

EMPLOYMENT	<b>Allen Institute for AI</b> <i>Research Scientist</i> • Advisor: Yejin Choi	Seattle, US 2023 - 2024
	<b>Allen Institute for AI</b> <i>Postdoctoral Fellow</i> • Advisor: Yejin Choi	Seattle, US 2022 - 2023
	<b>Mila – Quebec Artificial Intelligence Institute / McGill University</b> <i>Visiting Scholar</i> • Advisor: Siva Reddy	Montreal, CA 2021 - 2022
	<b>Google DeepMind</b> <i>Student Researcher</i> • Advisors: Tal Linzen, David Reitter, Hannah Rashkin	NYC, US 2020 - 2022
	<b>Microsoft Research</b> <i>Research Intern</i> • Advisors: Alessandro Sordoni, Goeff Gordon	Montreal, CA 2019 - 2020
	<b>Google DeepMind</b> <i>Research Intern</i> • Advisors: Diyi Yang, Tom Kwiatkowski	NYC, US 2019
EDUCATION	<b>Ph.D. Computing Science, University of Alberta</b> <i>Thesis: Mitigating Hallucinations in Conversational LLMs.</i> • Advisor: Osmar Zaiane, GPA: 4.00/4.00	Edmonton, Canada 2018 - 2022
	<b>MSc. Computing Science, University of Alberta</b> <i>Thesis: Response Generation For An Open-Ended Conversational Agent</i> • Advisor: Osmar Zaiane, GPA: 4.00/4.00	Edmonton, Canada 2016 - 2018
PREPRINTS	You can find an exhaustive list of my publications in my Google Scholar profile.	
	1. <b>RL Grokking Recipe: How Does RL Unlock and Transfer New Algorithms in LLMs?.</b> Yiyu Sun, Yuhao Cao, Pohao Huang, Haoyue Bai, Hannaneh Hajishirzi, Dawn Song*, <u>Nouha Dziri</u> * (* = equal advising role) <i>Arxiv 2025 (in submission to ICLR)</i>	
	2. <b>Climbing the Ladder of Reasoning: What LLMs Can-and Still Can't-Solve after SFT?.</b> Yiyu Sun, Shawn Hu, Georgia Zhou, Ken Zheng, Hannaneh Hajishirzi, <u>Nouha Dziri</u> , Dawn Song. <i>Arxiv 2025</i>	
	3. <b>OpenAgentSafety: A Comprehensive Framework for Evaluating Real-World AI Agent Safety.</b> Sanidhya Vijayvargiya, Aditya Bharat Soni, Xuhui Zhou, Zora Zhiruo Wang, <u>Nouha Dziri</u> , Graham Neubig, Maarten Sap. <i>Arxiv 2025 (in submission to ICLR)</i>	

4. **OMEGA: Can LLMs Reason Outside the Box in Math? Evaluating Exploratory, Compositional, and Transformative Generalization.** Yiyu Sun, Georgia Zhou, Hao Wang, Dacheng Li, Dawn Song\*, Nouha Dziri\* (\* = equal advising role)  
*NeurIPS 2025*
5. **Artificial Hivemind: The Open-Ended Homogeneity of Language Models (and Beyond)** Liwei Jiang, Yuanjun Chai, Margaret Li, Mickel Liu, Raymond Fok, Maarten Sap, Yulia Tsvetkov, Nouha Dziri, Yejin Choi.  
*NeurIPS 2025*  
**Oral presentation (Top 1.4% among 25k submissions)**
6. **Why and How LLMs Hallucinate: Connecting the Dots with Subsequence Associations** Yiyu Sun, Yu Gai, Lijie Chen, Abhi Ravichander, Yejin Choi, Nouha Dziri, Dawn Song.  
*NeurIPS 2025*
7. **2 OLMo 2 Furious.** Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, ..., Nouha Dziri, Noah A. Smith, Hannaneh Hajishirzi.  
*COLM 2025*  
**2M downloads on HuggingFace**
8. **TÜLU 3: Pushing Frontiers in Open Language Model Post-Training.** Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V. Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, Yuling Gu, Saumya Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Chris Wilhelm, Luca Soldaini, Noah A. Smith, Yizhong Wang, Pradeep Dasigi and Hannaneh Hajishirzi.  
*COLM 2025*  
**300K downloads on HuggingFace**
9. **SafetyAnalyst: Interpretable, Transparent, and Steerable LLM Safety Moderation** Jing-Jing Li, Valentina Pyatkin, Max Kleiman-Weiner, Liwei Jiang, Nouha Dziri, Anne GE Collins, Jana Schaich Borg, Maarten Sap, Yejin Choi, Sydney Levine.  
*ICML 2025*
10. **AI as Humanity’s Salieri: Quantifying Linguistic Creativity of Language Models via Systematic Attribution of Machine Text against Web Text.** Ximing Lu, Melanie Sclar, Skyler Hallinan, Niloofar Mireshghallah, Jiacheng Liu, Seungju Han, Allyson Ettinger, Liwei Jiang, Khyathi Chandu, Nouha Dziri, Yejin Choi.  
*ICLR 2025*  
**Oral presentation (Top 1.8% among 12k submissions)**  
**Featured in Science Magazine**
11. **WildBench: Benchmarking LLMs with Challenging Tasks from Real Users in the Wild.** Bill Yuchen Lin, Yuntian Deng, Khyathi Chandu, Faeze Brahman, Abhilasha Ravichander, Valentina Pyatkin, Nouha Dziri, Ronan Le Bras, Yejin Choi.  
*ICLR 2025*
12. **Steering Masked Discrete Diffusion Models Via Discrete Denoising Posterior Prediction** Jarrod Rector-Brooks, Mohsin Hasan, Zhangzhi Peng, Zachary Quinn, Chenghao Liu, Sarthak Mittal, Nouha Dziri, Michael Bronstein, Yoshua Bengio, Pranam Chatterjee, Alexander Tong, Avishek Joey Bose.  
*ICLR 2025*
13. **RewardBench: Evaluating Reward Models for Language Modeling.** Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, Noah A Smith, Hannaneh Hajishirzi.  
*NAACL 2025*
14. **Rel-AI: An Interaction-Centered Approach To Measuring Human-LM Reliance.** Kaitlyn Zhou, Jena D Hwang, Xiang Ren, Nouha Dziri, Dan Jurafsky, Maarten Sap.  
*NAACL 2025*  
**Best Paper Runner-Up 2025 (Top 0.2%)**
15. **To Err is AI: A Case Study Informing LLM Flaw Reporting Practices** Sean McGregor, Allyson Ettinger, Nick Judd, Paul Albee, Liwei Jiang, Kavel Rao, Will Smith, Shayne Longpre, Avijit Ghosh, Christopher Fiorelli, Michelle Hoang, Sven Cattell, Nouha Dziri.  
*AAAI 2025*

16. **WildTeaming at Scale: From In-the-Wild Jailbreaks to (Adversarially) Safer Language Models.** Liwei Jiang, Kavel Rao, Seungju Han, Allyson Ettinger, Faeze Brahman, Sachin Kumar, Niloofar Mireshghallah, Maarten Sap, Yejin Choi, Nouha Dziri.  
*NeurIPS 2024*
17. **WildGuard: Open One-Stop Moderation Tools For Safety Risks, Jailbreaks, and Refusals of LLMs.** Seungju Han, Kavel Rao, Allyson Ettinger, Liwei Jiang, Bill Yuchen Lin, Nathan Lambert, Yejin Choi, and Nouha Dziri.  
*NeurIPS 2024*  
**215K downloads on HuggingFace**
18. **Multi-Attribute Constraint Satisfaction via Language Model Rewriting.** Ashutosh Baheti, Debanjana Chakraborty, Faeze Brahman, Ronan Le Bras, Ximing Lu, Nouha Dziri, Yejin Choi, Mark Riedl, Maarten Sap.  
*TMLR 2024*
19. **A Roadmap to Pluralistic Alignment.** Taylor Sorensen, Jared Moore, Jillian Fisher, Mitchell Gordon, Niloofar Mireshghallah, Christopher Michael Rytting, Andre Ye, Liwei Jiang, Ximing Lu, Nouha Dziri, Tim Althoff, Yejin Choi.  
*ICML 2024*
20. **Elastic Weight Removal For Faithful and Abstractive Dialogue generation.** Nico Daheim, Nouha Dziri, Mrinmaya Sachan, Iryna Gurevych, Edoardo Ponti.  
*NAACL 2024*
21. **Faith and Fate: Limits of Transformers on Compositionality.** Nouha Dziri, Ximing Lu, Melanie Sclar, Xiang Lorraine Li, Liwei Jian, Bill Yuchen Lin, Peter West, Chandra Bhagavatula, Ronan Le Bras, Jena D Hwang, Soumya Sanyal, Sean Welleck, Xiang Ren, Allyson Ettinger, Zaid Harchaoui, Yejin Choi.  
*NeurIPS 2023*  
**Spotlight (Top 2% among 12k submissions)**  
**Featured in Quanta Magazine, Science News, Montreal AI Ethics Institute**  
**15 invited talks, 1 podcast, 2 guest lectures**
22. **Fine-Grained Human Feedback Gives Better Rewards for Language Model Training.** Zeqiu Wu, Yushi Hu, Weijia Shi, Nouha Dziri, Alane Suhr, Prithviraj Ammanabrolu, Noah A Smith, Mari Ostendorf, Hannaneh Hajishirzi.  
*NeurIPS 2023*  
**Spotlight (Top 2% among 12k submissions)**
23. **Self-Refine: Iterative Refinement with Self-Feedback.** Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Sean Welleck, Bodhisattwa Prasad Majumder, Shashank Gupta, Amir Yazdanbakhsh, Peter Clark.  
*NeurIPS 2023*
24. **The Generative AI Paradox:” What It Can Create, It May Not Understand”. Peter West\*, Ximing Lu\*, Nouha Dziri\*, Faeze Brahman\*, Linjie Li\*, Jena D Hwang, Liwei Jiang, Jillian Fisher, Abhilasha Ravichander, Khyathi Chandu, Benjamin Newman, Pang Wei Koh, Allyson Ettinger, Yejin Choi. (\* = equal contribution)**  
*ICLR 2024*
25. **Phenomenal Yet Puzzling: Testing Inductive Reasoning Capabilities of Language Models with Hypothesis Refinement.** Linlu Qiu, Liwei Jiang, Ximing Lu, Melanie Sclar, Valentina Pyatkin, Chandra Bhagavatula, Bailin Wang, Yoon Kim, Yejin Choi, Nouha Dziri\*, Xiang Ren\*. (\* = equal advising role)  
*ICLR 2024*  
**Oral presentation (Top 1.8% among 7k submissions)**
26. **The Unlocking Spell on Base LLMs: Rethinking Alignment via In-Context Learning.** Bill Yuchen Lin, Abhilasha Ravichander, Ximing Lu, Nouha Dziri, Melanie Sclar, Khyathi Chandu, Chandra Bhagavatula, Yejin Choi.  
*ICLR 2024*

27. **Value Kaleidoscope: Engaging AI with pluralistic human values, rights, and duties**  
Taylor Sorensen, Liwei Jiang, Jena Hwang, Sydney Levine, Valentina Pyatkin, Peter West, Nouha Dziri, Ximing Lu, Kavel Rao, Chandra Bhagavatula, Maarten Sap, John Tasioulas, Yejin Choi.  
*AAAI 2024*
28. **Culture-Gen: Revealing global cultural perception in language models through natural language prompting**  
Huihan Li, Liwei Jiang, Jena D Hwang, Hyunwoo Kim, Sebastin Santy, Taylor Sorensen, Bill Yuchen Lin, Nouha Dziri, Xiang Ren, Yejin Choi.  
*COLM 2024*
29. **What Makes it Ok to Set a Fire? Iterative Self-distillation of Contexts and Rationales for Disambiguating Defeasible Social and Moral Situations** Kavel Rao, Liwei Jiang, Valentina Pyatkin, Yuling Gu, Niket Tandon, Nouha Dziri, Faeze Brahman, Yejin Choi.  
*EMNLP 2023*
30. **Evaluating Open-Domain Question Answering in the Era of Large Language Models.**  
Ehsan Kamalloo, Nouha Dziri, Charles LA Clarke, Davood Rafiei.  
*ACL 2023*
31. **CHAMPAGNE: Learning Real-world Conversation from Large-Scale Web Videos.**  
Seungju Han, Jack Hessel, Nouha Dziri, Yejin Choi, Youngjae Yu.  
*ICCV 2023*
32. **On the Origin of Hallucinations in Conversational Models: Is it the Datasets or the Models?** Nouha Dziri, Sivan Milton, Mo Yu, Osmar Zaiane, Siva Reddy.  
*NAACL 2022*
33. **FaithDial: A Faithful Benchmark for Information-Seeking Dialogue.**  
Nouha Dziri, Ehsan Kamalloo, Sivan Milton, Osmar Zaiane, Mo Yu, Edoardo M Ponti, Siva Reddy.  
*TACL 2022*
34. **Neural Path Hunter: Reducing Hallucination in Dialogue Systems via Path Grounding.**  
Nouha Dziri, Andrea Madotto, Osmar Zaiane, Avishek Joey Bose.  
*EMNLP 2021*
35. **Decomposed Mutual Information Estimation For Contrastive Representation Learning.**  
Alessandro Sordani\*, Nouha Dziri\*, Hannes Schulz\*, Geoff Gordon, Philip Bachman, Remi Tachet Des Combes. (\* = equal contribution)  
*ICML 2021*
36. **Evaluating Coherence in Dialogue Systems Using entailment.**  
Nouha Dziri\*, Ehsan Kamalloo\*, Kory W Mathewson, Osmar Zaiane. (\* = equal contribution)  
*NAACL 2019*

**Computer Use Agents**

*ICML workshop, Vancouver*

July 2025

**Panelists:** Graham Neubig, Ruslan Salakhutdinov, Yu Su, Victor Zhong, Alexandre Drouin

**Data in Generative Models**

*ICML workshop, Vancouver*

July 2025

**Panelists:** Tatsunori Hashimoto, Eric Wong, Serena Booth, Pin-Yu Chen, Ivan Evtimov

**AI Reasoning**

*Cross Future AI & Technology Summit, Vancouver*

July 2025

**Panelists:** Wenhu Chen, Ken Perlin

**Trustworthy Foundation Models**

*International Symposium on Trustworthy Foundation Models, UAE*

May 2025

**Panelists:** Tomas Mikolov, Hakim Hacid, Tongliang Liu

INVITED PANELS	<b>Making the Right Career Call across Academia and Industry</b> <i>ICLR Workshop Women in ML, Singapore</i> <b>Panelists: Claire Vernade, Reyhane Askari, Katherine Driscoll</b>	April 2025
	<b>Meta-Generation Algorithms for Large Language Models</b> <i>NeurIPS Tutorial, Vancouver</i> <b>Panelists: Noam Brown, Jakob Foerster, Rishabh Agarwal, Beidi Chen</b>	Dec 2024
	<b>Analyzing the Plasticity of LLMs During Language Interactions.</b> <i>NeurIPS Language Gamification Workshop, Vancouver</i> <b>Panelists: Aaron Courville, Alane Suhr, Tom Schaul, Marc Lanctot, Tom Griffiths, and Sam Devlin</b>	Dec 2024
	<b>Towards A Global Standard for AI Product Safety</b> <i>MLCommons, San Francisco</i> <b>Panelists: Ion Stoica, April Chen, Wan Sie Lee</b>	Nov 2024
INVITED TALKS	<b>Reasoning in LLMs</b> <i>Workshop on Cognitive Basis of Reasoning/IVADO, Montreal</i>	Jan 2026
	<b>Autonomous LLM Agents: Risks and Scientific Challenges</b> <i>IVADO/Mila - Quebec Artificial Intelligence Institute, Montreal</i>	Nov 2025
	<b>Faith and Fate: Limits of Transformers in Reasoning</b> <i>D.E.Shaw Research, NYC</i>	Sept 2025
	<b>LLM Reasoning: Advanced Inference-time Strategies</b> <i>Armenia LLM Summer School, Armenia</i>	July 2025
	<b>Can LLMs Reason Outside the Box in Math?</b> <i>Apple Workshop on Reasoning and Planning, Cupertino</i>	July 2025
	<b>Faith and Fate: Limits of Transformers in Reasoning</b> <i>Cross Future AI &amp; Technology Summit, Vancouver</i>	July 2025
	<b>Trustworthy LLMs: How Data Quality Shapes Performance and Where It Falls Short?</b> <i>ICML 2025 Workshop on Data in Generative Models, Vancouver</i>	July 2025
	<b>OpenAgentSafety: A Comprehensive Framework for Evaluating AI Agent Safety</b> <i>ICML 2025 Workshop on Computer-Use Agents, Vancouver</i>	July 2025
	<b>Does Scaling Guarantee Trustworthy LLMs</b> <i>International Symposium on Trustworthy Foundation Models, UAE</i>	May 2025
	<b>In-Context Learning in LLMs: Potential and Limits</b> <i>Causality in the Era of Foundation Models Workshop, Barbados</i>	Feb 2025
	<b>Red-teaming and Safeguarding LLMs</b> <i>International Association for Safe and Ethical AI, Paris</i>	Feb 2025
	<b>In-Context Learning in LLMs: Potential and Limits</b> <i>NeurIPS Language Gamification Workshop, Vancouver</i>	Dec 2024
	<b>What it can create, it may not understand: Studying the Limits of Transformers</b> <i>University of Cambridge</i>	May 2024
	<b>Limits of Generative AI Models and their Societal Implications.</b> <i>Princeton University</i>	Dec 2023
	<b>Faith and Fate: Limits of Transformers on Compositionality</b> <i>The Alan Turing Institute, UK</i>	Nov 2023

INVITED TALKS	<b>Faith and Fate: Limits of Transformers on Compositionality.</b> <i>University of Edinburgh</i>	Nov 2023
	<b>Faith and Fate: Limits of Transformers on Compositionality</b> <i>SAIL workshop on fundamental limits of LLMs, Germany</i>	Oct 2023
	<b>Faith and Fate: Limits of Transformers on Compositionality</b> <i>University of Pittsburgh</i>	Oct 2023
	<b>Faith and Fate: Limits of Transformers on Compositionality</b> <i>Formal Languages and Neural Networks Seminar, US</i>	Sep 2023
	<b>Towards Building Hallucination-Free Conversational Models</b> <i>Stanford University</i>	Aug 2022
	<b>FaithDial: A Faithful Benchmark for Information-Seeking Dialogue</b> <i>Google Research, NYC</i>	May 2022
	<b>FaithDial: A Faithful Benchmark for Information-Seeking Dialogue</b> <i>Amazon Research, Seattle</i>	June 2022
	<b>Evaluating Coherence in Dialogue Systems Using Entailment.</b> <i>Google DeepMind, Montreal</i>	Dec 2019
GUEST LECTURER	<b>LLM Summer School 2025</b> <i>Organizers: ServiceNow Research/ Nvidia/ Meta</i> <i>Yerevan, Armenia</i>	Summer 2025
	<b>Lecture: LLM Reasoning: Advanced Inference-Time Strategies</b>	
	<b>11-430/830 Ethics, Safety, and Social Impact in NLP and LLMs</b> <i>Instructor: Maarten Sap</i> <i>Carnegie Mellon University</i>	Winter 2025
	<b>Lecture: Redteaming and Safeguarding in LLMs</b>	
	<b>Generative AI Seminar Course</b> <i>Instructor: Adji Bousso Dieng</i> <i>Princeton University</i>	Fall 2023
	<b>Lecture: Limits of Generative AI Models and their Societal Implications.</b>	
	<b>CSE 599 D1: (Grad) Exploration on Language, Knowledge, and Reasoning</b> <i>Instructor: Yejin Choi</i> <i>University of Washington</i>	Winter 2023
	<b>Lecture: Can LLMs reason?</b>	
TEACHING ASSISTANT	<b>CMPUT 650: (Grad) Computational Semantics</b> <i>Instructor: Greg Kondrak</i> <i>University of Alberta</i>	Winter 2018
	<b>Lecture: Challenges in Conversational LLMs</b>	
	<b>CMPUT 101 - Introduction to Computing</b> <i>University of Alberta</i>	2018
	<b>Instructor: Marianne Morris</b>	

STUDENTS MENTORING	• <b>Ximing Lu</b> , Predoctoral student at Ai2 → PhD student at UW <i>ICLR 2025 Oral (top 1%), NeurIPS 2023 Spotlight (top 2%)</i>	2022-2024
	• <b>Liwei Jiang</b> , PhD student at UW <i>NeurIPS 2025 Oral (top 1%), NeurIPS 2024</i>	2023-2025
	• <b>Seungju Han</b> , BSc at Seoul National University → PhD at Stanford <i>NeurIPS 2024, ICCV 2023</i>	2023-2024
	• <b>Kavel Rao</b> , MSc student at UW → Researcher at Jane Street <i>NeurIPS 2024, EMNLP 2023</i>	2023-2024
	• <b>Linlu Qiu</b> , PhD student at MIT <i>ICLR 2024 Oral (top 1%)</i>	2023
	• <b>Huihan Li</b> , PhD student at University of Southern California <i>COLM 2024</i>	2023-2024
	• <b>Yiyu Sun</b> , Postdoc at University of California, Berkeley <i>2 papers at NeurIPS 2025</i>	2025-current
	• <b>Shrey Jain</b> , MSc student at Carnegie Mellon University	2025-current
WORKSHOP CO- ORGANIZER	• The first workshop for Research on Agent Language Models (REALM)	ACL 2025
	• The Reasoning and Planning for Large Language Models	ICLR 2025
	• System 2 Reasoning at Scale	NeurIPS 2024
	• The third workshop on Document-Grounded Dialogue Systems (DialDoc)	ACL 2023
ACADEMIC SERVICES	<b>Demo Chair:</b> NAACL 2025	
	<b>Senior Area Chair:</b> ACL 2025 in the area of Ethics, Bias, and Fairness	
	<b>Area Chair:</b> EMNLP 2023, COLM (2024-2025)	
	<b>Reviewer:</b> NeurIPS (2022-2024), ICLR (2022-2024), ACL (2018-2023), EMNLP (2018-2023), NAACL (2018-2023), EACL (2018-2022)	
PRESS	<ul style="list-style-type: none"> <li>• <b>Quanta Magazine:</b> <i>Chatbot Software Begins to Face Fundamental Limitations</i>, 2025.</li> <li>• <b>Science Magazine:</b> <i>AI writing is improving, but it can't match human creativity</i>, 2024.</li> <li>• <b>LeMonde:</b> <i>Should we be concerned about the 'hallucinations' of AI like ChatGPT or Gemini?</i>, 2024.</li> <li>• <b>ScienceNews:</b> <i>AI's reasoning skills can't be assessed by current tests</i>, 2024.</li> <li>• <b>TechCrunch:</b> <i>Treating a chatbot nicely might boost its performance — here's why</i>, 2024.</li> <li>• <b>Podcast (Women in AI Research):</b> <i>Limits of Transformers</i>, 2025.</li> </ul>	
AWARDS AND HONORS	• <b>Runner-up Best Paper</b> at NAACL 2025 (top 0.1%)	2025
	• <b>Outstanding Reviewer</b> at ACL 2021 (top 1%)	2021
	• <b>Alberta Doctoral Recruitment Scholarship</b> (\$10,000)	2018
	• <b>Mitacs Globalink PhD Graduate Fellowship</b> (\$44,000)	2018
	• <b>National Scholarship for MSc and PhD studies</b> (\$70,000)	2016
	• <b>Best Poster Award</b> , ACM Canadian Celebration of Women in Computing (\$400)	2017
	• <b>Mitacs Globalink MSc Graduate Fellowship</b> (\$15,000)	2016
	• <b>DAAD Scholarship</b> for Research Internship, Leipzig, Germany (€10,000)	2015
	• <b>Erasmus Mundus Exchange Scholarship</b> for BSc studies (€50,000)	2015