

# Thomas Jiralerspong

---

Université de Montréal  
Mila  
Montreal, Canada

[thomasjiralerspong@gmail.com](mailto:thomasjiralerspong@gmail.com)  
[superkaiba.github.io](https://superkaiba.github.io)  
+1 (514) 625-9308

[Google Scholar](#)  
[LinkedIn](#)  
[GitHub](#)

## Education

### Université de Montréal

PhD - Computer Science

Supervisors: [Yoshua Bengio](#) & [Guillaume Lajoie](#)

Vanier Canada Graduate Scholarship Scholarship (150 000\$)

FRQNT Scholarship (40 000\$) (Rank #1 among all applicants in category)

NSERC Canada Graduate Scholarship (17 500\$)

Hydro-Québec Excellence Scholarship (10 000\$)

Arbour Scholarship (7 500\$)

*In progress*

### Massachusetts Institute of Technology

Brains, Minds, and Machines Summer Course

2024

### McGill University

B.Sc., Honours Computer Science

2023

Supervisors: [Blake Richards](#) & [Doina Precup](#)

GPA: 4.00/4.00

Exchange semester at the **National University of Singapore**

J.W. McConnell Major Entrance Scholarship (9 000\$)

## Refereed Conferences

Thomas Jiralerspong, Trenton Bricken. "Cross-Architecture Model Differing With Cross-coders." Under review. 2026.

Luca Scimeca\*, Thomas Jiralerspong\*, Berton Earnshaw, Jason Hartford, Yoshua Bengio. "Learning What Matters: Steering Diffusion via Spectrally Anisotropic Forward Noise." Under review. 2026.

Eric Elmoznino\*, Thomas Jiralerspong\*, Yoshua Bengio, Guillaume Lajoie. "A Complexity-Based Theory of Compositionality." In *Forty-Second International Conference on Machine Learning (ICML)*. 2025.

Jin Hwa Lee\*, Thomas Jiralerspong\*, Lei Yu, Emily Cheng. "Geometric Signatures of Compositionality Across a Language Model's Lifetime." In *The 63rd Annual Meeting of the Association for Computational Linguistics (ACL)*. 2025.

Ezekiel Williams, Avery Ryoo\*, Thomas Jiralerspong\*, Matt Perich, Guillaume Lajoie. "Expressivity of neural networks with random weights and learned biases." In *The 13th International Conference on Learned Representations (ICLR)*. 2025.

Jean-Pierre Falet, Hae Beom Lee, Nikolay Malkin, Chen Sun, Dragos Secrieru, **Thomas Jiralerpong**, Dinghuai Zhang, Guillaume Lajoie, Yoshua Bengio. “Delta-AI: Local Objectives for Amortized Inference in Sparse Graphical Models” In *Twelfth International Conference on Learning Representations (ICLR)*. 2024.

Chen Sun, Wannan Yang, **Thomas Jiralerpong**, Dane Malenfant, Benjamin Alsbury-Nealy, Yoshua Bengio, Blake Richards. “Contrastive Retrospection: honing in on critical steps for rapid learning and generalization in RL.” In *Thirty-seventh Annual Conference on Neural Information Processing Systems (NeurIPS)*. 2023.

Flemming Kondrup\*, **Thomas Jiralerpong\***, Elaine Lau, Nathan de Lara, Jacob Shkrob, My Duc Tran, Doina Precup, Sumana Basu. “Towards Safe Mechanical Ventilation Treatment Using Deep Offline Reinforcement Learning.” In *Thirty-seventh AAAI Conference on Artificial Intelligence (AAAI)*. 2023.

Marshall Wang, John Willes, **Thomas Jiralerpong**, Matin Moezzi. “A Comparison of Classical and Deep Reinforcement Learning Methods for HVAC Control.” In *20th IEEE International Conference on Ubiquitous Intelligence and Computing (UIC)*. 2023.

## Refereed Workshops

**Thomas Jiralerpong**, Berton Earnshaw, Jason Hartford, Yoshua Bengio, Luca Scimeca. “Shaping Inductive Bias in Diffusion Models through Frequency-Based Noise Control” In *The ICLR Workshop on Deep Generative Models in Machine Learning: Theory, Principle and Efficacy (DeLTA)*. 2025.

Marco Jiralerpong, **Thomas Jiralerpong**, Vedant Shah, Dhanya Sridhar, Gauthier Gidel. “General Causal Imputation via Synthetic Interventions.” In *The Causal Representation Learning Workshop at NeurIPS*. 2024.

**Thomas Jiralerpong\***, Xiaoyin Chen\*, Yash More, Vedant Shah, Yoshua Bengio. “Efficient Causal Graph Discovery Using Large Language Models.” In *How Far Are We From AGI? Workshop at ICLR*. 2024.

**Thomas Jiralerpong\***, Flemming Kondrup\*, Doina Precup, Khimya Khetarpal. “Forecaster: Towards Temporally Abstract Tree-Search Planning from Pixels.” In *Seventh Workshop on Generalization in Planning at NeurIPS*. 2023.

Flemming Kondrup\*, **Thomas Jiralerpong\***, Elaine Lau, Nathan de Lara, Jacob Shkrob, My Duc Tran, Doina Precup, Sumana Basu. “Deep Conservative Reinforcement Learning for Personalization of Mechanical Ventilation Treatment.” In *The Multi-disciplinary Conference on Reinforcement Learning and Decision Making (RLDM)*. 2022.

## Research Experience

### Astra Fellowship - Google DeepMind Stream

*Researcher*

*Jan 2026 - Present*

**Project:** Evaluation of in-context chain of thought awareness and obfuscation

---

\* Equal Contribution

	<b>LawZero</b> <i>Researcher</i> <b>Project:</b> Training a LLM explicitly to output the truth	<i>Jan 2026 - Present</i>
	<b>Anthropic</b> <i>Research Fellow</i> <i>Mentored by Trenton Bricken</i> <b>Project:</b> Mechanistic interpretability to discover behavioral differences between models	<i>Apr 2025 - Oct 2025</i>
	<b>Occam AI</b> <i>Research Scientist</i> <b>Projects:</b> Optimization of interactions between network of LLM agents, automated SQL query generation using LLMs	<i>Jun 2024 - Apr 2025</i>
	<b>Waabi</b> <i>Deep Learning Research Intern</i> <i>Mentored by Kelvin Wong and Chris Zhang</i> <b>Project:</b> Realistic and controllable traffic simulation using a transformer based variational autoencoder	<i>Jun 2023 – Aug 2023</i>
	<b>Reasoning and Learning Lab, Mila/McGill University</b> <i>Research Intern</i> <i>Supervised by Prof. Doina Precup</i> <b>Project:</b> Model-based reinforcement learning with affordance aware tree-search planning directly from pixels	<i>Jan 2022 – Aug 2023</i>
	<b>Learning in Neural Circuits Lab, Mila/McGill University</b> <i>Research Intern</i> <i>Supervised by Prof. Blake Richards</i> <b>Project:</b> Contrastive learning to discover critical states for reinforcement learning in sparse reward environments	<i>Sep 2022 – Aug 2023</i>
	<b>Vector Institute for A.I.</b> <i>Machine Learning Research Intern</i> <i>Mentored by John Willes and Marshall Wang</i> <b>Project:</b> Model-based reinforcement learning for HVAC control	<i>Sep 2022 – Dec 2022</i>
	<b>Project X, Machine Learning Research Competition</b> <i>Co-leader of McGill's Team</i> <i>Received the highest score out of 25 submitted papers</i> <b>Project:</b> Deep offline conservative reinforcement learning for mechanical ventilation treatment	<i>Jun 2021 – Feb 2022</i>
Industry Experience	<b>Amazon Web Services (AWS) – S3 Team</b> <i>Software Development Engineer Intern</i> <b>Project:</b> JavaScript/Python tool to automate the Incremental Backup recovery system for AWS S3 (stores ~14 trillion objects)	<i>May 2022 – Jul 2022</i>

	<b>Square Enix</b> <i>Software Development Intern</i> <b>Project:</b> Localization system to allow a MOBA game to be translated into over 10 languages	<i>May 2021 – Aug 2021</i>
	<b>Expedia</b> <i>Software Development Intern</i> <b>Project:</b> React/TypeScript tool to identify which elements of a webpage are broken and conveniently display them to developers	<i>Jun 2019 – Aug 2019</i>
Teaching	<b>Université de Montréal</b> Teaching Assistant, Representation Learning	2023
	<b>McGill A.I. Society</b> Organizer/Teaching Assistant, Accelerated Intro to ML	2021 – 2023
	<b>McGill University</b> Teaching Assistant, Software Systems Guest Lecturer, Theory of Machine Learning	2021 – 2022 2022
Honors	Vanier Canada Graduate Scholarship (150000\$) FRQNT Master's Scholarship (40000\$) (Rank #1 among all applicants in category) 2024 Arbour Scholarship (7500\$) Hydro-Québec Excellence Scholarship (10000\$) Chosen to attend the 10th Heidelberg Laureate Forum NSERC Canada Graduate Scholarship (17500\$) University of Montreal Master's Scholarship (5000\$) McGill Mobility Bursary for Exchanges (6000\$) Winner of UofT AI's Project X competition (25000\$) J.W. McConnell Major Entrance Scholarship (9000\$) CIBPA Foundation Bursary (1000\$, 2500\$, 1000\$) Marianopolis College Valedictorian Governor General of Canada's Academic Medal	2025 2024 2024 2023 2023 2023 2022 2022 2020 – 2022 2021, 2022, 2023 2020 2020
Invited Talks	Canadian Undergraduate Conference on AI (CUCAI) University of Toronto AI Conference McGill AI Society Learnathon	2022 2022 2022
Professional Activities	<b>Mila</b> Chairman of Lab Representatives Chairman of Social Committee Executive Member of Recruitment Committee	2023 – Present 2023 – Present 2023 – Present
	<b>McGill AI Society</b> Senior Advisor	2023 – Present

	Technical Project Manager	2021 – 2023
	<b>Montreal AI &amp; Neuroscience Conference</b> Organizer – Introduction to deep learning with PyTorch workshop	2022
	<b>McGill NeuroTech</b> Machine Learning Developer	2021 – 2022
	<b>McGill Robotics</b> Software Developer	2020 – 2021
Languages	<b>Native:</b> English, French <b>Advanced:</b> Italian, Spanish <b>Beginner:</b> Mandarin, Japanese	
Skills	<b>Programming Languages:</b> Python, Java, JavaScript, R, C, C++, C#, OCaml, SQL, HTML, CSS  <b>Machine Learning Libraries:</b> PyTorch, TensorFlow, Keras, Pandas, NumPy, Matplotlib  <b>Other:</b> L <sup>A</sup> T <sub>E</sub> X, Slurm, Jupyter Notebooks, Perforce, GitHub, Jira, Unity	
Press	<b>SciLogs.</b> Nina Beier. Jan 24, 2024. <a href="#">What Do Food and Research Have in Common? More Than You Might Think.</a>  <b>The McGill Tribune.</b> Mikaela Shadick. March 15, 2022. <a href="#">Six McGill undergrads win UofT international artificial intelligence competition.</a>  <b>McGill Reporter.</b> Richard Deschamps. March 1, 2022. <a href="#">Undergrad team uses machine learning to create a better hospital ventilator.</a>	
Advanced Coursework	<b>Université de Montréal</b> Representation Learning Reinforcement Learning & Optimal Control Scaling Laws Causal Inference & Machine Learning Probabilistic Graphical Models  <b>McGill University</b> Reinforcement Learning Brain Inspired Artificial Intelligence Honours Math for Machine Learning Probabilistic Programming Network Science	
	<b>National University of Singapore</b>	

Quantum Computing  
Information Theory