

Andrew Lee

Contact Information	E-mail: ajyl@umich.edu Website: https://ajyl.github.io Google Scholar: Link	
Research Interests	Interpretability, Language Models, Dialogue, Natural Language Processing, Deep Learning	
Education	University of Michigan	Ann Arbor, MI
	Ph.D. Candidate in Computer Science	2020 - Present
	Advisor: Rada Mihalcea	
	University of Michigan	Ann Arbor, MI
	Master's in Computer Science	2015
	Northwestern University	Evanston, IL
	Bachelor of Science in Computer Science	2013
Publications	<p>Andrew Lee, Xiaoyan Bai, Itamar Pres, Martin Wattenberg, Jonathan K. Kummerfeld, Rada Mihalcea. 2024. A Mechanistic Understanding of Alignment Algorithms: A Case Study on DPO and Toxicity. <i>Pre-print</i>.</p> <p>Andrew Lee, Jonathan K. Kummerfeld, Larry An, Rada Mihalcea. 2024. A Comparative Multidimensional Analysis of Empathetic Systems. <i>European Chapter of the Association for Computational Linguistics (EACL)</i>.</p> <p>Oana Ignat, Zhijing Jin, Artem Abzaliev, Laura Biester, Santiago Castro, Naihao Deng, Xinyi Gao, Aylin Gunal, Jacky He, Ashkan Kazemi, Muhammad Khalifa, Namho Koh, Andrew Lee, Siyang Liu, Do June Min, Shinka Mori, Joan Nwatu, Veronica Perez-Rosas, Siqi Shen, Zekun Wang, Winston Wu, Rada Mihalcea. 2024. Has It All Been Solved? Open NLP Research Questions Not Solved by Large Language Models. In <i>The 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING)</i>.</p> <p>Neel Nanda*, Andrew Lee*, Martin Wattenberg. 2023. Emergent Linear Representations in World Models of Self-Supervised Sequence Models. In <i>Proceedings of the 6th BlackboxNLP Workshop: Analyzing and Interpreting Neural Networks for NLP (BlackboxNLP)</i>. Honorable Mention for Best-Paper.</p> <p>Jing Xu, Andrew Lee, Sainbayar Sukhbaatar, Jason Weston. 2023. Some things are more CRINGE than others: Preference Optimization with the Pairwise Cringe Loss. <i>Pre-print</i></p> <p>Andrew Lee, Jonathan K. Kummerfeld, Larry An, Rada Mihalcea. 2023. Empathy Identification Systems are not Accurately Accounting for Context. <i>European Chapter of the Association for Computational Linguistics (EACL)</i>.</p> <p>Andrew Lee, David Wu, Emily Dinan, Michael Lewis. 2022. Improving Chess Commentaries by Combining Language Models with Symbolic Reasoning Engines. <i>Pre-print</i></p> <p>Andrew Lee, Zhenguo Chen, Kevin Leach, Jonathan K. Kummerfeld. 2022. Augmenting Task-Oriented Dialogue Systems with Relation Extraction. In <i>Proceedings of the 10th Dialog System Technology Challenges (AAAI Workshop)</i>.</p>	

* Equal Contribution

Andrew Lee, Jonathan K. Kummerfeld, Larry An, Rada Mihalcea. 2021. Micromodels for Efficient, Explainable, and Reusable Systems: A Case Study on Mental Health. *Findings of Empirical Methods in Natural Language Processing (Findings of EMNLP)*.

Stefan Larson, Anish Mahendran, Joseph J. Peper, Chris Clarke, **Andrew Lee**, Parker Hill, Jonathan K. Kummerfeld, Kevin Leach, Michael Laurenzano, Lingjia Tang and Jason Mars. 2019. An Evaluation Dataset for Intent Classification and Out-of-Scope Prediction. *Empirical Methods in Natural Language Processing (EMNLP)*.

Stefan Larson, Anish Mahendran, **Andrew Lee**, Jonathan K. Kummerfeld, Parker Hill, Michael Laurenzano, Johann Hauswald, Lingjia Tang, Jason Mars. 2019. Outlier Detection for Improved Data Quality and Diversity in Dialog Systems. *North American Chapter of the Association for Computational Linguistics (NAACL)*.

Fellowships **OpenAI Superalignment Student Fellowship** 2024
Awarded to 50 out of 2,700 applicants.

Employments **Meta AI Research** New York, NY
Research Intern - RAM (Reasoning, Attention, Memory) Team May 2023 - October 2023
Advisors: Jing Xu, Sainbayar Sukhbaatar, Jason Weston

Meta AI Research New York, NY
Research Intern - Diplomacy Team May 2022 - December 2022
Advisors: Emily Dinan, Mike Lewis

Microsoft Research Redmond, WA
Research Intern - KTX Team May 2021 - August 2021
Advisor: Silviu-Petru Cucerzan

Clinic, Inc. Ann Arbor, MI
Core AI R&D - Senior Software Engineer, Team Lead June 2019 - August 2020
Core AI R&D - Software Engineer June 2017 - June 2019

Ford Motor Company Dearborn, MI
Software Engineer March 2016 - June 2017

Invited Talks **University of Texas - Austin: Social Applications and Impact of NLP** 2024
A Mechanistic Understanding of Alignment Algorithms

University of Cambridge 2024
A Mechanistic Understanding of Alignment Algorithms

Patents Systems and methods for slot relation extraction for machine learning task-oriented dialogue systems.
Andrew Lee, Zhenguo Chen, Jonathan K. Kummerfeld.
US Patent 11,734,519. 2023.

Systems and methods for constructing an artificially diverse corpus of training data samples for training a contextually-biased model for a machine learning-based dialogue system.
Andrew Lee, Stefan Larson, Chris Clarke, Kevin Leach, Jonathan K. Kummerfeld, Parker Hill, Johann Hauswald, Michael Laurenzano, Lingjia Tang, Jason Mars.
US Patent 10,796,104. 2020.

Systems and methods for automatically configuring training data for training machine learning models of a machine learning-based dialogue system including seeding training samples or curating a corpus of training data based on instances of training data identified as anomalous. Stefan Larson, Anish Mahendran, **Andrew Lee**, Jonathan K. Kummerfeld, Parker Hill, Michael Laurenzano, Johann Hauswald, Lingjia Tang, Jason Mars.
US Patent 10,679,150. 2020.

Awards

University of Michigan CSE Excellence in Climate, Diversity, Equity, and Inclusion (2023)
Award given to 5 students in the Computer Science and Engineering department for one's commitment and action in increasing diversity, equity, and inclusion in the division.

University of Michigan CSE Research Honors AI Lab Nominee (2023)

Professional Services

ACL - Association for Computational Linguistics
2021, 2022 - ACL Rolling Review
CLPsych - Workshop on Computational Linguistics and Clinical Psychology (@NAACL)
2022 - Organizing Committee
2021 - Volunteer
CBO - International Symposium on Code Generation and Optimization
2019 - Program Committee