

ATTICUS GEIGER

Stanford University
atticusg@stanford.edu

EDUCATION

Ph.D. Linguistics

Stanford University, Stanford, CA

September 2019 - December 2023

M.S. Computer Science

Stanford University, Stanford, CA

September 2016 - June 2019

B.S. Symbolic Systems

Stanford University, Stanford, CA

September 2015- June 2019

CAREER HISTORY

Principal Investigator, Practical AI Alignment and Interpretability Research Group

October 2023 -

- Acquired grant from Open Philanthropy to start a remote, non-profit research group focused on analyzing and interpreting AI models
- Support students internationally with remote research internship opportunities
- Conduct interpretability research grounded in theories of causality and abstraction

PhD Research, Stanford University

September 2019 - December 2023

- Advised by Christopher Potts and Thomas Icard (Also mentored by Noah Goodman and Mike Frank)
- Publications at Psychological Review, NeurIPS, ICML, EMNLP, ACL, NAACL, CLear, CogSci and BlackBoxNLP
- Collaborated with and mentored several graduate students and undergraduates

Honors Thesis Research, Stanford University

June 2018 - September 2018

- Acquired grant for self led natural language inference research project advised by Chris Potts, Thomas Icard, and Lauri Karttunen
- Constructed artificial natural language inference dataset using logic models
- Designed task specific neural model with standout performance on the generated datasets

Symbolic Systems Research Intern, Stanford University

June 2017 - September 2017

- Worked with Lauri Karttunen and Ignacio Cases to create a Natural Language Inference dataset focused on implicatives
- Implemented neural network models for natural language inference in Tensor Flow

Software Engineer Intern, Alaska Satellite Facility

June 2016 - September 2016

- Created a AWS cloud processing infrastructure and website interface
- Provided a service that automatically processes incoming satellite data
- Accomplished significant, largely self led software implementation

SKILLS

Programming: Python, C++, C, TensorFlow, PyTorch, AWS

Social Skills: I love collaborating and find great joy in mentoring!

AWARDS AND HONORS

Firestone Medal Award for B.S. honors thesis titled *Can Natural Language Inference Models Perform Natural Logic Reasoning?* and advised by Chris Potts and Thomas Icard.

GRANTS

PI, Fair Adversarial Tasks for Natural Language Understanding. Facebook Research: Robust Deep Learning for Natural Language Processing. Co-PI Christopher Potts. 2019-2020

PI, Practical AI Alignment and Interpretability Research Group. Open Philanthropy. 2023-2025

1. Amir Zur, Elisa Kreiss, Karel D'Oosterlinck, Christopher Potts, and **Atticus Geiger**. Updating clip to prefer descriptions over captions, 2024
2. Zhengxuan Wu*, Aryaman Arora*, Zheng Wang, **Atticus Geiger**, Dan Jurafsky, Christopher D. Manning, and Christopher Potts. Reft: Representation finetuning for language models. *CoRR*, abs/2404.03592, 2024
3. **Atticus Geiger***, Zhengxuan Wu*, Christopher Potts, Thomas Icard, and Noah D. Goodman. Finding alignments between interpretable causal variables and distributed neural representations. In Francesco Locatello and Vanessa Didelez, editors, *Causal Learning and Reasoning, 1-3 April 2024, Los Angeles, California, USA*, volume 236 of *Proceedings of Machine Learning Research*, pages 160–187. PMLR, 2024
4. Zhengxuan Wu, **Atticus Geiger**, Aryaman Arora, Jing Huang, Zheng Wang, Noah D. Goodman, Christopher D. Manning, and Christopher Potts. pyvene: A library for understanding and improving pytorch models via interventions. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT 2024)*, page To appear, Mexico City, Mexico, June 2024. Association for Computational Linguistics. Demonstrations Track
5. Jing Huang, Zhengxuan Wu, Christopher Potts, Mor Geva, and **Atticus Geiger**. RAVEL: evaluating interpretability methods on disentangling language model representations. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL 2024)*, Bangkok, Thailand, August 2024. Association for Computational Linguistics. To appear
6. Zhengxuan Wu, **Atticus Geiger**, Jing Huang, Aryaman Arora, Thomas Icard, Christopher Potts, and Noah D. Goodman. A reply to Makelov et al. (2023)’s “interpretability illusion” arguments. arXiv:2401.12631, 2024
7. Curt Tigges, Oskar John Hollinsworth, **Atticus Geiger**, and Neel Nanda. Linear representations of sentiment in large language models, 2023
8. Zhengxuan Wu*, **Atticus Geiger***, Christopher Potts, and Noah D. Goodman. Interpretability at scale: Identifying causal mechanisms in Alpaca. In *Advances in Neural Information Processing Systems*, 2023
9. Jing Huang, **Atticus Geiger**, Karel D'Oosterlinck, Zhengxuan Wu, and Christopher Potts. Rigorously assessing natural language explanations of neurons. In *Proceedings of the Sixth BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*. Association for Computational Linguistics, December 2023
10. Jingyuan S. She, Sam R. Bowman, Christopher Potts, and **Atticus Geiger**. Scone: Benchmarking negation reasoning in language models with fine-tuning and in-context learning. In *Proceedings of the 2023 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Toronto, Canada, July 2022. Association for Computational Linguistics
11. Angela Cao*, **Atticus Geiger***, Elisa Kreiss*, Thomas Icard, and Tobias Gerstenberg. A semantics for causing, enabling, and preventing verbs using structural causal models. In *Proceedings of the Cognitive Science Society*, 2023
12. Riccardo Massidda, **Atticus Geiger**, Thomas Icard, and Davide Bacciu. Causal abstraction with soft interventions. In *2nd conference on Causal Learning and Reasoning (CLear)*, Tbingen, Germany, 2023
13. **Atticus Geiger**, Alexandra Carstensen, Michael C Frank, and Christopher Potts. Relational reasoning and generalization using nonsymbolic neural networks. *Psychol. Rev.*, 130(2):308–333, March 2023
14. Zhengxuan Wu*, Karel D'Oosterlinck*, **Atticus Geiger***, Amir Zur, and Christopher Potts. Causal Proxy Models for concept-based model explanations. In *Proceedings of the 40th International Conference on Machine Learning*, 2023. arXiv:2209.14279
15. Eldar David Abraham* and Karel D'Oosterlinck* and Amir Feder* and Yair Ori Gat* and **Atticus Geiger*** and Christopher Potts* and Roi Reichart* and Zhengxuan Wu*. Cebab: Estimating the causal effects of real-world concepts on NLP model behavior. *CoRR*, abs/2205.14140, 2022

16. Zhengxuan Wu*, **Atticus Geiger***, Josh Rozner, Elisa Kreiss, Hanson Lu, Thomas Icard, Christopher Potts, and Noah D. Goodman. Causal distillation for language models. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4288–4295, Seattle, United States, July 2022. Association for Computational Linguistics
17. **Atticus Geiger***, Zhengxuan Wu*, Hanson Lu*, Josh Rozner, Elisa Kreiss, Thomas Icard, Noah Goodman, and Christopher Potts. Inducing causal structure for interpretable neural networks. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 7324–7338. PMLR, 17–23 Jul 2022
18. **Atticus Geiger***, Hanson Lu*, Thomas Icard, and Christopher Potts. Causal abstractions of neural networks. In *Advances in Neural Information Processing Systems*, 2021
19. Christopher Potts, Zhengxuan Wu, **Atticus Geiger**, and Douwe Kiela. DynaSent: A dynamic benchmark for sentiment analysis. In *Proceedings of the Association for Computational Linguistics*, 2021
20. Douwe Kiela, Max Bartolo, Yixin Nie, Divyansh Kaushik, **Atticus Geiger**, Zhengxuan Wu, Bertie Vidgen, Grusha Prasad, Amanpreet Singh, Pratik Ringshia, Zhiyi Ma, Tristan Thrush, Sebastian Riedel, Zeerak Waseem, Pontus Stenetorp, Robin Jia, Mohit Bansal, Christopher Potts, and Adina Williams. Dynabench: Rethinking benchmarking in NLP. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4110–4124, Online, June 2021. Association for Computational Linguistics
21. **Atticus Geiger**, Kyle Richardson, and Christopher Potts. Neural natural language inference models partially embed theories of lexical entailment and negation. In *Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 163–173, Online, November 2020. Association for Computational Linguistics
22. **Atticus Geiger**, Ignacio Cases, Lauri Karttunen, and Christopher Potts. Posing fair generalization tasks for natural language inference. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4475–4485, Stroudsburg, PA, November 2019. Association for Computational Linguistics
23. **Atticus Geiger**. Can natural language inference models perform natural logic reasoning? B.s. thesis, Stanford University, 2019
24. Ignacio Cases, Clemens Rosenbaum, Matthew Riemer, **Atticus Geiger**, Tim Klinger, Alex Tamkin, Olivia Li, Sandhini Agarwal, Joshua D. Greene, Dan Jurafsky, Christopher Potts, and Lauri Karttunen. Recursive routing networks: Learning to compose modules for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Stroudsburg, PA, June 2019. Association for Computational Linguistics
25. **Atticus Geiger**, Ignacio Cases, Lauri Karttunen, and Christopher Potts. Stress-testing neural models of natural language inference with multiply-quantified sentences. Ms., Stanford University. arXiv 1810.13033, 2018
26. Steve A. Arko, Rose Hogenson, **Atticus Geiger**, Jake Herrmann, Brian Buechler, and Kirk Hogenson. Sentinel-1 Archive and Processing in the Cloud using the Hybrid Pluggable Processing Pipeline (HyP3) at the ASF DAAC. In *AGU Fall Meeting Abstracts*, volume 2016, pages G43A–1040, December 2016
27. Kirk. Hogenson, Steve A. Arko, Brian Buechler, Rose Hogenson, Jake Herrmann, and **Atticus Geiger**. Hybrid Pluggable Processing Pipeline (HyP3): A cloud-based infrastructure for generic processing of SAR data. In *AGU Fall Meeting Abstracts*, volume 2016, pages IN21B–1740, December 2016

INVITED TALKS

1. 2023. AI Alignment Workshop at NeurIPS in New Orleans. Theories and Tools for Mechanistic Interpretability via Causal Abstraction. December.
2. 2023. AI safety and mechanistic interpretability conference at MIT. Theories and Tools for Mechanistic Interpretability via Causal Abstraction. May 6-7.

3. 2023. Finding Alignments Between Interpretable Causal Variables and Distributed Neural Representations. Workshop on Causal Representation Learning (CRL 2023) at the Max Planck Institute for Intelligent Systems in Tbingen. April 17-19th.
4. 2022. Causal Abstraction and Computational Explanations in Artificial Intelligence. Seminar Series in Stanford Psychology for cognitive and neuroscience areas. January 21.
5. 2021. Causal Abstractions and Interchange Intervention Training. Allen AI Aristo Group. September 23.
6. 2021. Causal Abstraction for Neural Network Analysis. McDonnell Plenary. July 8.
7. 2020. Modular Representation in Neural Natural Language Inference Models. Allen AI Aristo Group. June 19.

TEACHING

1. 2021 (Fall). Probabalistic Pragmatics. Guest Lecture for LINGUIST 245 *Psycholinguistics* at Stanford. (Graduate and Undergraduate)
2. 2020 (Spring). Causal Abstraction with Applications to Computational Implementation and Neural Network Analysis. Guest Lecture for PHIL 359 *Logic Seminar* at Stanford. (Graduate)
3. 2020 (Spring). Natural Logic. Guest Lecture for LINGUIST 130B *Lexical Semantics* at Stanford. (Undergraduate)
4. 2019 (Spring). Evaluating NLU Models with Harder Generalization Tasks. Guest Lecture for CS 224U *Natural Language Understanding* at Stanford. (Graduate)
5. Teaching Assistant Positions: CS 224U *Natural Language Understanding* at Stanford 2019 (Spring), LINGUIST 130B *Lexical Semantics* at Stanford 2020 (Spring) , and LINGUIST 245 *Psycholinguistics* at Stanford 2021 (Fall).