

JONATHAN NÖTHER | Curriculum Vitae

✉ s8jonoet@stud.uni-saarland.de • Ottweilerstraße 38 / St. Wendel • 📞 01773 5507105
🐙 GitHub

INTERESTS

Secure Machine Learning, Attacks against ML Models, AI Safety, Reinforcement Learning

EDUCATION

SAARLAND UNIVERSITY

M.Sc. in Data Science and Artificial Intelligence
ECTS: 1.7

12/2022 - Ongoing

Saarbrücken, Germany

SAARLAND UNIVERSITY

B.Sc. in Data Science and Artificial Intelligence
ECTS: 1.7

10/2019 - 11/2022

Saarbrücken, Germany

EXPERIENCE

RESEARCH ASSISTANT

Machine Teaching Group

Conducted research projects and presented my work and related papers

08/2022-07/2024

MPI-SWS

TEACHING EXPERIENCE

TEACHING ASSISTANT FOR THE SEMINAR "TRUSTWORTHINESS OF FOUNDATION MODELS"

Multi-Agent Systems Group

Prepare the seminar project on red-teaming and watermarking of foundation models

Summer 2024

MPI-SWS

TEACHING ASSISTANT FOR THE LECTURE "STATISTICS LAB"

Modeling and Simulation Group

Explained course topics to students and graded tests and exams

Summer 2022

Saarland University

TEACHING ASSISTANT FOR THE LECTURE "ARTIFICIAL INTELLIGENCE"

Foundations of Artificial Intelligence Group

Prepared Exercises, explained Course topics to students, and graded tests and exams

Summer 2022

Saarland University

TEACHING ASSISTANT FOR "PROGRAMMING 1"

Reactive Systems Group

Prepared Exercises, explained course topics to students, and graded tests and exams

Winter 2019/2020

Saarland University

SKILLS

PROGRAMMING LANGUAGE CONCEPTS

Experienced: Python

Familiar: C++

Experienced: Machine Learning | LLMs | Reinforcement Learning | Adversarial ML

Familiar: Cybersecurity | Computer Vision | Diffusion Models

LIBRARIES LANGUAGES

matplotlib | Pytorch | numpy

Native: German | Fluent: English (C1)

PUBLICATIONS

TEXT-DIFFUSION RED-TEAMING OF LARGE LANGUAGE MODELS: UNVEILING HARMFUL BEHAVIORS WITH PROXIMITY CONSTRAINTS

Under Review at COLM 2024

Safety of LLMs

DEFENDING AGAINST UNKNOWN CORRUPTED AGENTS: REINFORCEMENT LEARNING OF ADVERSARIALLY ROBUST NASH EQUILIBRIA

Under Review at TMLR

Robust Reinforcement Learning

IMPLICIT POISONING ATTACKS IN TWO-AGENT REINFORCEMENT LEARNING: ADVERSARIAL POLICIES FOR TRAINING-TIME ATTACKS

AAMAS 2023

PDF
Adversarial Reinforcement Learning

PROJECTS

INTERVIEW PERFORMANCE PREDICTION AND LIE DETECTION

Project Seminar Data Science and Artificial Intelligence

Implementation of model that evaluated the performance and detected lies of a participant of mock-job interviews. Grade : 1.0

Type: Computer Vision, NLP

INPAINTING DETECTION

High Level Computer Vision Course

Combine automatic segmentation with inpainting to automatically create edited images.

Additionally experimented with detecting these faked images. Grade : 1.0

Type: Computer Vision, Generative AI



REINFORCEMENT LEARNING PROJECT

Reinforcement Learning Course

Implementation of a RL-agent that solves the gridworld. Grade(course) : 1.0

Type: Reinforcement Learning



SAFE STREETS

AI for the Social Good Seminar

Extend pedestrian route recommendation by taking into account the safety of the route

(e.g. lights, open shops). Grade : 1.0

Type: Data Science, Geospatial Data