcal{H}}(x, u_{\mathcal{H}}, \theta, u_{\mathcal{R}}))$$



$$b_{t+1}(\theta) \propto b_t(\theta) \cdot p(u_{\mathcal{H}}|x_t, \theta, u_{\mathcal{R}})$$

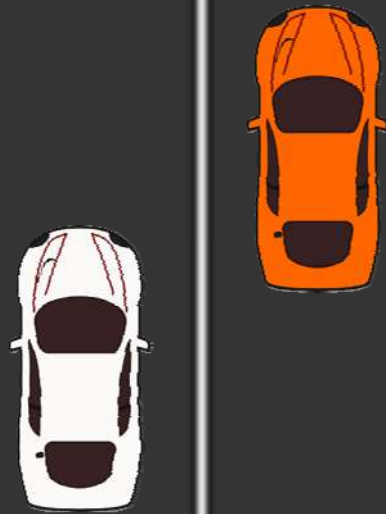
Info Gathering

$$R_{\mathcal{R}}(x, u_{\mathcal{H}}, \theta, u_{\mathcal{R}}) = \underbrace{\mathbb{H}(b_t) - \mathbb{H}(b_{t+1})}_{\text{Info Gathering}} + \underbrace{\lambda \cdot R_{goal}(x, u_{\mathcal{H}}, \theta, u_{\mathcal{R}})}_{\text{Goal}}$$

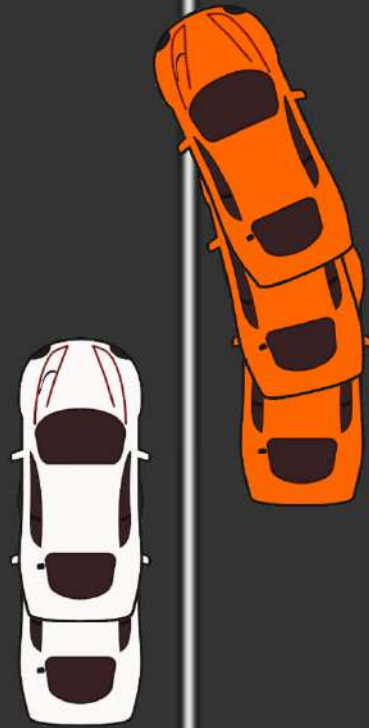
Goal

$$u_{\mathcal{R}} = \operatorname{argmax}_{u_{\mathcal{R}}} \mathbb{E}_{\theta} [R_{\mathcal{R}}]$$

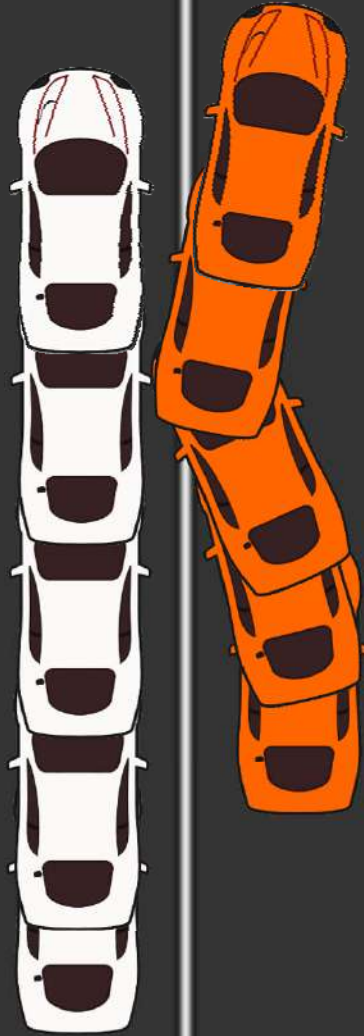
Nudging in for Active Info Gathering



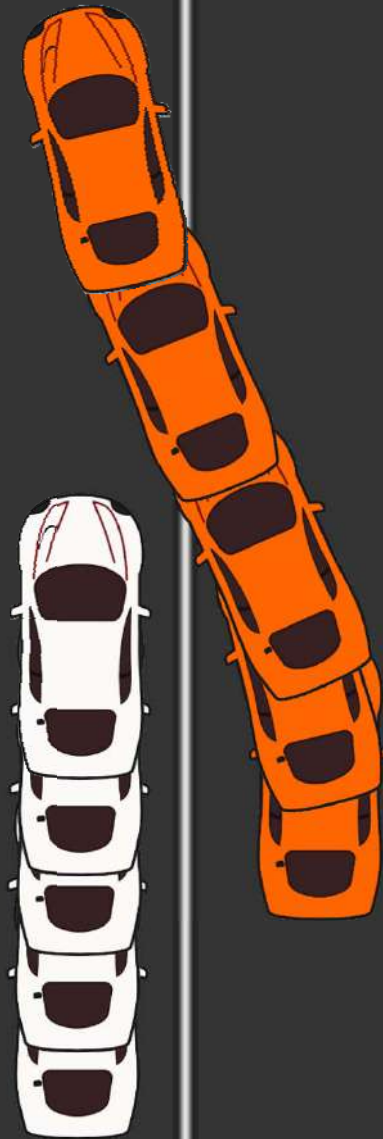
Nudging in for Active Info Gathering



Distracted Human Driver



Attentive Human Driver



Key Ideas:

Actively gather data about the environment (human) by *affecting* the environment's behavior

Learn environment (human) model from data, update online

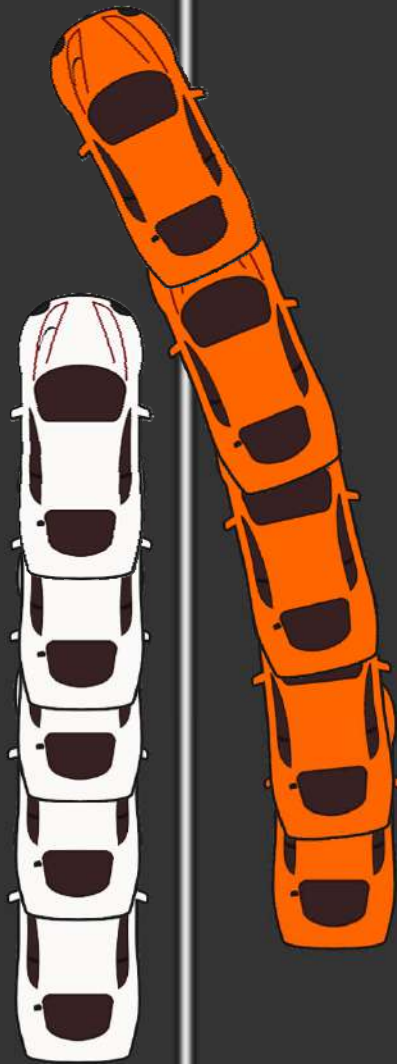
Questions:

- How to verify such human-robot systems?
- What are more realistic human models? (e.g. “bounded rationality”)

Talk Outline

- Environment Modeling Challenge
 - Interaction-Aware Control for Human-CPS
- Specification (& Verification) Challenge
 - Verifying Robustness (of Interaction-Aware Controller)
 - Falsification for Deep Learning based CPS
- Conclusions and Future Directions
 - Towards a New Design Methodology for AI-based Systems

More Efficient, but is it Safe?



Verifying Temporal Logic Requirements

Signal Temporal Logic (STL) [Maler & Nickovic, '04]

Predicates over continuous signals, Propositional Formulas φ (\wedge, \vee, \neg of the predicates), Temporal Operators (G, F, X, U), real-time interval τ .

$G_{\tau} \varphi$	φ is true at all future moments in τ .
$F_{\tau} \varphi$	φ is true in some future moment in τ .
$\varphi_1 U_{\tau} \varphi_2$	φ_1 is true until φ_2 becomes true in τ .

Safety (invariance): Vehicle maintains specified distance from obstacles.

$$G_{[0,\tau]} [\text{dist}(\text{vehicle}, \text{obstacle}) > \Delta]$$

From Logical Formulas to Objective Functions

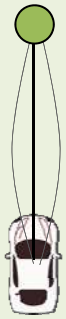
- STL formula has both
 - Boolean semantics: true/false
 - Quantitative semantics: value in \mathbb{R}

- Example:

$$G_{[0,\tau]}(\text{dist}(\text{vehicle}, \text{obstacle}) > \Delta)$$



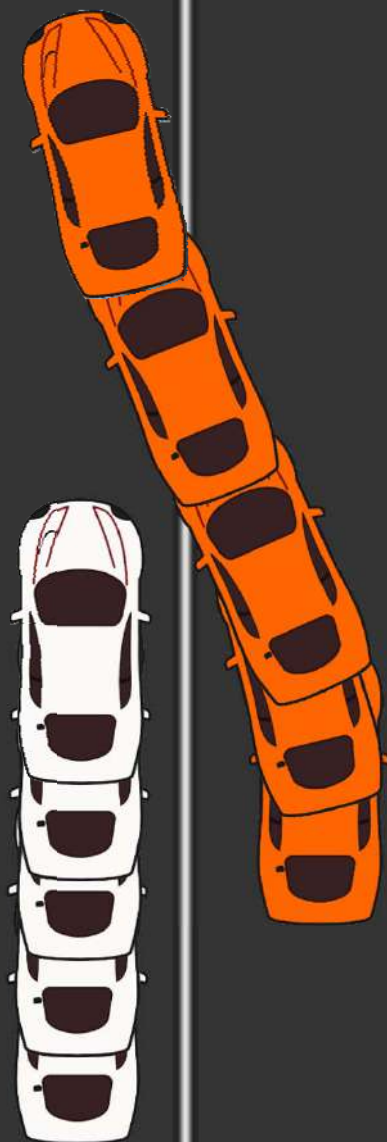
$$\inf_{[0,\tau]} [\text{dist}(\text{vehicle}, \text{obstacle}) - \Delta]$$



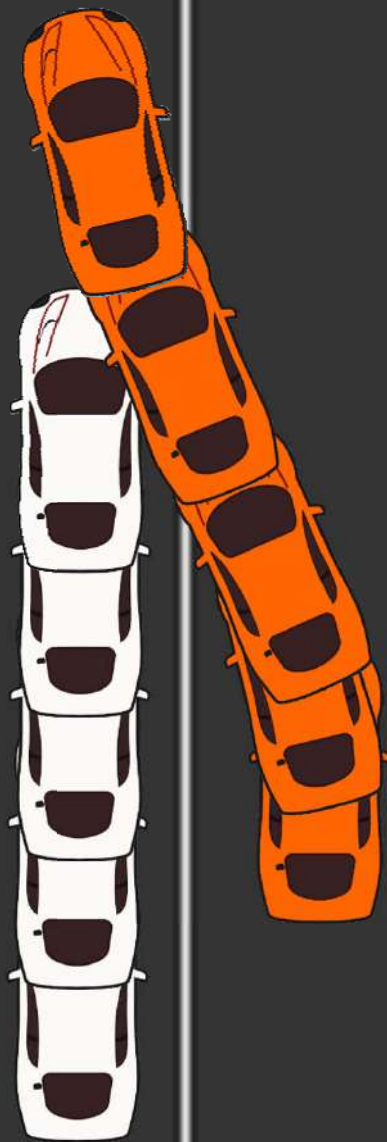
$$u_{\mathcal{H}}^* = \operatorname{argmax}_{u_{\mathcal{H}}} R_{\mathcal{H}}(x, u_{\mathcal{R}}, u_{\mathcal{H}}^*)$$



$$u_{\mathcal{R}}^* = \operatorname{argmax}_{u_{\mathcal{R}}} R_{\mathcal{R}}(x, u_{\mathcal{R}}, u_{\mathcal{H}}^*(x, u_{\mathcal{R}}))$$



$$R_{\mathcal{H}}(x, u_{\mathcal{H}}, u_{\mathcal{R}})$$



$$R_{\mathcal{H}}^{\dagger}(x, u_{\mathcal{H}}, u_{\mathcal{R}})$$

$$|R_{\mathcal{H}}^{\dagger} - R_{\mathcal{H}}| < \delta$$

How **robust** is the learning-based controller?

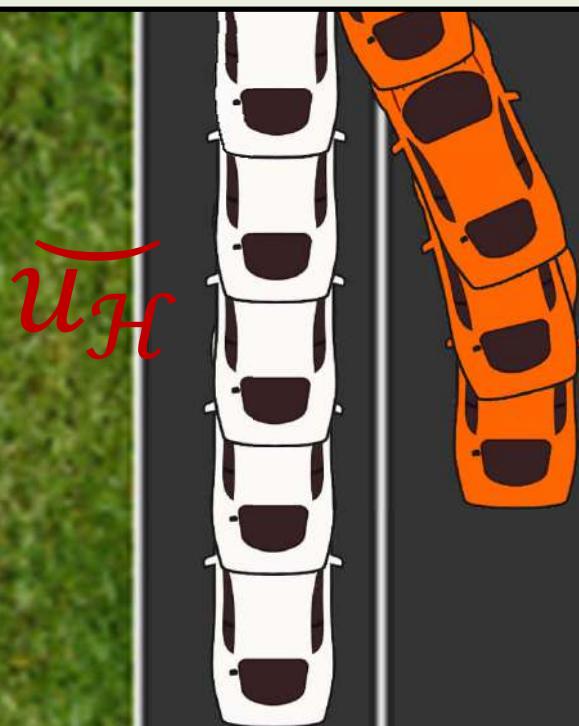
How to algorithmically find **falsifying** actions by the human?

$$\widetilde{u}_{\mathcal{H}} = \arg \min_{u_{\mathcal{H}}} R_{\mathcal{R}}(x, u_{\mathcal{R}}^*, u_{\mathcal{H}}) \quad \text{Falsifying actions}$$

$$\text{s. t. } \exists R_{\mathcal{H}}^{\dagger} : u_{\mathcal{H}} = \arg \max_{\widehat{u}_{\mathcal{H}}} R_{\mathcal{H}}^{\dagger}(x, u_{\mathcal{R}}^*, \widehat{u}_{\mathcal{H}})$$

$$|R_{\mathcal{H}}^{\dagger} - R_{\mathcal{H}}| < \delta$$

Optimizing a perturbed version of the learned reward function.



Theorem:

$$\begin{aligned} \widetilde{u}_{\mathcal{H}} &= \arg \min_{u_{\mathcal{H}}} R_{\mathcal{R}}(x, u_{\mathcal{R}}^*, u_{\mathcal{H}}) & |R_{\mathcal{H}}^{\dagger} - R_{\mathcal{H}}| < \delta \\ \text{s. t. } \exists R_{\mathcal{H}}^{\dagger} : u_{\mathcal{H}} &= \arg \max_{\widehat{u}_{\mathcal{H}}} R_{\mathcal{H}}^{\dagger}(x, u_{\mathcal{R}}^*, \widehat{u}_{\mathcal{H}}) \end{aligned}$$

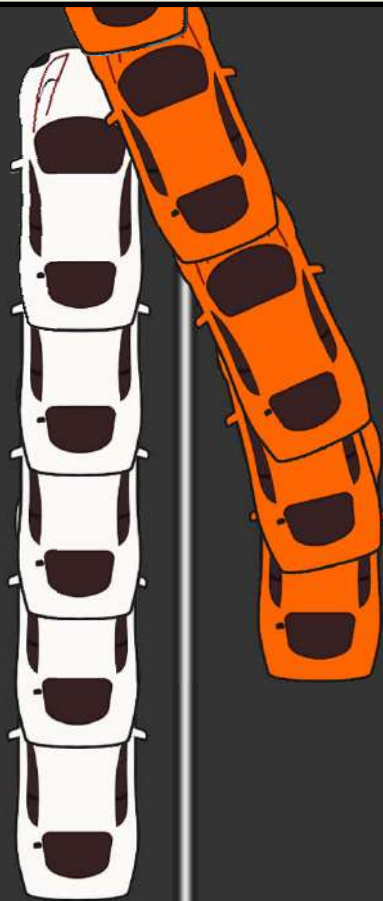


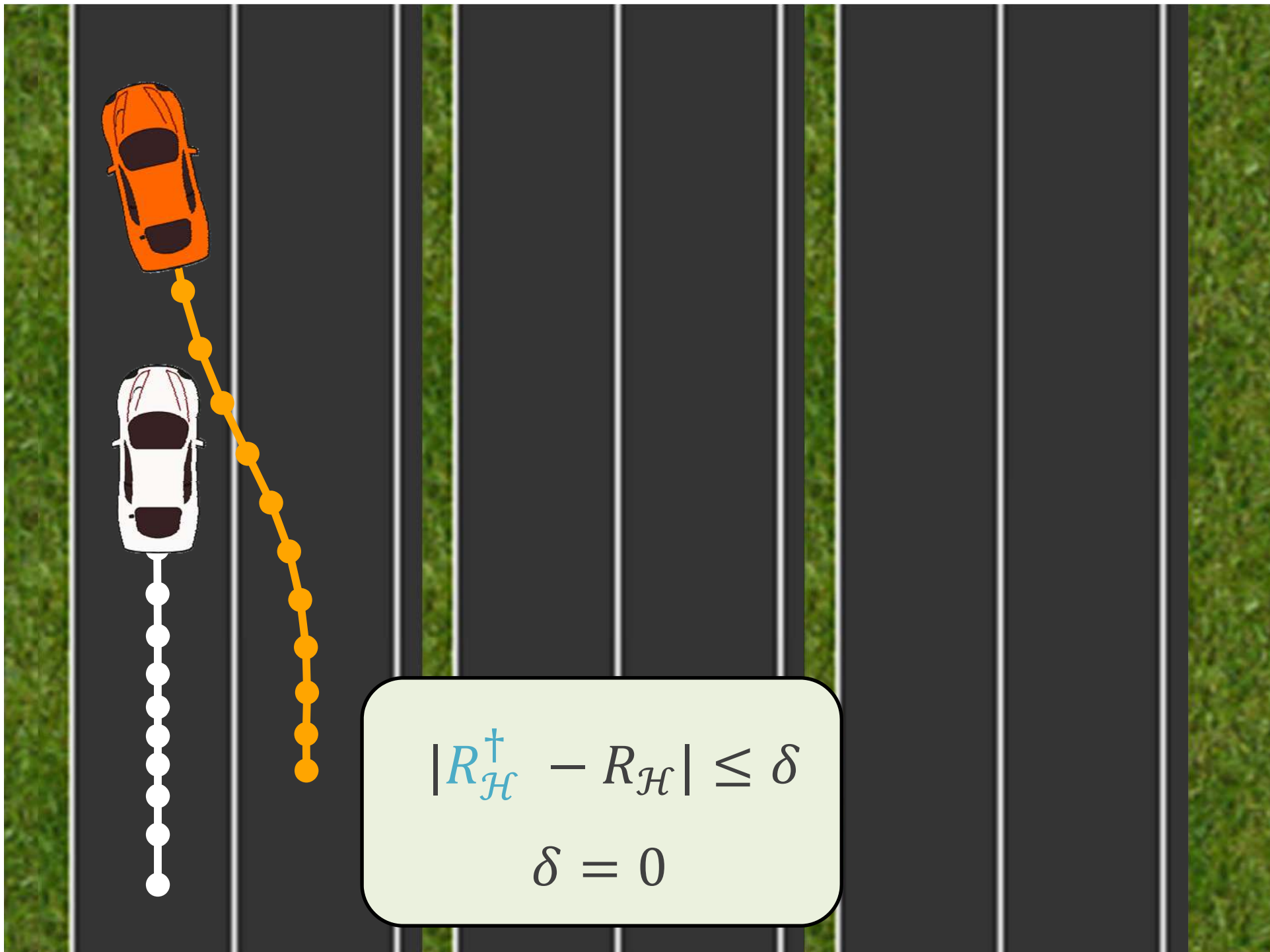
Reduction

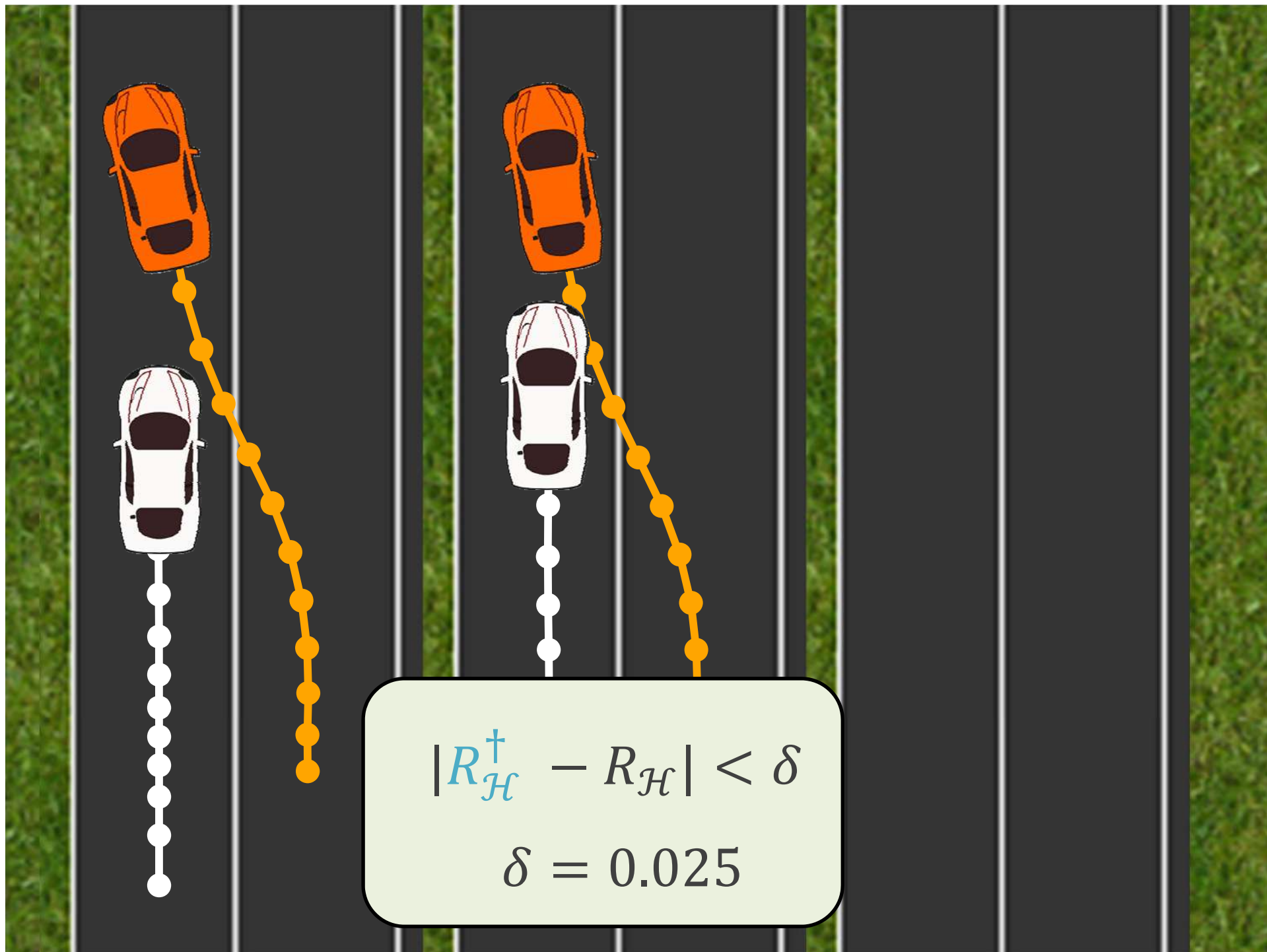
$$\begin{aligned} \widetilde{u}_{\mathcal{H}} &= \arg \min_{u_{\mathcal{H}}} R_{\mathcal{R}}(x, u_{\mathcal{R}}^*, u_{\mathcal{H}}) \\ \text{s. t. } R_{\mathcal{H}}(x, u_{\mathcal{R}}^*, u_{\mathcal{H}}) &> R_{\mathcal{H}}(x, u_{\mathcal{R}}^*, u_{\mathcal{H}}^*) - 2\delta \end{aligned}$$

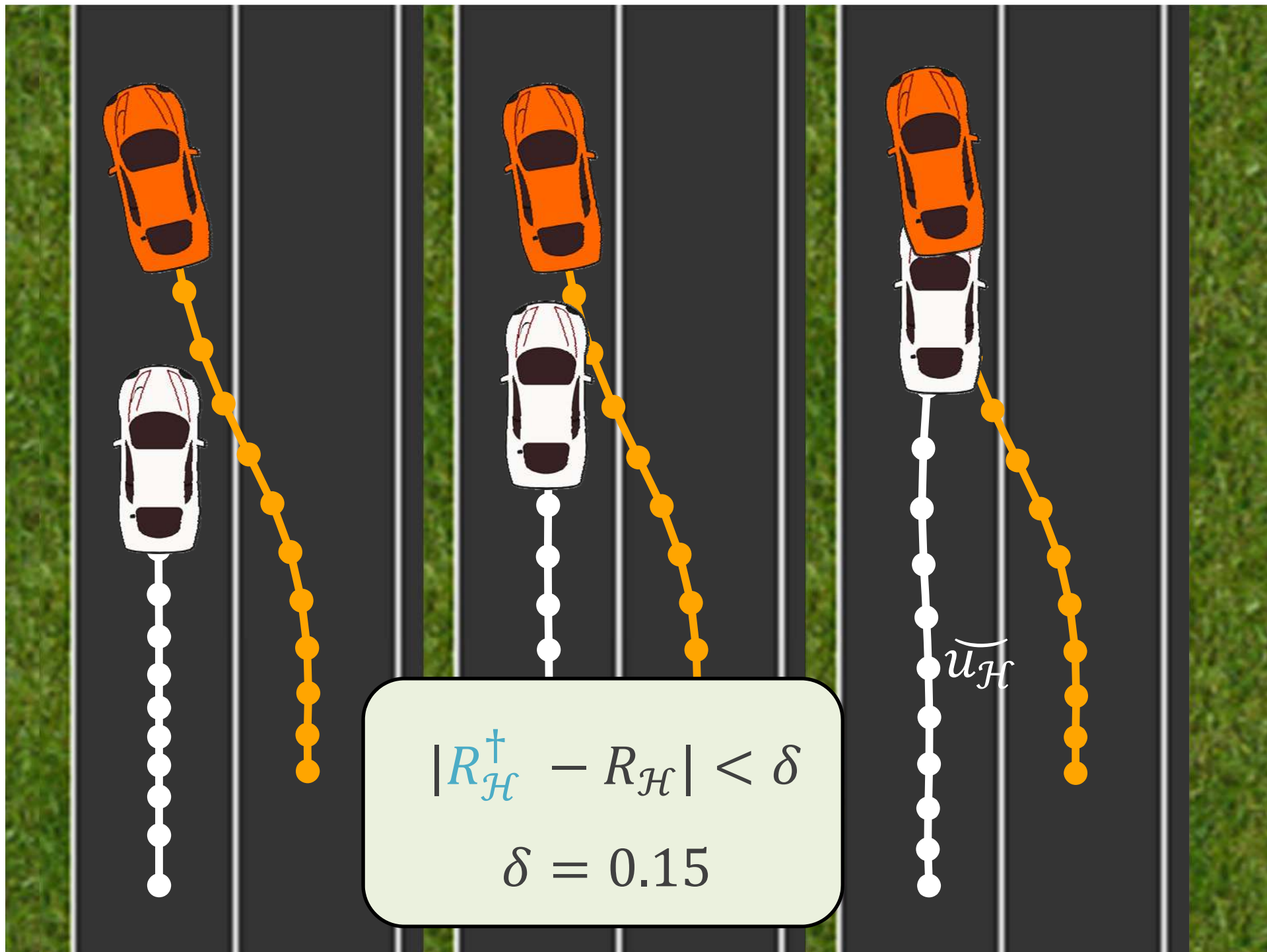
$$\widetilde{u}_{\mathcal{H}} = \arg \min_{u_{\mathcal{H}}} R_{\mathcal{R}}(x, u_{\mathcal{R}}^*, u_{\mathcal{H}})$$

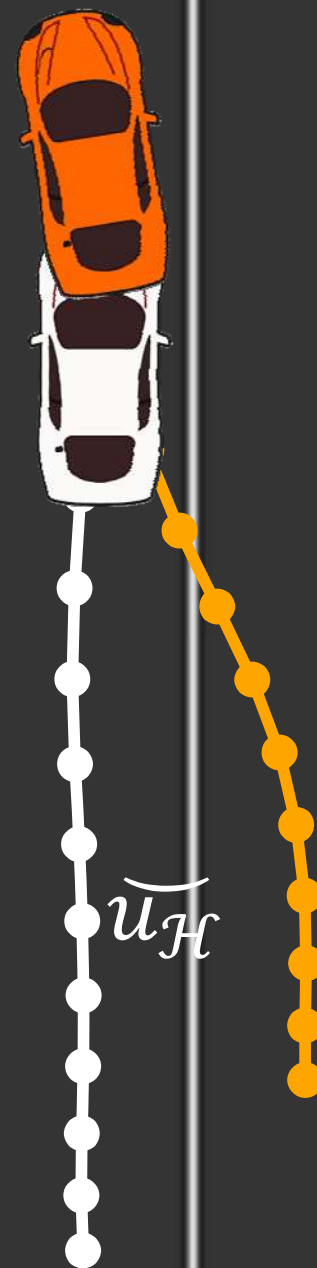
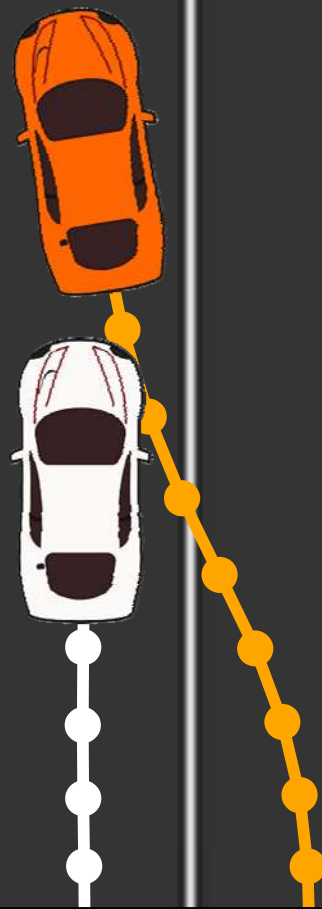
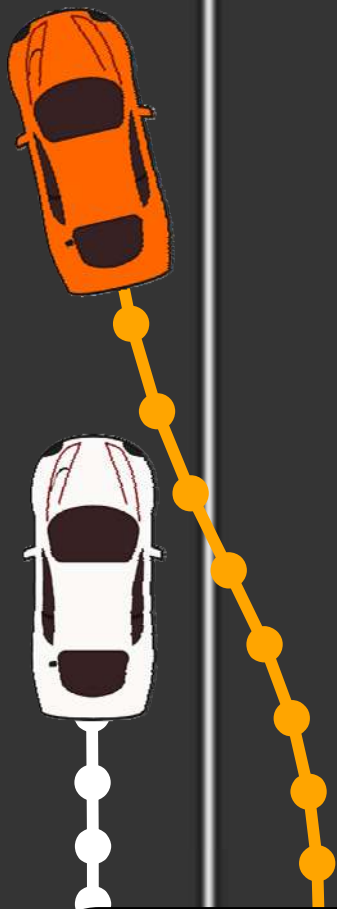
$$\text{s. t. } R_{\mathcal{H}}(x, u_{\mathcal{R}}^*, u_{\mathcal{H}}) > R_{\mathcal{H}}(x, u_{\mathcal{R}}^*, u_{\mathcal{H}}^*) - 2\delta$$











$$\Pr(|\textcolor{teal}{R}_{\mathcal{H}}^{\dagger} - R_{\mathcal{H}}| < \delta) > 0.9$$

$$\delta = 0.15$$

Key Ideas:

Turn Verification (falsification) into **Optimization**

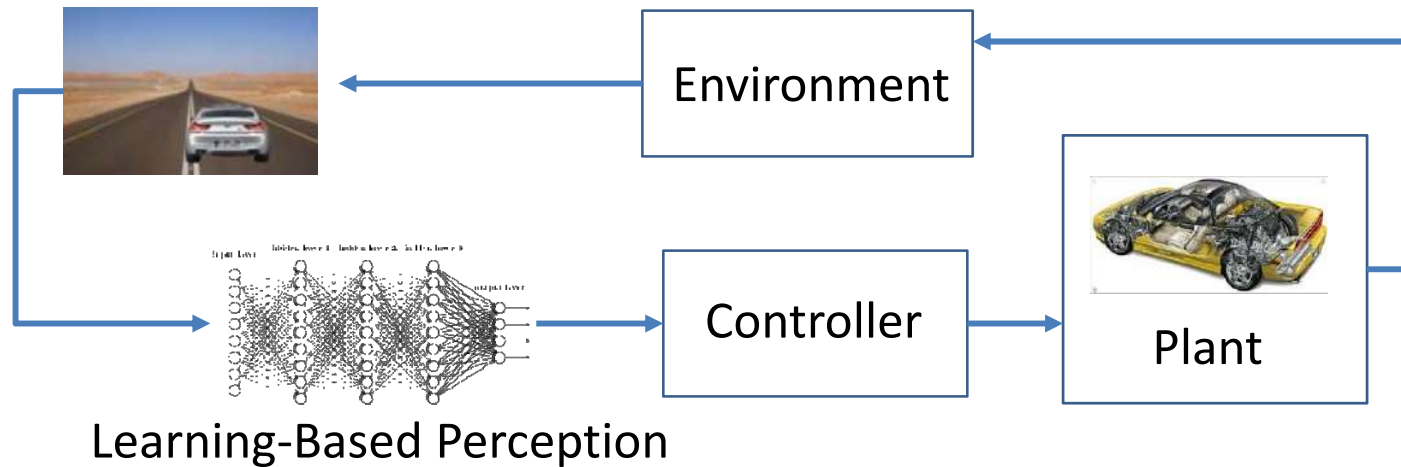
Important Property: **Robustness** of AI/Learning-based system to small perturbations in data/learned function

**Challenge: Specification, Verification,
Training/Testing for Learning Systems**

Falsification of Cyber-Physical Systems with Machine Learning Components

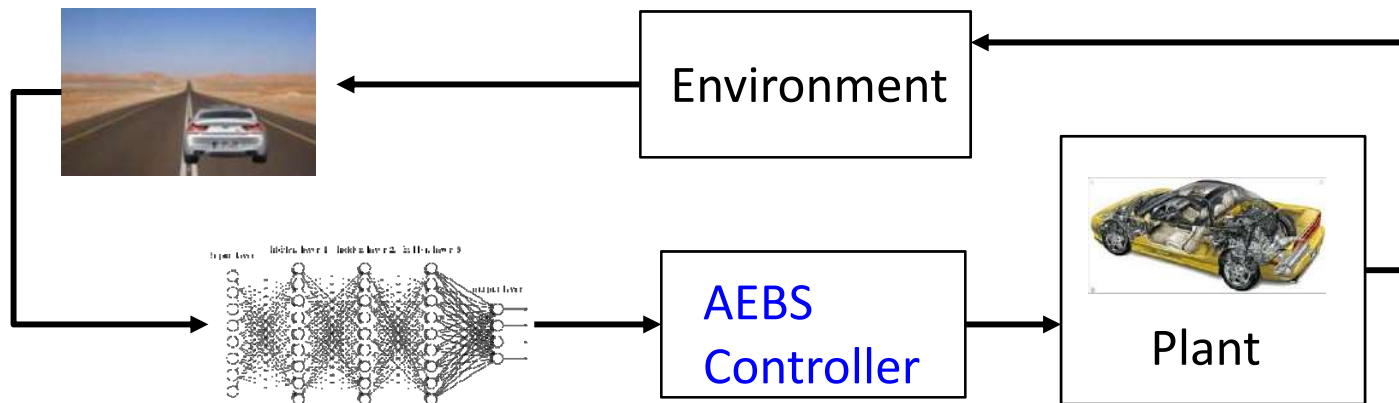
T. Dreossi, A. Donze, and S. A. Seshia. *Compositional Falsification of Cyber-Physical Systems with Machine Learning Components*, In NASA Formal Methods Symposium, May 2017.

Problem: Verify Automotive System (CPS) that uses ML-based Perception



- Focus:
 - **Falsification**: finding scenarios that violate safety properties
 - **Test (Data) Generation**: generate “interesting” data for training / testing → improve accuracy
 - **Deep Neural Networks**, given the increasing interest and use in the automotive context.

Automatic Emergency Braking System (AEBS)

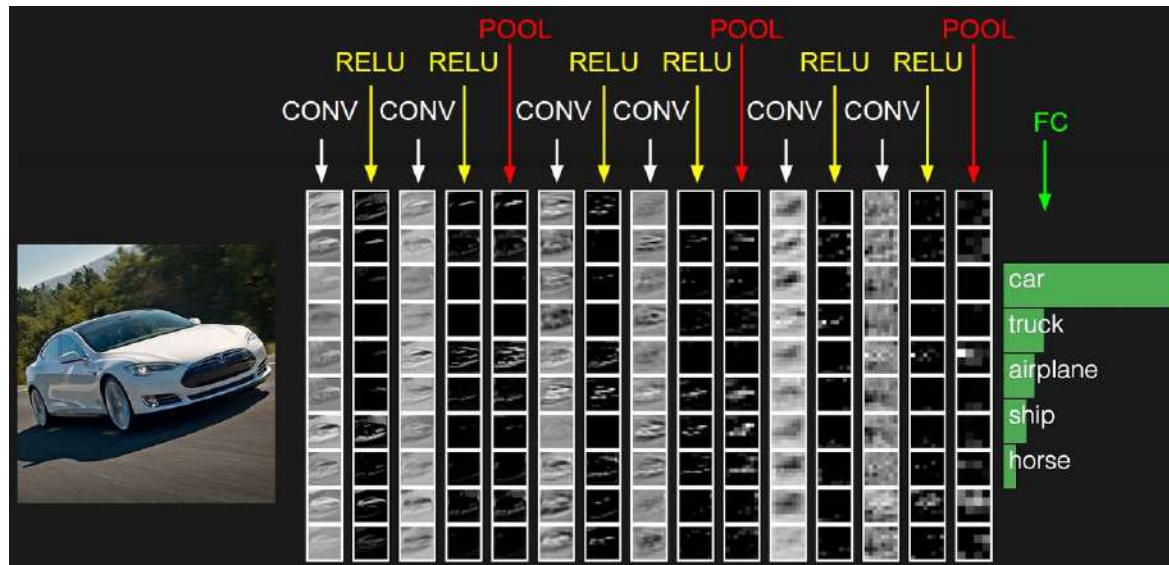


Deep Learning-Based Object Detection

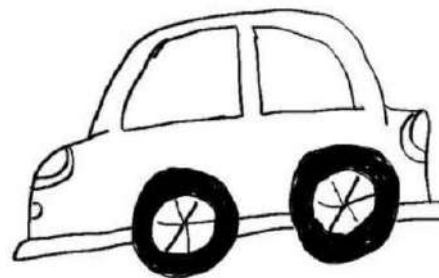
- Goal: Brake when an obstacle is near, to maintain a minimum safety distance
 - Controller, Plant, Env models in Matlab/Simulink
- Object detection/classification system based on deep neural networks
 - Inception-v3, AlexNet, ... trained on ImageNet

What's the Specification for Perception Tasks?

Convolutional Neural Network trained to recognize cars



How do you formally specify “a car”?



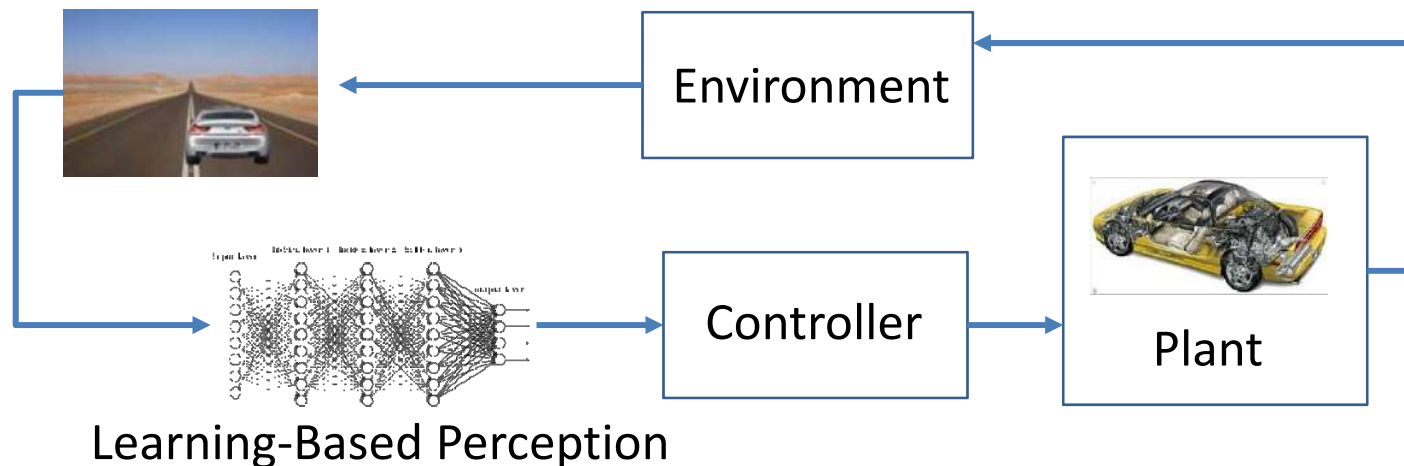
Idea: Use a **System-Level** Specification

✗ “Verify the Deep Neural Network”

✓ “Verify the System containing the Deep Neural Network”

Formally Specify the *End-to-End Behavior* of the System

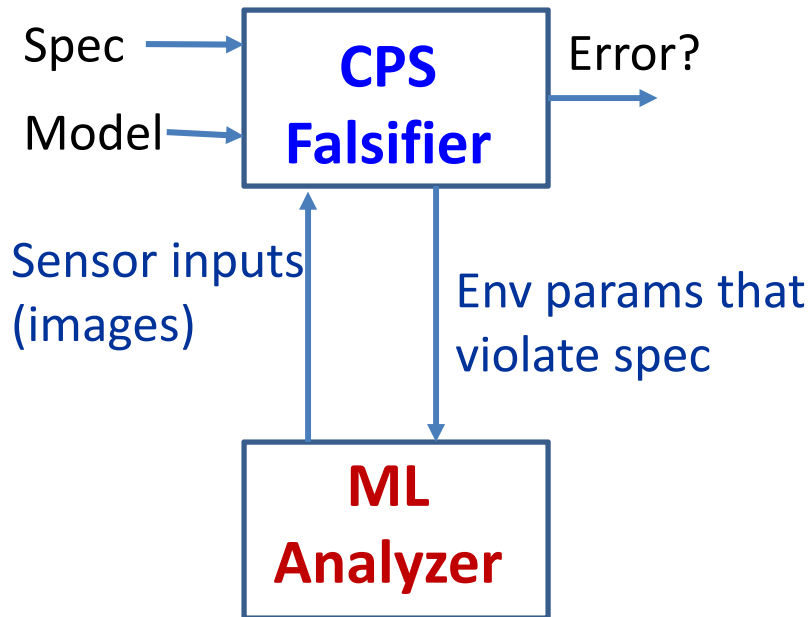
STL Formula: **G** ($dist(ego\ vehicle, env\ object) > \Delta$)



Tool: Simulation-Based Falsification of Signal Temporal Logic for CPS

- STL has quantitative semantics
 - Logical formula \rightarrow Cost Function ρ
 - Quantifies “how much” a trace satisfies a property
- *Advantage:* Finding a bug (property violation) corresponds to minimizing the function ρ and checking if the value falls below 0.
 - This view of “verification as optimization” underlies the Breach toolkit and similar tools

Our Approach: Combine Temporal Logic CPS Falsifier with ML Analyzer

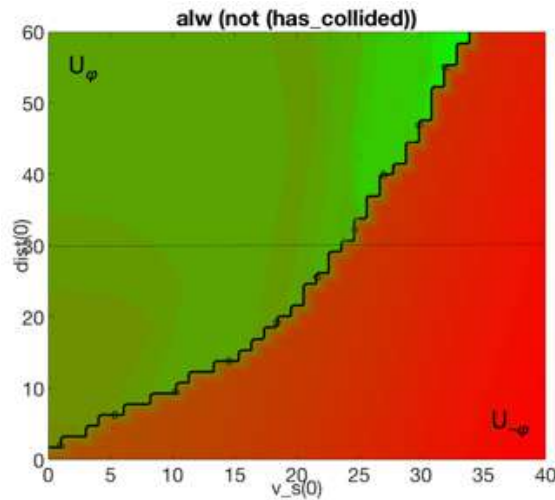


- CPS Falsifier uses **abstraction** of ML component
 - **Optimistic analysis**: assume ML classifier is always correct
 - **Pessimistic analysis**: assume classifier is always wrong
- Difference is the **region of interest** where output of the ML component “matters”

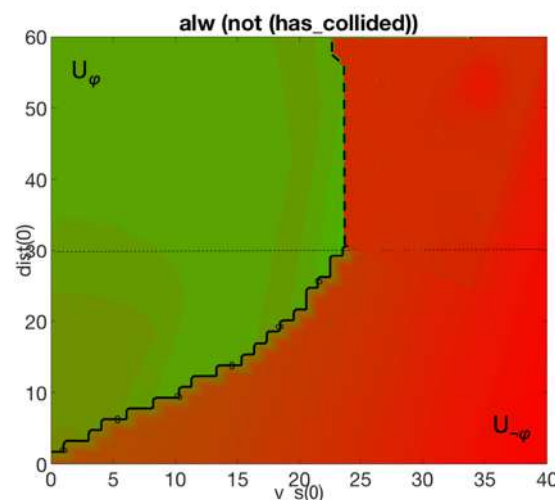
Compositional:

CPS Falsifier and ML Analyzer can be designed and run independently (& communicate)!

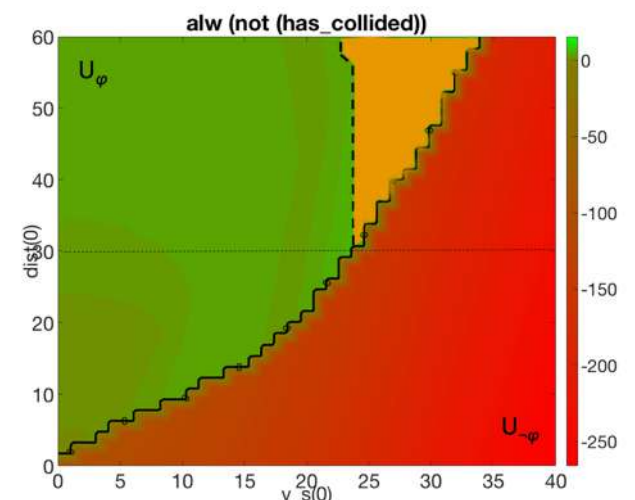
Identifying Region of Interest for Automatic Emergency Braking System



ML always correct



ML always wrong

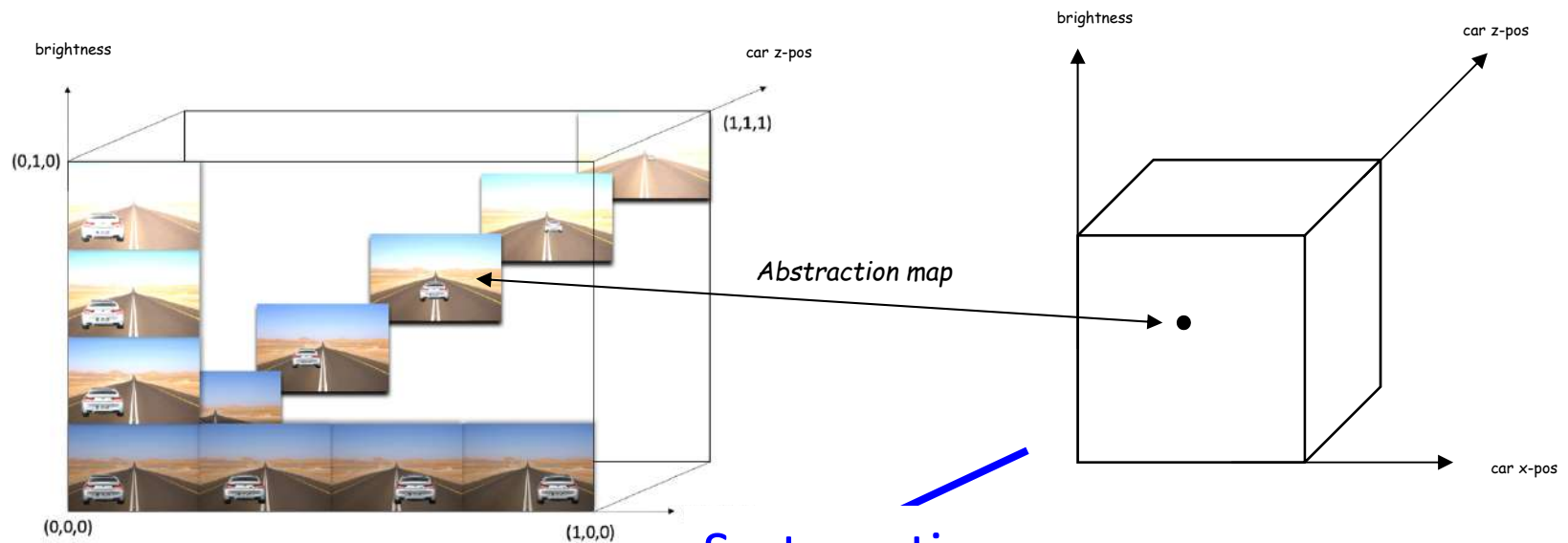


*Potentially unsafe region
depending on ML
component (yellow)*

Perform Optimistic and Pessimistic Analyses on the Deep Neural Network

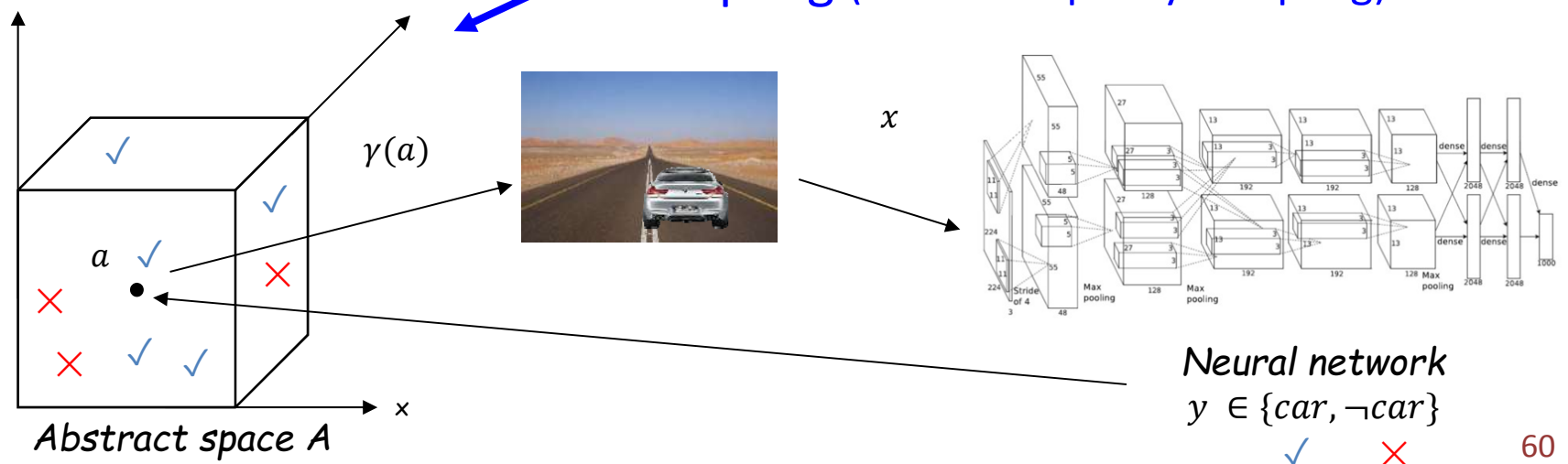
Machine Learning Analyzer

Systematically Explore Region of Interest in the Image (Sensor) Space



Feature space \tilde{X}

Systematic
Sampling (low-discrepancy sampling)

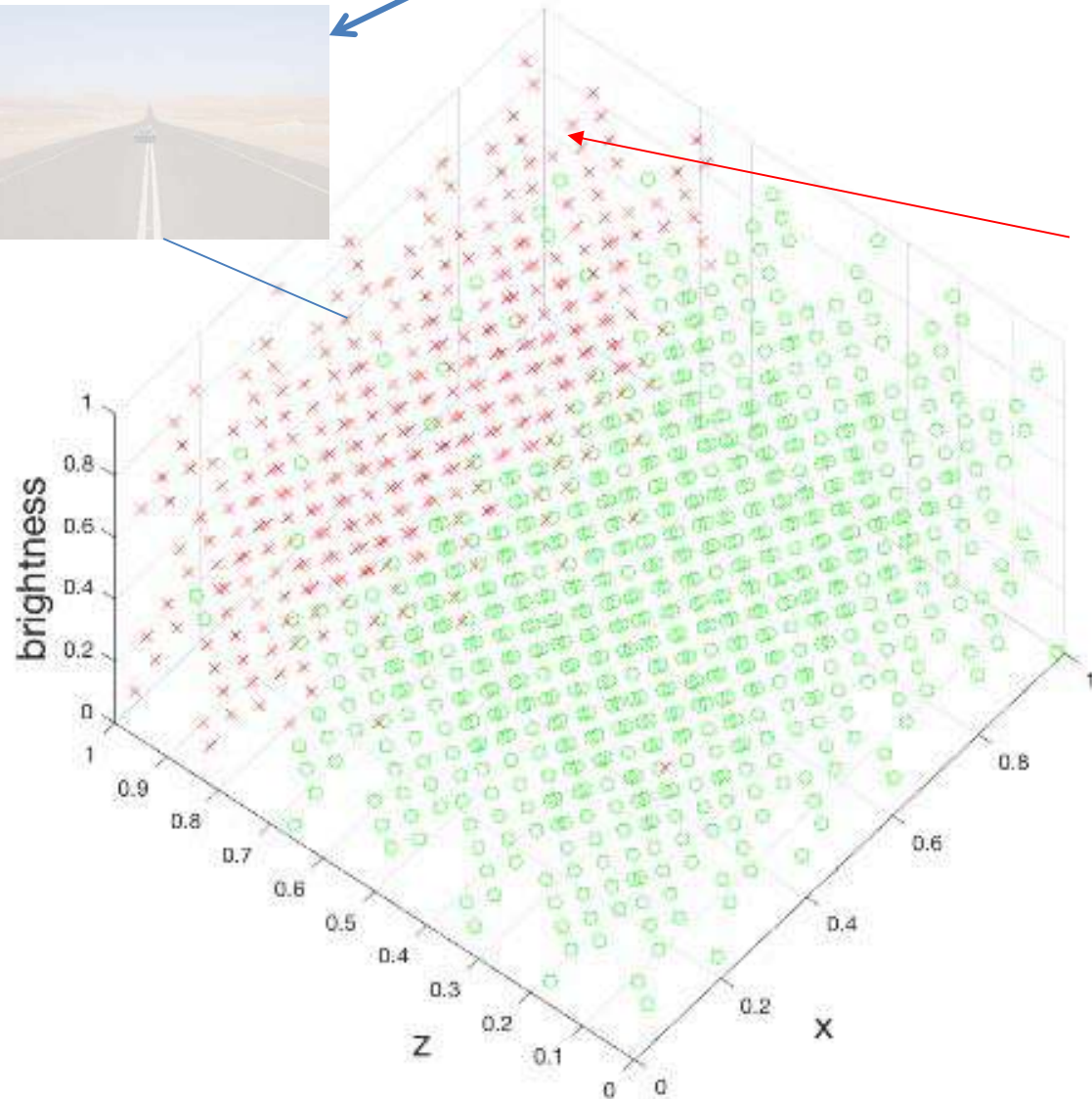


Sample Result



This misclassification may not be of concern

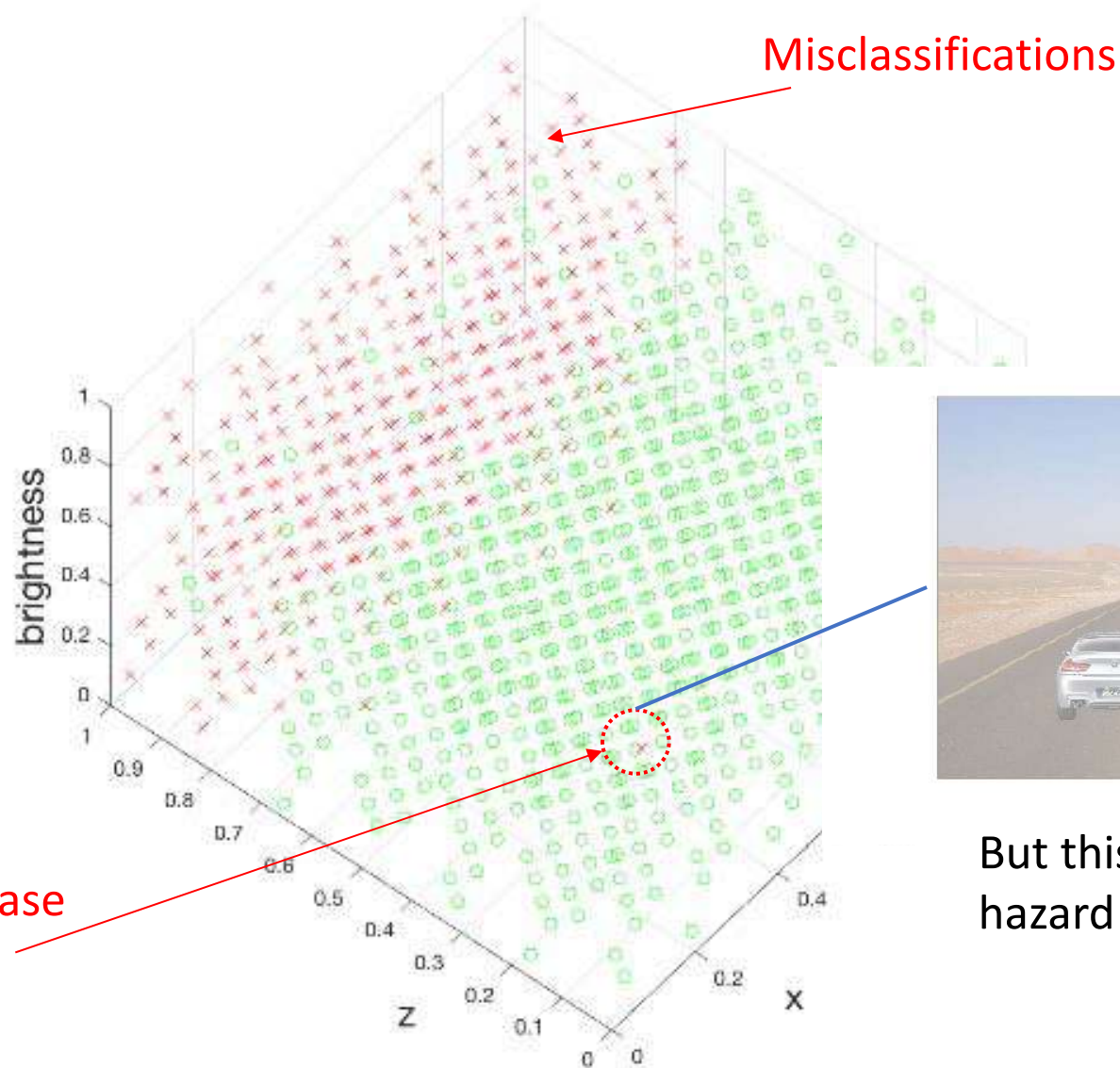
Inception-v3
Neural
Network
(pre-trained on
ImageNet using
TensorFlow)



Misclassifications

Sample Result

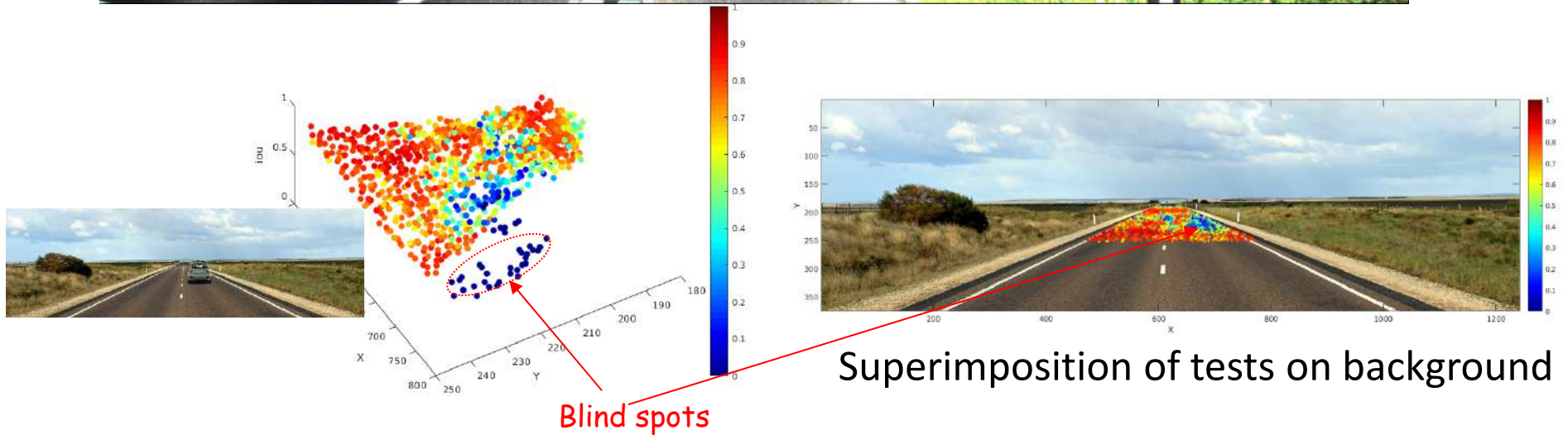
Inception-v3
Neural
Network
(pre-trained on
ImageNet using
TensorFlow)



But this one is a real hazard!

Newer Results

[Dreossi, Ghosh, et al., ICML 2017 workshop]



Summary of Key ideas

- Generate adversarial examples that violate *system-level specification*
- *Compositional* Approach blends the strengths of the CPS Falsifier with a Machine Learning Analyzer
- Counterexample images can be added to the training set to improve ML accuracy (“right” data vs. “big” data)
- Ongoing/Future Work:
 - Improving ML analyzer
 - New benchmarks (datasets and networks)
 - Evaluating training/test accuracy improvements

Concluding Thoughts

Towards Verified Artificial Intelligence

Challenges

Principles

1. Environment (incl. Human) Modeling	→	Data-Driven, Introspective Environment Modeling
2. Specification	→	System-Level Specification; Robustness/Quantitative Spec.
3. Learning Systems Complexity	→	Abstract & Explain
4. Efficient Training, Testing, Verification	→	Verification-Guided, Adversarial Analysis and Improvisation
5. Design for Correctness	→	Formal Inductive Synthesis

S. A. Seshia, D. Sadigh, S. S. Sastry. *Towards Verified Artificial Intelligence*. July 2016. <https://arxiv.org/abs/1606.08514>.

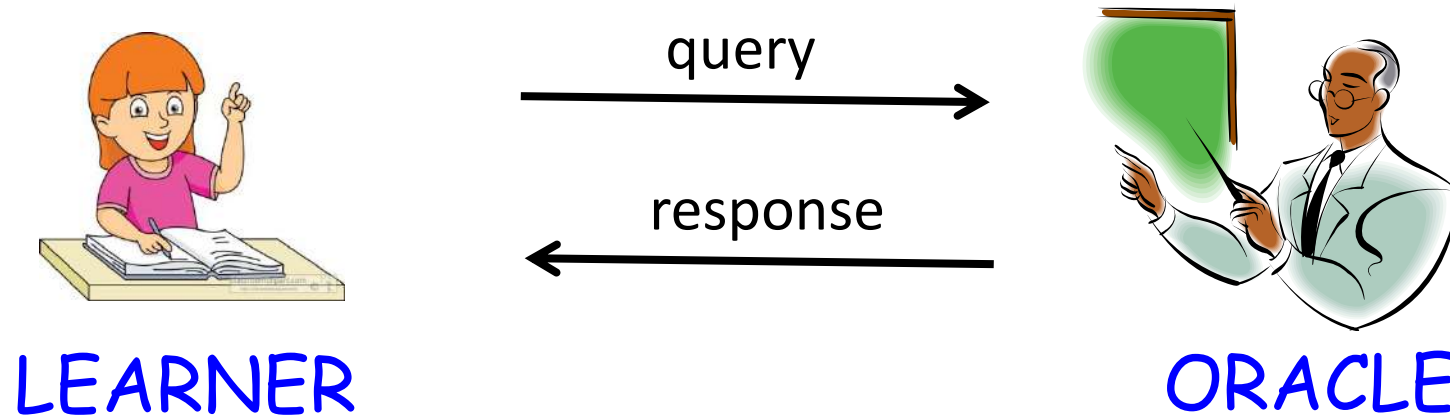
Correct-by-Construction Design with Formal Inductive Synthesis

Inductive Synthesis: Learning from Examples (ML)

Formal Inductive Synthesis: Learn from Examples *while satisfying a Formal Specification*

Key Idea: **Oracle-Guided Learning**

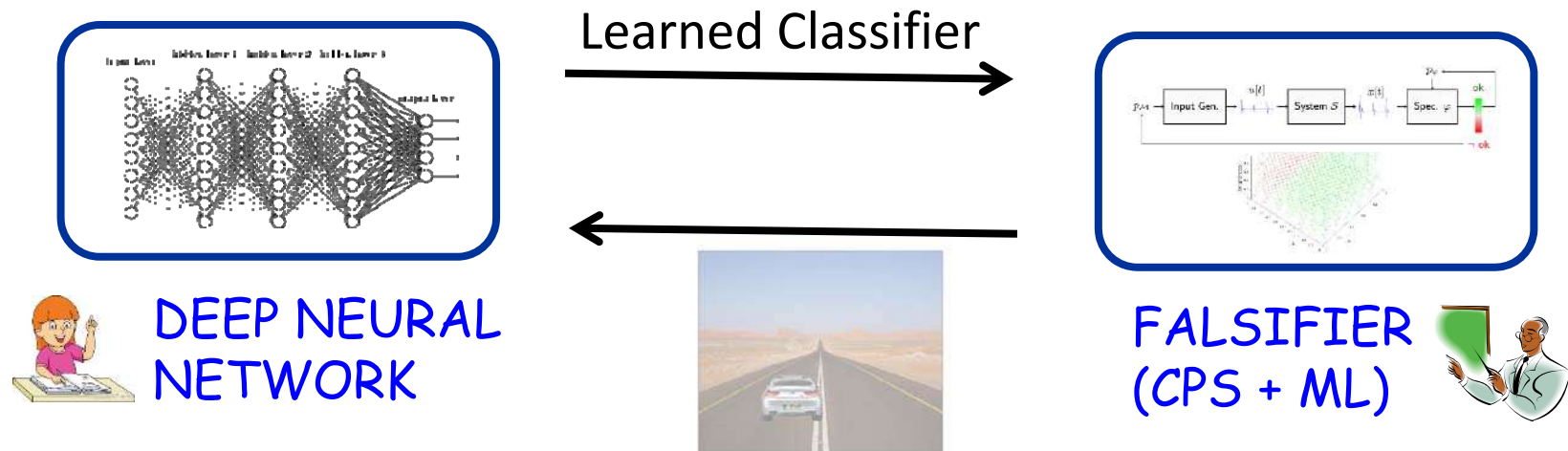
Combine Learner with Oracle (e.g., Verifier) that answers Learner's Queries



[Jha & Seshia, “A Theory of Formal Synthesis via Inductive Learning”, 2015, Acta Informatica 2017.]

Verifier-Guided Training of Deep Neural Networks

- Instance of Oracle-Guided Inductive Synthesis
- Oracle is Verifier (CPSML Falsifier) used to perform counterexample-guided training of DNNs
- Substantially increase accuracy with only few additional examples



Towards Verified Artificial Intelligence

Challenges

1. Environment (incl. Human) Modeling
2. Specification
3. Learning Systems Complexity
4. Efficient Training, Testing, Verification
5. Design for Correctness

Principles

- Data-Driven, Introspective Environment Modeling
- System-Level Specification; Robustness/Quantitative Spec.
- Abstract & Explain
- Verification-Guided, Adversarial Analysis and Improvisation
- Formal Inductive Synthesis

Exciting Times Ahead!!! Thank you!