Edoardo **Debenedetti**

PHD STUDENT IN COMPUTER SCIENCE @ ETH ZÜRICH

🛛 (+41) 76 699 43 27 📋 🔽 edebenedetti@inf.ethz.ch 📔 🌴 https://edoardo.science 📔 🖸 dedeswim 📔 🖑 Edoardo Debenedetti

Education

ETH Zürich - Federal Institute of Technology Zürich Zürich, Switzerland PhD in Computer Science 08/2022 - 12/2026 (exp.) Focus: Real-world machine learning security and privacy, advised by Prof. Florian Tramèr in the SPY Lab. • IT Coordinator for the group: managing the GPU servers and hardware resources. • Fully funded by the CYD Doctoral Fellowship, awarded by the Armasuisse Cyber-Defense Campus. EPFL - Federal Institute of Technology Lausanne Lausanne, Switzerland MSC IN COMPUTER SCIENCE 09/2019 - 04/2022 • GPA 5.63/6, focus on Machine Learning ∩ Security ∩ Privacy. • Master's Thesis about the adversarial robustness of Vision Transformers supervised by Princeton University's Prof. Mittal. Politecnico di Torino Turin. Italv BSC IN COMPUTER ENGINEERING 09/2016 - 07/2019 • GPA 28.4/30, graduation mark 110/110, top 9%. • Exchange year at 同济大学 (Tongji University), in Shanghai (China), supported by a full scholarship granted to the top 31% applicants. Industry experience **Bloomberg LP** London, United Kingdom SOFTWARE ENGINEERING INTERN 07/2021 - 09/2021 • Worked in the Multi Asset Risk System team, on the re-design and implementation of the configuration of a distributed logging library. Move the configuration of a distributed logging library from an internal technology to a centralized SQL DB, using a cache and a C++ service. • The configuration is checked ~1M times per minute, and the usage of the cache gave a ~23x speed improvement w.r.t. querying the DB.

Armasuisse Cyber-Defence Campus

Research Intern

- Worked on Machine Unlearning and Membership Inference Attacks against Generative Models, supervised by Prof. Mathias Humbert.
- Adapt the MIA technique proposed by the GAN-Leaks work (by Chen et al.), to work after the removal some datapoints from the training set.
- The technique achieved **promising results** when attacking DCGAN trained on the CelebA dataset

Publications_

Conference proceedings

- Debenedetti, E., Carlini, N., Tramèr, F., "Evading Black-box Classifiers Without Breaking Eggs", 2nd IEEE Conference on Secure and Trustworthy Machine Learning, 2024, Distinguished Paper Award Runner-up.
- Debenedetti, E., Sehwag, V., Mittal, P., "A Light Recipe to Train Robust Vision Transformers", 1st IEEE Conference on Secure and Trustworthy Machine Learning, 2023.
- Croce, F.*, Andriushchenko, M.*, Sehwag, V.*, **Debenedetti, E.***, Flammarion, N., Chiang, M., Mittal, P., Hein, M., "*RobustBench: a standardized adversarial robustness benchmark*", Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track, 2021.

Workshop papers

• Edoardo Debenedetti, Zishen Wan, Maksym Andriushchenko, Vikash Sehwag, Kshitij Bhardwaj, Bhavya Kailkhura, "Scaling Compute Is Not All You Need for Adversarial Robustness", ICLR 2024 Workshop on Reliable and Responsible Foundation Models.

Manuscripts

- Chao, P.*, **Debenedetti, E.***, Robey, A.*, Andriushchenko, M.*, Croce, F., Sehwag, V., Dobriban, E., Flammarion, N., Pappas, G., Tramèr, F, Hassani, H., Wong, E., *"JailbreakBench: An Open Robustness Benchmark for Jailbreaking Language Models"*, arXiv ePrint 2404.01318.
- Qi, X., Huang, Y., Zeng, Y., **Debenedetti, E.**, Geiping, J., He, L., Huang, K., Madhushani Sehwag, U., Sehwag, V., Shi, W., Wei, B., Xie, T., Chen, D., Chen, P., Ding, J., Jia, R., Ma, J., Narayanan, A., Su, W., Wang, M., Xiao, C., Li, B., Song, D., Henderson, P., Mittal, P., *"Al Risk Management Should Incorporate Both Safety and Security"*, under review.
- Debenedetti, E.*, Severi, G.*, Carlini, N., Choquette-Choo, C. A., Jagielski, M., Nasr, M., Wallace, E., Tramèr, F., "Privacy Side Channels in Machine Learning Systems", arXiv ePrint 2309.05610.

* denotes equal contribution.

Honors and Awards

- 2024 Distinguished Paper Award Runner-up IEEE SaTML, Top 2 out of 34 accepted papers.
- 2023 Oral presentation ICML AdvML Frontiers Wokshop, Top 10% accepted papers.
- 2023 CYD Doctoral Fellowship, full PhD funding for 4 years, worth USD 536'000 (CHF 461'000), from Armasuisse CYD Campus and EPFL.
- 2021 Google TPU Research Cloud Program, extensive hardware support for 8 months to work on the Master's Thesis.
- **2021** Best Paper Honorable Mention ICLR Workshop on Security and Safety in ML Systems, top 2 out of 50 accepted papers.

Teaching

- Information Security Lab ETH Zürich: 2022, 2023 (Teaching Assistant)
- Large Language Models ETH Zürich: 2023, 2024 (Teaching Assistant)

Lausanne, Switzerland

08/2020 - 01/2021

Professional Service

Reviewer

- NeurIPS Datasets and Benchmarks Track: 2022, 2023
- CCS AlSec workshop: 2023

Conference service

- Competition organizer at SaTML 2024: lead organizer of the Large Language Models Capture-the-Flag. More than 400 users and 140 teams signed up and more than 70 defenses were submitted.
- Volunteer at NeurIPS 2021: helped with monitoring the website and technical issues.

Open Source Maintainer

- RobustBench: adversarial robustness benchmarking library and model zoo.
 - More than 150 models spanning 3 datasets and 3 threat models.
 - 569 stars, with 202 unique cloners in 2 weeks (measured in January 2024).
 - Refactored the code to improve the extensibility of the library.

Repository at https://github.com/RobustBench/robustbench.

Student supervision

- **Yize Cheng** (BSc project): universal black-box adversarial examples, *2023*.
- Qianjun Zheng (MSc project): white-box membership inference attacks, 2023.
- Kazuki Egashira (MSc project): membership inference attacks via memorization localization and neural functional networks, 2023-24.
- Fredrik Nestaas (MSc thesis): LLM SEO attacks via indirect prompt injection, in progress.
- Joshua Freeman (MSc project): Black-box LLM training data extraction, in progress.

Invited talks

- ACL SIGSEC Privacy Side-channels in Machine Learning Systems, 2023.
- TU Graz EfficientML Reading Group Privacy Side-channels in Machine Learning Systems, 2023.