# JONATHAN NÖTHER | Curriculum Vitae

✉ jnoether@mpi-sws.org     • 🐙 GitHub     • 🔗 LinkedIn

## INTERESTS

Secure Machine Learning, Attacks against ML Models, Reinforcement Learning, Agentic Systems

## EDUCATION

**MAX PLANCK INSTITUTE FOR SOFTWARE SYSTEMS**
PhD in Computer Science

10/2024 - Ongoing
Saarbrücken, Germany

**SAARLAND UNIVERSITY**
M.Sc. in Data Science and Artificial Intelligence
ECTS: 1.3

12/2022 - 08/2024
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 COURSE "GENERATIVE AI"**
Machine Teaching Group
Prepare exercise sheets, answered student's questions and graded exercises and the exam

Winter 2024/2025
MPI-SWS

**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 2020/2021
Saarland University

# SKILLS

| | |
|---|---|
| PROGRAMMING LANGUAGES | **Experienced:** Python    **Familiar:** C++ |
| CONCEPTS | **Experienced:** Machine Learning \| Reinforcement Learning \| Adversarial ML \| Agentic Systems \| Large Language Models |
| | **Familiar:** Cybersecurity \| Computer Vision \| Diffusion Models |
| LIBRARIES | matplotlib \| Pytorch \| numpy \| AutoGen \| TRL \| transformers |
| LANGUAGES | **Native:** German \| **Fluent:** English (C1) |

# PUBLICATIONS

### MAMA: A GAME-THEORETIC APPROACH FOR DESIGNING SAFE AGENTIC SYSTEMS
PDF
Preprint, under Review                                                                            Safety of Agentic Systems
**TL;DR:** Automatic Design of Safe Agentic Systems using a two-player game between a system designer and an attacker

### AGENTICRED: OPTIMIZING AGENTIC SYSTEMS FOR AUTOMATED RED-TEAMING
PDF
Preprint, under Review                                                                                    Jailbreaking of LLM's
**TL;DR:** Automatically design red-teaming workflows without human intervention

### BENCHMARKING THE ROBUSTNESS OF AGENTIC SYSTEMS TO ADVERSARIALLY-INDUCED HARMFUL ACTIONS
PDF
Preprint, under Review                                                                            Safety of Agentic Systems
**TL;DR:** Benchmark for testing the robustness of LLM-based agents against adversaries that aim to manipulate them into performing dangerous actions

### TEXT-DIFFUSION RED-TEAMING OF LARGE LANGUAGE MODELS: UNVEILING HARMFUL BEHAVIORS WITH PROXIMITY CONSTRAINTS
PDF
AAAI 2025 (Oral)                                                                                              Safety of LLMs
**TL;DR:** Applying text-diffusion models to red-teaming to satisfy proximity constraints with regards to a reference prompt

### POLICY TEACHING VIA DATA POISONING IN LEARNING FROM HUMAN PREFERENCES
PDF
AISTATS 2025                                                                                                  Safety of LLMs
**TL;DR:** Forcing an LLM to adapt a target policy by synthesizing preference data

### DEFENDING AGAINST UNKNOWN CORRUPTED AGENTS: REINFORCEMENT LEARNING OF ADVERSARIALLY ROBUST NASH EQUILIBRIA
PDF
TMLR 08/2024                                                                            Robust Reinforcement Learning
**TL;DR:** Training robust agents in an MARL setting where an attacker can arbitrarily corrupt a subset of peer agents of a given cardinality

### IMPLICIT POISONING ATTACKS IN TWO-AGENT REINFORCEMENT LEARNING: ADVERSARIAL POLICIES FOR TRAINING-TIME ATTACKS
PDF
AAMAS 2023                                                                            Adversarial Reinforcement Learning
**TL;DR:** Attacking an agent by poisoning the policy of a peer agent during training

# AWARDS

### GÜNTER-HOTZ MEDAILLE
Award for the top computer science graduates of Saarland University

### TOP-REVIEWER
Neurips 2025

# REVIEWING

AAAI 2024 \| AAAI 2025 \| NeurIPS 2025 \| ICLR 2026