

# Chris Cundy

## Experience

- June 2024– **Research Scientist**, *FAR.AI*, Berkeley, California, USA  
Characterising and mitigating catastrophic risks from frontier AI systems. My main responsibilities are:  
o Leading research projects and determining the direction for teams of researchers and engineers.  
o Designing and implementing training and evaluation schemes with state-of-the-art models, conducting empirical studies, and communicating results through papers and presentations.  
o As a research team lead, ensuring that research aims are consistent with FAR's mission and that my direct reports are supported in their output and professional development.
- June 2022– **Research Scientist Intern**, *Technical AI Safety Team*, DeepMind, London, UK
- September 2022 Investigating robust and reliable machine learning in theory and at scale  
o Investigated susceptibility of autoregressive models to 'delusions', where unobserved latent variables lead to incorrect probabilistic judgments.  
o Developed a theoretical model for delusions; investigated delusions at scale by analysing performance of DeepMind's Gato (a large generalist, multi-task autoregressive model) on custom environments.
- October 2022 **Visiting Scholar**, *Future of Humanity Institute*, University of Oxford, Oxford, UK
- 2017–January 2018 Developing algorithms to predict human judgments. Supervised by Owain Evans and Andreas Stuhlmüller  
o Designed algorithms to collate quick, noisy human judgments to predict the answer to complicated tasks which would typically require deliberation.
- June– September 2017 **Visiting Scholar**, *Centre for Human-Compatible AI*, University of California, Berkeley, US  
Supervised by Daniel Filan & Stuart Russell, researching topics in AI safety  
o Extended previous work on inverse reinforcement learning to hierarchical setting. Formalized the problem, derived theoretical results, performed experiments on data and presented at an ICML workshop.

## Education

- 2018–2024 **PhD - Computer Science**, *Stanford University*, Stanford, California, USA  
o Advised by Stefano Ermon.  
o Investigating topics in inverse reinforcement learning, sequence modelling and variational inference.  
o Thesis: *Beyond Maximum Likelihood: Distribution-Aware Machine Learning*.
- 2016–2017 **MEng - Computer Science**, *University of Cambridge*, Cambridge, UK  
o Grade: Distinction. Modules include: Data Science, Probabilistic Machine Learning, Network Analytics.  
o Thesis: *Investigating Variational Gaussian Process State-Space Models with Gaussian Likelihood*.  
Supervised by Carl E. Rasmussen.
- 2013–2016 **BA - Natural Sciences (Physics)**, *University of Cambridge*, Cambridge, UK  
o Grade: 1st. Modules: Physics, Maths, Chemistry, Computer Science.

## Selected Publications

- 2025 **Preference Learning with Lie Detectors can Induce Honesty or Evasion**,  
**Chris Cundy, Adam Gleave**, NeurIPS 2025
- 2024 **SequenceMatch: Imitation Learning for Autoregressive Sequence Modelling with Backtracking**,  
**Chris Cundy, Stefano Ermon**, ICLR 2024
- 2024 **Privacy-Constrained Policies via Mutual Information Regularized Policy Gradients**,  
**Chris Cundy, Rishi Desai, Stefano Ermon**, AISTATS 2024
- 2022 **LMPriors: Pre-Trained Language Models as Task-Specific Priors**,  
**Kristy Choi\*, Chris Cundy\*, Sanjari Srivasta, Stefano Ermon**, First Workshop on Foundation Models for Decision Making, NeurIPS 2022

- 2021 **BCD Nets: Scalable Variational Approaches for Bayesian Causal Discovery**,  
*Chris Cundy, Aditya Grover, Stefano Ermon, NeurIPS 2021*
- 2020 **Flexible Approximate Inference via Stratified Normalizing Flows**,  
*Chris Cundy, Stefano Ermon, UAI 2020*
- 2018 **Parallelizing Linear Recurrent Neural Nets over Sequence Length**,  
*Eric Martin, Chris Cundy, ICLR 2018*

---

## Additional Publications

- 2025 **Sharpe Ratio-Guided Active Learning for Preference Optimization in RLHF**,  
*Syrine Belakaria, Joshua Kazdan, Charles Marx, Chris Cundy, Willie Neiswanger, Sanmi Koyejo, Barbara E Engelhardt, Stefano Ermon, CoLM 2025*
- 2024 **A physics-informed machine learning model for the prediction of drop breakup in two-phase flows**,  
*Chris Cundy, Shahab Mirjalili, Charlélie Laurent, Stefano Ermon, Gianluca Iaccarino, Ali Mani, International Journal of Multiphase Flow*
- 2023 **Neural Networks and the Chomsky Hierarchy**,  
*Grégoire Delétang, Anian Ruoss, Jordi Grau-Moya, Tim Genewein, Li Kevin Wenliang, Elliot Catt, Chris Cundy, Marcus Hutter, Shane Legg, Joel Veness, Pedro A. Ortega, ICLR 2023*
- 2022 **Towards a foundation model for geospatial artificial intelligence**,  
*Gengchen Mai, Chris Cundy, Kristy Choi, Yingjie Hu, Ni Lao, Stefano Ermon, Proceedings of the 30th International Conference on Advances in Geographic Information Systems*
- 2021 **IQ-Learn: Inverse soft-Q Learning for Imitation**,  
*Divyansh Garg, Shuvam Chakraborty, Chris Cundy, Jiaming Song, Stefano Ermon, NeurIPS 2021*
- 2018 **Exploring Hierarchy-Aware Inverse Reinforcement Learning**,  
*Chris Cundy, Daniel Filan, Workshop on Goal Specifications for Reinforcement Learning, ICML 2018*
- 2017 **Predicting Slow Judgment**,  
*Owain Evans, Andreas Stuhlmüller, Ryan Carey, Neal Jean, Andrew Schreiber, Girish Sastry, Chris Cundy, Aligned Artificial Intelligence Workshop, NeurIPS 2017*

---

## Service

- 2025 **Participant, EU AI Act Code of Practice Working Groups 2 and 4**  
Participated, as an independent expert, in working groups 2 and 4 for the development of the EU AI Act Code of Practice (CoP). I advocated, via written and oral presentation, for the importance of pre- and post-mitigation model evaluations, outlined in an earlier position paper I authored.
- 2023 **Teaching Assistant–CS228 (Probabilistic Graphical Models)**, *Stanford University*
- 2022 **Head Teaching Assistant–CS228 (Probabilistic Graphical Models)**, *Stanford University*  
Received award for excellence (awarded to top 5% of Teaching Assistants).
- 2023 **Project Supervisor**, *Supervised Project for Alignment Research (SPAR)*, Stanford AI Alignment  
Supervised five undergraduates on a project finding scaling laws in prompt injections.  
Presented work at the 7th Center for Human-Compatible AI workshop.
- 2021 **Project Supervisor**, *Undergraduate Research Program*, Stanford Existential Risk Initiative  
Served as supervisor for an undergraduate project on forecasting AI progress.
- 2020– **Reviewer**  
Reviewed for the following conferences: UAI (2020,2022,2025), ICML (2019,2020,2023,2025), ICLR (2021–2026), NeurIPS (2021–2025), AAAI-(Safe and Robust AI track) (2023–2024).

---

## Relevant Awards

- March 2024 **Winner, OpenAI Preparedness Challenge**  
○ One of the top ten submissions for the OpenAI Preparedness Challenge, for submitting *the most unique, while still being probable, potentially catastrophic misuse of the [OpenAI API]*.  
○ Developed proof-of-concept showing how GPT4-V, and speech-to-text with GPT4, could be used to parse vast amounts of unlabelled surveillance data, finding actionable insights for blackmail or insider trading.  
○ Prize: \$25,000 in OpenAI credits.