

Master's Thesis

# Lightweight Attack-Firewall Classifier for Computer Vision Models

Neural networks are gaining more attention in the field of autonomous driving. This is especially true for computer vision-based applications. Since these algorithms are increasingly used in safety-relevant scenarios, it must be ensured that the prediction does not lead to any malfunction of the system. Despite the good generalization capability of convolutional neural networks on unseen input data, a minimal perturbation caused by an adversarial attack can disturb the model and impair its prediction. In the context of this work, an efficient agent is to be learned that detects potential attacks on the underlying CNN and provides appropriate countermeasures to ensure trouble-free operation.

## Prerequisites

### Prerequisites

To successfully complete this project, you should have the following skills and experiences:

- Good programming skills in Python and Tensorflow
- Good knowledge of neural networks, basic in adversarial attacks

The student is expected to be highly motivated and independent.

## Contact

### Nael Fafous

Department of Electrical and Computer Engineering  
Chair of Integrated Systems

**Phone:** +49.89.289.23858

**Building:** N1 (Theresienstr. 90)

**Room:** [N2116](#)

**Email:** [nael.fafous@tum.de](mailto:nael.fafous@tum.de)

## Advisors

Nael Yousef Abdullah Al-Fafous, Alexander Frickenstein