responsible foundation models [4, HAIC@ICLR25 (workshop)] for educational LLM applications [7, Frontiers AI (journal)].

Doctoral Researcher at UIUC

- Studied foundational mechanistic interpretability, including...
 - defining and evaluating the reliability of leading causal probing methods [1, IAI@NeurIPS24 (workshop oral)].
 - introducing a general causal probing framework for LLM interpretation and analvsis and new causal probing methods based on adversarial machine learning

1

Doctoral Researcher at INVITE

RESEARCH EXPERIENCE

PhD Research Intern at Microsoft Research

- (conference)].
- Studied LLM-based learner agent simulation for educational AI [4, AAAI25
- (conference oral)].
- Studied LLM steering for distribution-shift robustness and bias mitigation [6, ICML25

• Worked with domain-area experts to define and operationalize principles of socially

- Studying how LLMs internally represent the latent structure of symbolic reasoning tasks in order to predict and improve their compositional generalization.

University of Illinois Urbana-Champaign , Urbana, IL	08/2021 - Present
Ph.D. in Computer Science (anticipated graduation May 2026)	1
University of Utah, Salt Lake City, Utah	08/2016 - 05/2021
B.S. in Computer Science (May 2021, cum laude)	

B.S. in Cognitive Science (May 2021, cum laude)

EDUCATION

PhD candidate at UIUC (University of Illinois Urbana-Champaign), advised by Profs. ChengXiang Zhai and Julia Hockenmaier.

<u>Research areas:</u> natural language processing, (mechanistic) interpretability, distribution-shift robustness, causal machine learning, synthetic data, and multimodal representation learning.

Adam Davies

adavies4@illinois.edu +1 (801) 357-9217 https://ahdavies6.github.io/

05/2025 - Present

05/2024 - 05/2025

08/2022 - 05/2024

- [3, IAI@NeurIPS24 (workshop)].
- surveying the history of interpretability and its parallels with cognitive science, up through current categories of interpretability methods and associated goals, key assumptions, and shared challenges [2, preprint].
- Evaluated the abstract shape recognition abilities of vision-language models by synthesizing benchmarks using conditional generative models [5, NeurIPS24 (conference)], and studied how synthetic data from text-to-image models can improve distribution-shift robustness of image classifiers [9, ICML24 (conference)].

Doctoral Researcher at \underline{NCSA}

08/2021 - 08/2022

• Researched intersection of NLP, data mining, and computational social science for studying social construction using "big data" historical newspaper collections [10, **JCSS** (journal)] and [8, **PASC** (conference oral)].

PUBLICATIONS

- [1] Marc Canby*, <u>Adam Davies*</u>, Chirag Rastogi, and Julia Hockenmaier. Measuring the reliability of causal probing methods: Tradeoffs, limitations, and the plight of nullifying interventions. In *NeurIPS 2024 Workshop on Interpretable AI*, 2024. URL https: //openreview.net/forum?id=tmpMQLxVHh.
- [2] <u>Adam Davies</u> and Ashkan Khakzar. The cognitive revolution in interpretability: From explaining behavior to interpreting representations and algorithms. arXiv preprint arXiv:2408.05859, 2024. URL https://arxiv.org/abs/2408.05859.
- [3] <u>Adam Davies</u>, Jize Jiang, and ChengXiang Zhai. Competence-based analysis of language models. In *NeurIPS 2024 Workshop on Interpretable AI*, 2024. URL https: //openreview.net/forum?id=x6ZM5Is2Po.
- [4] <u>Adam Davies</u>, Elisa Nguyen, Michael Simeone, Erik Johnston, and Martin Gubri. Social science is necessary for operationalizing socially responsible foundation models. In *ICLR 2025 Workshop on Human-AI Coevolution*, 2025. URL https://openreview. net/forum?id=zbB2vjAq7X.
- [5] Arshia Hemmat, <u>Adam Davies</u>, Tom A. Lamb, Jianhao Yuan, Philip Torr, Ashkan Khakzar, and Francesco Pinto. Hidden in plain sight: Evaluating abstract shape recognition in vision-language models. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang, editors, *Advances in Neural Information Processing Systems*, volume 37, pages 88527–88556. Curran Associates, Inc., 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/a13ff984831deea39e6132bafdfdd6d5-Paper-Datasets_and_Benchmarks_Track.pdf.

- [6] Tom A. Lamb, <u>Adam Davies</u>, Alasdair Paren, Philip Torr, and Francesco Pinto. Focus on this, not that! steering LLMs with adaptive feature specification. In *Forty-second International Conference on Machine Learning*, 2025. URL https://openreview.net/ forum?id=rbI5mOUA8Z.
- [7] Amogh Mannekote, <u>Adam Davies</u>, Juan D Pinto, Shan Zhang, Daniel Olds, Noah L Schroeder, Blair Lehman, Diego Zapata-Rivera, and ChengXiang Zhai. Large language models for whole-learner support: opportunities and challenges. *Frontiers in Artificial Intelligence*, 7:1460364, 2024. URL https://www.frontiersin.org/journals/artificial-intelligence/articles/10.3389/frai.2024.1460364/full.
- [8] Sandeep Puthanveetil Satheesan, <u>Adam Davies</u>, Alan B Craig, Yu Zhang, and ChengXiang Zhai. Toward a big data analysis system for historical newspaper collections research. In *Proceedings of the Platform for Advanced Scientific Computing Conference*, pages 1–11, 2022. URL https://doi.org/10.1145/3539781.3539795.
- [9] Jianhao Yuan*, Francesco Pinto*, <u>Adam Davies*</u>, and Philip Torr. Not just pretty pictures: Toward interventional data augmentation using text-to-image generators. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp, editors, *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 57924–57952. PMLR, 21–27 Jul 2024. URL https://proceedings.mlr. press/v235/yuan24e.html.
- [10] Yu Zhang, <u>Adam Davies</u>, and ChengXiang Zhai. Understanding the social construction of juvenile delinquency: insights from semantic analysis of big-data historical newspaper collections. *Journal of Computational Social Science*, pages 1–43, 2024. URL https://link.springer.com/article/10.1007/s42001-024-00254-x.

TECHNICAL SKILLS

- Deep Learning in Python: PyTorch, TensorFlow, Keras, 😕 Transformers
- Data Science & Machine Learning in Python: NumPy, SciPy, scikit-learn, Pandas, 😕 Datasets
- Classic NLP in Python: spaCy, NLTK, CoreNLP, WordNet, gensim
- Scientific Visualization in Python: Matplotlib, Seaborn, Plotly, Jupyter
- Collaboration and Publishing: Git, LATEX, Overleaf, and Markdown.

TALKS

• Measuring the Reliability of Causal Probing Methods (Oral. NeurIPS24 Workshop on Interpretable AI)	12/2024
• Cognitive Interpretability in the Era of LLMs (Guest Lecture, UIUC Seminar in Psychology)	10/2024
 Causal Probing for Language Model Interpretability and Analysis (Tutorial, University of Oxford) 	09/2023
• Computational Social Science with Historical Text Analysis (Oral, Platform for Advanced Scientific Computing Conference)	06/2022

TEACHING AND MENTORSHIP

Research Supervision and Mentoring

Advised the following undergraduate students:

• <u>Chirag Rastogi</u> (UIUC BS) • <u>Publication</u> [1] (topic: <i>evaluating interpretability methods</i>)	07/2023 - 10/2024
 <u>Jize Jiang</u> (UIUC BS → MS) Ondergraduate thesis (topic: formal reasoning with LLMs) First publication [3] (topic: language model interpretability) 	01/2023 - 05/2023
Co-advised the following undergraduate students:	
 <u>Arshia Hemmat</u> (Oxford internship) First conference publication [5] (topic: evaluating abstract shift) 	01/2024 - 08/2024 ape recognition)
 <u>Jianhao Yuan</u> (Oxford BS → PhD) o Undergraduate thesis [9] (topic: synthetic data for distribution 	10/2022 - 05/2023 m-shift robustness)
Teaching Assistant at <u>UIUC</u>	08/2023 - 05/2024

- Applied Machine Learning (Spring 2024)
- Natural Language Processing (Fall 2023)