

Ilias Tsingenopoulos

Postdoctoral Researcher | AI Safety & Security
Department of Computer Science | KU Leuven, Belgium

+32486319227 | ✉ ilias.tsingenopoulos@cs.kuleuven.be | 🌐 [Webpage](#)
🌐 [LinkedIn](#) | 📄 [itsiggen](#) | 📖 [Google Scholar](#)

ABOUT ME

I work on the intersection of AI and Computer Security, specializing in Adversarial Machine Learning and Reinforcement Learning. My research spans the theoretical and practical aspects of adversarial attacks and defenses across a broad range of AI systems and modalities: from bot detection like reCaptcha, to hardening commercial antivirus against adversarial malware, and more generally on adaptation and optimization of attacks and defenses against each other under the competitive game they form.

I investigate the fundamentals of robust learning and its adversarial and counterfactual aspects, essential components towards achieving trustworthy and safe AI. As human decision-making increasingly becomes automated, I am developing new techniques and methodologies for training models resilient to evasion and other failure modes. Currently, I focus on safeguarding LLMs through adversarial training, alignment, and activation engineering to ensure safe and correct outputs in real-world deployments.

PROFESSIONAL & RESEARCH EXPERIENCE

DistriNet, KU Leuven

Postdoctoral Researcher in AI Safety & Security

- Developing principled approaches for robust LLM generation and classification.
- Investigating distinct approaches in model hardening: problem- and feature-space based.

Leuven, Belgium

October 2024 – Present

DistriNet, KU Leuven

Doctoral Researcher in Robust & Secure AI

- Adversarial attack adaptation and optimization.
- Model hardening and active defense development.
- Domains: image classification, bot & malware detection.

Leuven, Belgium

May 2019 – September 2024

S2Lab, University College London

Visiting Scholar

- Rendering a commercial antivirus robust to evasive malware.

London, UK

January, 2023 – April 2023

DistriNet, KU Leuven

Research Assistant

- Optimization of black-box adversarial attacks.
- Preventing abusive DNS registrations in the .eu domain.

Leuven, Belgium

May, 2018 – April 2019

Information Technologies Institute, CERTH

Research Assistant

- Research and implementation in several Horizon 2020 projects.
- Formulating and writing research proposals.

Thessaloniki, Greece

January, 2017 – April 2018

EDUCATION

Aristotle University of Thessaloniki

Diploma (5-year degree, M.Eng equivalent) in Electrical & Computer Engineering

Dissertation: Fuzzy Clustering Algorithms in Feature Subspaces

Thessaloniki, Greece

Graduated Jul. 2016

Technical University of Berlin

Exchange Semester, Department of Electrical Engineering & Computer Science

Berlin, Germany

Spring 2012

SELECTED PUBLICATIONS

- **I. Tsingenopoulos**, J. Cortellazzi, B. Bosansky, S. Aonzo, D. Preuveneers, W. Joosen, F. Pierazzi, and L. Cavallaro. How to Train your Antivirus: RL-based Hardening through the Problem-Space. In: 27th International Symposium on Research in Attacks, Intrusions and Defenses (RAID 2024).
- **I. Tsingenopoulos**, V. Rimmer, D. Preuveneers, F. Pierazzi, L. Cavallaro, and W. Joosen. On Adaptive Decision-Based Attacks and Defenses. In: DLSP Workshop, IEEE S&P 2024.
- **I. Tsingenopoulos**, D. Preuveneers, L. Desmet, and W. Joosen. Captcha me if you can: Imitation Games with Reinforcement Learning. In: 2022 IEEE 7th European Symposium on Security and Privacy (EuroS&P 2022).
- **I. Tsingenopoulos**, A. M. Shafiei, L. Desmet, D. Preuveneers, and W. Joosen. Adaptive Malware Control: Decision-Based Attacks in the Problem Space of Dynamic Analysis. In: 1st Workshop on Robust Malware Analysis (WoRMA 2022).
- C. J. Hernández-Castro, Z. Liu, A. Serban, **I. Tsingenopoulos**, and W. Joosen. Adversarial Machine Learning. In: Security and Artificial Intelligence: A Crossdisciplinary Approach. Springer, 2022.

KNOWLEDGE & TECHNICAL SKILLS

Domain Expertise:

- **Machine Learning:** Clustering, SVMs, Decision Trees/Random Forests, XGBoost
- **Deep Learning:** Convolutional/Graph/Recurrent Neural Networks, GANs
- **Reinforcement Learning:** Q-Learning, Policy Optimization, POMDPs, Multi-agent
- **Adversarial ML:** White/Black-box Attacks, Adversarial Training
- **Foundation Models:** Alignment, Guardrails, Activation Engineering, Agentic workflows

Development and Tools:

- **Programming Languages:** Python, C, C++, Java, Matlab
- **Deep Learning:** PyTorch, Tensorflow, Keras
- **Reinforcement Learning:** Gym/Gymnasium, Stable Baselines, RLlib
- **Others:** Git, LaTeX

ONLINE COURSES & SUMMER SCHOOLS

- **Advanced LLM Agents:** MOOC Spring 2025
- **LLM Agents:** MOOC Fall 2024
- **CS 285: Deep Reinforcement Learning** UC Berkeley, Fall 2021
- **M2L 2023: Mediterranean Machine Learning Summer School**
- **Security & Privacy in the Age of AI:** Organizer and participant in editions 2022-2024

LANGUAGES

Greek: Native | English: Excellent - C2 | German: Very Good - B2/C1 | Dutch: Good - B1

OTHER INTERESTS

An inquisitive and driven spirit, in my free time I enjoy science fiction, creative writing, and table-top/computer games; I am also a capoeirista and water polo player.