

USMAN ANWAR

 usmananwar391@gmail.com

 www.linkedin.com/in/uzman-anwar

 <https://uzman-anwar.github.io/>

 <https://scholar.google.com/uanwar>

EDUCATION

University of Cambridge, Cambridge, UK

PhD In Engineering (Machine Learning)

Advisors: David Krueger and Yashar Ahmadian

Awards: Open Phil AI Fellowship and Future of Life Fellowship in AI Existential Safety

October 2022 - July 2026

Information Technology University, Lahore

MS Data Science (Rank: 2nd, Dean's Honours List)

Thesis Paper: Inverse Constrained Reinforcement Learning (*published at ICML*)

September 2019 – August 2021

University of Engineering and Technology, Lahore

BS Electrical Engineering

Undergraduate Thesis/FYP: Single Channel Acoustic Source Separation And Speech Enhancement

August 2015 – May 2019

AWARDS & HONOURS

Open Phil AI Fellowship by Open Philanthropy Foundation (2022-2026). 

Vitalik Buterin PhD Fellowship by Future Of Life Institute (2022-2026). 

Free Registration Award at virtual MLSS 2021 Taipei.

Graduate Student Fellowship, ITU, Lahore.

Merit Scholarship, ITU, Lahore.

SELECTED PUBLICATIONS (GOOGLE SCHOLAR)

* denotes equal contribution

Machine Learning and AI Safety Research

U. Anwar, J. von Oswald, L. Kirsch, D. Krueger, and S. Frei. Understanding in-context learning of linear models in transformers through an adversarial lens. *Transactions on Machine Learning Research*, 2025. **(Featured Certification)** 

T. Bush, **U. Anwar***, S. Chung*, A. Garriga-Alonso, and D. Krueger. Interpreting emergent planning in model-free reinforcement learning. In *The Thirteenth International Conference on Learning Representations*, 2025. **(Oral)** 

M. Farrugia-Roberts, K. A. Abdel Sadek*, H. Erlebach, C. Schroeder de Witt, D. Krueger, **U. Anwar**, and M. D. Dennis. Mitigating goal misgeneralization via minimax regret. In *Reinforcement Learning Conference*, 2025. 

T. Coste, **U. Anwar**, R. Kirk, and D. Krueger. Reward model ensembles help mitigate overoptimization. In *The Twelfth International Conference on Learning Representations*, 2024. 

D. Papadimitriou, **U. Anwar**, and D. S. Brown. Bayesian methods for constraint inference in reinforcement learning. *Transactions on Machine Learning Research*, 2022. ISSN 2835-8856. 

U. Anwar*, S. Malik*, A. Aghasi, and A. Ahmed. Inverse constrained reinforcement learning. In *Proceedings of the 38th International Conference on Machine Learning (ICML)*, 2021. **(Spotlight)** 

U. Anwar*, A. Pandian*, J. Wan, D. Krueger, and J. N. Foerster. Noisy ZSC: Breaking the common knowledge assumption in zero-shot coordination games. In *NeurIPS 2023 Workshop on Open Ended Learning*. [🔗](#)

J. M. L. Rangel, **U. Anwar**, S. Schoepf, J. Foster, and D. Krueger. Learning to forget using diffusion hypernetworks. In *The Third Workshop on New Frontiers in Adversarial Machine Learning*, 2024. [🔗](#)

M. Brumley, J. Kwon, D. Krueger, D. Krasheninnikov, and **U. Anwar**. Comparing bottom-up and top-down steering approaches on in-context learning tasks. In *MINT: Foundation Model Interventions*, 2024. [🔗](#)

S. Casper, X. Davies, ... **U. Anwar**, et al. Open problems and fundamental limitations of reinforcement learning from human feedback. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. (**Survey Certification, Featured Certification**) [🔗](#)

AI Policy and Sociotechnical Research

U. Anwar, A. Saparov, J. Rando, D. Paleka, et al. Foundational challenges in assuring alignment and safety of large language models. *Transactions on Machine Learning Research*, 2024. (**Survey Certification**) [🔗](#)

A. Peppin, A. Reuel, S. Casper, E. Jones, A. Strait, **U. Anwar***, et al. The reality of ai and biorisk. In *Proceedings of the 2025 ACM Conference on Fairness, Accountability, and Transparency*, 2025. [🔗](#)

A. Chan, N. Kolt, P. Wills, **U. Anwar**, C. S. de Witt, N. Rajkumar, L. Hammond, D. Krueger, L. Heim, and M. Anderljung. IDs for AI systems. In *NeurIPS 2024 Workshop on Regulatable ML*, 2024. [🔗](#)

WORK EXPERIENCE

Senior Machine Learning Engineer February 2020 – March 2022
Scientific Computing Department
NetSol Technologies, Lahore

Research Assistant & Graduate Student Fellow July 2019 – June 2020
Center of Artificial Intelligence and Computational Science
Information Technology University, Lahore

Research Intern July – September 2018
Internet of Things Laboratory
Khwarizmi Institute of Computer Science, Lahore

Junior Data Scientist June – August 2017
ADDO AI, Lahore

TEACHING EXPERIENCE

Teaching Assistant - Discrete Mathematics September 2019 – January 2020
Department Of Computer Science
Information Technology University, Lahore

(SELECTED) TALKS

Foundational Challenges in LLM Alignment and Safety - SERI MATS, ACM Summer School on Responsible & Safe AI, GovAI, Chalmers AI Ethics Seminar and others.

Reward Modelling for AGI Safety - Future of Life Seminar on AI Safety, October 2022.

PROFESSIONAL ACTIVITIES & SERVICES

Lead Organizer: Workshop on Socially Responsible Language Modelling Research (SoLaR) at NeurIPS 2023, 2024 and COLM 2025

Peer Reviewer: ICML 2023, NeurIPS 2023, ICLR 2025, NeurIPS 2025

Grant Reviewer: Vitalik Buterin Fellowship for AI Existential Safety (2023)

MENTORING & SUPERVISIONS

Karim Abdul Sadek *for* Mitigating Goal Misgeneralization via Minimax Regret at University of Cambridge, UK (resulted in ICLR 2025 submission).

Thomas Bush *for* Interpreting Emergent Planning in Model-Free Reinforcement Learning at University of Cambridge, UK (resulted in ICLR 2025 submission).

Miguel Lara *for* Learning to Forget Using Diffusion Hypernetworks at University of Cambridge, UK (resulted in NeurIPS workshop paper).

Thomas Coste *for* Conservative Agency For Safe Optimization of Learned Reward Models at University of Cambridge, UK (resulted in ICLR 2024 paper).

Yawen Duan *for* Red Teaming (Learned) Reward Models at University of Cambridge, UK.

Jason Brown *for* Likert Scale Feedback to Improve Robustness of Learning From Human Preferences at University of Cambridge, UK.

Abdul Rehman & Arslan Malik *for* Privacy Preserving Recommendation System at ITU, Pakistan.

VOLUNTEER WORK

Managing Director & Co-Founder Spectra Magazine

April 2017 – May 2020

Spectra Magazine is a student-powered online magazine aiming to enhance public understanding of science and shape the narrative of science journalism in Pakistan. Under my leadership, we published more than 215 articles and mentored more than 50 high school and undergraduate students in science writing, editing and design.

NON-DEGREE STUDIES

Wilson Center Pathways to AI Policy Fellowship (2025)

Cooperative AI Summer School (*July 2023*)

Eastern European Machine Learning School (*July 2021*)

SKILLS

- Python (Numpy, Scipy, Matplotlib) • Jax • Pytorch • Tensorflow • C • SQL • NoSQL