



Felippe S. Roza¹, Simon Hadwiger², Ingo Thorn², Karsten Roscher¹

Towards Safety Assurance of Uncertainty-Aware Reinforcement Learning Agents

AI/RL capabilities

Motivation



RL is not the choice for safety-critical applications

Motivation

- Industrial players are historically very conservative
- High complexity of DNNs
- Standard verification and validation are not compatible to Deep RL
- **Unreliable under distributional shifts**



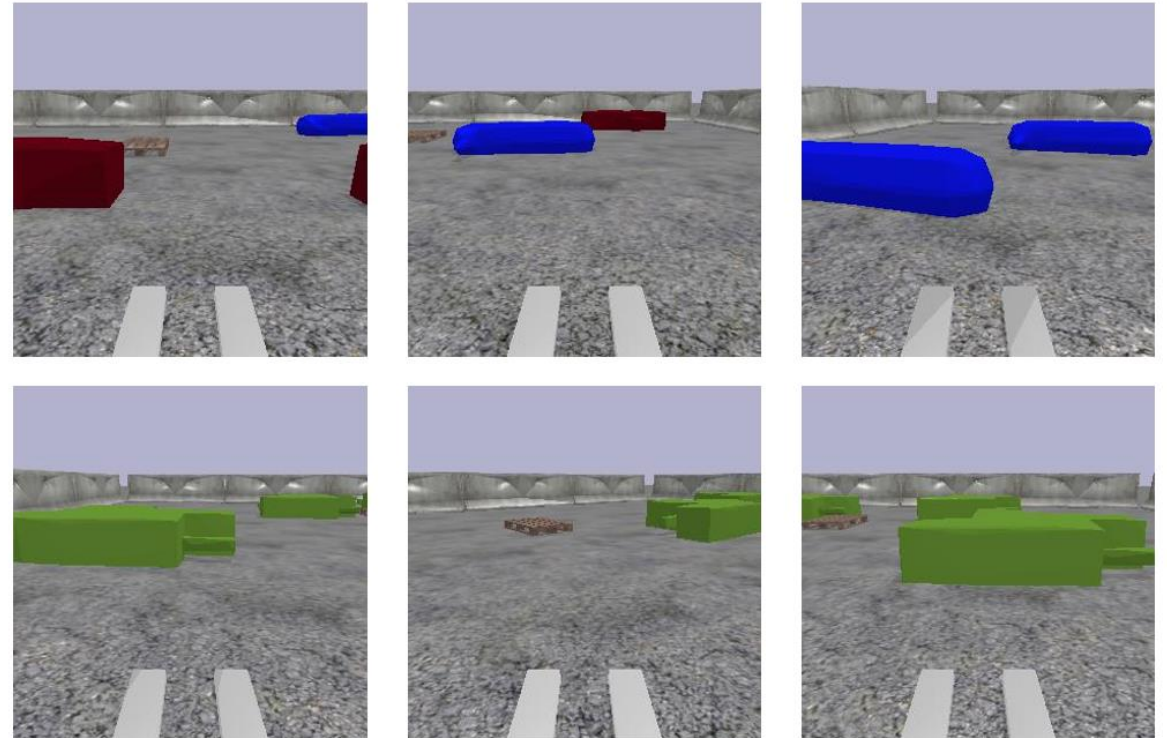
Contribution

As a position paper:

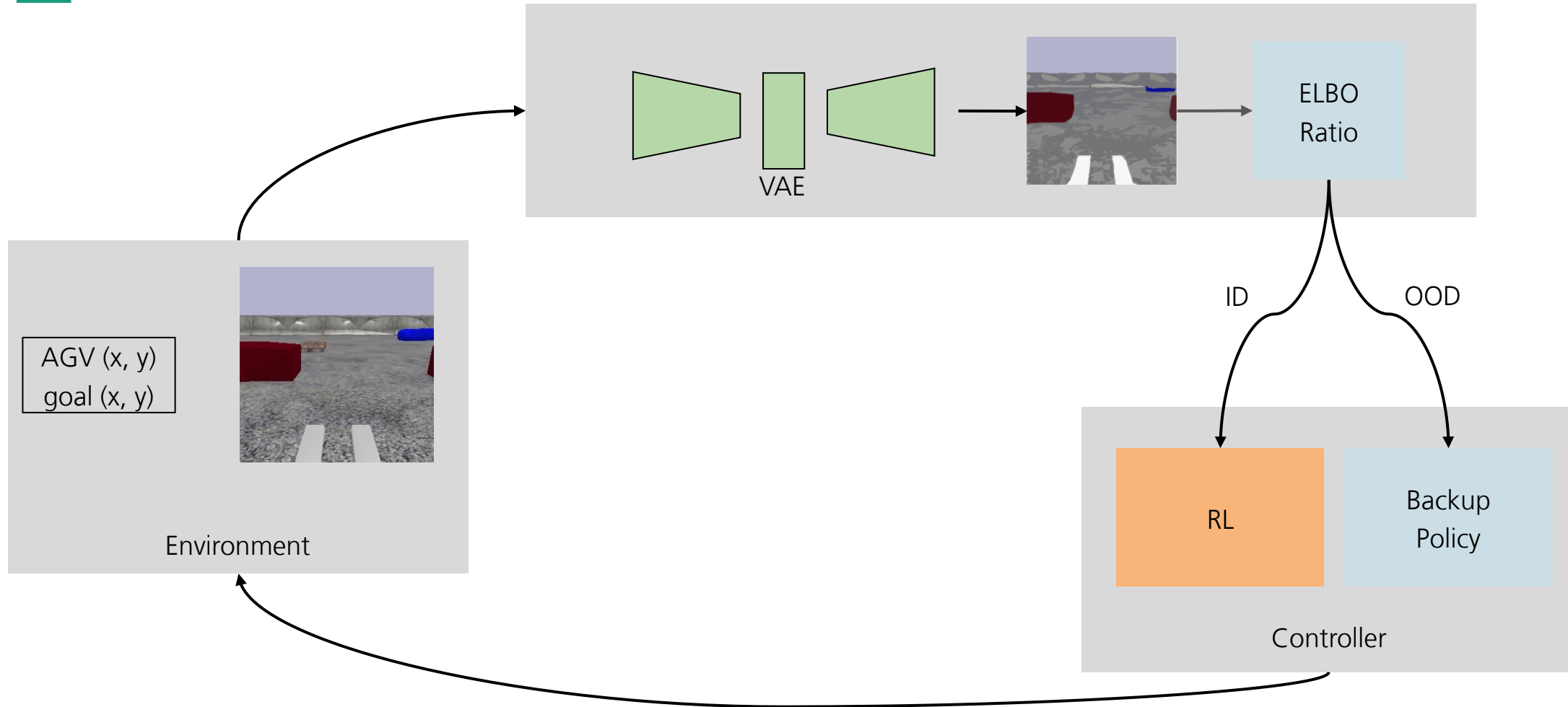
1. Propose an architecture that integrates RL and uncertainty estimation
2. Provide evidence that uncertainty-based OOD detection can help identifying unsafe states
3. Outline future steps towards building a safety case for uncertainty-aware RL systems

Use case

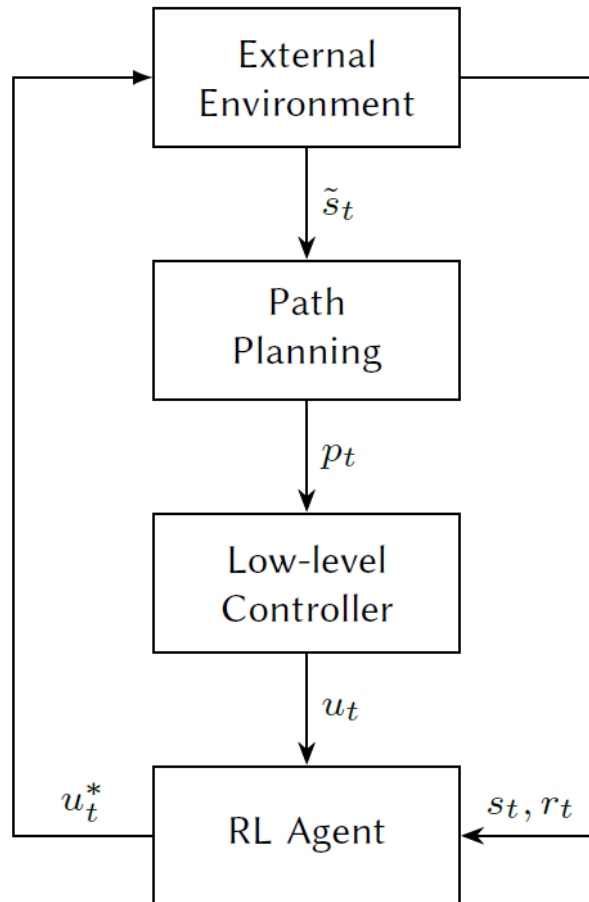
- Moving goods in a warehouse
- Goal: reach the pallet while avoiding obstacles
- Custom AGV/forklift environment
- ID vs OOD obstacles: different colors and shapes



Uncertainty-aware RL



RL-based controller

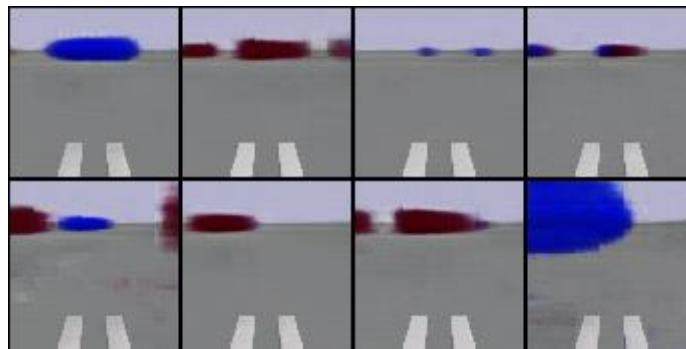
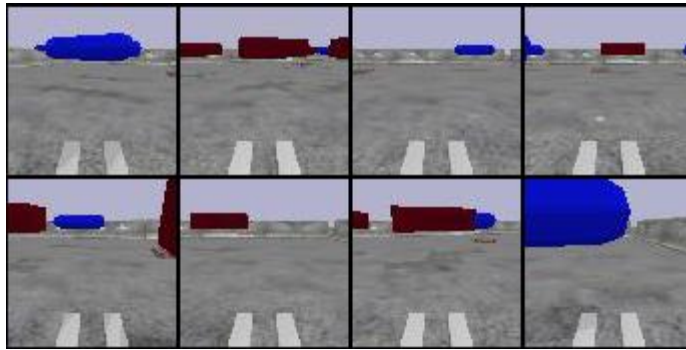


- Path planner calculates optimal trajectory
- Low-level controller calculates control actions
- RL - obstacle avoidance
- Better fit for industrial-grade applications

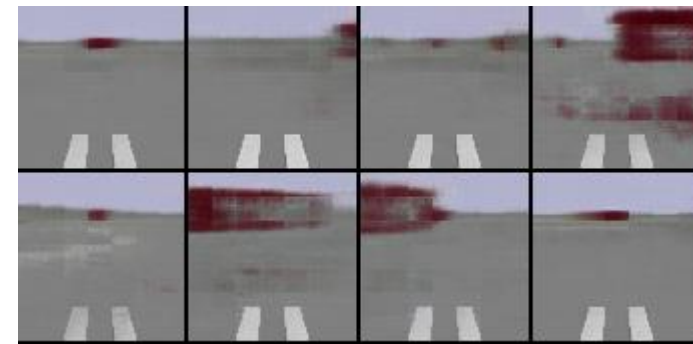
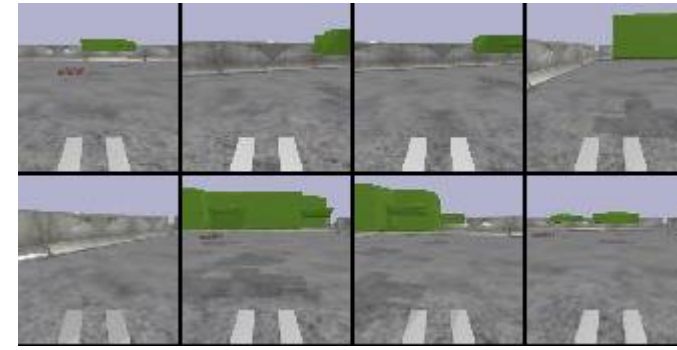
VAE reconstruction

Preliminary results

ID

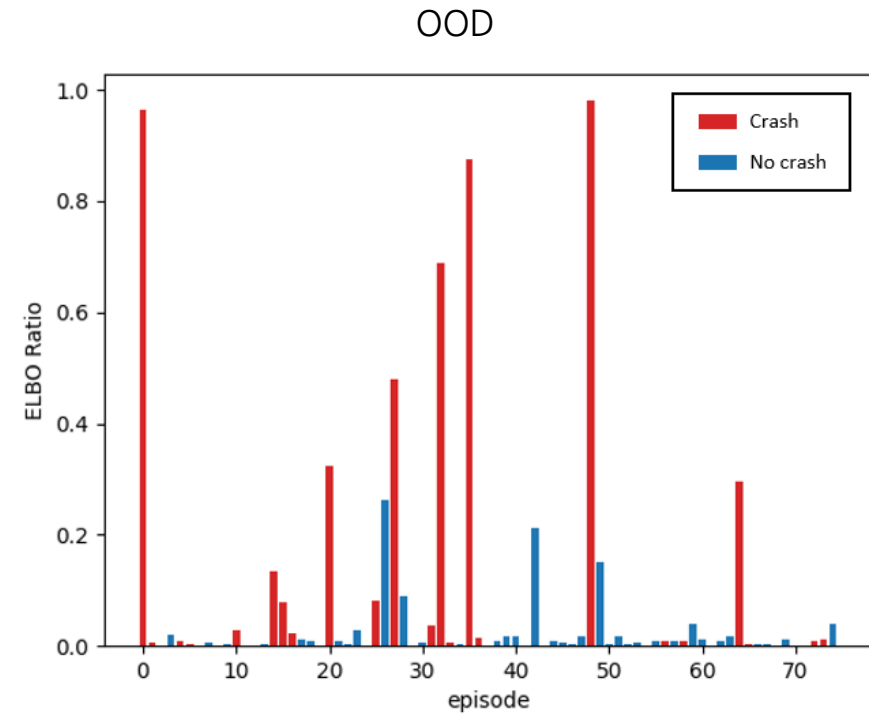
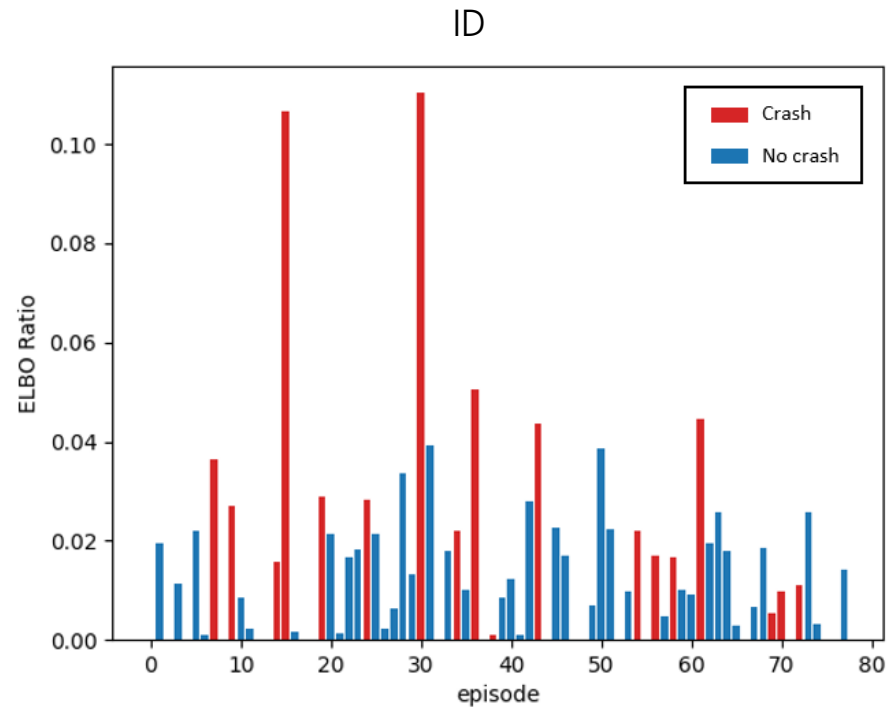


OOD



Uncertainty scores in termination states

Preliminary results



Future perspective

Shortcomings

- Only covers a single aspect of safety in RL
- Limited testing scenarios
- Naïve/simplistic OOD scenarios
- Limitations of VAE as OOD detector

Future steps

- Extensive experimentation
- Reason behind FPs and FNs
- Impact of OOD's residual error
- Failure rate calibration
- SOTIF



SIEMENS

We'd love to hear from you!

Felippe Schmoeller da Roza
felippe.schmoeller.da.roza@iks.fraunhofer.de

<https://iks.fraunhofer.de>
<https://safe-intelligence.fraunhofer.de>

RL reward function

Encourage following optimal trajectory:

$$r = 0.0 \text{ if } u_t^* \cong u_t$$
$$r = -0.1 \text{ otherwise}$$

Accomplish task

$$r = 100 \text{ if reaching the goal}$$
$$r = -10 \text{ if timing out}$$

Obstacle avoidance

$$r = -100 \text{ if crashing}$$