

# Jonathan Colaço Carr

jonathan.colacocarr@mail.mcgill.ca | j-c-carr.github.io | linkedin.com/in/jonathan-colaco-carr

## Education

- M.Sc in Computer Science, McGill University** 2023 – 2025 (Expected)
- Supervisors: Prakash Panangaden and Doina Precup. GPA: 4.0/4.0.
- Visiting Graduate Student, Stanford University** 2024 – 2025
- Supervisor: Benjamin Van Roy
- B.Sc in Honours Mathematics, McGill University** 2018 – 2023
- Minors in Computer Science and Physics. GPA: 3.8/4.0.

## Research Experience

- Visiting Student Researcher, Stanford University** October 2024 – March 2025
- Supervisor: Benjamin Van Roy
  - Extended the maximal lottery social choice rule to sequential decision-making problems, and designed a new game theory algorithm that provably converges to the solution. Showed that the algorithm scales to function approximation.
- Visiting Student Researcher, UC Berkeley** June 2024 – September 2024
- Supervisors: Cameron Allen and Stuart Russell
  - Developed a simple improvement to three unsupervised reinforcement learning algorithms using state abstraction, tripling their exploration efficiency and enabling them to automatically avoid negative side effects in downstream tasks.
- Student Researcher, Mila** May 2022 – present
- Supervisors: Prakash Panangaden and Doina Precup
  - Audited the Anthropic Helpful and Harmless preference dataset for label quality and topic coverage. Fine-tuned nine LLMs on this dataset, showing how demographic imbalances caused disparities in model safety.
  - Proved it was possible for sequential decision-making agents to solve tasks that can't be expressed by expected utility.
- Student Research Engineer, McGill Space Institute** May 2021 – December 2021
- Supervisor: Adrian Liu
  - Implemented a deep computer vision algorithm (a modified U-Net) to recover signals from noisy telescope data.

## Manuscripts in Preparation

- [1] **J. Colaço Carr**, “Sequential decision making on the basis of maximally preferred alternatives,” M.S. thesis, McGill University, 2025, To be submitted for the Canadian AI Association's Best Master's Thesis competition.
- [2] **J. Colaço Carr**, D. Precup, P. Panangaden, and B. Van Roy, *Markov decision contests: Competitive sequential decision making with social choice guarantees*, 2025, To be submitted for the International Conference on Machine Learning (ICML).

## Peer-reviewed Conference Proceedings

\* denotes co-first authorship.

- [3] K. Chehbouni<sup>\*</sup>, **J. Colaço Carr**<sup>\*</sup>, Y. More, J. C. Cheung, and G. Farnadi, “Beyond the safety bundle: Auditing the helpful and harmless dataset,” in *Proceedings of the Annual Conference of the Nations of the Americas Chapter of the ACL (NAACL)*, 2025, **Selected for an oral presentation**.
- [4] **J. Colaço Carr**, P. Panangaden, and D. Precup, “Conditions on preference relations that guarantee the existence of optimal policies,” in *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2024.

## Journal Articles

- [5] **J. Colaço Carr**, Q. Sun, and C. Allen, “Focused skill discovery: Learning to control specific state variables while minimizing side effects,” in *Reinforcement Learning Journal (RLJ)*, 2025.

[6] J. Kennedy, **J. Colaço Carr**, S. Gagnon-Hartman, A. Liu, J. Mirocha, and Y. Cui, "Machine-learning recovery of foreground wedge-removed 21-cm light cones for high- $z$  galaxy mapping," in *Monthly Notices of the Royal Astronomical Society (MNRAS)*, 2024.

## Work Experience

- Lead Software Developer, Hortus AI** December 2024 – present
- Led a team of four software developers with startup founder, Thomas Krendl Gilbert, to create a platform that helps government entities manage AI alignment and safety concerns during the procurement process.
  - Improved the user interface through iterative design, integrating user feedback from focus groups and usability studies with the GovAI coalition.
- Teaching Assistant, McGill University** January 2024 – April 2024
- Held weekly office hours and created assignments for a graduate-level course in machine learning (COMP 551).
- Associate Machine Learning Developer, AltaML** January 2022 – April 2022
- Implemented decision trees and clustering algorithms to identify cross-selling opportunities and enhance customer segmentation. Presented technical insights to the client's sales and marketing teams.

## Extracurricular Experience

- Software Contributor, Center for Human-Compatible AI** July 2024
- Built a docker container for the *DignityIndex*, a tool that detects disrespectful language in political discourse.
- Co-founder, McGill AI Safety x Law Group** April 2022 – December 2023
- Co-founded and taught an 8-week course for McGill Law students interested in AI regulation. Organized a panel discussion with speakers from AI Governance & Safety Canada.
- VP Technology, McGill Research and Sustainability Network** April 2021 – April 2022
- Organized monthly panel discussions with students, professors and professionals interested in sustainability research.
- Team Lead, RBC x Microsoft AI for Social Impact Challenge** September 2019 – April 2020
- Led a team of four students in an eight-month competition to design a robotic arm for automated textile recycling. Placed 2nd out of 167 teams from across Canada.

## Awards

- McGill Graduate Mobility Award** October 2024
- Research award for McGill graduate students studying abroad, determined by academic standing.
- MITACS Globalink Research Award** October 2024
- Competitive research award for Canadian graduate students studying abroad.
- NSERC Canada Graduate Scholarship – Master's** May 2023
- Graduate scholarship determined by academic excellence, research potential and interpersonal skills.
- McGill Computer Science Undergraduate Research Award** May 2022
- Summer research award determined by academic record and research aptitude.
- McGill Space Institute Summer Undergraduate Research Award** May 2021
- Summer research award determined by academic record and extra-curricular leadership.
- RBC x Microsoft AI for Social Impact Challenge, 2nd Place Prize** April 2020
- National AI competition prize determined by originality, technical excellence, team diversity, impact, and feasibility.