

# Andrew Lee

---

<b>Contact Information</b>	<b>E-mail:</b> <a href="mailto:ajyl@umich.edu">ajyl@umich.edu</a> <b>Website:</b> <a href="https://ajyl.github.io">https://ajyl.github.io</a> <b>Google Scholar:</b> <a href="#">Link</a>	
<b>Research Interests</b>	Interpretability, Language Models, Dialogue, Natural Language Processing, Deep Learning	
<b>Education</b>	<b>University of Michigan</b>	Ann Arbor, MI
	Ph.D. Candidate in Computer Science	2020 - Present
	Advisor: Rada Mihalcea	
	<b>University of Michigan</b>	Ann Arbor, MI
	Master's in Computer Science	2015
	<b>Northwestern University</b>	Evanston, IL
	Bachelor of Science in Computer Science	2013
<b>Publications</b>	Core Francisco Park, Maya Okawa, <b>Andrew Lee</b> , Ekdeep Singh Lubana, and Hidenori Tanaka. 2024. Emergence of Hidden Capabilities: Exploring Learning Dynamics in Concept Space. <i>Workshop on High-dimensional Learning Dynamics (ICML)</i> .	
	<b>Andrew Lee</b> , 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>The International Conference on Machine Learning (ICML)</i> . <b>Oral presentation (Top 1.5% of submissions).</b>	
	<b>Andrew Lee</b> , 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> .	
	Oana Ignat, Zhijing Jin, Artem Abzaliev, Laura Biester, Santiago Castro, Naihao Deng, Xinyi Gao, Aylin Gunal, Jacky He, Ashkan Kazemi, Muhammad Khalifa, Namho Koh, <b>Andrew Lee</b> , 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> .	
	Neel Nanda*, <b>Andrew Lee*</b> , 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> . <b>Honorable Mention for Best-Paper.</b>	
	Jing Xu, <b>Andrew Lee</b> , Sainbayar Sukhbaatar, Jason Weston. 2023. Some things are more CRINGE than others: Preference Optimization with the Pairwise Cringe Loss. <i>Pre-print</i>	
	<b>Andrew Lee</b> , 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> .	

\* Equal Contribution

**Andrew Lee**, David Wu, Emily Dinan, Michael Lewis. 2022. Improving Chess Commentaries by Combining Language Models with Symbolic Reasoning Engines. *Pre-print*

**Andrew Lee**, Zhenguo Chen, Kevin Leach, Jonathan K. Kummerfeld. 2022. Augmenting Task-Oriented Dialogue Systems with Relation Extraction. In *Proceedings of the 10th Dialog System Technology Challenges (AAAI Workshop)*.

**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)*.

<b>Fellowships</b>	<b>OpenAI Superalignment Student Fellowship</b> Awarded to 50 out of 2,700 applicants.	2024
<b>Employments</b>	<b>Meta AI Research</b> Research Intern - RAM (Reasoning, Attention, Memory) Team Advisors: Jing Xu, Sainbayar Sukhbaatar, Jason Weston	New York, NY May 2023 - October 2023
	<b>Meta AI Research</b> Research Intern - Diplomacy Team Advisors: Emily Dinan, Mike Lewis	New York, NY May 2022 - December 2022
	<b>Microsoft Research</b> Research Intern - KTX Team Advisor: Silviu-Petru Cucerzan	Redmond, WA May 2021 - August 2021
	<b>Clinic, Inc.</b> Core AI R&D - Senior Software Engineer, Team Lead Core AI R&D - Software Engineer	Ann Arbor, MI June 2019 - August 2020 June 2017 - June 2019
	<b>Ford Motor Company</b> Software Engineer	Dearborn, MI March 2016 - June 2017
<b>Invited Talks</b>	<b>University of Texas - Austin: Social Applications and Impact of NLP</b> A Mechanistic Understanding of Alignment Algorithms	2024
	<b>University of Cambridge</b> A Mechanistic Understanding of Alignment Algorithms	2024
<b>Patents</b>	Systems and methods for slot relation extraction for machine learning task-oriented dialogue systems. Andrew Lee, Zhenguo Chen, Jonathan K. Kummerfeld. <i>US Patent 11,734,519</i> . 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