# CPS: Frontier: VeHICaL: Verified Human Interfaces, Control, and Learning for Semi-Autonomous Systems

S. A. Seshia, R. Bajcsy, B. Hartmann, S. S. Sastry, C. Tomlin (Berkeley),
R. Murray (Caltech), T. Griffiths (Princeton), C. Sturton (UNC Chapel Hill)

VeHICaL  http://vehical.org

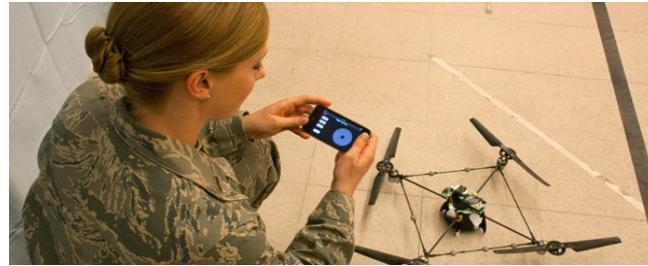Caltech  Berkeley UNIVERSITY OF CALIFORNIA  PRINCETON UNIVERSITY  THE UNIVERSITY of NORTH CAROLINA at CHAPEL HILL

# Design of Human Cyber-Physical Systems (h-CPS)

## CPS that operate in concert with humans



Semi-Autonomous Driving

UAVs with Human Operators

Robotic Surgery & Medicine

...and other applications.

*Project Goal: To develop a **science of verified co-design** of controllers for semi-autonomous cyber-physical systems and interfaces between humans and cyber-physical components*
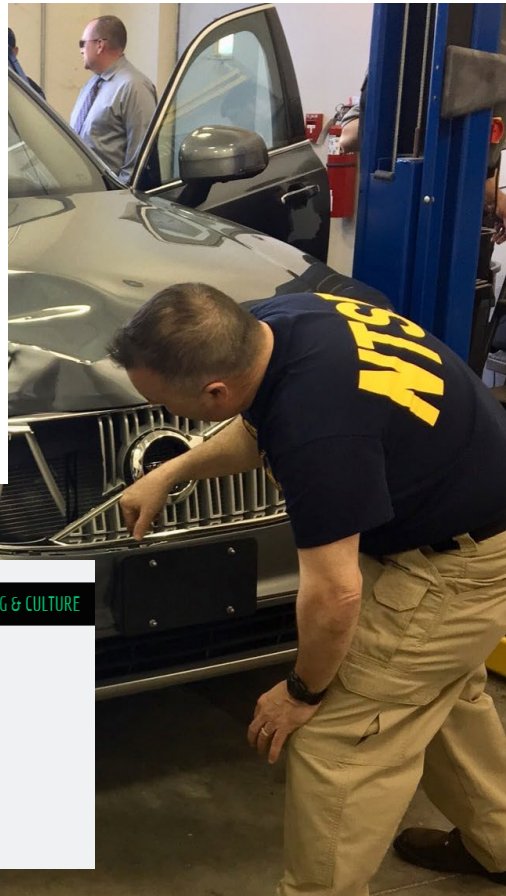
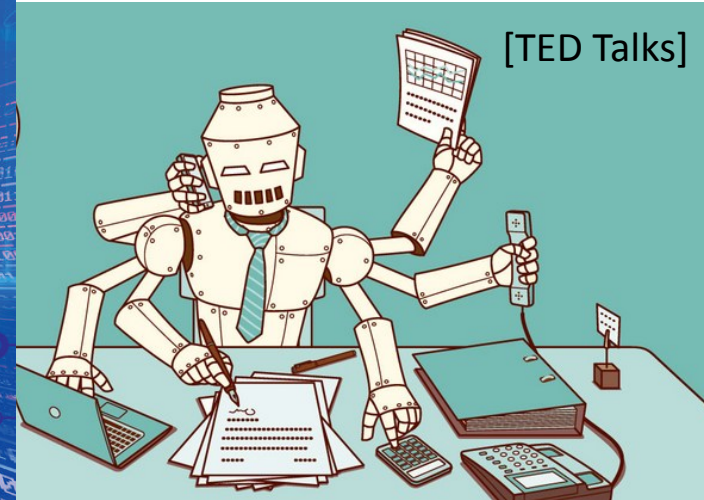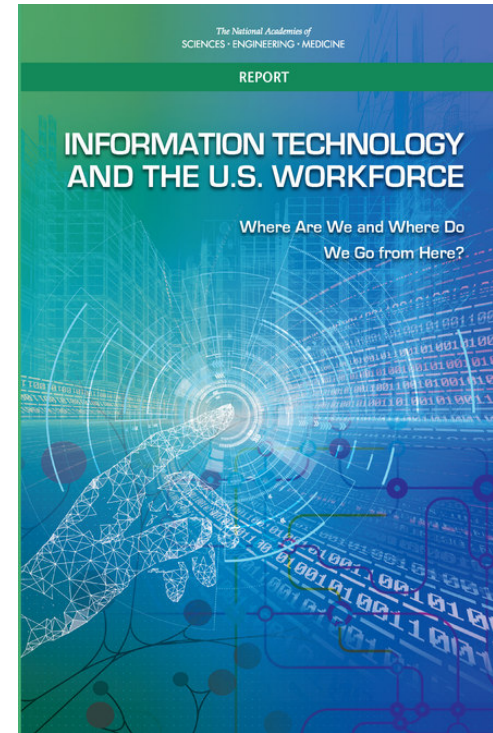# Why is this Important?

## SAFETY-CRITICAL & MISSION-CRITICAL

**Tesla driver dies in first fatal autonomous car crash in US**

Hands-off driving faces tough questions
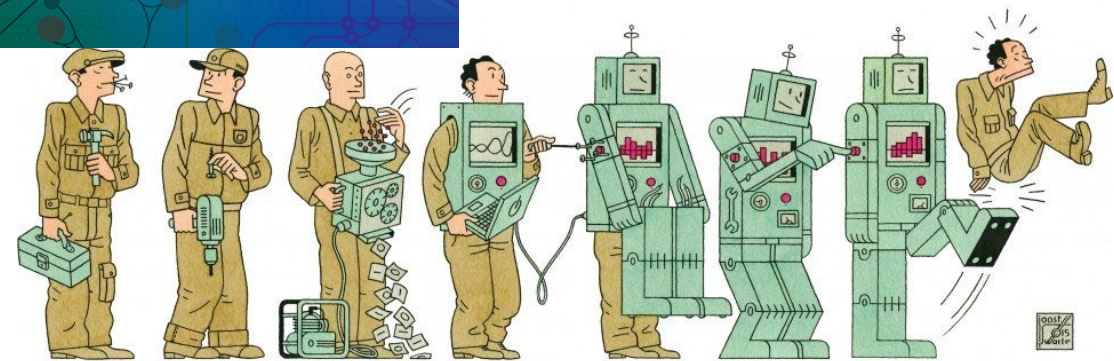Beck Diefenbach/Reuters

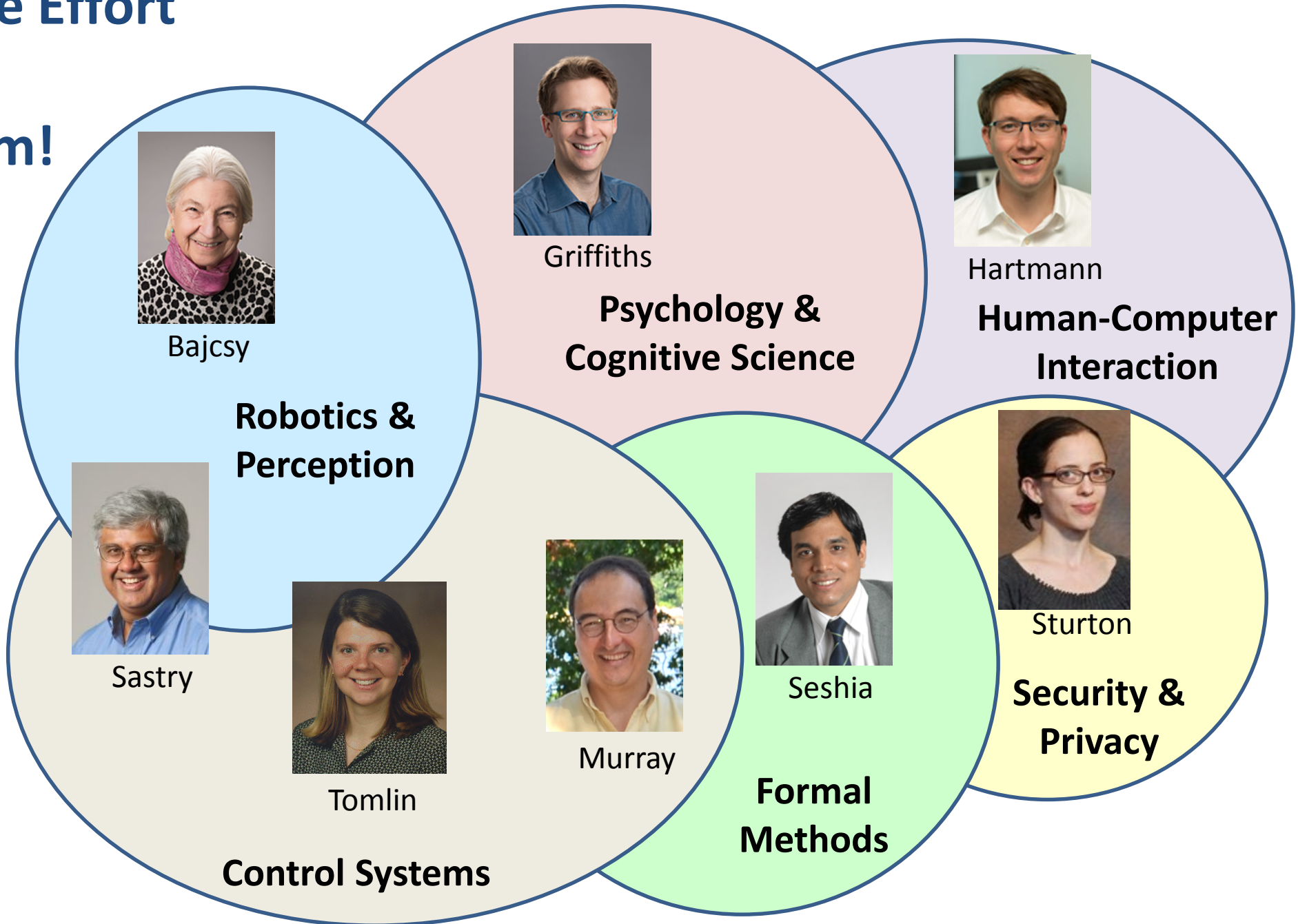**ars TECHNICA**   BIZ & IT   TECH   SCIENCE   POLICY   CARS   GAMING & CULTURE

*DRIVERLESS CAR SAFETY —*
Report: Software bug led to death in Uber's self-driving crash

Sensors detected Elaine Herzberg, but software reportedly decided to ignore her.

TIMOTHY B. LEE - 5/7/2018, 3:12 PM

## IMPACT OF AUTOMATION ON WORK/JOBS

The National Academies of
SCIENCES · ENGINEERING · MEDICINE
REPORT

**INFORMATION TECHNOLOGY AND THE U.S. WORKFORCE**

Where Are We and Where Do
We Go from Here?

[TED Talks]

[MIT Technology Review]

3

# Key Envisioned Contributions to CPS Science

- Developing a Science of Co-Design of Human Interfaces and Control
  - Turning design of h-CPS from an art to a science by systematic design and verification of human interfaces
- Making Uncertainty a first-class citizen in Verification and Control
  - New algorithms and models to deal with uncertainty in CPS dynamics and CPS design
- Bridging the Schism between Model-Based Design and Data-Driven Methods
  - A new design methodology for CPS that blends data-driven learning with formal modeling and proof engines

# Design for
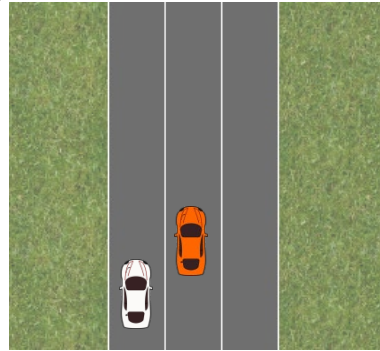# Effective Communication between Humans and Automation
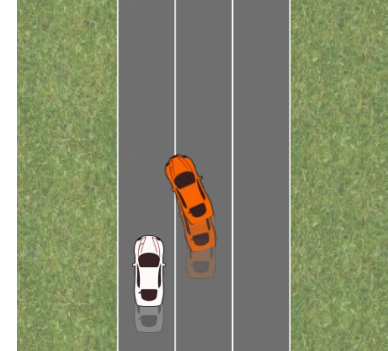
# Interaction-Aware Control

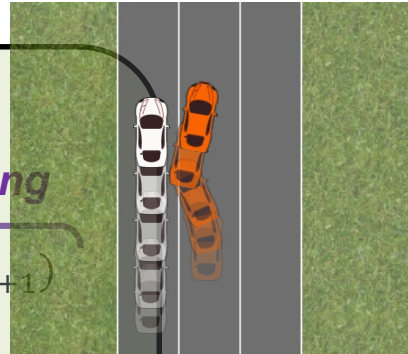Leverage human *responses* to estimate human internal state, and learn human model.

$$p(u_H|x, \theta, u_R) \propto \exp(R_H(x, u_H, \theta, u_R))$$

$$b_{t+1}(\theta) \propto b_t(\theta) \cdot p(u_H|x_t, \theta, u_R)$$

**Info Gathering**

$$r_R(x_t, u_H, \theta, u_R) = \underbrace{\mathcal{H}(b_t) - \mathcal{H}(b_{t+1})}_{} + \underbrace{\lambda \cdot r_{goal}(x_t, u_H, \theta, u_R)}_{\textbf{Goal}}$$

$$u_R = \operatorname*{argmax}_{u_R} \mathbb{E}_\theta[R_R]$$

## SAMPLE RESULT

**Lane Change**

**Nudging In**

**Distracted Human**

**Attentive Human**

## VERIFYING ROBUSTNESS

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

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

$$|R_{\mathcal{H}}^\dagger - R_{\mathcal{H}}| < \delta \quad \text{Optimizing a perturbed version of the learned reward function.}$$

$\widetilde{u_{\mathcal{H}}}$

$\widetilde{u_{\mathcal{H}}}$

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

$\delta = 0$   $\delta = 0.025$   $\delta = 0.15$

[Sadigh, Sastry, Seshia, Dragan; RSS, IROS '16]

[Sadigh, Sastry, Seshia: CPHS '18]

# Learning and Teaching (Multiple) Task Specifications

Good Communication is Crucial

Demonstrations,
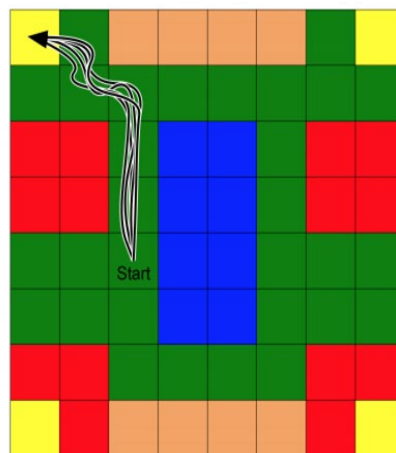Natural Language
…

Cost Functions,
Logical Specs.
…



Humans and machines must coordinate actions and processing

Boolean (logic) specifications:
- Composable
- Non-Markovian tasks
- Leverage formal methods

How can we hand off control reliably and intuitively?

Learning Boolean Specifications from Demonstrations

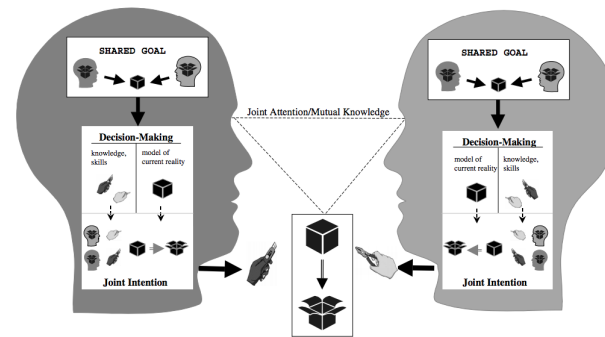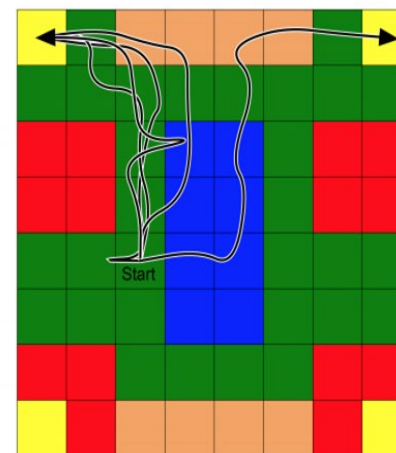On the Utility of Learning about Humans for Human-AI Coordination

Target Specification:

Go to a yellow tile without going on a red tile. If a blue tile is steped on, step on a brown tile before stepping on a yellow tile
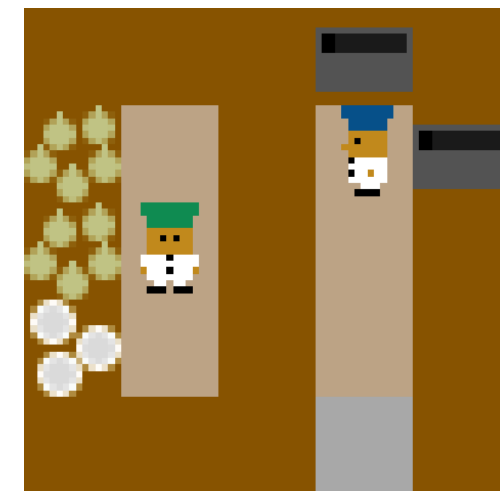
Doing Task

Communicating Task



[Vazquez-Chanlatte, Jha, Ho, et al., NeurIPS'18; CPHS'18]
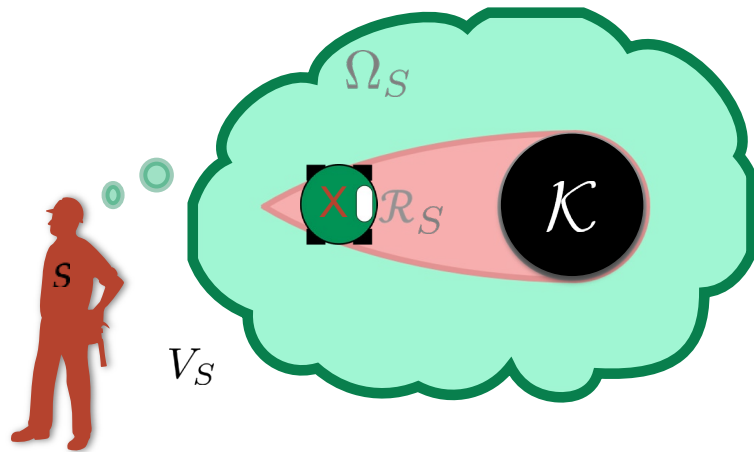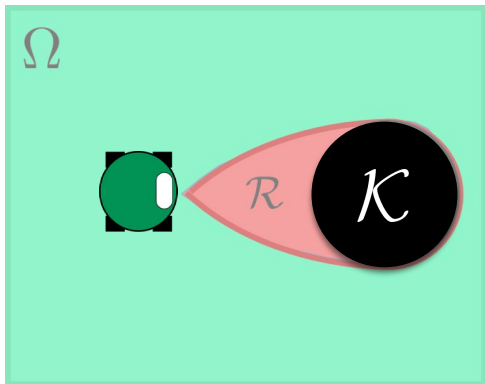
[Carroll, Shah, Ho, et al., NeurIPS'19]

11

# Inferring Supervisor Safe Sets for Human-Robot Teams

**Standard Reachability Safe Set:**

$\Omega$

$\mathcal{R}$  $\mathcal{K}$

**Human's Perceived Safe Set:**

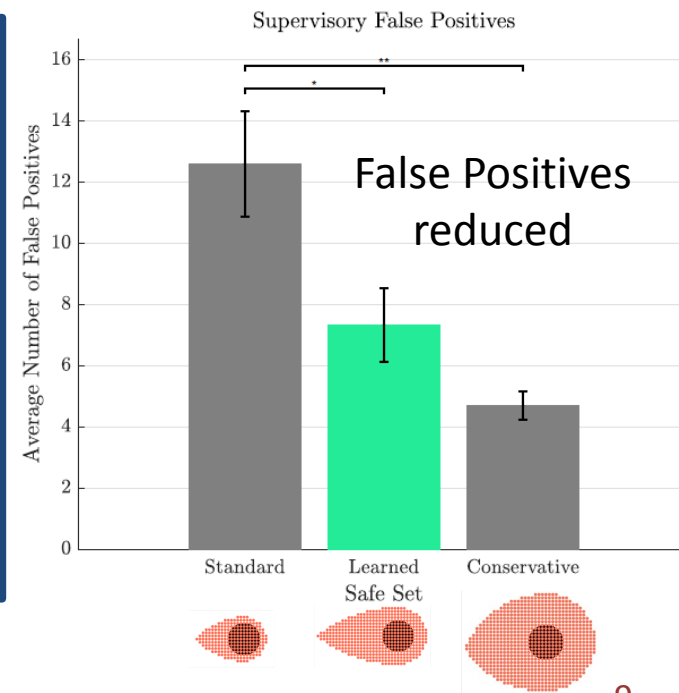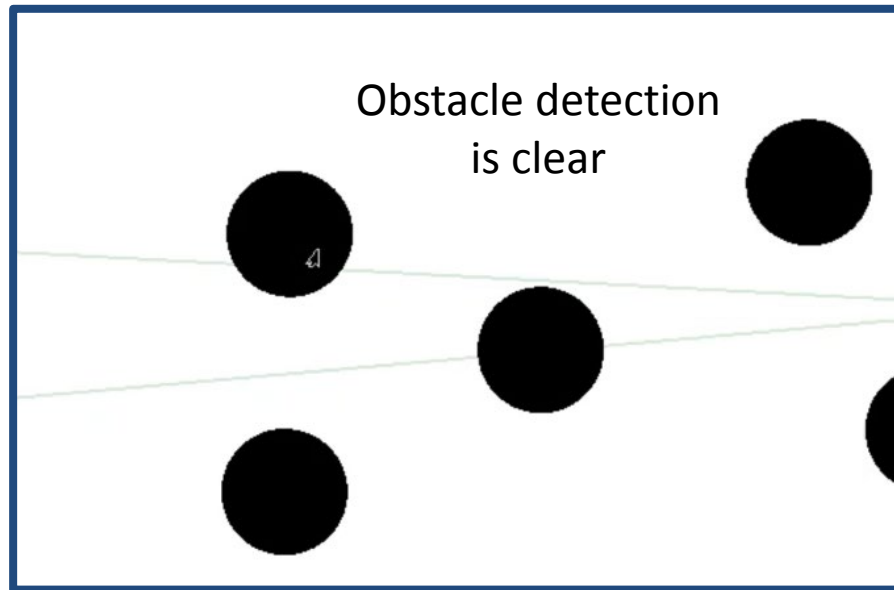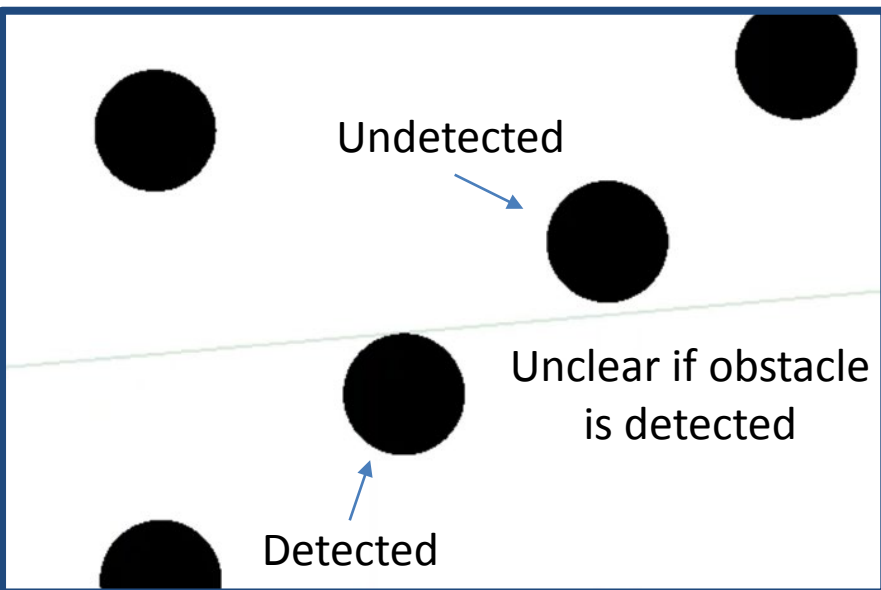$\Omega_S$

$S$

$V_S$

$\mathcal{R}_S$  $\mathcal{K}$

- Reachability analysis can compute safe sets
- Humans may perceive safety differently (from each other and from the system designer)
- Our technique can learn these individual safety preferences
- Avoid obstacles using learned safe sets to clearly communicate obstacle detection
- False positives (human thinks obstacle is undetected) are reduced

**Avoidance with Standard:**

Undetected

Unclear if obstacle is detected

Detected

**Avoidance with Learned:**

Obstacle detection is clear

**Supervisory False Positives**

False Positives reduced

Standard | Learned Safe Set | Conservative

D. McPherson, D. Scobee, J. Menke, A. Yang and S. Sastry, "Modeling Supervisor Safe Sets for Improving Collaboration in Human-Robot Teams." IROS 2018
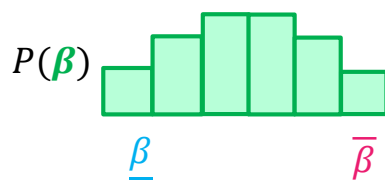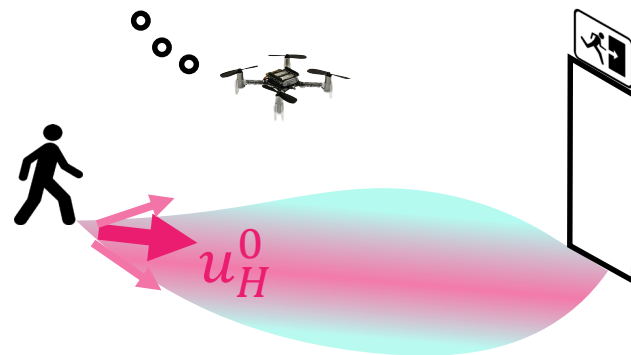
9

# Probabilistically Safe Motion Planning Around People

Use Human Models not as Ground Truth, but to Inform *Confidence in Predictions*

**Prediction**:
Confidence-aware Human
Prediction w/ Boltzmann Model

$+$

**Planning & Control:**
Fast and Safe Tracking (FaSTrack)

$$P(u_H^0 \mid x_H^0; \theta, \beta) \propto e^{\beta Q(x_H, u_H; \theta)}$$

$$P(\text{Crash}(x_R^\tau)) = \mathbb{E}_{\beta, \theta} \int_{\mathcal{H}_\varepsilon(x_R^\tau)} dP(x_H^\tau \mid x_H^t; \beta, \theta)$$

$u_H^0$

$P(\beta)$

$\underline{\beta}$       $\overline{\beta}$

$x_R^\tau$

$x_H^0$

$x_H^t$

[Andrea Bajcsy, Sylvia Herbert, David Fridovich-Keil, Jaime Fisac, Claire Tomlin 2018]

10

# Probabilistically Safe Motion Planning Around People
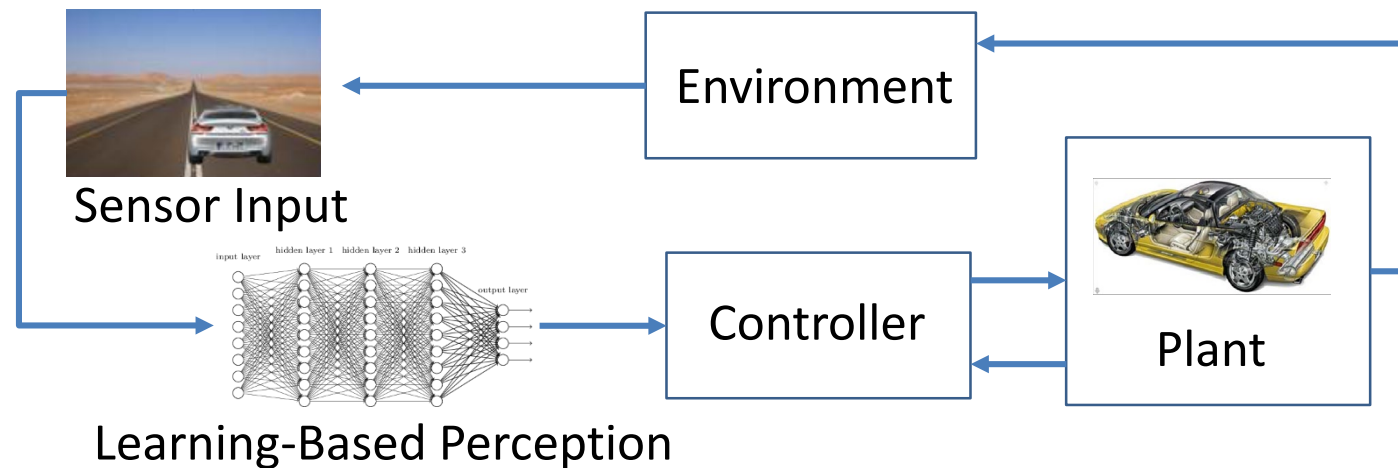
**Large Scale Simulation**

**Top-down View**

**Hardware Experiment**

[A Scalable Framework for Real-time, Multi-Robot, Multi-Human Collision Avoidance, ICRA 2019]

# A Semantic Approach to the Design of High-Assurance Learning-Based CPS

# SCENIC: Scenario Description Language

- *Scenic* is a probabilistic programming language defining *distributions over scenes*
- *Use cases:* data generation, test generation, verification, debugging, design exploration, etc.

```
from gta import Car, curb, roadDirection

ego = Car

spot = OrientedPoint on visible curb
badAngle = Uniform(1.0, -1.0) * (10, 20) deg
Car left of (spot offset by -0.5 @ 0),
    facing badAngle relative to roadDirection
```

Platoons

Images created with GTA-V

Bumper-to-bumper

[D. Fremont et al., "Scenic: A Language for Scenario Specification and Scene Generation", TR 2018, PLDI 2019.]

S. A. Seshia

# Some Applications of Scenic

- ## Data Generation, (Re)-Training
  - More controllable, interpretable
  - Improves performance significantly
  - Rare scenarios, controlled distributions, etc.
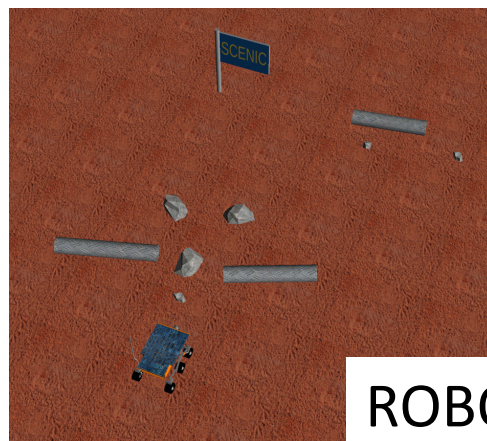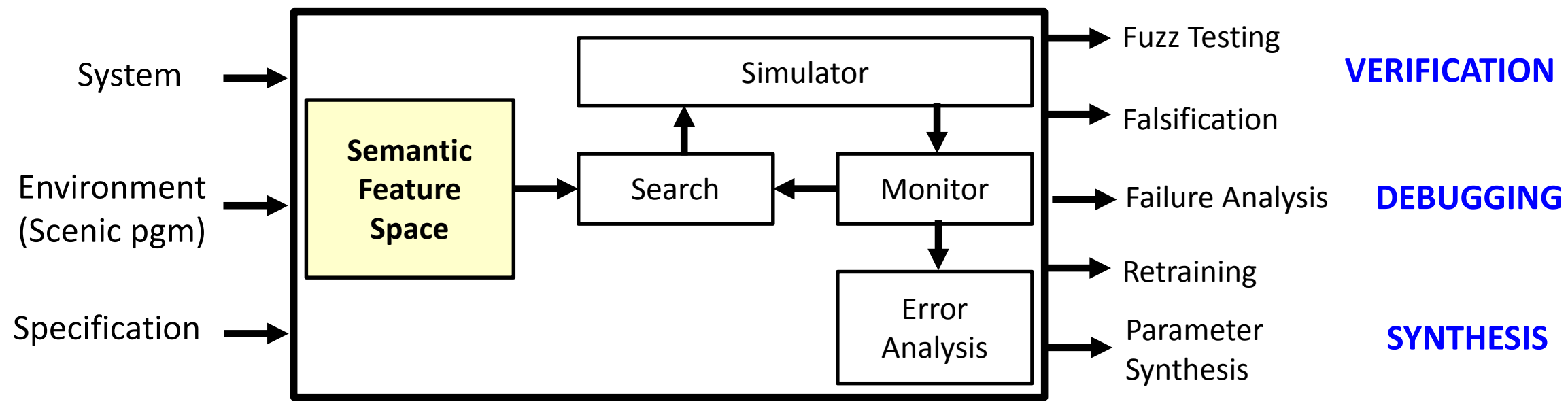


Car detection with occlusions

- ## Debugging Failures
  - Vary scenarios systematically
  - Explain failures of ML



- ## Design Space Exploration

Test Hypothesis: does the car model lead to a mis-detection?

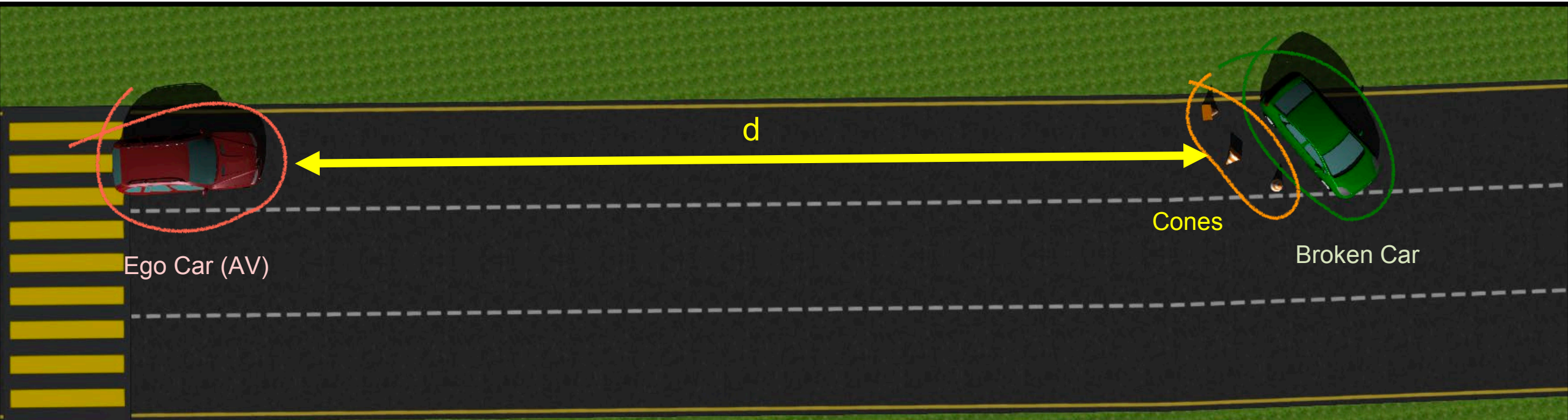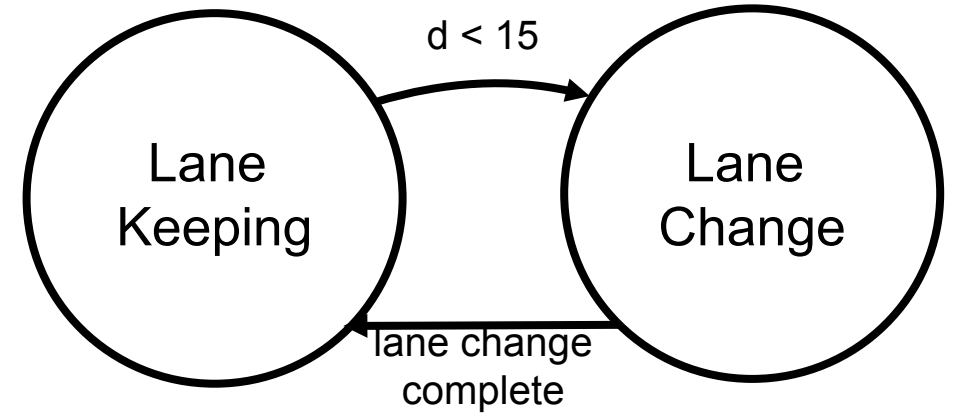# VERIFAI: A Toolkit for the Design and Analysis of AI-Based Systems [CAV 2019]

https://github.com/BerkeleyLearnVerify/VerifAI



ROBOTICS

AUTONOMOUS DRIVING

AIRCRAFT

# Case Study for Temporal Logic Falsification with VerifAI: Navigation around an Accident Scenario



Lane Keeping → Lane Change : $d < 15$

Lane Change → Lane Keeping : lane change complete

d

Ego Car (AV)

Cones

Broken Car

# Modeling Accident Scenario in the SCENIC Language



```
# Pick location for blockage randomly along curb
blockageSite = OrientedPoint on curb

# Place traffic cones
spot1 = OrientedPoint left of blockageSite by (0.3, 1)
cone1 = TrafficCone at spot1,
                   facing (0, 360) deg

...

# Place disabled car ahead of cones
SmallCar ahead of spot2 by (-1, 0.5) @ (4, 10),
        facing (0, 360) deg
```

Fremont et al., *Scenic: A Language for Scenario Specification and Scene Generation*, PLDI 2019.
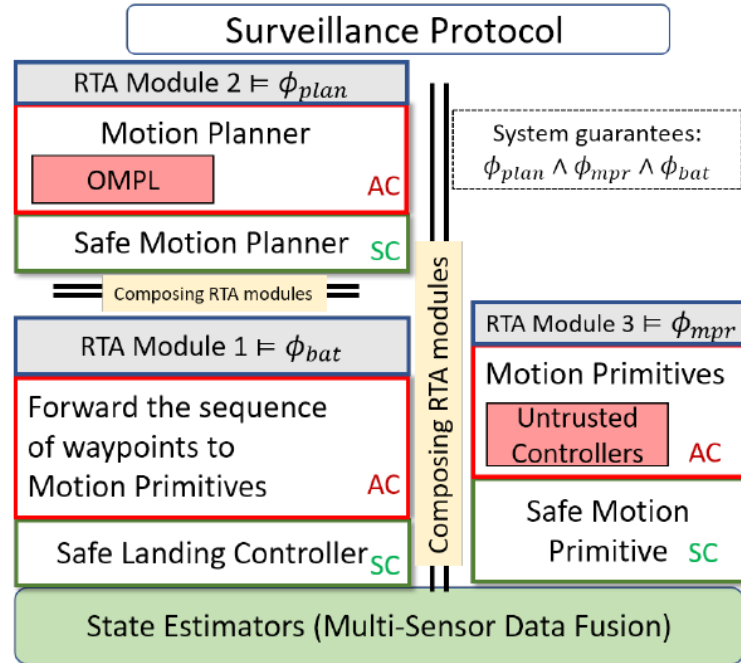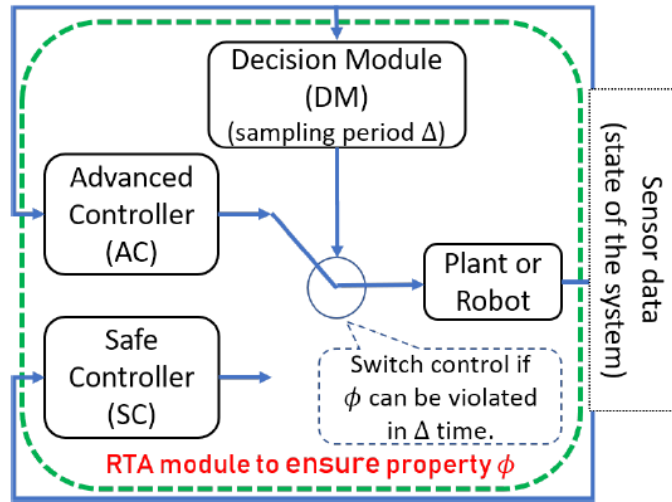
Temporal Logic Falsification

# From Models to Real World: Bridging the Gap

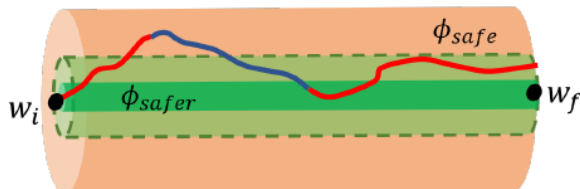**SPEC Specification Across Distance Between Behavior** [HSCC 2019]

Models                                                     Real World

## SOTER: Programming Framework for Run-Time Assurance [DSN 2019]



The

ther

There

Mo

Theorem: The following is an inductive invariant:

$$\text{Mode} = \text{SC} \land s \in \varphi_{safe}$$
$$\lor$$
$$\text{Mode} = \text{AC} \land \text{Reach}(s, *, \Delta) \in \varphi_{safe}$$

Surveillance Protocol

RTA Module 2 $\vDash \phi_{plan}$

Motion Planner

OMPL                          AC

Safe Motion Planner   SC

Composing RTA modules

RTA Module 1 $\vDash \phi_{bat}$

Forward the sequence of waypoints to Motion Primitives          AC

Safe Landing Controller SC

RTA Module 3 $\vDash \phi_{mpr}$

Motion Primitives

Untrusted Controllers    AC

Safe Motion Primitive  SC

System guarantees:
$\phi_{plan} \land \phi_{mpr} \land \phi_{bat}$

State Estimators (Multi-Sensor Data Fusion)

$\phi_{safe}$

$w_i$  $\phi_{safer}$  $w_f$   `module system = RTAModule1 || RTAModule2 || RTAModule3;`

# Formal Models Key to Co-Design

**Formal, Semantic and Predictive Models of Human and Environment Behavior**

**Verified Control Design**

**Verified Human-Machine Interaction Design**
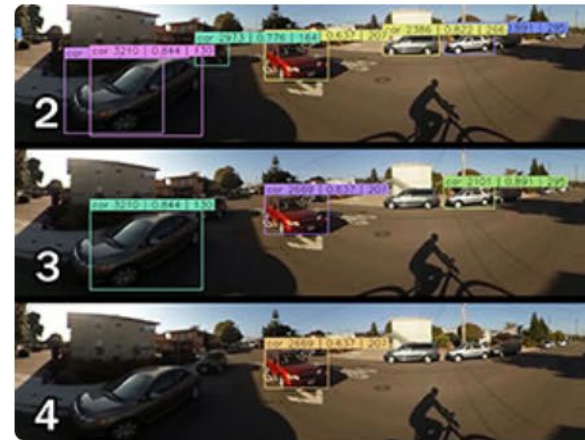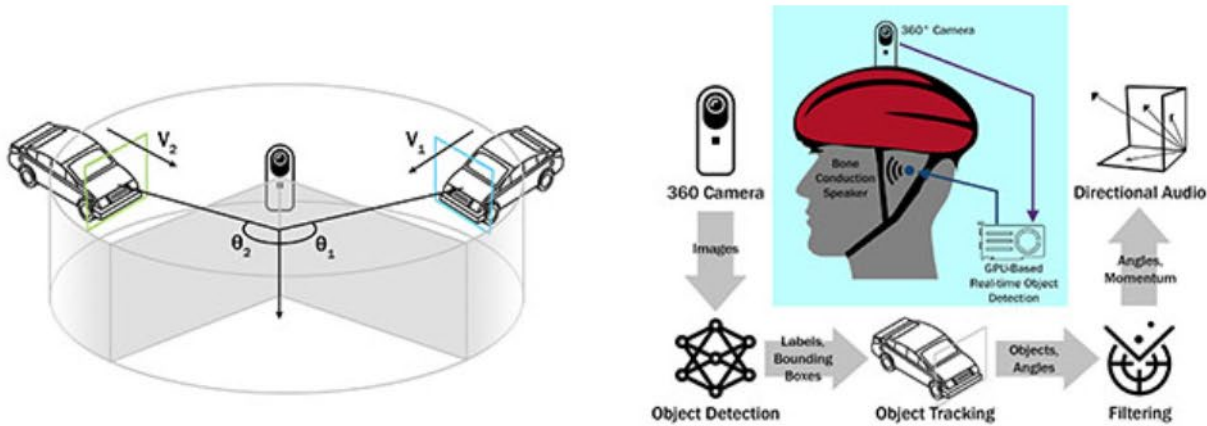
# A Selection of Other Results from VeHICaL

**Abstractions for Neural Network Analysis**

**Vibro-Acoustical Approach to Driver Interfaces**

**Drowsy Driver Detection**

**HindSight: Bicyclist Assistance Systems**



*HindSight* increases the environmental awareness of cyclists by warning them of vehicles approaching from outside their visual field. A panoramic camera mounted on a bicycle helmet streams real-time, 360-degree video to a laptop running YOLOv2, a neural object detector designed for real-time use. Detected vehicles are passed through a filter bank to find the most relevant.

# Industrial Impact

- Several workshops with strong industry participation

- Open-Source Tools and Datasets
  - VerifAI, Scenic, …
  - Drowsy Driver Dataset, Visual-Acoustic Vehicle Dataset, …

- Tools/ideas being adopted by Industry

- Working with AAA & LG on AV scenario specification and testing at GoMentum test facility

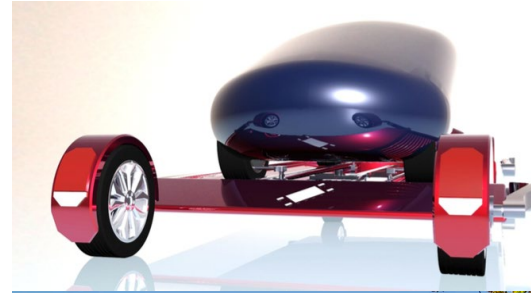- Advice to NHTSA project on AV Test Cases







A Framework for Automated Driving System Testable Cases and Scenarios

# Impact on Graduate and Undergraduate Education

- Several courses impacted by VeHICaL
- Reimagining Mobility – collaboration with Ford Greenfield Labs
  - at the Jacobs Institute of Design Innovation @ Berkeley
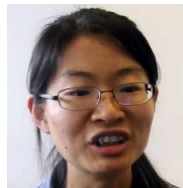- Academic/industry positions for graduates from VeHICaL project

UIUC  Stanford  UCSC  Princeton

TRI  Waymo  … and more

Design Project 3:
The Goods Delivery Interface Between Humans and Autonomous Vehicles

See-Thru:
Towards Minimally Obstructive Eye-Controlled Wheelchair Interfaces
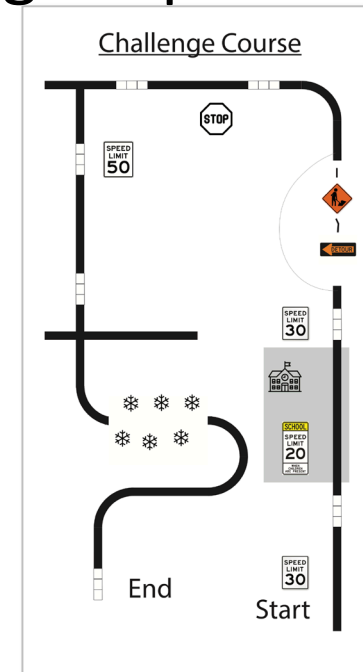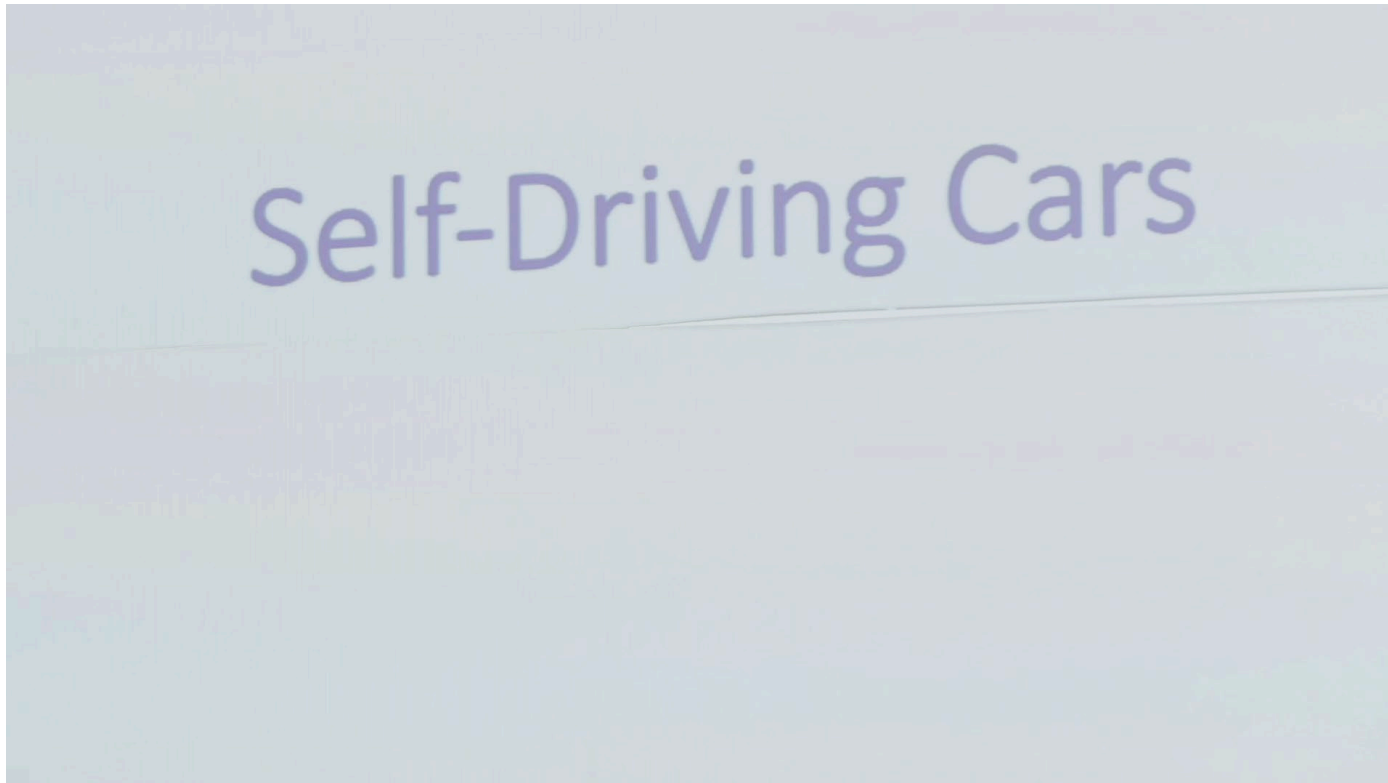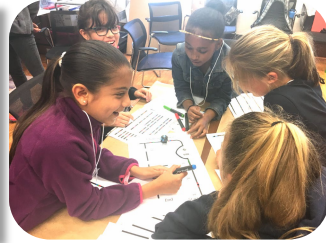
# Broader Impacts – Girls in Engineering (GiE) VeHICaL modules

- Summer program for middle-school girls at Berkeley
- VeHICaL provided instructors/mentors, funding, content
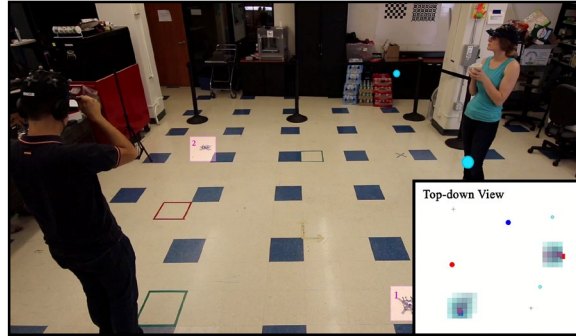- Modules on self-driving car technology using simple Ozobot platform

# VeHICaL: Verified Human Interfaces, Control, and Learning for Semi-Autonomous Systems



## Challenge:

- *Co-design human interfaces and control* for human-cyber-physical systems with *provable guarantees*
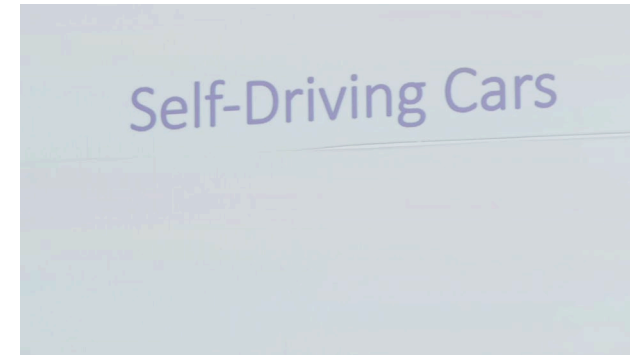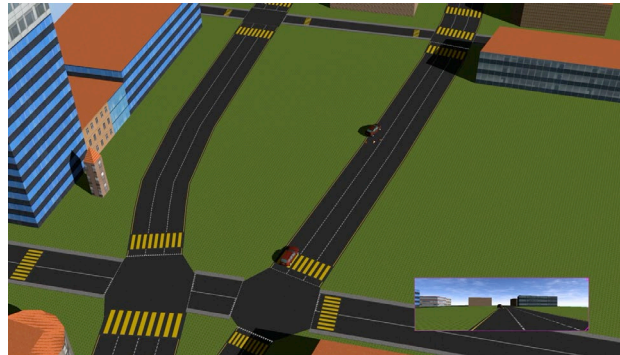- Apply to semi-autonomous vehicles (ground and air)

## Solution:

- Integrate Learning, Verification and Control
- Data-Driven Resource Rational Human Modeling
- Prototype Controllers & Interfaces, Evaluate on Testbed

## Scientific Impact:

- Developing a Science of Co-Design of Human Interfaces and Control
- Bridging Model-Based and Data-Driven Design of CPS

## Broader Impact:

- Significantly improve safety, security, and performance of systems where humans interact closely with automation
- Involve middle/high-school and undergraduate students in VeHICaL activities

```
from gta import Car, curb, roadDirection

ego = Car

spot = OrientedPoint on visible curb
badAngle = Uniform(1.0, -1.0) * (10, 20) deg
Car left of (spot offset by -0.5 @ 0),
    facing badAngle relative to roadDirection
```

## THANK YOU!