

Diego DORN

Research Engineer

Diego works on the mitigation of **systemic risks** from **general-purpose artificial intelligence** systems.

✉ cv@ddorn.fr
🌐 cozyfractal.com
🐙 github.com/ddorn

He has extensive expertise in **software engineering** and experience in **teaching**, **leadership** and **communication** from his volunteering.
He finishes his master in Communication Systems in August 2024.

WORK EXPERIENCE

- Paris 🇫🇷 **Research engineer at EffiSciences**
Feb. 2024 – present *Design of a benchmark to evaluate monitoring systems of LLM agents, including detection of out-of-distribution failure modes by monitoring systems.*
- Berlin 🇩🇪 **Head Teacher for two ML4Good, a summer school on systemic AI risk**
Aug. 2023 *Delivery and improvement of 10 days of technical and conceptual content for ~20 participants, covering threat modeling, technical safety and AI policy.*
March 2024
- Cambridge 🇬🇧 **Research assistant, Machine Learning Group, Cambridge University**
July – Sep. 2023 *Research on goal misgeneralisation in Reinforcement Learning (RL) with N. Alex and D. Krueger. Published “Goal Misgeneralization as Implicit Goal Conditioning” in the GCRL workshop at NeurIPS 2023*
- Lausanne 🇨🇭 **Lead developer for the startup SPRIG (sprigproofs.org)**
Jan. 22 – May 23 *Developing a distributed platform to increase confidence in mathematical proofs.*

EDUCATION

- Lausanne 🇨🇭 **Master’s in Communication Systems, Ecole Polytechnique Fédérale de Lausanne (EPFL)**
Sep. 2021 – present *Focus on artificial intelligence, formal verification and advanced algorithms. Minor in Mathematics.*
- Interlaken 🇨🇭 **Summer school “Science and Policy – How to bridge the gap?”**
July 2023 *5 days on science for policy, science communication, open science and the Swiss policy landscape.*
- London 🇬🇧 **ARENA, Alignment Research Engineer Accelerator (arena.education)**
May – June 2023 *6 weeks intensive training on interpretability, RL and training at scale.*
- Lausanne 🇨🇭 **Bachelor’s in Mathematics at EPFL**
Sep. 18 – July 2021 *Passed with a 5.42/6 average and top 5/100 of my year.*

VOLUNTEERING

- Lausanne 🇨🇭 **Founder and President of the Safe AI Lausanne student association**
Sep. 22 – March 24 *Led a team of 8 through the design of a strategy, resulting in a 10-day winter school on systemic AI risks, 3 talks and 2 panel discussions with a total of 10 experts, and a talk for TEDxEcublens.*
- Lausanne 🇨🇭 **President of CQFD, the mathematics students’ association of EPFL**
Sep. 20 – Sep. 21 *Management of a team of 14 people, dialogue with the direction of the faculty.*
- Many places 🇫🇷 **National organisation committee of the french tournament of young mathematicians, TFJM²**
Sep. 20 – May 21 *Coordination of 9 events across France with 600 participants, development of a new online infrastructure and communication.*

AWARDS & EXTRA

- Bruxelles 🇧🇪 **1st place in the hackathon the “Digital Services Act RAG Race”**
February 2024 *Creation of a Q&A system for questions on the DSA based on open-source models, in a team of 3, during a 7 hours hackathon organised by the PEReN and the European Commission.*
- Earth 🌍 **Game development, tool design, websites (cozyfractal.com/showcase)**
2014 – present *Creation of 10+ small games under strong time constraints and pressure for game jams, a 2D EsoLang (Asciiidots), multiple software tools and websites. Teamwork and sprints.*

HARD SKILLS

Python (pytorch, huggingface, streamlit, click, mpy, pytest...) 6000h
JavaScript / CSS / HTML (VueJS, TailwindCSS) 500h
Rust, C++, Scala, LaTeX 300h each
System Administration (Git, Docker, Bash, remote machines...) 200h

SOFT SKILLS

- Training in Non-Violent Communication
- Public speaking
- Native in French (C2)
- Fluent in English (C1)