

Robustness Certification of Deep Learning

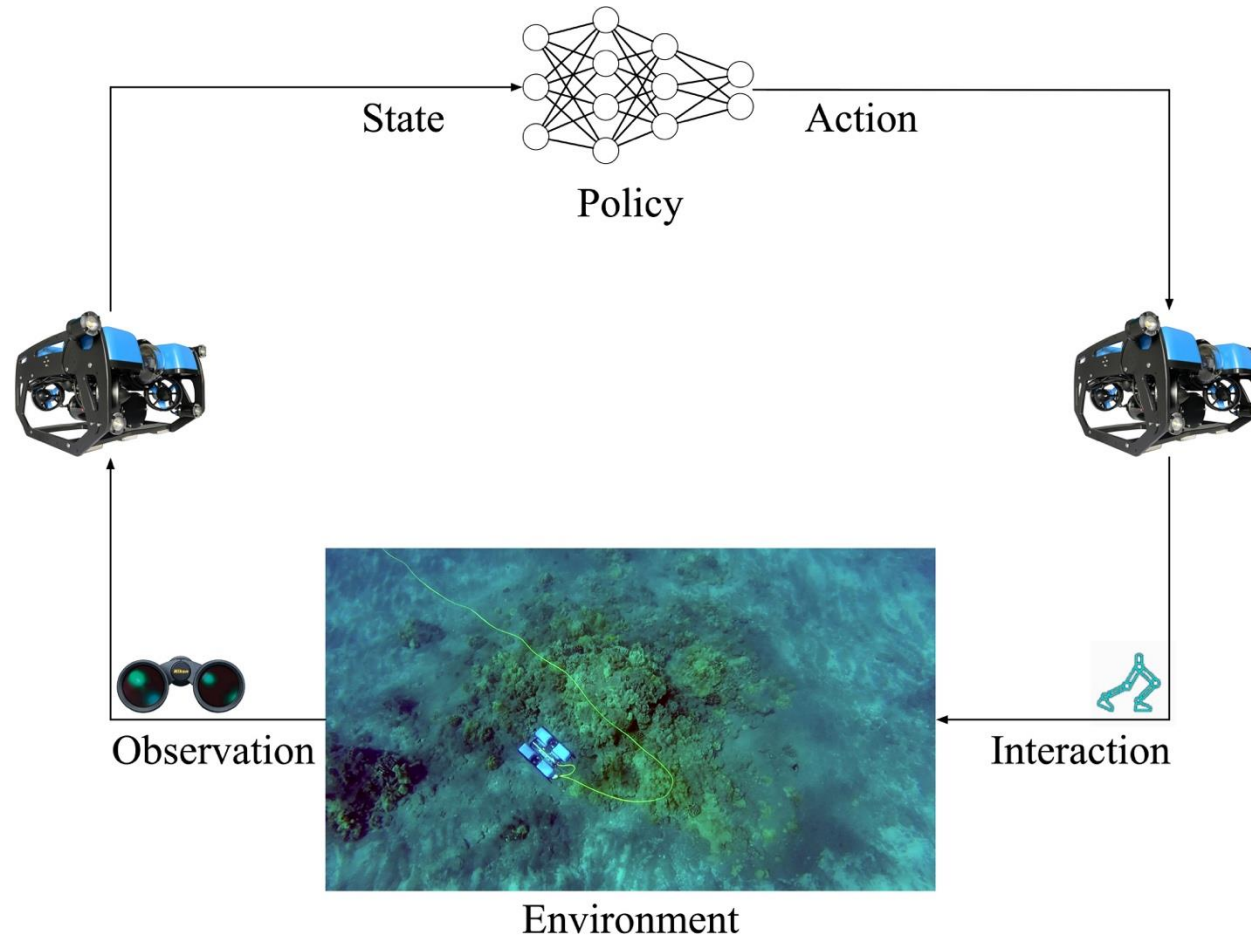
09/06/2025

Dr. Yi Dong & Prof. Xiaowei Huang

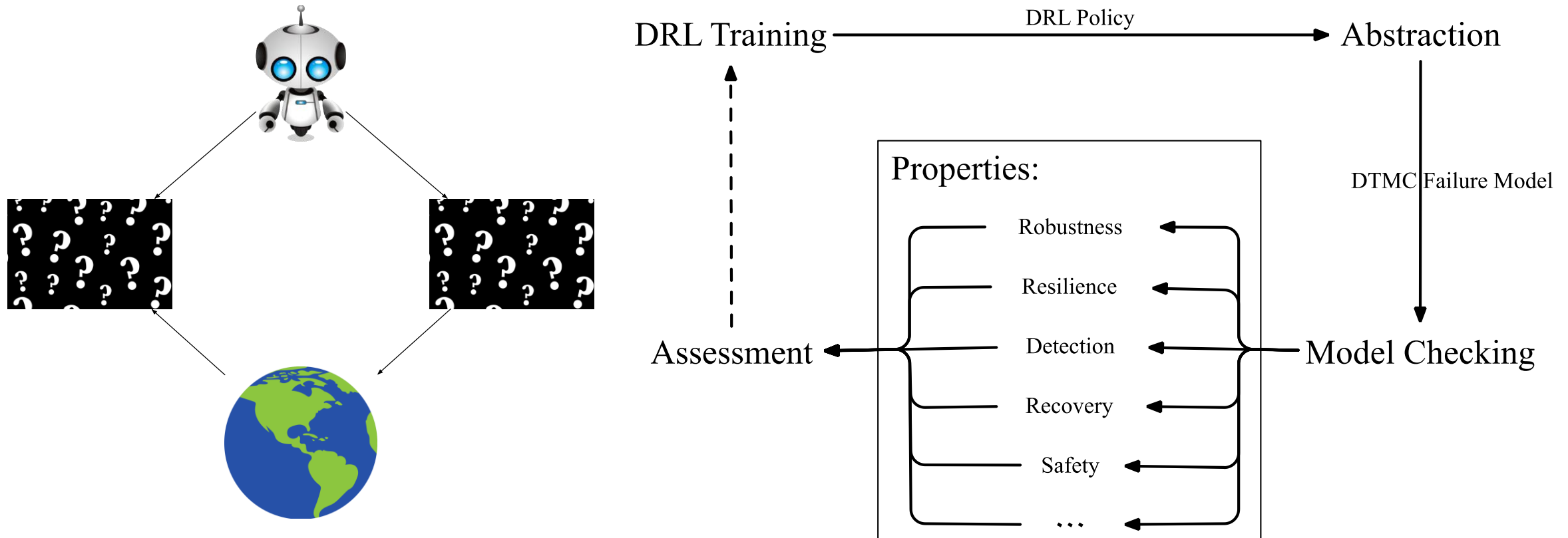
Outline

- Probabilistic Verification of Deep Reinforcement Learning
- Reachability Verification of Deep Reinforcement Learning
- Privacy-preserving Distributed Learning

Deep Reinforcement Learning



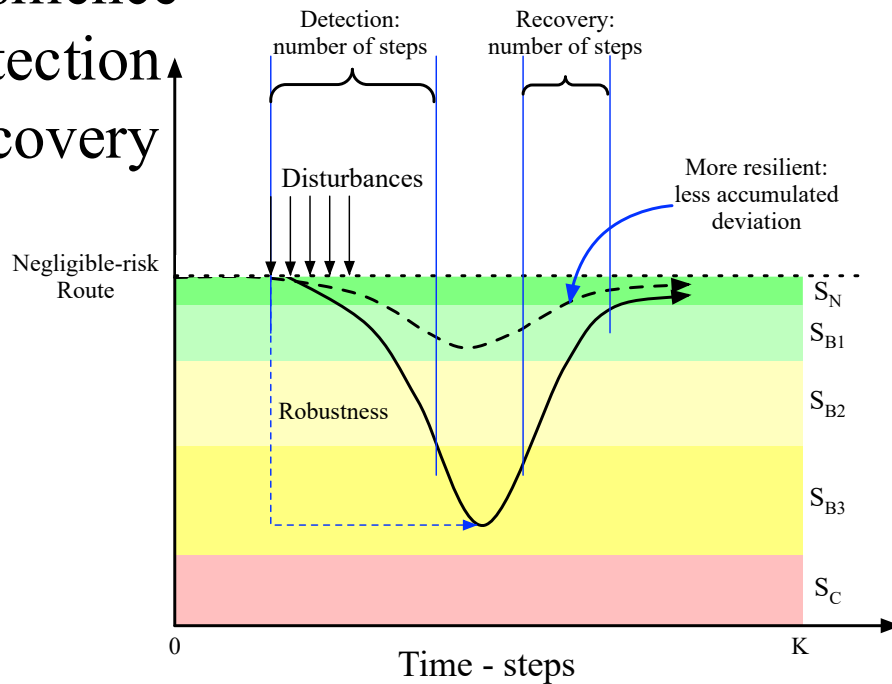
Case 1: Unknown Environment



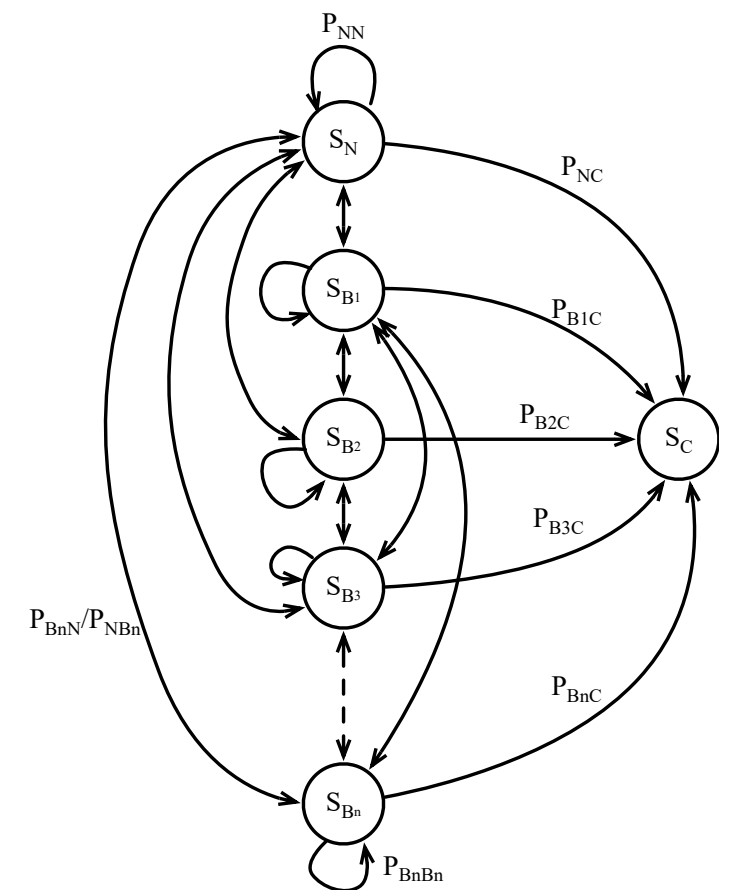
Our Solutions

1. Safety Properties:

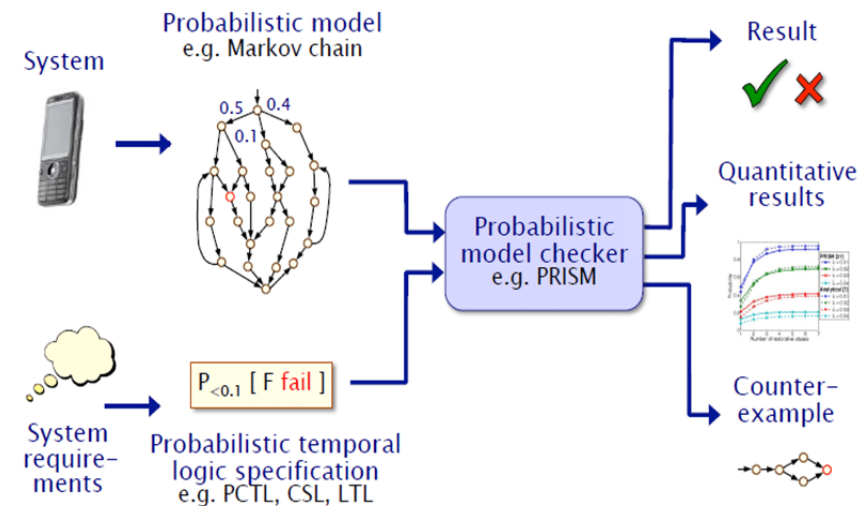
- Robustness
- Resilience
- Detection
- Recovery
- ...



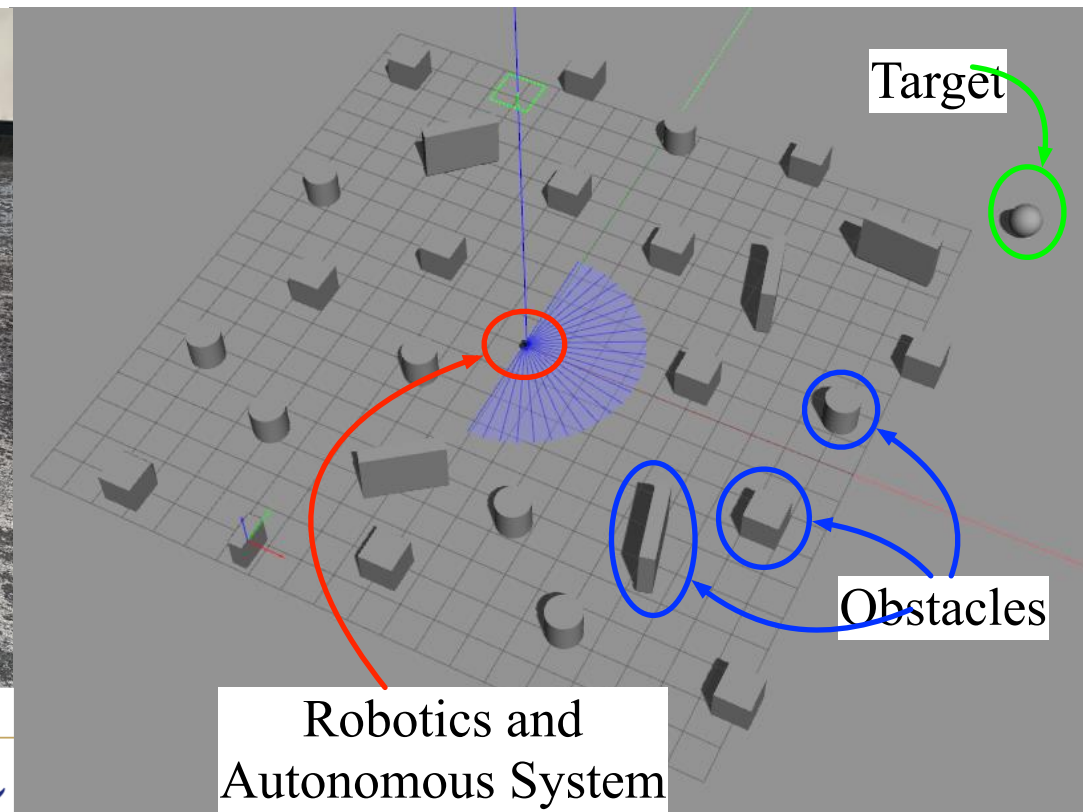
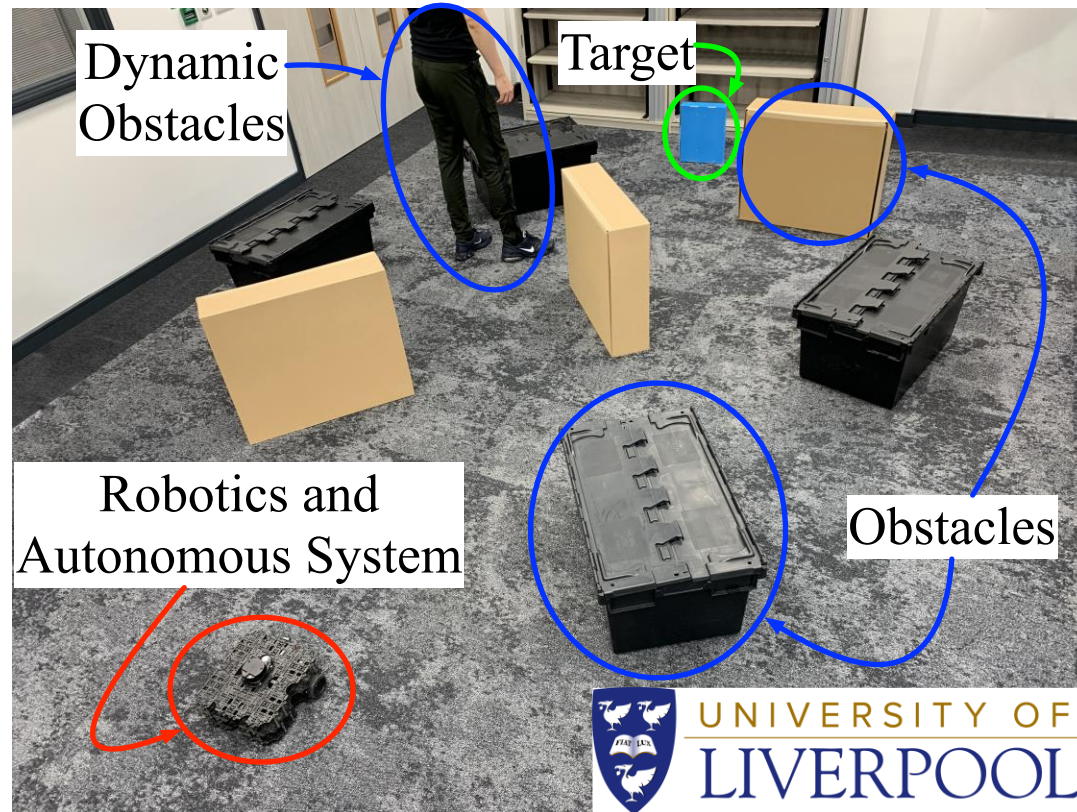
2. DTMC model:



3. Probabilistic Model Checker:

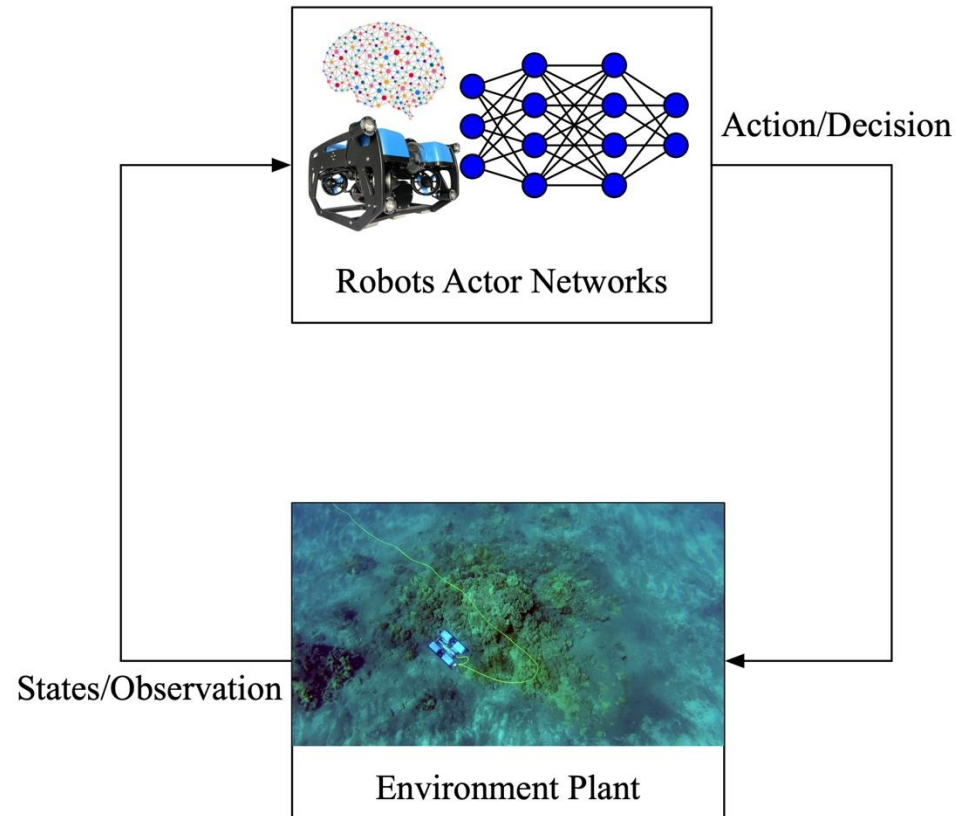


Experiments



Case 2: White-box Verification

Target: Probability of crash per random initial-condition (pci)



Robot Actor Networks:

- Number of Layers
- Number of Neurons
- Weights and Bias

Environment Plant:

- Transition Probability
- Kinetic Model



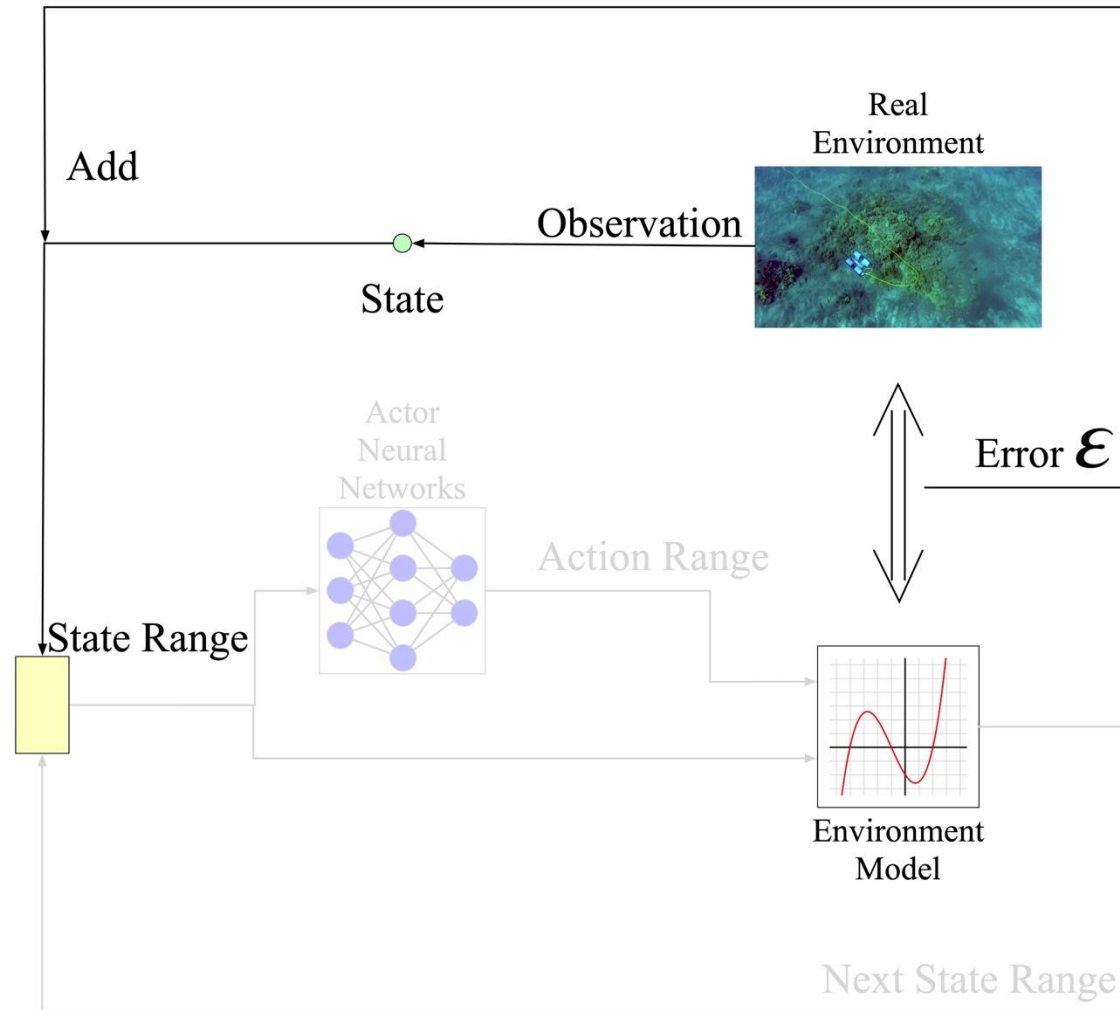
UNIVERSITY OF
LIVERPOOL

[2] Y. Dong, X. Zhao, S. Wang and X. Huang, "Reachability Verification Based Reliability Assessment for Deep Reinforcement Learning Controlled Robotics and Autonomous Systems," in *IEEE Robotics and Automation Letters*, vol. 9, no. 4, pp. 3299-3306, April 2024

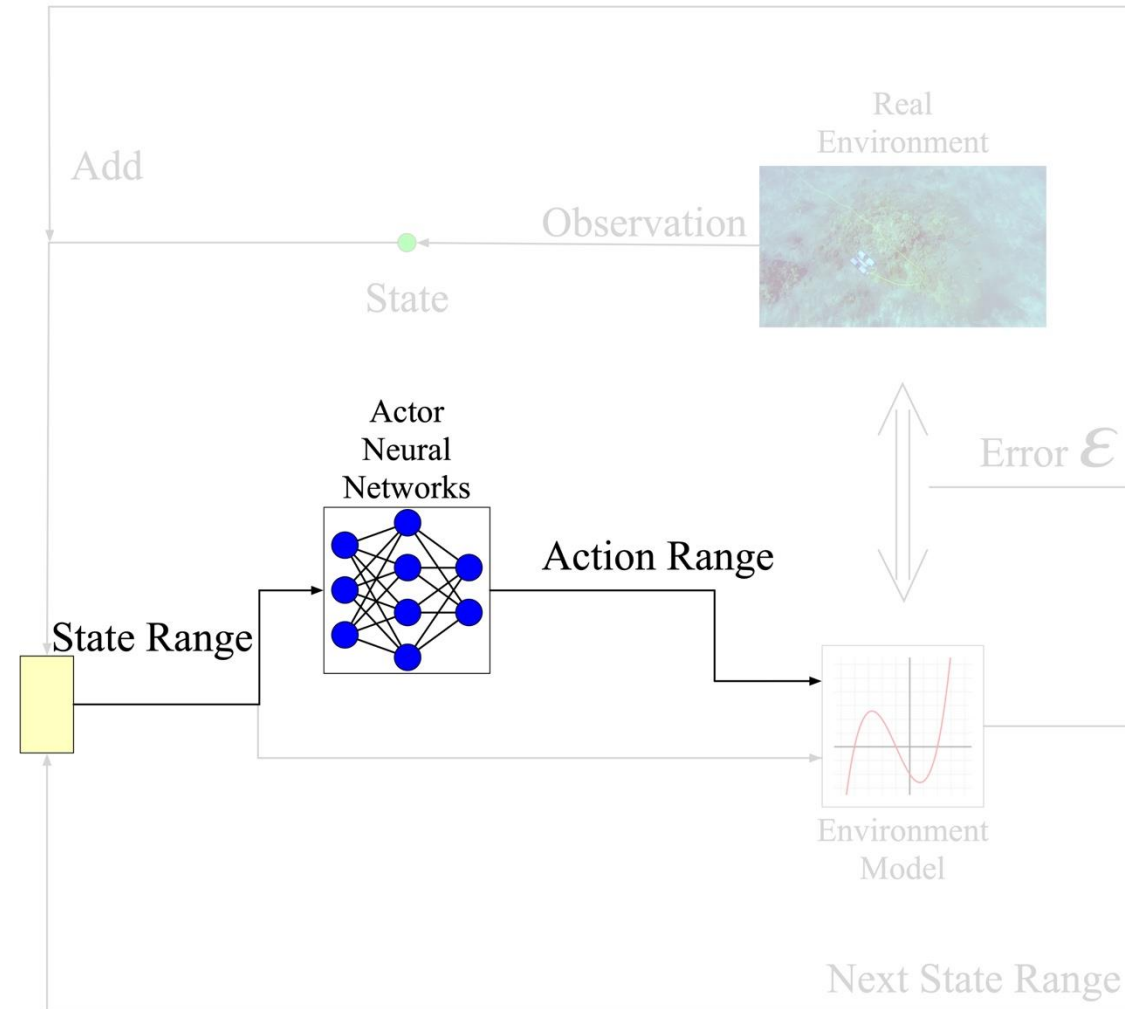


EnnCore

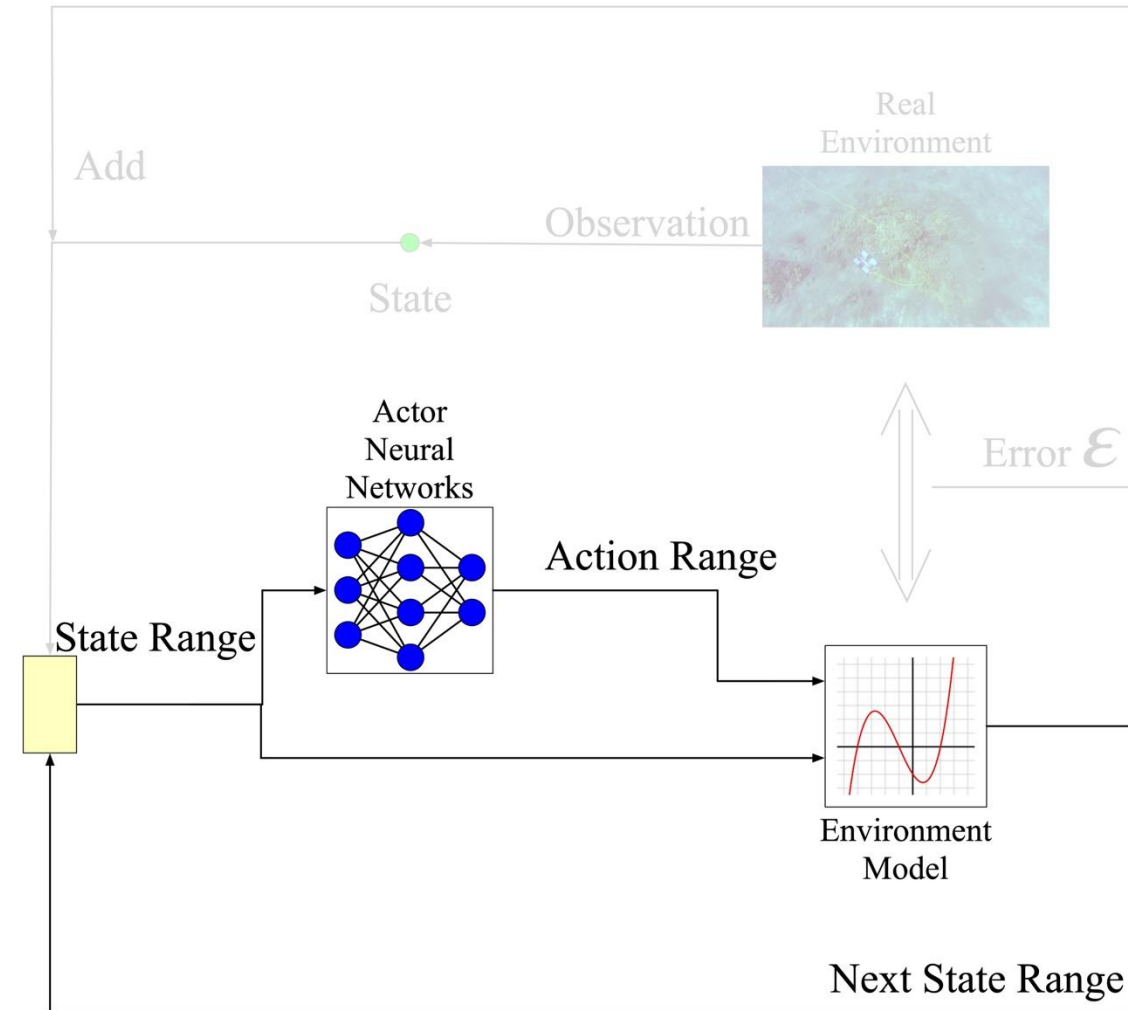
Reachability Verification



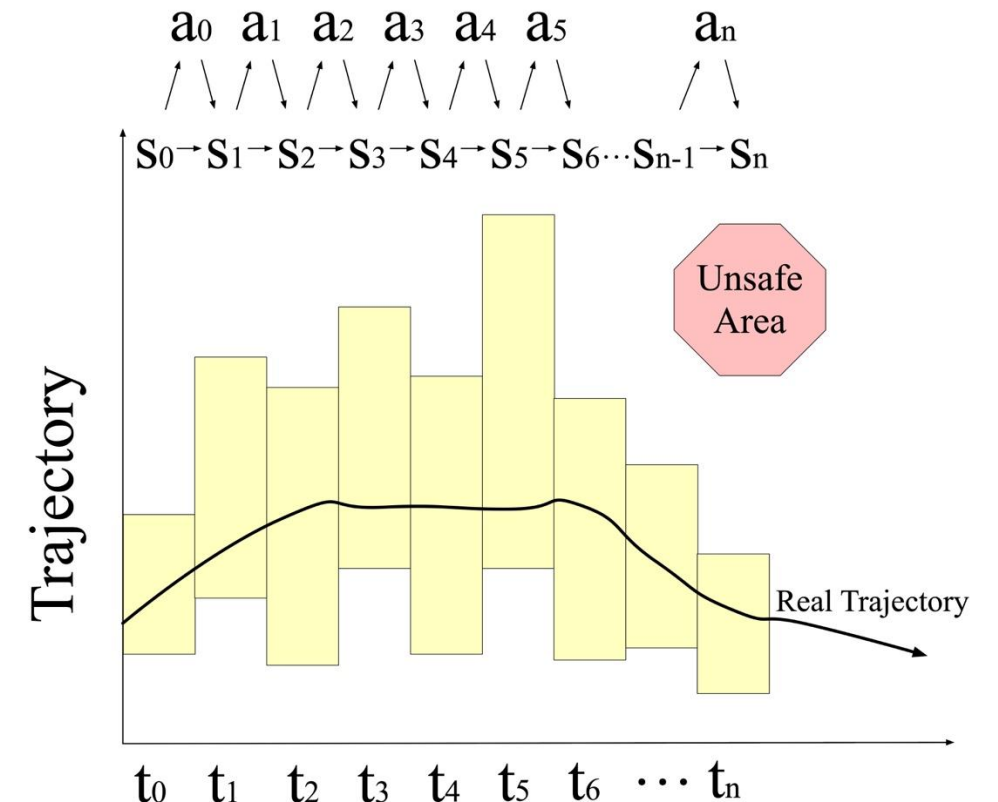
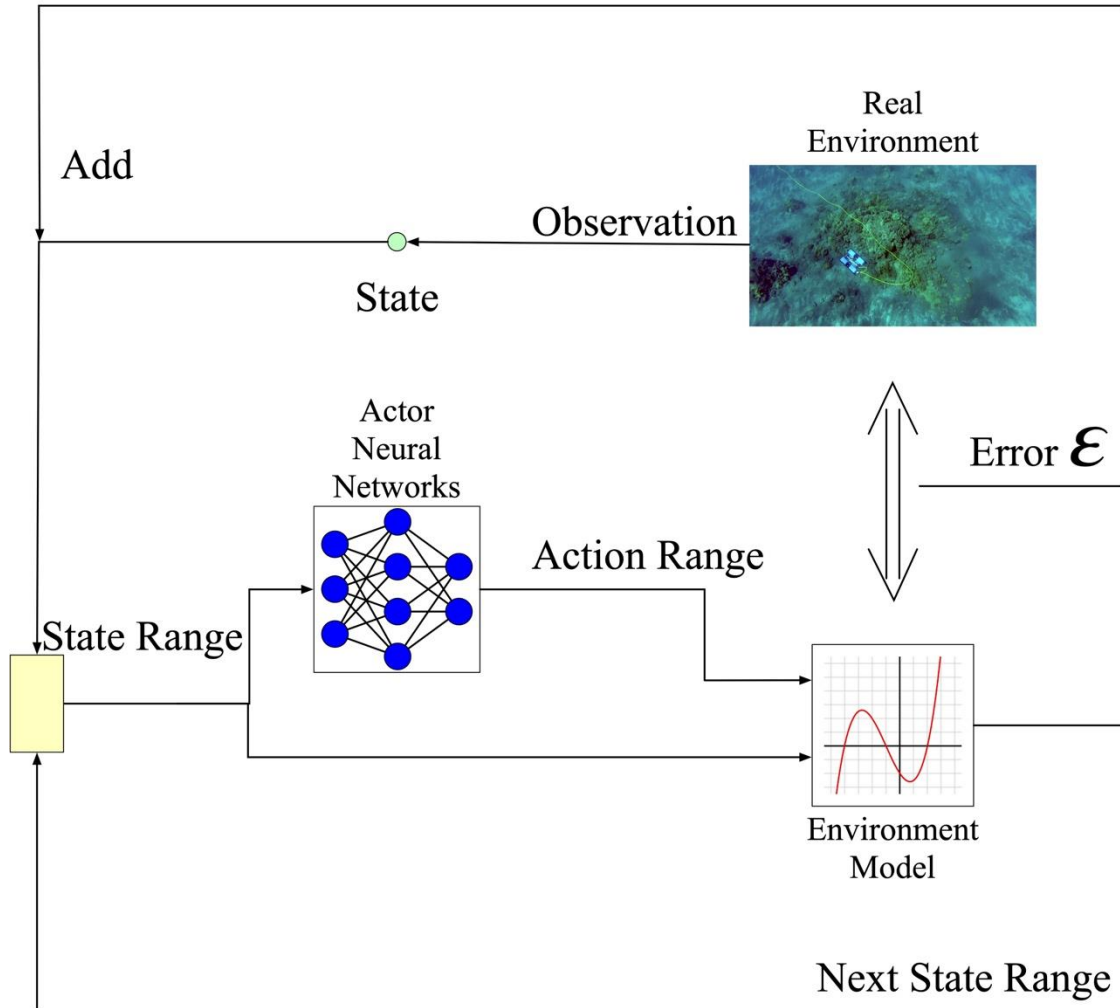
Reachability Verification



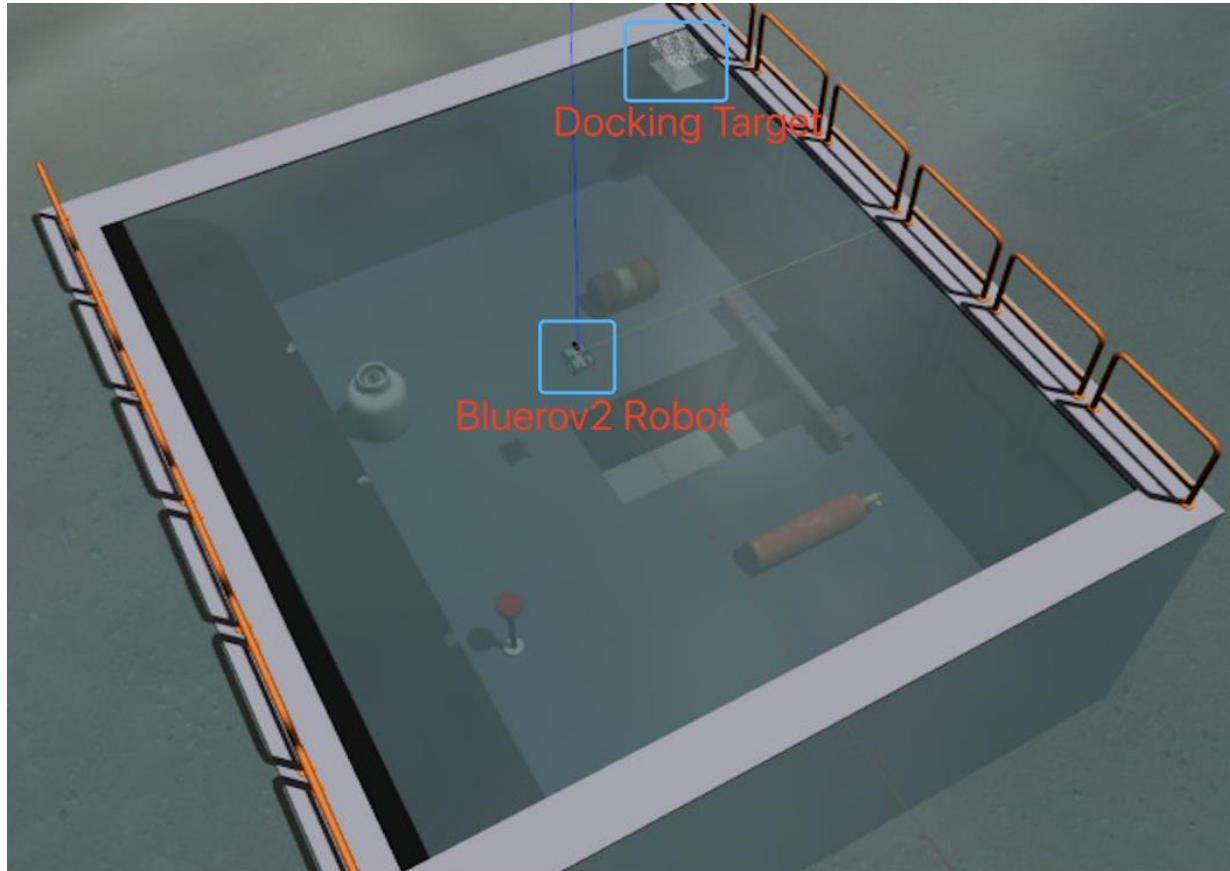
Reachability Verification



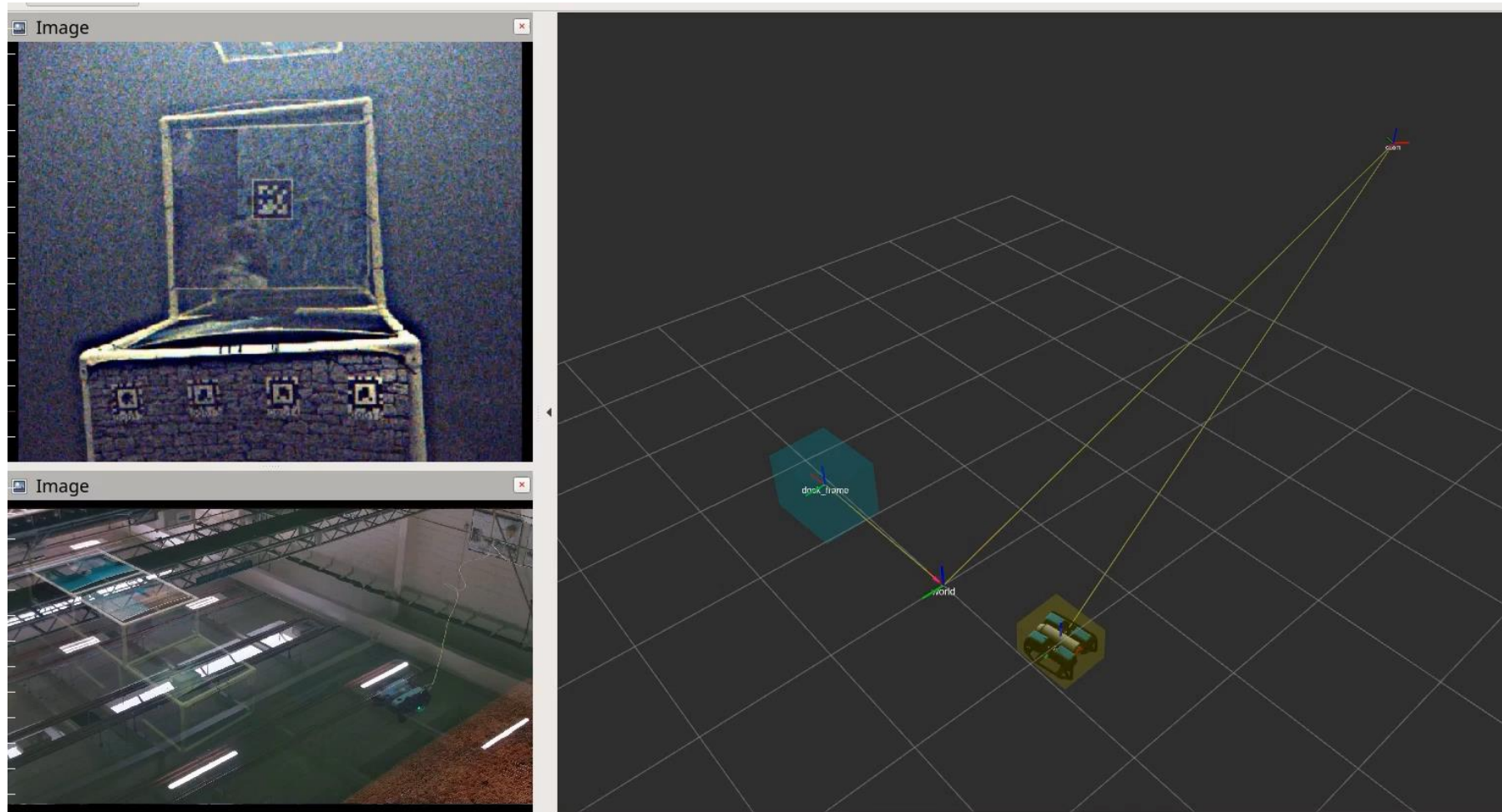
Reachability Verification



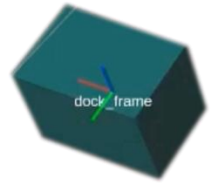
Environment Setup



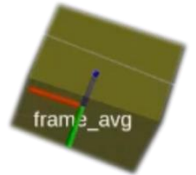
Demo a Safe Route



- Docking Cage



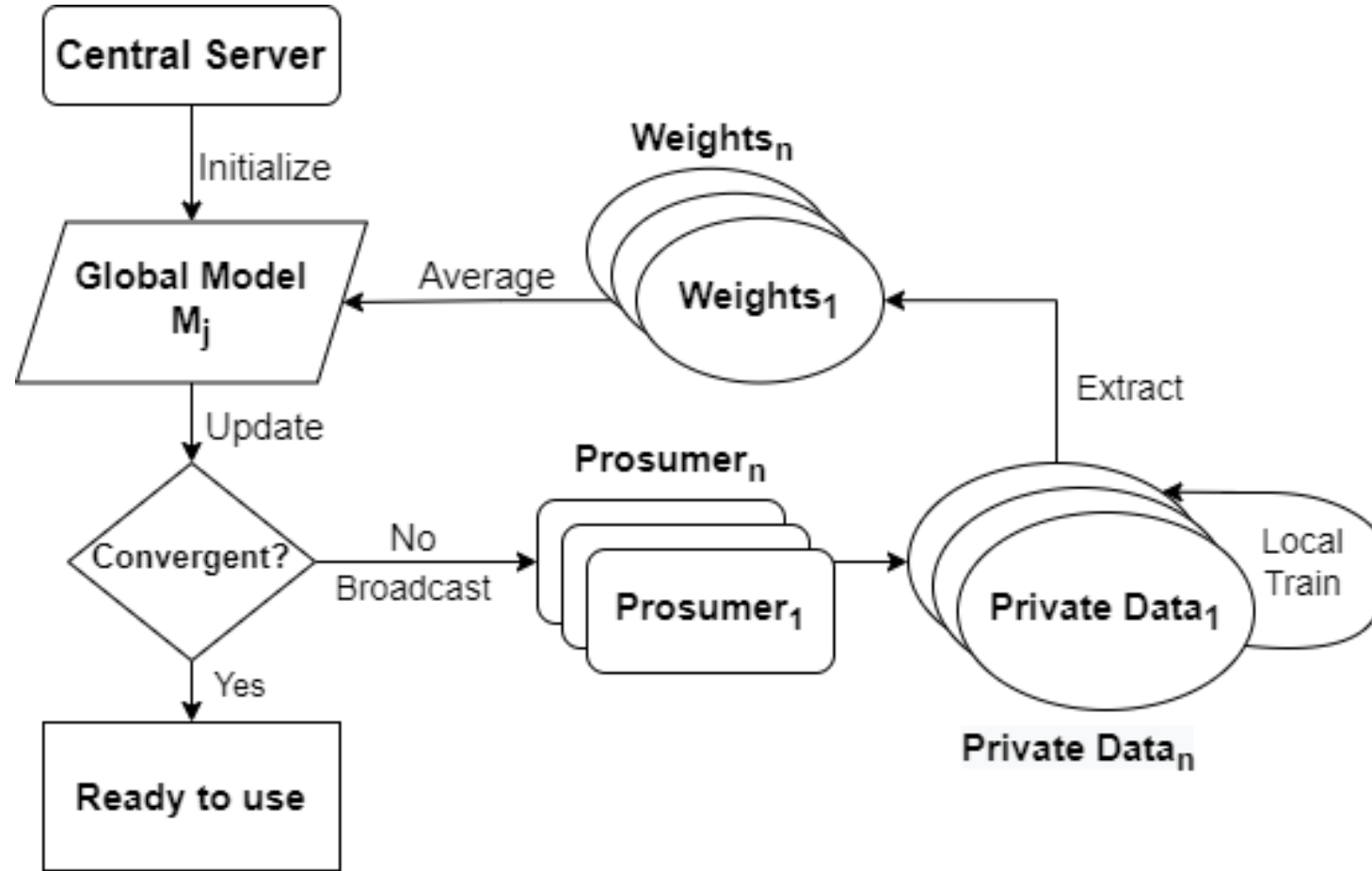
- Reachable Range of UUV



- World Frame (Static)



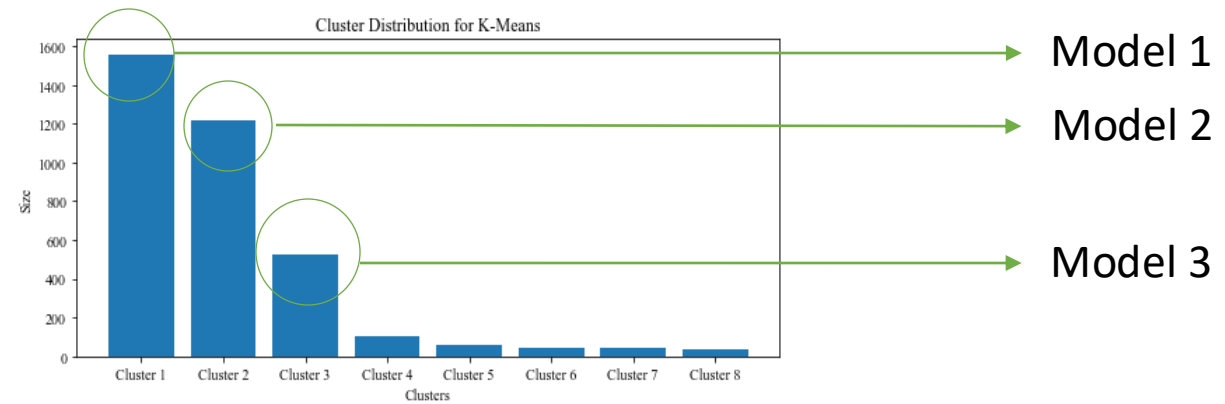
Case 3: Privacy-preserving Distributed Learning



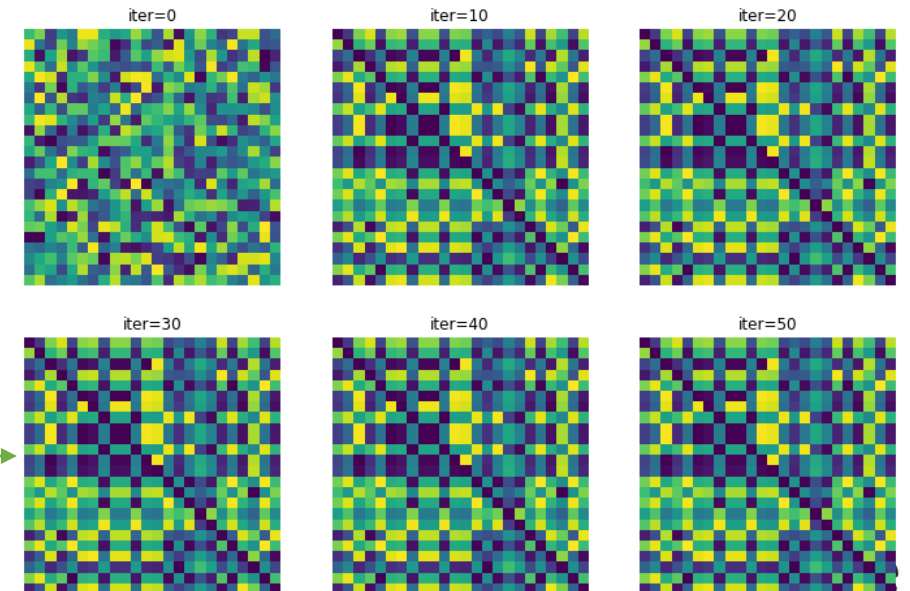
Research Challenges:

- Scalability
- Accuracy
- Robustness
- Privacy

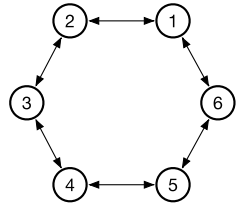
Our Solution



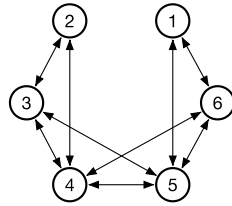
- Scalability: K-means clustering to handle the non-IID dataset.
- Accuracy: 4 types of NN models (DNN, LSTM, CNN, WaveNet)
- Robustness: Damaged training data.
- Privacy: Deep Leakage from Gradients.



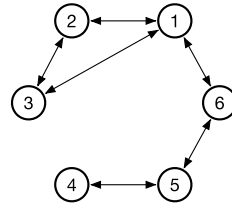
Our Solution



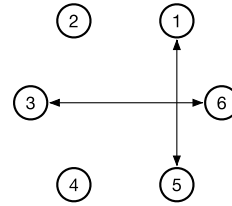
Ring Topology



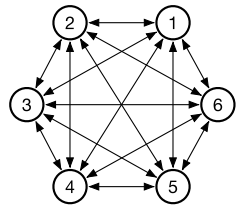
$t = 0$



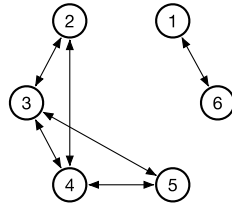
$t = 1$



$t = 2$

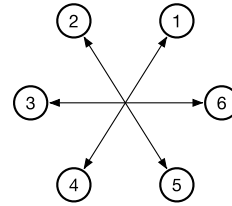


Fully Connected Topology



$t = 3$

...

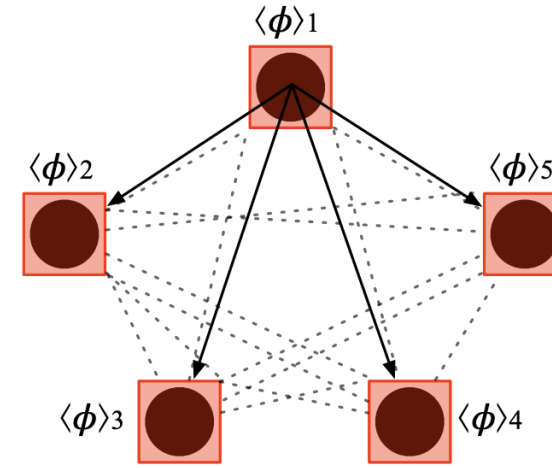


$t = n$

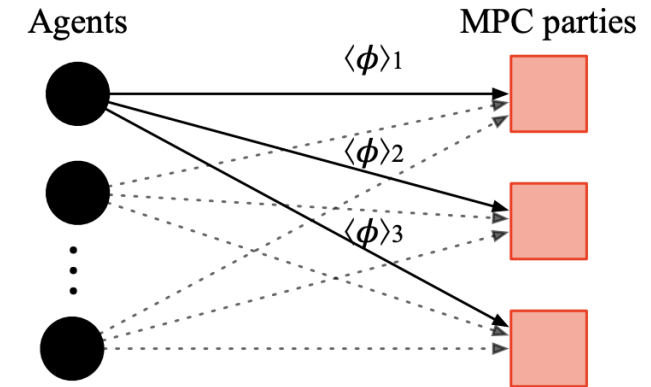
Markovian Switching Topologies

Distributed Frameworks

Agents/MPC parties



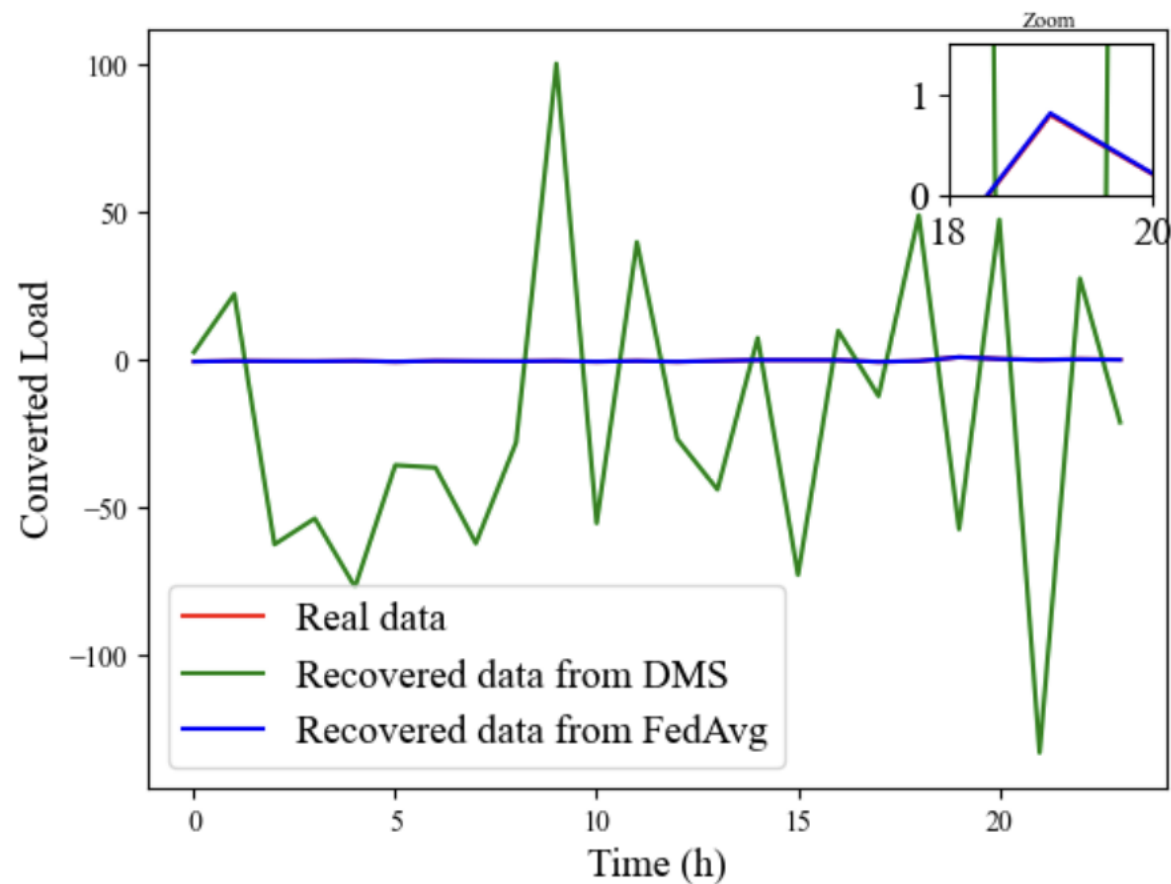
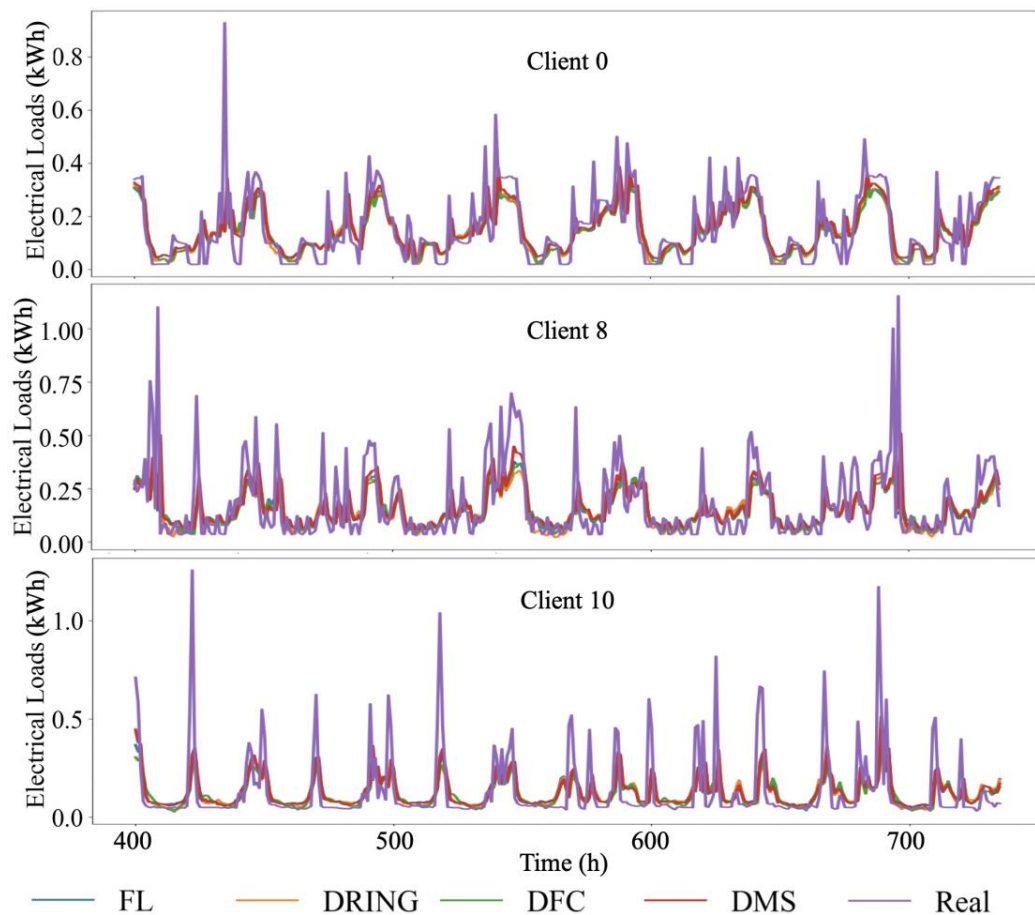
(a) MPC parties played by the agents.



(b) MPC parties played by external servers.

Agent secret sharing a gradient ϕ with the MPC parties.

Experiment Results



Recovered data from DLG attack.

Thank You

- Please refer to our EnnCore project website for more details:
- <https://enncore.github.io/>
- yi.dong@liverpool.ac.uk