

Evžen Wybitul

Education

Oxford University Incoming student, Oct 2025
DPhil in Technical AI Governance and Safety

(Incoming student) Joint DPhil with the Oxford Martin AI Governance Initiative and the EPSRC CDT in Autonomous Intelligent Machines & Systems.

ETH Zurich Sep 2022 – (Aug 2025)
MSc in Data Science

Courses include: Causality, Large Language Models, Natural Language Processing, Reliable and Trustworthy Artificial Intelligence, Data Science in Law and Policy.

Charles University Sep 2021 – Aug 2022
Auditing courses in MSc in Artificial Intelligence

Best possible grades in all courses I took.

Courses include: Reinforcement Learning, Evolutionary Algorithms, Probability Theory, Artificial Intelligence Theory.

Charles University Sep 2018 – Aug 2021
BSc in Bioinformatics

Top student in the programme in all three years, graduated with honors.

Courses include: Deep Learning, Mathematical Analysis, Linear Algebra, Data Structures.

Publications

Gradient Routing: Masking Gradients to Localize Computation in Neural Networks 2024
[ArXiv](#), joint first author; MATS 6

A modification of backpropagation for learning specific capabilities in specific modules in the network, which can be used for unlearning and steering. Mentored by **Alex Turner** (Google DeepMind).

ViSta Dataset: Do Vision-language Models Understand Sequential Tasks? 2024
[ArXiv](#), first author; MATS 5

A dataset of 4,000+ videos of sequential tasks with descriptions. We use ViSta to evaluate if visual-language models could serve as task supervisors in reinforcement learning. Mentored by **David Lindner** (Google DeepMind).

Fellowships

Talos AI Governance Fellowship, EU Track 2025
Weekly discussion groups on AI governance, with a small research project focused on helping the EU lead international efforts in AI safety.

Other Research Experience

Assesing Vurneabilities in LLMs 2024
[GitHub](#), course project

Evaluated the safety of Large Language Model (LLM) agents, with a specific emphasis on prompt injections. Mentored by **Florian Tramèr**.

Training Steering Vectors 2024
[PDF](#), course project

Produced first steering vectors for GPT-2 small. Explored the usage of sparse auto-encoder features for steering. Supervised by **Elliott Ash**.

Measuring Emotion in Political Language 2024
[PDF](#), course project

Mapped how emotionality in political speeches developed over time. Supervised by **Elliott Ash**.

Certifying Robustness of Neural Networks 2023
[GitHub](#), course project

Formulated an algorithm based on DeepPoly to certify neural network robustness against input perturbations.

Teaching Experience

ETH Zurich Feb 2024 – Aug 2024
Teaching Assistant, Large Language Models

Taught a tutorial on the intuitions behind the transformer architecture and parameter-efficient fine-tuning methods.

Haviřov Grammar School Sep 2020 – July 2022
Haskell curriculum developer & instructor

Developed and taught an introductory course in functional programming.

Work Experience

IOCB Prague June 2020 – July 2021
Assistant in a bioinformatics research group

1. Improved effectivity of a lengthy manual procedure that identifies cysteine bonds in proteins by partially automating it (thesis project).
2. [GitHub](#). Built a web application for managing internal experiment requests (full stack web development). Enhanced the reusability and accessibility of experimental data.

MSD Sep 2019 – Feb 2020
Junior data scientist in a pharmaceutical company

Contributed to cost reduction and increased drug yields by optimizing a complex drug preparation process using classical ML on time-series data.

Selected Software Projects

Technologies: Python (PyTorch), R, Julia, Haskell, Purescript, React, Typescript, PostgreSQL, Docker.

Hate Speech Detection in Online Comments

[GitHub](#). Fine-tuned a BERT-based model for research and industry use in hate speech detection.

Racket Language Extension for VS Code

[GitHub](#). The most popular Racket extension for VS Code, with over 200 stars and 60,000 downloads.

Optimizing Exam Schedule

[GitHub](#). A program designed to assist students in optimizing their exam preparation schedules.

Selected Awards and Achievements

Most Active Student Award
Bakala Foundation, 2024

Long-Term Future Fund Grant, \$40 000
EA Funds, 2024

AI Safety Grant, \$21 000
AI Safety Support, 2024

Scholarship, \$33 000
Bakala Foundation, 2022

Scholarship for Outstanding Academic Achievement
Charles University, 2019

Scholarship for Outstanding Academic Achievement
Charles University, 2018

Team Debating League Finalist
National Team Debating League, 2017

Best A3 speaker
National Team Debating League, 2017

Contact

E-mail: wybitul.evzen@gmail.com

GitHub: github.com/Eugleo