



The AAIL's Workshop on  
Artificial Intelligence Safety

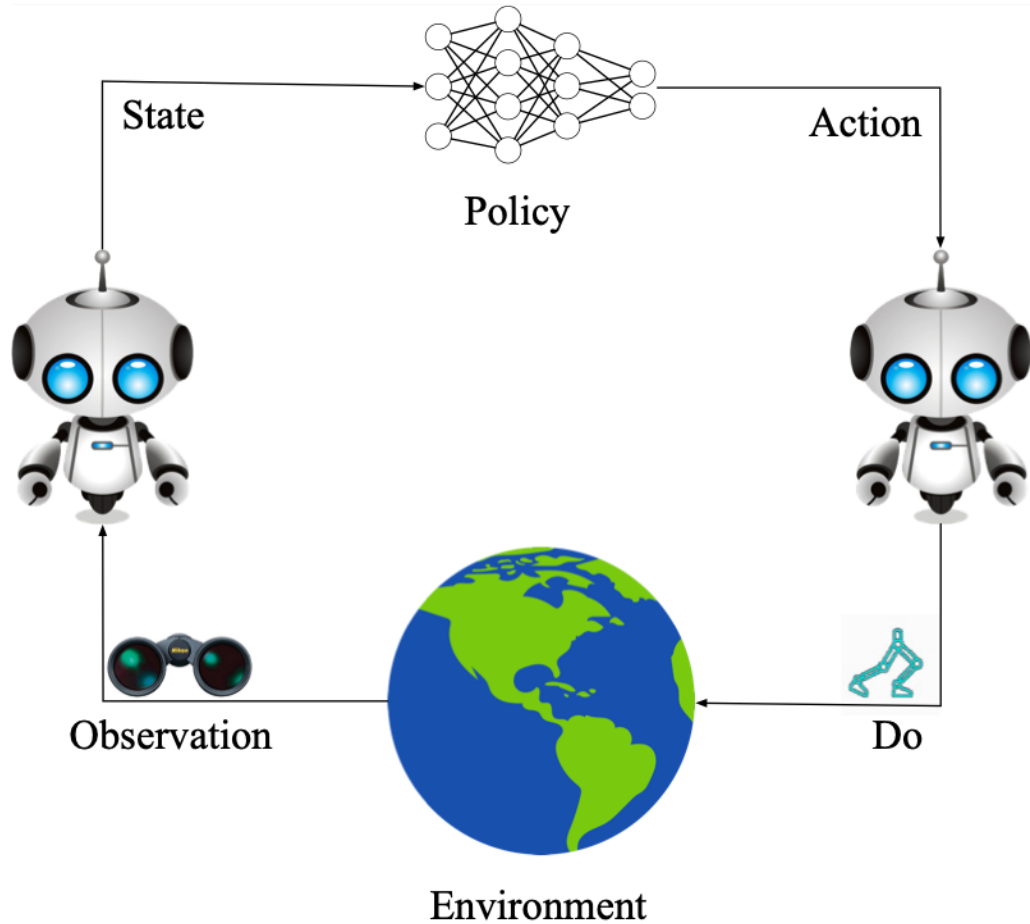


# EnnCore: Safety Verification of Deep Reinforcement Learning

Yi Dong, Xingyu Zhao, and Xiaowei Huang

Department of Computer Science, University of Liverpool

# Deep Reinforcement Learning

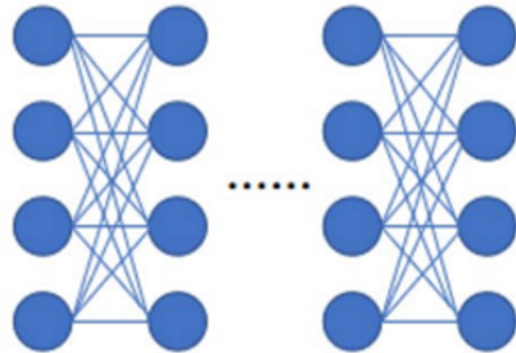


- Black-box
- White-box

# DRL verification vs CNN verification

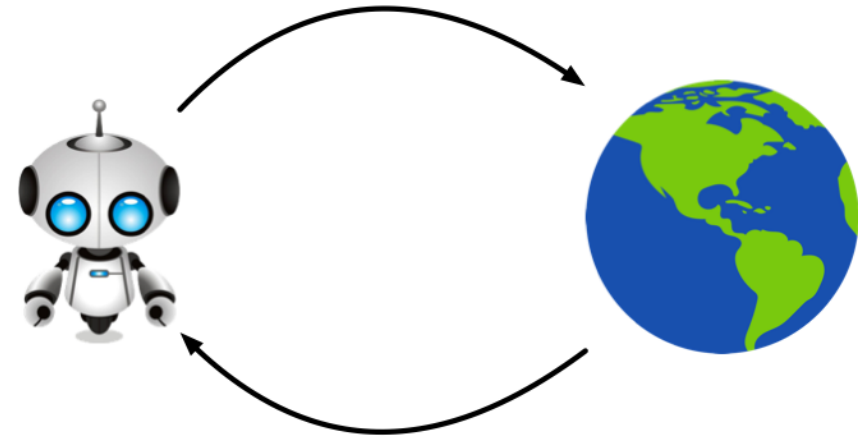


adversarial example



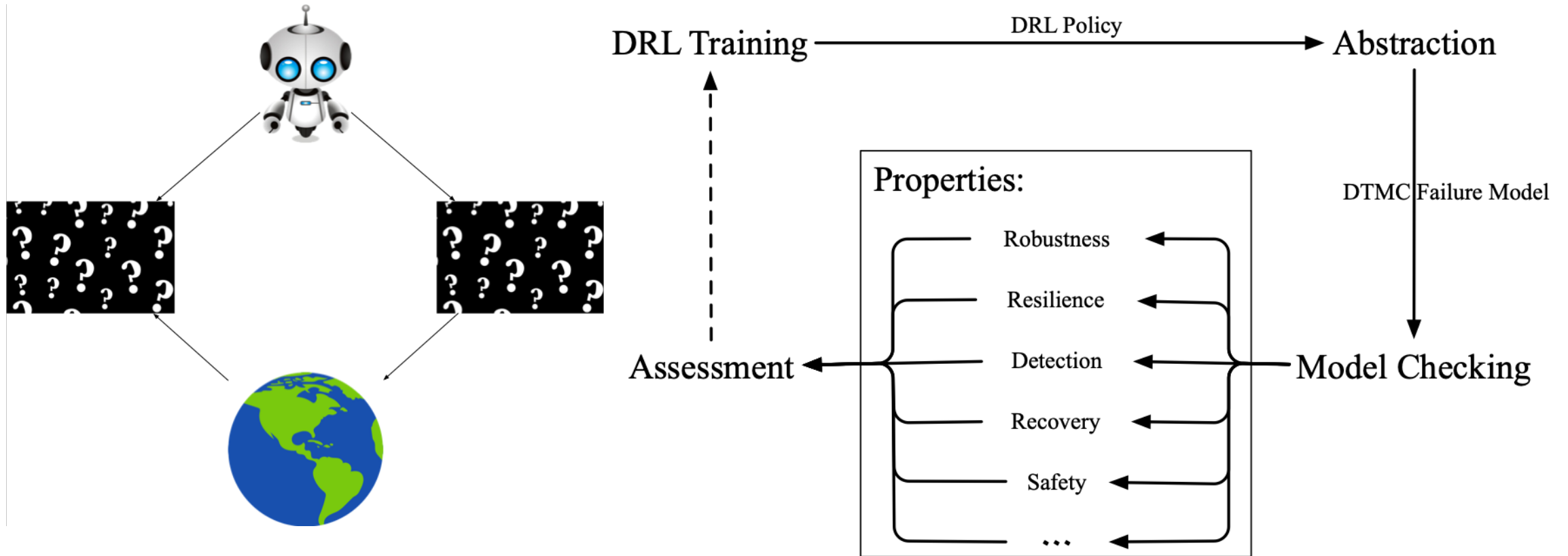
DNN model

CNN Verification



DRL Verification

# Black-box Verification



# Black-box Verification

Technical Problem:

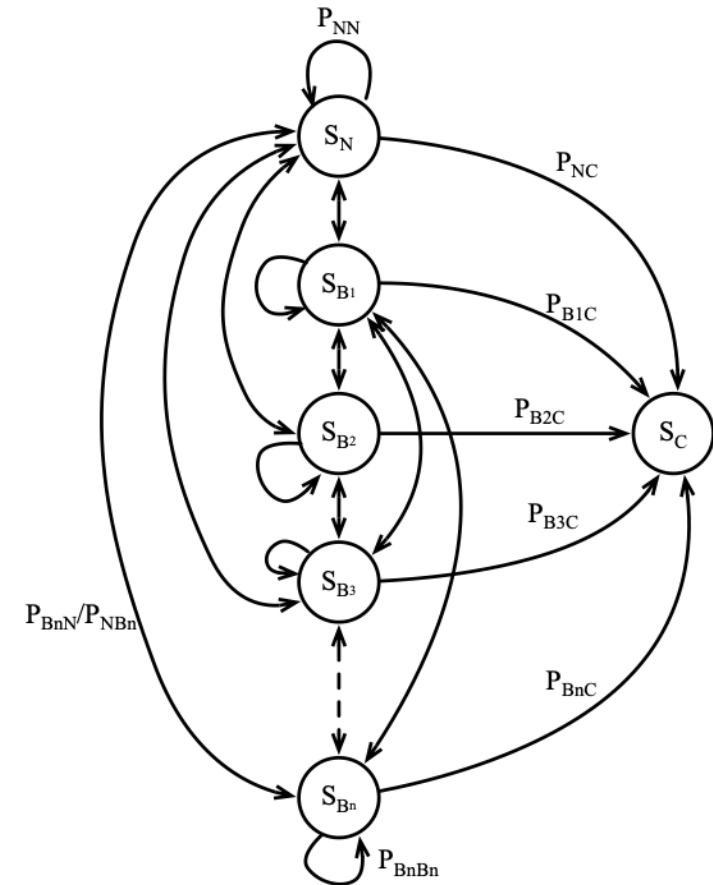
- How to synthesize the Markov model (DTMC, MDP)?
- How to evaluate the dependability of the DRL policy?
- How do properties perform in DRL algorithms ?

Technical Solution:

- Defined different properties
- Construct DTMC model
- Probabilistic model checking

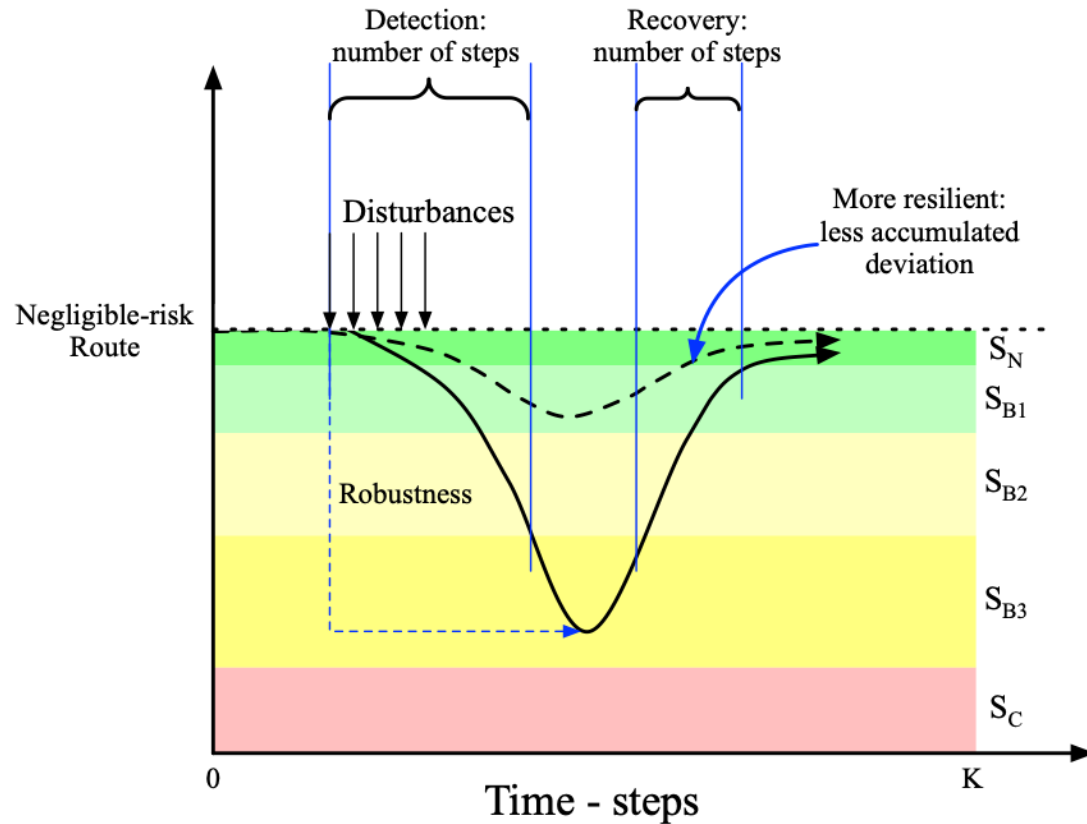
Technical Observation:

- The dependability analysis are insensitive to the sample size
- trade-offs between different properties



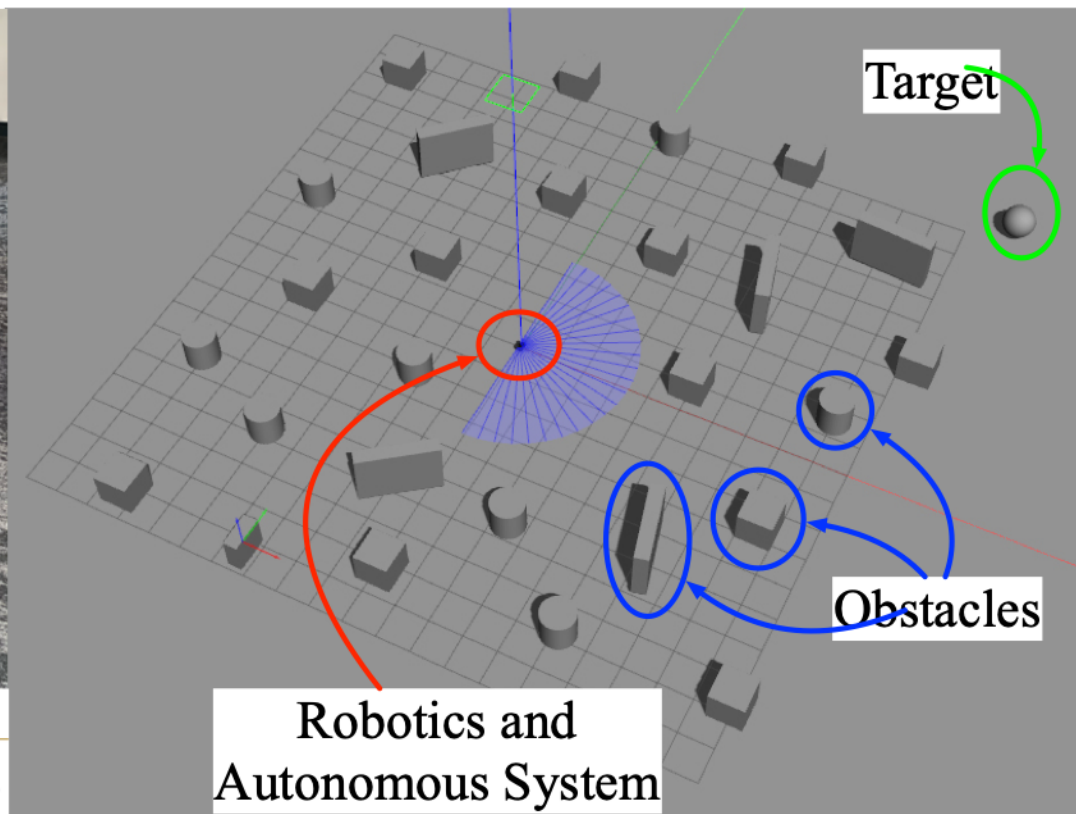
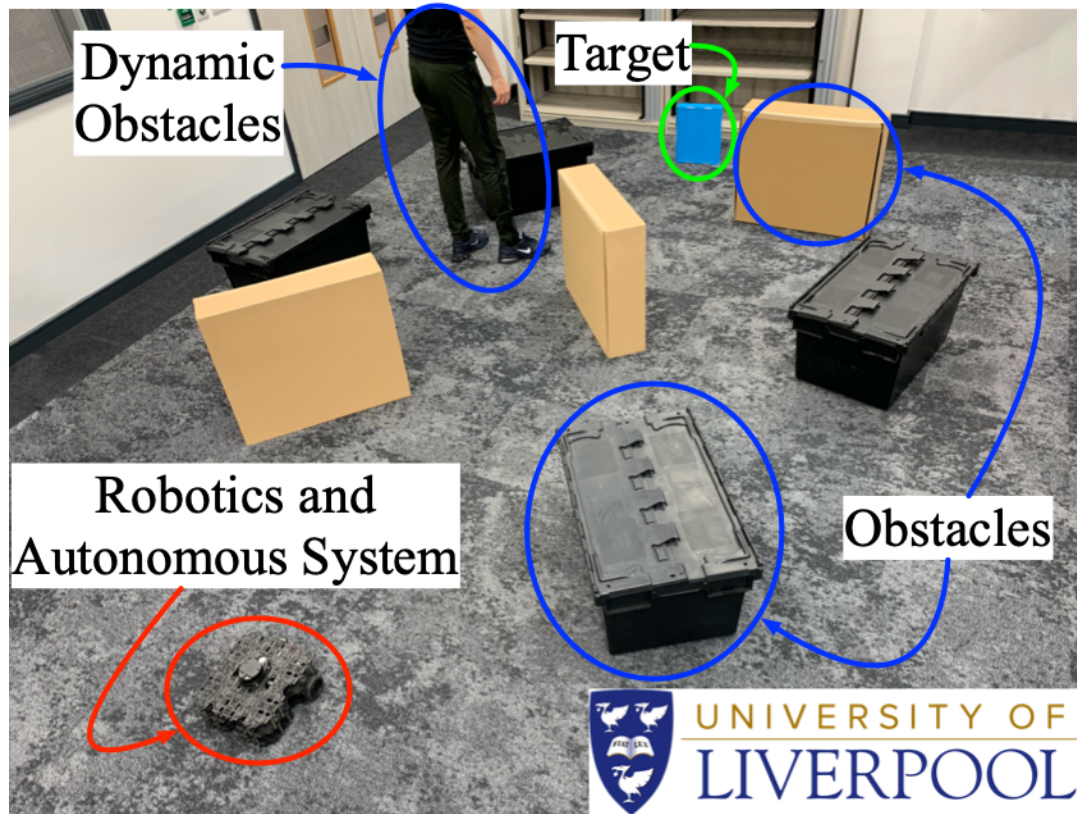
# Black-box Verification

## Safety Properties

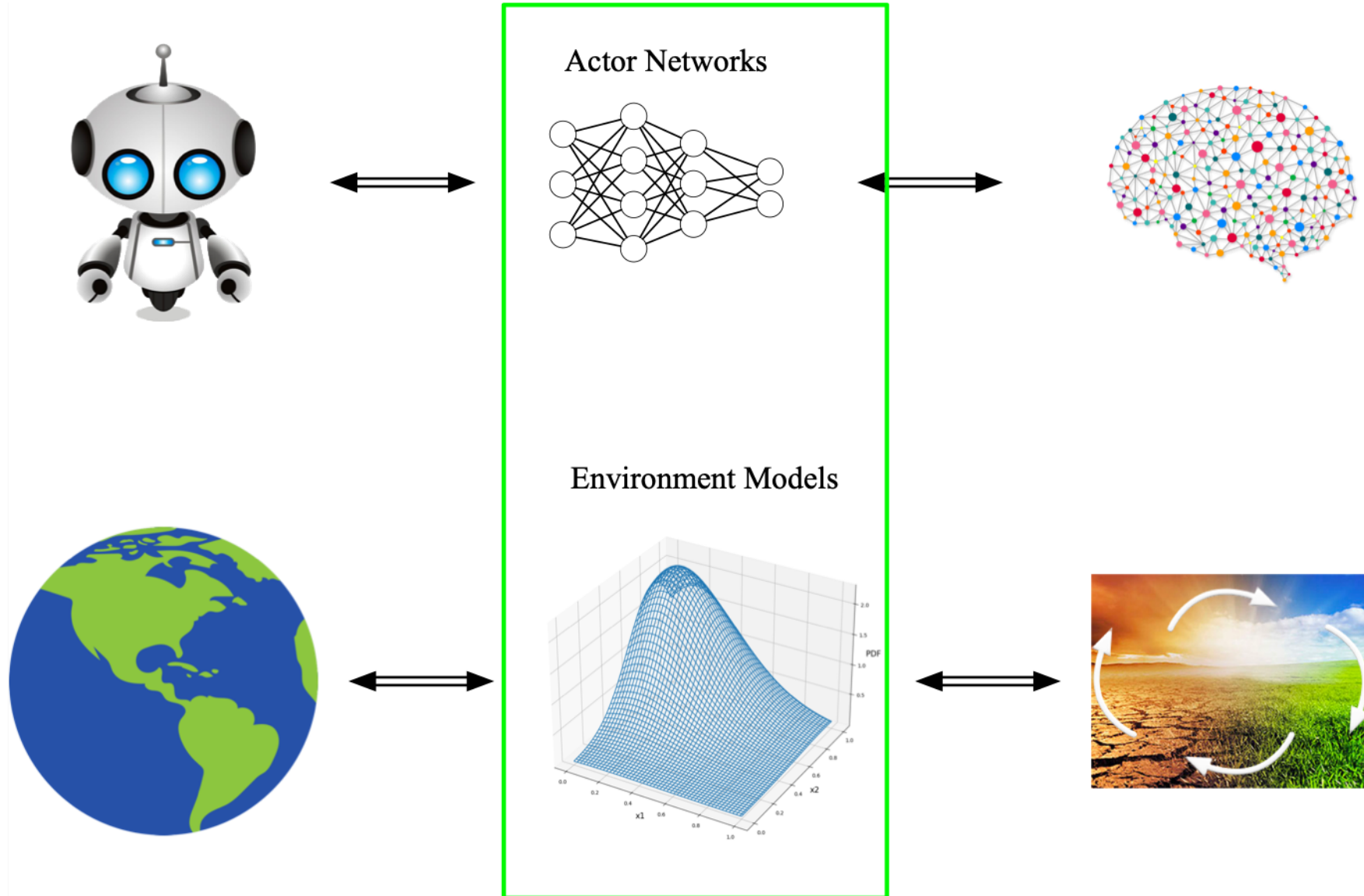


- Safety
- Robustness
- Resilience
- Detection
- Recovery
- ...

# Black-box Verification



# White-box Verification





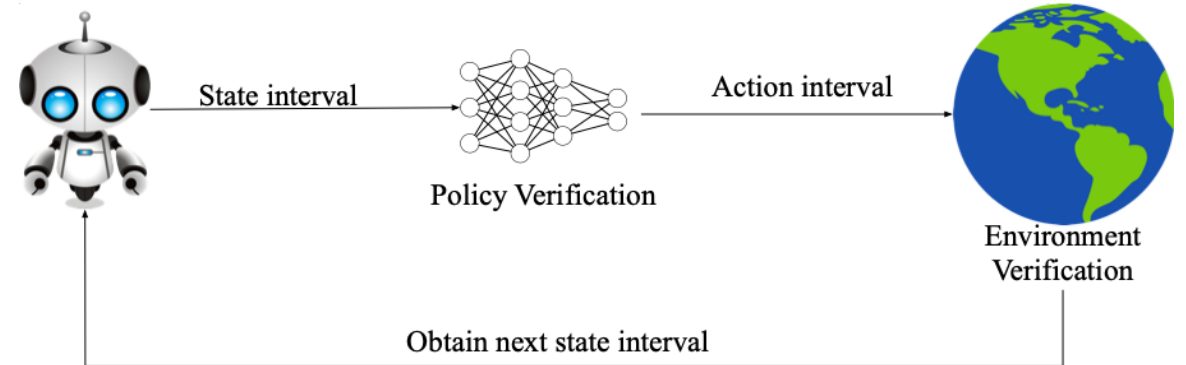
# White-box Verification

## Two-level verification

**Low-level:** for a given actor network, calculate reachable set of actions.

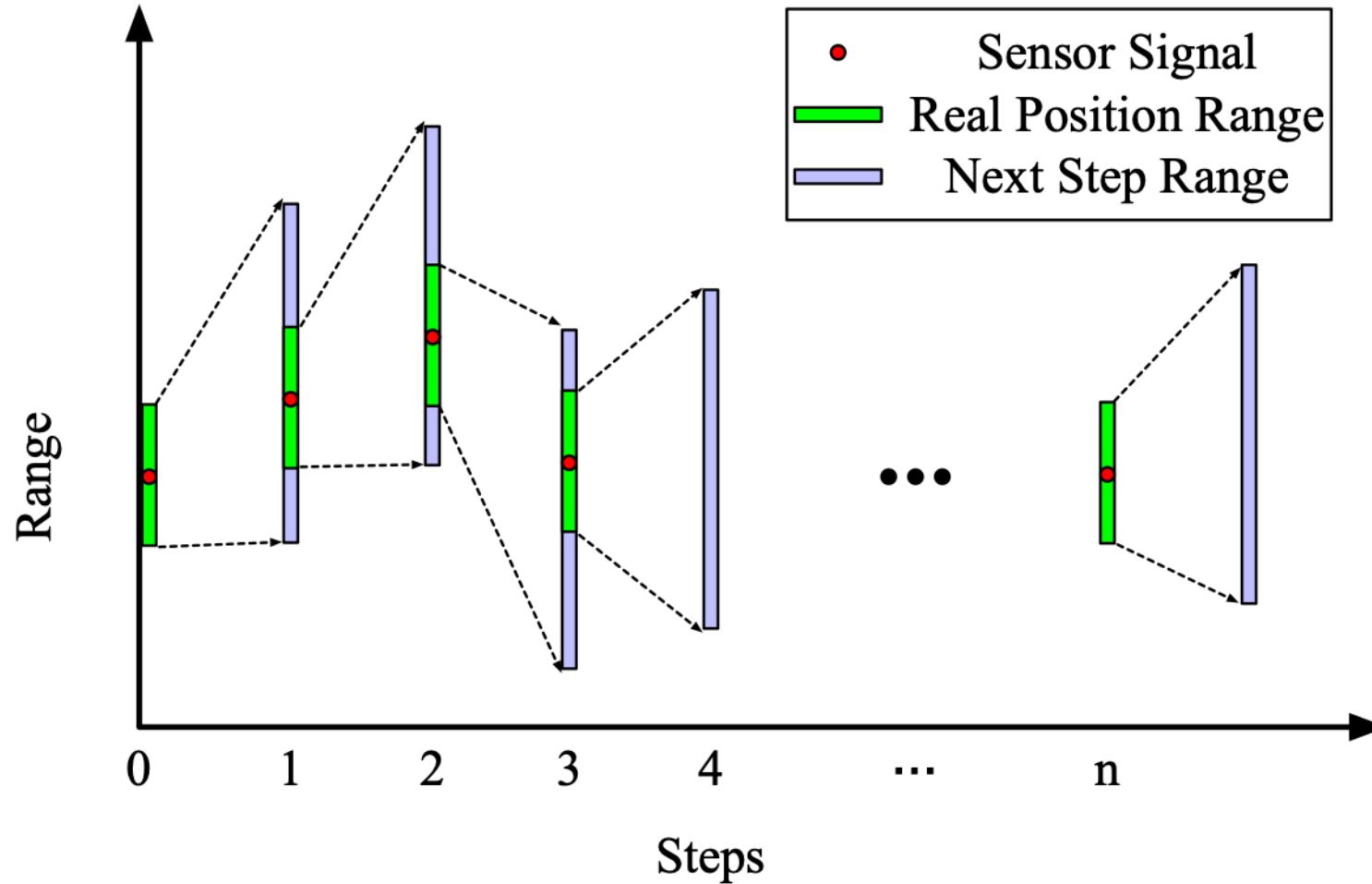
**High-level:** similar to black-box verification methods.

**Connection:** Do high-level verification based on low-level verification results.

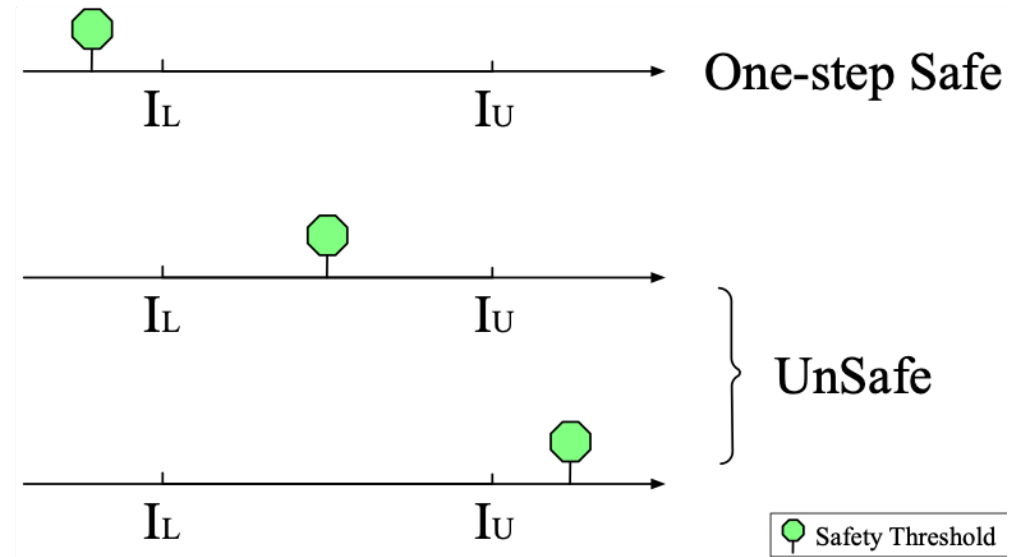
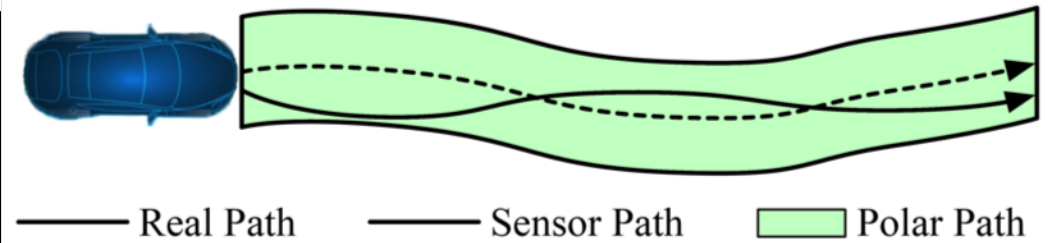
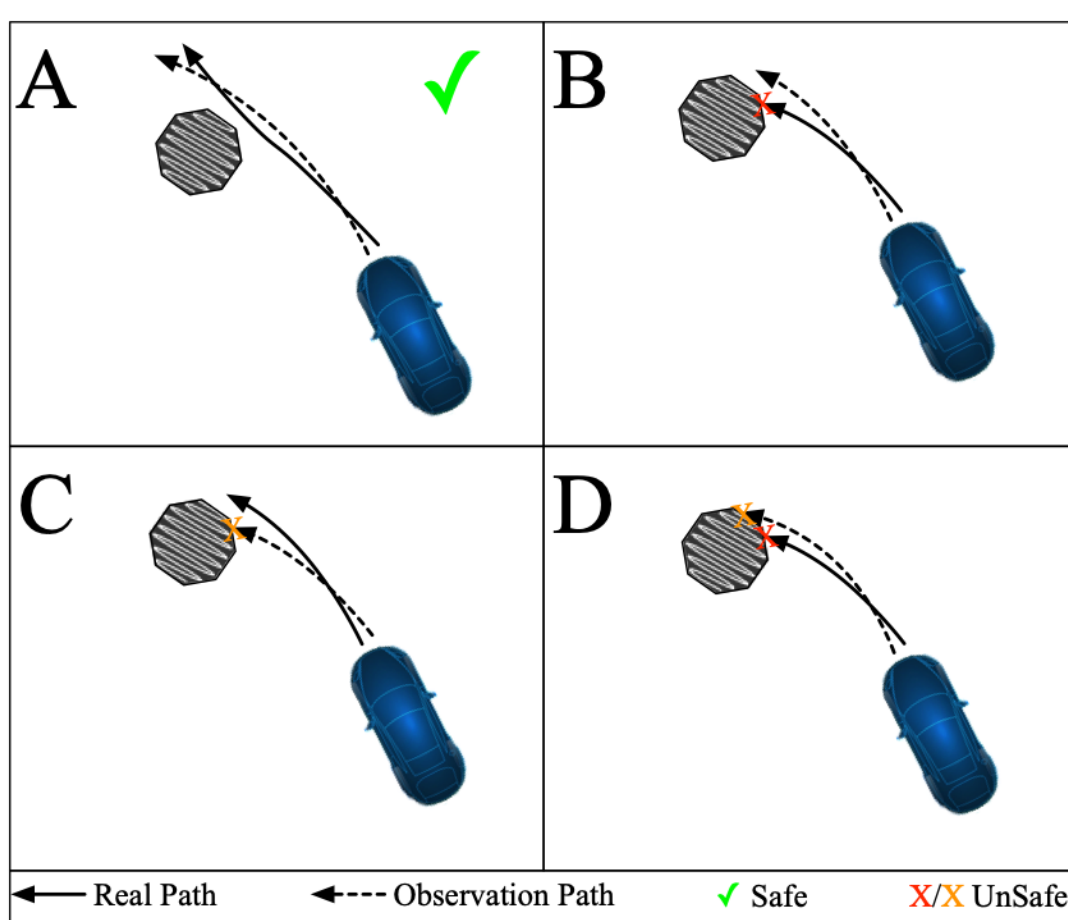


# White-box Verification

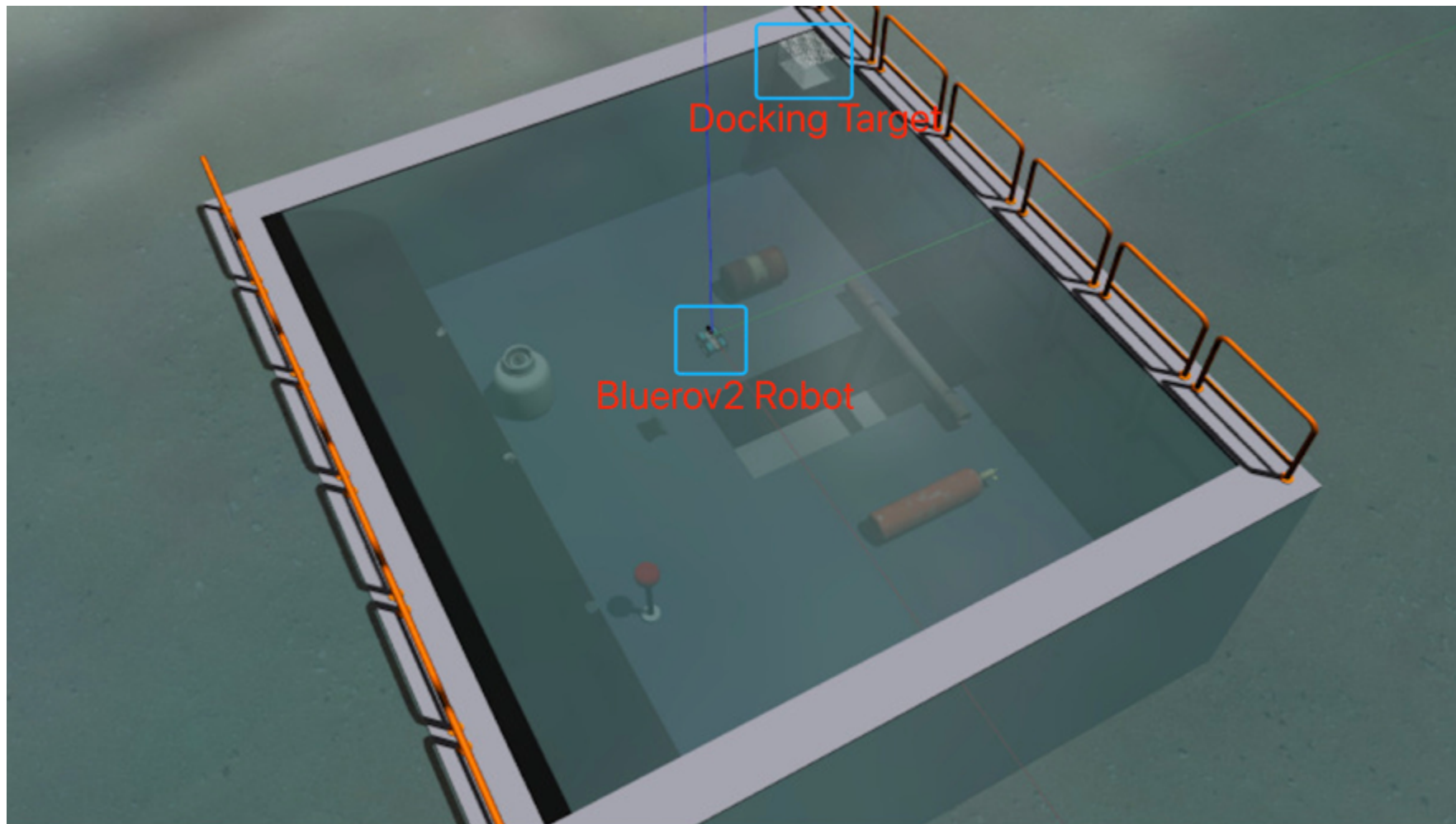
Temporal path



# White-box Verification



# White-box Verification



# Challenges

- **DTMC models**
  - build accurate DTMC with less samples
- **Verification of DRL Models**
  - A key problem remains on the **scalability and time-efficiency**: trade-off between **soundness** and **completeness**
- **Unknown Environment**
  - Verify of Neural network model working in an **unknown environment**

# EnnCore Program at SafeAI 2022

(<https://safeai.webs.upv.es/>)

Time (UTC)	Description
14:00-14:10	Welcome, overview, Lucas Cordeiro (University of Manchester, UK)
14:10-14:25	Verifying Quantized Neural Networks using SMT-Based Model Checking, Edoardo Manino (University of Manchester, UK)
14:25-14:40	Explainability and Inference Controls, André Freitas (University of Manchester UK & Idiap Research Institute, Switzerland)
14:40-14:55	Safety Verification of Deep Reinforcement Learning, Yi Dong (University of Liverpool, UK)
14:55-15:10	Privacy Friendly Energy Consumption Prediction: Real Case-Studies, Mustafa A. Mustafa (University of Manchester, UK / KU Leuven, Belgium)
15:10-15:30	Closed-loop Safety of Bayesian Neural Networks and Stochastic Control Systems, Mathias Lechner, IST Austria