Towards Verified Artificial Intelligence

Sanjit A. Seshia

Professor

EECS Department, UC Berkeley

Joint work with:



Ankush Desai, Tommaso Dreossi, Daniel Fremont, Shromona Ghosh, Edward Kim, Hadi Ravanbakhsh, Marcell Vazquez-Chanlatte, Xiangyu Yue, Dorsa Sadigh, Wenchao Li, Alexandre Donze, Alberto Sangiovanni-Vincentelli, S. Shankar Sastry, Natarajan Shankar, Ashish Tiwari

Vet-JiCal

MLSE International Symposium Tokyo, October 18, 2019

Growing Use of Machine Learning/Artificial Intelligence in Safety-Critical Cyber-Physical Systems



How do we make AI/ML-based Autonomous Systems Safe & Dependable?

- Numerous papers showing that Deep Neural Networks can be easily fooled
- *Fatal accidents* involving potential failure of object detection/classification systems in self-driving cars

Challenges for Verified Al

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

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



Need Design Principles for Verified AI

Challenges

- 1. Environment (incl. -Human) Modeling
- 2. Formal Specification
- 3. Learning Systems Representation
- 4. Scalable Training, Testing, Verification
- 5. Design for Correctness

Principles

Verified AI Project VerifAI software toolkit

https://github.com/BerkeleyLearnVerify/VerifAI

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

Outline

- Motivating Example
- Specification/Requirements
- Verification: Compositionality and Abstraction
- Environment Modeling ODDs, Test Scenarios, etc.
- Synthesis: Data-Set Design, Run-Time Assurance, etc.
- The VerifAI toolkit & Conclusion

Example: Automatic Emergency Braking System (AEBS) using Deep Learning for Perception



- Goal: Brake when an obstacle is near, to maintain a minimum safety distance
- Models: Controller, Plant, Env modeled in a software-in-the-loop simulator
 - Matlab/Simulink, Udacity, Webots, CARLA, ...
- <u>Perception</u>: Object detection/classification system based on deep neural networks
 - Inception-v3, AlexNet, ... trained on ImageNet
 - squeezeDet, Yolo, ... trained on KITTI

[Dreossi, Donze, Seshia, "Compositional Falsification of Cyber-Physical Systems with Machine Learning Components", NASA Formal Methods (NFM), May 2017.] ₆



Challenge: Formal Specification

Principle: Start at the System Level (i.e. Specify Semantic Behavior of the Overall System)

Our Approach: Start with a System-Level Specification



"Verify the System containing the Deep Neural Network"

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



Property does not mention inputs/outputs of the neural network

Do We Need to Formally Specify ML Component Requirements?

Sometimes...

- 1. When Component Specifications are Meaningful
 - ML-based Controllers
 - Semantic Robustness*, Input-Output Relations, Monotonicity, Fairness, etc.
- 2. For Compositional Analysis
 - Derive component specifications from system-level specifications

[S. A. Seshia, et al., "Formal Specification for Deep Neural Networks", ATVA 2018.]

Challenge: Scalability of Verification

Principle: Compositional Simulation-Based Verification (Falsification)

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. (Extended version: Journal of Automated Reasoning, 2019.)

Automotive system with ML-based perception (CPSML)



*s*_p

Our Approach: Three Key Ideas

- 1. Reduce CPSML falsification problem to combination of CPS falsification and ML analysis by *abstraction*
- 2. Simulation-based *temporal logic falsification* for CPS model
 - Scalable technology already used on production automotive systems (powertrain)
- 3. Semantic feature space analysis of ML component
 - Derive constraints on input feature space of neural network (pixels) from semantic constraints on environment model

[Dreossi et al., NFM'17, JAR'19; Dreossi et al., CAV'18]

Compositional Approach: Combine Temporal Logic CPS Falsifier with ML Analyzer



Identifying Component-Level Input Constraints (ROU) for Automatic Emergency Braking System

Green \rightarrow environments where system-level safety property is satisfied



[Dreossi et al., NFM'17, JAR'19; Dreossi et al., CAV'18]

Result on AEBS Example

This misclassification not of concern



0 0

Sample image

Result on AEBS Example



Extending to Image Streams

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



Results on squeezeDet NN and KITTI dataset for autonomous driving

<u>Challenge:</u> How to (Re)Design Learning Components

Principle: Use

Oracle-Guided Inductive Synthesis (OGIS)

Standard Machine Learning



LEARNER

labeled/unlabeled data





SOURCE OF DATA

Learned Model only as good as the Data!

Correct-by-Construction Design with Formal (Oracle-Guided) Inductive Synthesis/Learning

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.]

Counterexample-Guided Training of Deep Neural Networks

- Instance of Oracle-Guided Inductive Synthesis
- Oracle is Verifier (CPSML Falsifier) used to find counterexample inputs to DNN
- Substantially increase accuracy with relatively few additional examples



"Counterexample-Guided Data Augmentation", T. Dreossi, S. Ghosh, X. Yue, K. Keutzer, A. Sangiovanni-Vincentelli, S. A. Seshia, IJCAI 2018.



Challenge: Environment Modeling

Principles: Data-Driven, Probabilistic, Introspective, Modeling

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

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.
 - Example scenario: a badly-parked car

S. A. Seshia

```
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
```



Images created with GTA-V

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

SCENIC: Scenario Description Language

Scenic makes it possible to specify broad scenarios with complex structure, then generate many concrete instances from them automatically:



Platoons

Bumper-to-Bumper Traffic



Use Case: Debugging a Known Failure



Use Case: Debugging a Known Failure

Scenic makes it easy to vary a scenario along different dimensions:



Add noise



Change car model



Change global position



VERIFAI: A Toolkit for the Design and Analysis of Al-Based Systems <u>https://github.com/BerkeleyLearnVerify/VerifAI</u>



Case Study for Temporal Logic Falsification with VerifAI: Navigation around an accident scenario







Modeling Case Study in the SCENIC Language

 $\cdot \cdot \cdot \cdot$

```
# 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.

Using Scenic to Generate Initial Scenes



Using Scenic to Generate Initial Scenes



Using Scenic to Generate Initial Scenes



Falsification

TAXABLE PARTY



Analyzing the failure

Fix the controller: Update model assumptions and re-design controller Retrain the perception module: Collect the counter-example images and retrain the network [IJCAI'18]



Different Case Study: Automated Taxiing

- NN-based research prototype from Boeing
- Specification: track centerline within 1.5 m



Falsifying Run Found with VerifAl/Scenic



Links and Ongoing Work

- Open-Source Releases: the Verified AI toolkit [CAV 2019]
 - VerifAI: <u>https://github.com/BerkeleyLearnVerify/VerifAI</u>
 - Scenic: <u>https://github.com/BerkeleyLearnVerify/Scenic</u>
 - Operate with any simulator
- Other results/ongoing projects:
 - Bridging Simulation and Real World
 - Domain adaptation to produce "real" data from simulated data
 - Quantifying distance between simulated and real behaviors [HSCC 2019]
 - More Complex Sensors: Video + LiDAR + RADAR + ...
 - Counterexample-Guided Retraining/Data Set Design [IJCAI 2018]

Conclusion: Towards Verified AI/ML based CPS

Challenges			Core Principles
1.	Environment (incl. Human) Modeling	\rightarrow	Data-Driven, Introspective, Probabilistic Modeling
2.	Specification		Start with System-Level Specification, then Component Spec (robustness,)
3.	Learning Systems Complexity	\longrightarrow	Abstraction, Semantic Representation, and Explanations
4.	Efficient Training, Testing, Verification	\longrightarrow	Compositional Analysis and Semantics- directed Search/Training
5.	Design for Correctness		Oracle-Guided Inductive Synthesis; Run-Time Assurance

Exciting Times Ahead!!! Thank you!

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