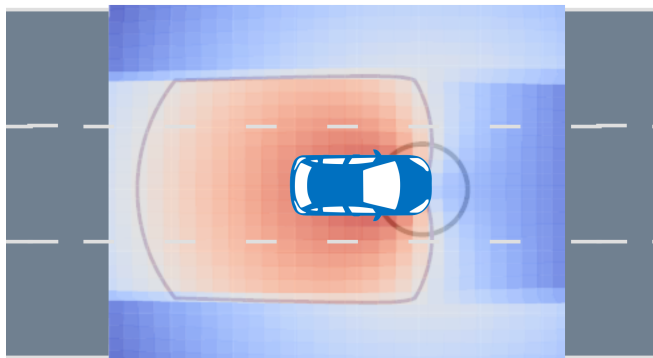**Master Thesis / Bachelor Thesis**

# Safe Deep Reinforcement Learning based on Responsibil Sensitive-Safety

Deep Reinforcement learning is a promising approach for decision-making for autonomous vehicles. However, it is difficult to guarantee safety for such systems. Incorporating safety constraints into the reward function is difficult because then there is always a trade-off between safety and the other objectives, depending on the weights. Moreover, since the policy is learned, there are no guarantees on the output of the neural network. One approach to verify safety of autonomous vehicles is the Responsibility-Sensitive-Safety (RSS) model [1].



Safety index of a car approaching another vehicle

To incorporate safety constraints into reinforcement learning, one line of research uses a safety index that captures safe states [2,3,4]. Then a Lagrangian is optimized that maximizes the reward while not exceeding a threshold on constraints violations. This work should investigate how the typically used safety index can be adapted for the automated driving domain based on the RSS model.

This sounds exciting? Then apply to us! Methods and scope of the thesis can be adapted to your interests and previous knowledge. The proposed thesis consists of the following parts:

+ Literature research about constrained reinforcement learning
+ Design of a domain-specific safety index
+ Training of the approach in different scenarios
+ Evaluation in different scenarios

I am happy to answer any questions you might have. Feel free to ask for an appointment or directly ask at my office!

## References

[1] Shalev-Shwartz, Shammah, and Shashua, *On a Formal Model of Safe and Scalable Self-driving Cars*, 2018
[2] Zhao, He, and Liu, "Model-Free Safe Control for Zero-Violation Reinforcement Learning", 2022
[3] Ma et al., *Joint Synthesis of Safety Certificate and Safe Control Policy Using Constrained Reinforcement Learning*, 2022
[4] Fischer et al., "Safety Reinforced Model Predictive Control (SRMPC): Improving MPC with Reinforcement Learning for Motion Planning in Autonomous Driving", 2023

---

**Institute of Measurement and Control Systems (MRT)**
Prof. Dr.-Ing. Christoph Stiller

**Advisor:**
Johannes Fischer, M.Sc.

**Programming language(s)[1]:**
Python or Julia — advanced

**System, Framework(s):**
Linux

**Required skills:**
- Solid mathematical foundations
- Work on your own

**Language(s):**
German, English

---

For more information please contact:

**Johannes Fischer**

Room: 039 → just come by!
Phone: +49 721 608-48760
Email: johannes.fischer@kit.edu

Or directly send in your application including your current grades as well as our questionnaire!

---

[1] **skill levels:**
*beginner* < 500 lines of code (LOC)
*advanced* 500 – 5000 LOC
*proficient* > 5000 LOC