

# HARRY MAYNE

[www.harrymayne.com](http://www.harrymayne.com)

✉ [harrymayne@gmail.com](mailto:harrymayne@gmail.com)   [in](#) LinkedIn   [GH](#) GitHub   ☎ +44(0)7425 889951   [TW](#) @HarryMayne5

I am a PhD candidate at the University of Oxford. My research evaluates the reliability of LLM self-explanations: natural language explanations LLMs give to justify their own decision-making. I measure the extent to which these explanations are faithful to the real internal reasoning and I develop new training incentives to improve self-explanation faithfulness. I'm motivated by explainability and AI safety. Further research details can be found on my [website](#). I've also produced the [LINGOLY](#) (NeurIPS 2024 Oral) and [LINGOLY-TOO](#) benchmarks, as well as forthcoming science of evals papers.

## EDUCATION

### University of Oxford

*DPhil LLM Explainability and Interpretability*

*Oct. 2023 – Present*

- Supervised by Prof. Adam Mahdi (Oxford Internet Institute) and Prof. Jakob Foerster (Engineering).
- Thesis on LLM explainability and interpretability.
- Fully funded by the Grand Union DTP (Economic and Social Research Council).

### University of Oxford

**Distinction (77%)**

*MSc Social Data Science*

*Oct. 2022 – Aug. 2023*

- Ranked 1<sup>st</sup> in cohort for overall exam performance and achieved the highest thesis mark (88%).
- Courses included Applied Machine Learning, Natural Language Processing, Data Analytics at Scale, Frontiers of Data Science and Applied Analytical Statistics.
- Fully funded by the Grand Union DTP (Economic and Social Research Council).
- **Thesis:** Unsupervised Learning Approaches to Intensive Care Reform: Opportunities and Challenges.

### University of Cambridge

**Double First Class**

*BA Economics*

*Oct. 2019 – Jun. 2022*

- Ranked 1<sup>st</sup>/155 and 2<sup>nd</sup>/155 in econometrics and microeconomics examinations, respectively.
- Awarded the Patrick Cross Prize, two Tripos Prizes and the Corfield Scholarship for academic excellence.
- Courses included Time Series Econometrics, Microeconometrics, Labour Economics and Industrial Economics.
- **Thesis:** A Feast in the Time of Plague: A Theoretical Model of the UK Housing Market During COVID-19.

### A Levels and GCSEs

*Sep. 2011 – Jul. 2019*

- 4 A\*s at A Level and 11 A\*s (or equivalent) at GCSE.

## SELECTED PUBLICATIONS

2025	<b>The Validity–Minimality Trade-Off: The Limits of Self-Generated Counterfactual Explanations.</b> Mayne, H., Kearns, R. O., Yang, Y., Delaney, E., Russell, C., Mahdi, A.	<i>EMNLP 2025</i>
	<a href="#">🔗</a> <b>Ablation is Not Enough to Emulate DPO: A Mechanistic Analysis of Toxicity Reduction.</b> Yang, Y., Sondej, F., Mayne, H., Lee, A., Mahdi, A.	<i>EMNLP 2025</i>
	<b>Measuring What Matters: Construct Validity in Large Language Model Benchmarks.</b> Bean, A., et al. (incl. Mayne, H.)	<i>Under review at NeurIPS 2025</i>

🔗 **LINGOLY-TOO: Disentangling Reasoning from Knowledge with Templatised Orthographic Obfuscation.**  
Khouja, J., Korgul, K., Hellsten, S., Yang, L., Neacsu, V., **Mayne, H.**, Kearns, R., Bean, A., Mahdi, A.

*Under review at  
NeurIPS 2025*

2024 🔗 **LingOly: A Benchmark of Olympiad-Level Linguistic Reasoning Puzzles in Low-Resource and Extinct Languages.**  
Bean, A., Hellsten, S., **Mayne, H.**, Magomere, J., Chi, E., Chi, R., Hale, S., Kirk, H.

*NeurIPS 2024,  
Oral, Top 0.5%*

🔗 **Can Sparse Autoencoders be used to Decompose and Interpret Steering Vectors?**  
**Mayne, H.**, Yang, Y., Mahdi, A.

*Interpretable AI,  
NeurIPS 2024*

🔗 **Ablation is Not Enough to Emulate DPO: A Mechanistic Analysis of Toxicity Reduction.**  
Yang, Y., Sondej, F., **Mayne, H.**, Mahdi, A.

*Interpretable AI,  
NeurIPS 2024*

🔗 **Large language models can help boost food production, but be mindful of their risks.**  
De Clercq, D., Nehring, E., **Mayne, H.**, Mahdi, A.

*Frontiers in  
Artificial  
Intelligence*

## TEACHING

### University of Oxford, Oxford Internet Institute

*Sep. 2023 – Present*

- Teaching Assistant for Applied Analytical Statistics (Social Data Science MSc).
- Taught over 50 Master's and PhD students. Voted the department's top TA of the year in 2024.
- Designed a new seminar syllabus to improve student understanding.

### Stanford University, Stanford House

*Nov. 2023 – Present*

- Employed by Stanford University to tutor computer science undergraduate students completing semesters abroad at the University of Oxford.
- Designed and taught 8-week personalised courses in the students' areas of interest (focusing on machine learning, AI and data science).

### King Saud University, Oxmedica/Mawhiba

*Jun. 2024*

- Designed and delivered a three-week course in introductory machine learning to 20 talented Saudi teenagers. The full course is available on my [website](#). Received excellent feedback.

## POSITIONS

### International Growth Centre (IGC)

*Jun. 2025 – Present*

- AI consultant in the Zambian Evidence Lab. Using AI to aid public service delivery. See my [IGC page](#).

### Bank of England

*Summer 2020 & 2021*

- **2020:** First-year intern in Insurance Supervision. Selected as one of 9 interns from a pool of 3,000 applicants; offered a place on the penultimate-year internship.
- **2021:** Penultimate-year intern in the General Insurance Actuarial team (PRA). Offered a place on the actuarial graduate programme.

## ADDITIONAL SKILLS AND INTERESTS

**Programming Languages:** Python (PyTorch, NumPy, Matplotlib, vLLM, Sklearn...etc), L<sup>A</sup>T<sub>E</sub>X.

**Conference Reviewing:** ICML 2024 DCML Workshop, NeurIPS 2024 MINT Workshop

**Soft Skills:** Public speaking (Grade 8, Distinction), organisation, leadership, presentation.

**Other Leadership Roles:** Selwyn College Boat Club Captain (2022-2023), Selwyn College JCR Entertainment Officer (2021-2022), Selwyn College JCR Freshers' Representative (2020-2021).

**Other Interests:** Travel, surfing, and music.