

# Akbir Khan

[akbir.dev](https://akbir.dev)

## Education

2021-2024	Ph.D. in Computer Science, University College London Advised by Ed Grefenstette & Tim Rocktäschel
2017-2018	MPhil in Advanced Computer Science, <i>with distinction</i> , University of Cambridge
2013-2017	MSci in Mathematics with Physics, <i>with 1<sup>st</sup> class honours</i> , University College London
2015-2016	Exchange student, <i>as Mathematics Specialist</i> , University of Toronto

## Experience

2023-	Research Analyst at <a href="#">Cooperative AI Foundation</a> , grant-making and field-building to mitigate risks posed by multi-polar AI outcomes
2021-2023	Senior Applied Researcher at <a href="#">Tractable AI</a> , built OCR pipeline which generates £8M ARR
2017-2020	CRO at <a href="#">Spherical Defence</a> , built Seq2seq models for anomalous web traffic detection
2016	Four month internship at <a href="#">Deutsche Bank</a> as a Software Engineer
2015	Research Intern at the <a href="#">Quantum Optics and Laser Group</a> , Imperial College London

## Selected Publications

2024	Debating with More Persuasive LLMs Leads to More Truthful Answers — <b>A Khan</b> , J Hughes, D Valentine, L Ruis, K Sachan, A Radhakrishnan, E Grefenstette, S Bowman, T Rocktäschel & E Perez. <b>Oral [top 1.5%]</b> at <i>International Conference on Machine Learning (ICML)</i>
2024	Scaling Opponent Shaping to High Dimensional Games — <b>A Khan</b> , T Willi, N Kwan, A Tachetti, C Lu, T Rocktäschel, E Grefenstette & J Foerstor. Oral at <i>The Autonomous Agents and Multi-Agent Systems (AAMAS)</i>
2023	The Goldilocks of Pragmatic Understanding: Fine-Tuning Strategy Matters for Implication Resolution by LLMs — L Ruis, <b>A Khan</b> , S Biderman, S Hooker, T Rocktäschel, & E Grefenstette. <b>Spotlight [top 3%]</b> at <i>The Neural Information Processing Systems (NeurIPS)</i>
2023	MAESTRO: Open-Ended Environment Design for Multi-Agent Reinforcement Learning — M Samvelyan, <b>A Khan</b> , M Dennis, M Jiang, J Parker-Holder, JN Foerster, R Raileanu, T Rocktäschel. Accepted at <i>International Conference on Learning Representations (ICLR)</i>

## Awards

2024	SuperAlignment Fellowship, OpenAI
2023	Astra Fellowship, Redwood Research
2020	Foundational Artificial Intelligence Scholarship, ESPRC

## Technical Projects

[Deep Equilibrium Models](#), a Haiku implementation of the NeurIPS 2019 paper, an implicit-depth differentiable architecture that simulates an infinitely deep network

[Bad Flamingo](#), a gamified collection of adversarial training examples; awarded 1<sup>st</sup> place at Hack Cambridge Ternary

Skills: Python [PyTorch, JAX ([contributor](#)), Pandas, Scikit-learn], Docker, GoLang