

Doctoral applicant aiming for academic career researching *computation and learning* with applications to *understanding risks from advanced artificial intelligence*. Academic website: <https://far.in.net>.

§ Education

Master of Computer Science (with Distinction)

The University of Melbourne

ETH Zürich (semester exchange, 2020)

Advanced coursework in computation and learning. Thesis in deep learning theory (4) supervised by Daniel Murfet, leading to two sole-author conference papers (1, under review; 2, NeurIPS 2023). Founded a reading group on AI safety.

Awards: *Dean's Honours List* (top 5% marks in Faculty of Engineering and IT). *Top thesis mark* since the degree was first conferred in 2021. Thesis mark in the 95%+ category, reserved for theses described as follows: “*Truly outstanding in every way. In an entire academic career such a student may be encountered only once or twice. The student would be welcome as a PhD candidate in the School and would be expected to succeed with a hands-off supervision style.*”

Part-time 2019–2022
Coursework average **98.8%**, thesis **95.5%**

GPA **5.9 / 6.0**

Bachelor of Science (Computing and Software Systems)

The University of Melbourne

Major in computer science and software engineering. Electives mainly in mathematics and physics.

Awards: *Dean's Honours List I, II, III* (top 1% marks in Faculty of Science in first, second, and third year). *Australian Computing Society Bachelor of Science Student Award* (top marks in third-year computer science coursework). *Australian Artificial Intelligence Institute Prize* (top marks in AI coursework). Top marks in many other computer science classes.

2014–2016

Average **93%**

Victorian Certificate of Education (secondary school)

Mount Lilydale Mercy College

Maths/Science Prefect (elected by peers). Initiated/presented mathematics exam revision lecture.

Awards: *Dux* (valedictorian). *Victorian Premier's Award (Physics)* (top 3 physics students, state). *Australian Student Prize* (top 500 students, national). *Australian Defence Force Long Tan Leadership and Teamwork Award* (recognising leadership and contribution to school community).

2013

National percentile **99.8th**

§ Publications

Machine Learning Theory

- (1) **Matthew Farrugia-Roberts**, 2023, “Computational complexity of detecting proximity to losslessly compressible neural network parameters”. Preprint [arXiv:2306.02834](https://arxiv.org/abs/2306.02834). Under review.
- (2) **Matthew Farrugia-Roberts**, 2023, “Functional equivalence and path connectivity of reducible hyperbolic tangent networks”. Preprint [arXiv:2305.05089](https://arxiv.org/abs/2305.05089). Conference paper, **NeurIPS 2023**.
- (3) Joar Skalse,⁽⁼⁾ **Matthew Farrugia-Roberts**,⁽⁼⁾ Alessandro Abate, Stuart Russell, and Adam Gleave, 2023, “Invariance in policy optimisation and partial identifiability in reward learning.” Preprint [arXiv:2203.07475](https://arxiv.org/abs/2203.07475). Conference paper, **ICML 2023**.
- (4) **Matthew Farrugia-Roberts**, 2022, *Structural Degeneracy in Neural Networks*, Master's thesis, School of Computing and Information Systems, the University of Melbourne. Available [online](#).

Computer Science Education

- (5) **Matthew Farrugia-Roberts**, Bryn Jeffries, and Harald Søndergaard, 2022, “Teaching simple constructive proofs with Haskell programs.” Extended abstract presented at TFPIE 2022, journal paper published in EPTCS. [doi:10.4204/EPTCS.363.4](https://doi.org/10.4204/EPTCS.363.4).
- (6) **Matthew Farrugia-Roberts**, Bryn Jeffries, and Harald Søndergaard, 2022, “Programming to learn: Logic and computation from a programming perspective.” Conference paper presented at ITiCSE 2022, published by ACM. [doi:10.1145/3502718.3524814](https://doi.org/10.1145/3502718.3524814).

§ Research Experience

- Research Assistant (AI alignment & reward hacking)** Sep 2023–present
Computational and Biological Learning Laboratory, University of Cambridge
 Working on understanding and mitigating goal misgeneralisation in deep reinforcement learning.
- Research Affiliate (Developmental interpretability)** Aug 2023–present
Timaeus (academic research institute)
 Working on understanding the emergence of in-context learning in transformers and other projects.
- Research Assistant (Human-agent interaction)** Jan 2023–Jul 2023
School of Computing and Information Systems, the University of Melbourne
 Contributed to ongoing explainable AI project, evaluating human understanding of automated decision-making systems. Automated the creation of dynamic surveys with thousands of variants.
- Virtual Research Intern** Jun 2021–Oct 2021
Center for Human-compatible AI, University of California (Berkeley)
 Project work leading to a paper on reward learning theory (3, ICML 2023). Initiated a virtual mini-conference for interns to share presentations about their projects.
- See also** [Master of Computer Science](#) (Master’s research project). Part-time Feb 2021–Oct 2022

§ Teaching Experience

- Tutor** 2021, 2023
Centre for AI and Digital Ethics, the University of Melbourne
 Facilitating classes in the ethics and governance of AI for technical graduate students.
- Sessional Subject Coordinator and Lecturer** Jan 2018–Jul 2018
School of Computing and Information Systems, the University of Melbourne
 Co-coordinated/lectured a summer intensive class on introductory programming (150 students).
 Co-coordinated a semester-long class on algorithms and data structures (400 students).
- Head Tutor** 2017–2021
School of Computing and Information Systems, the University of Melbourne
 Designed coursework/assessment for classes on algorithmics, theoretical computer science, and artificial intelligence (300–600 students/class). Coordinated tutor teams to assess students fairly. Pioneered digital teaching/assessment methods leading to two CS education publications (5; 6).
Awards: *School of Engineering Tutor Community Excellence Award (finalist). School of Engineering Most Innovative Academic (finalist). Head Tutor Special Commendation Award.*
- Tutor** 2016–2021
School of Computing and Information Systems, the University of Melbourne
 Taught above-listed classes plus classes in programming, operating systems, networks, and security.
Awards: *School of Computing and Information Systems Excellence in Tutoring Award, 2016.*
- Mathematics and Physics Presenter** Dec 2014–Nov 2016
ATAR Notes (Australian student support community)
 Created accessible lectures, subject notes, and webinars for secondary students.
- Private Tutor and School-based Tutor** Jan 2014–Jun 2017
Private; Mount Lilydale Mercy College; Scotch College (indigenous student support program)
 Mathematics and study-skills support for secondary students of diverse backgrounds. Volunteered to organise annual final exam revision classes, including delivering mathematics and physics classes and recruiting other high-achieving alumni to deliver classes on other topics.

§ Technical Proficiencies

- Programming:** Python (including NumPy, SciPy, matplotlib, PyTorch, PyTorch/XLA, einops); Haskell; C; JavaScript (including React).
- Markup:** \LaTeX (including TikZ, beamer, BibTeX styles); HTML & CSS; Markdown; pandoc.
- Other:** Unix-like operating systems (macos, Ubuntu Linux, Arch Linux); git & GitHub.