

Eric Wallace

EDUCATION	UC Berkeley Ph.D. in Computer Science Research Advisors: Dan Klein, Dawn Song Thesis: <i>Emerging Vulnerabilities of Large Language Models</i> University of Maryland B.S. in Computer Engineering Research Advisor: Jordan Boyd-Graber	2019–2024 2014–2018
WORK EXPERIENCE	OpenAI <i>Member of Technical Staff</i> Google Deepmind <i>Research Intern</i> Research Advisors: Dustin Tran, Denny Zhou, Xinyun Chen Facebook AI Research (FAIR) <i>Research Intern</i> Research Advisors: Robin Jia, Douwe Kiela Allen Institute for Artificial Intelligence (AI2) <i>Research Intern</i> Research Advisors: Matt Gardner, Sameer Singh	San Francisco, CA Nov 2023–Present Mountain View, CA June 2023–Aug 2023 Menlo Park, CA June 2021–Sep 2021 Irvine, CA Jan 2019–Aug 2019
SELECTED AWARDS	ICML Best Paper Award, 2024 Apple Fellowship in AI/ML, 2022–2024 Outstanding Paper Award, NeurIPS 2023 RegML Workshop PET Award Runner Up (test of time award for [27]), 2023 Best Poster, NeurIPS 2021 ENLSP Workshop First Superhuman Crossword AI, ACPT 2021 Best Demo Paper, EMNLP 2019 AI2 Intern of the Year, 2019	
PUBLICATIONS	<ul style="list-style-type: none">[1] OpenAI o3 and o4-mini System Card Eric Mitchell, ..., Eric Wallace[2] Trading Inference-Time Compute for Adversarial Robustness Wojciech Zaremba, Evgenia Nitishinskaya, Boaz Barak, Stephanie Lin, Sam Toyer, Yaodong Yu, Rachel Dias, Eric Wallace, Kai Xiao, Johannes Heidecke, Amelia Glaese <i>arXiv preprint</i>, 2025.[3] OpenAI Deep Research System Card Zhiqing Sun, Isa Fulford, ..., Eric Wallace[4] OpenAI o3 Mini System Card Hongyu Ren, ..., Eric Wallace[5] OpenAI o1 System Card Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, ..., Eric Wallace, et al. <i>arXiv preprint</i>, 2024.[6] Deliberative Alignment: Reasoning Enables Safer Language Models Melody Guan, Manas Joglekar, Eric Wallace, Saachi Jain, Boaz Barak, Alec Helyar, Rachel Dias, Andrea Vallone, Hongyu Ren, Jason Wei, Hyung Won Chung, Sam Toyer, Johannes Heidecke, Alex Beutel, Amelia Glaese <i>arXiv preprint</i>, 2024.[7] GPT-4o System Card Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, ..., Eric Wallace, et al. <i>arXiv preprint</i>, 2024.[8] The Instruction Hierarchy: Training LLMs to Prioritize Privileged Instructions Eric Wallace*, Kai Xiao*, Reimar Leike*, Lilian Weng, Johannes Heidecke, Alex Beutel <i>arXiv preprint</i>, 2024.	

- [9] Unfamiliar Finetuning Examples Control How Language Models Hallucinate
Katie Kang, **Eric Wallace**, Aviral Kumar, Claire Tomlin, Sergey Levine
NAACL, 2025.
- [10] Stealing Part of a Production Language Model
Nicholas Carlini, Krishnamurthy Dvijotham, Milad Nasr, A. Feder Cooper, Katherine Lee, Matthew Jagielski, Thomas Steinke, Daniel Paleka, Jonathan Hayase, Arthur Conmy, David Rolnick, Florian Tramér, **Eric Wallace**
International Conference on Machine Learning (ICML), 2024.
Best Paper Award
- [11] Covert Malicious Finetuning: Subverting LLM Safety Training Without Detection
Danny Halawi*, Alexander Wei*, **Eric Wallace**, Tony Tong Wang, Nika Haghtalab, Jacob Steinhardt
International Conference on Machine Learning (ICML), 2024.
- [12] What Evidence Do Language Models Find Convincing?
Alex Wan, **Eric Wallace**, Dan Klein
Association for Computational Linguistics (ACL), 2024.
- [13] Scalable Extraction of Training Data from (Production) Language Models
Milad Nasr, Nicholas Carlini, Jonathan Hayase, Matthew Jagielski, A Feder Cooper, Daphne Ippolito, Christopher A Choquette-Choo, **Eric Wallace**, Florian Tramér, Katherine Lee
arXiv preprint, 2023.
- [14] Privacy Side Channels in Machine Learning Systems
Edoardo Debenedetti, Giorgio Severi, Nicholas Carlini, Christopher A. Choquette-Choo, Matthew Jagielski, Milad Nasr, **Eric Wallace**, Florian Tramér
arXiv preprint, 2023.
- [15] SILO Language Models: Isolating Legal Risk In a Nonparametric Datastore
Sewon Min*, Suchin Gururangan*, **Eric Wallace**, Hannaneh Hajishirzi, Noah A. Smith, Luke Zettlemoyer
International Conference on Learning Representations (ICLR), 2024.
Spotlight Presentation (Top 5%)
- [16] The False Promise of Imitating Proprietary LLMs
Arnav Gudibande*, **Eric Wallace***, Charlie Snell*, Xinyang Geng, Hao Liu, Pieter Abbeel, Sergey Levine, Dawn Song
International Conference on Learning Representations (ICLR), 2024.
Spotlight Presentation (Top 5%)
- [17] Extracting Training Data from Diffusion Models
Nicholas Carlini, Jamie Hayes, Milad Nasr, Matthew Jagielski, Vikash Sehwal, Florian Tramér, Borja Balle, Daphne Ippolito, **Eric Wallace**
USENIX Security Symposium, 2023.
- [18] Poisoning Language Models During Instruction Tuning
Alexander Wan*, **Eric Wallace***, Sheng Shen, Dan Klein
International Conference on Machine Learning (ICML), 2023.
- [19] Large Language Models Struggle to Learn Long-Tail Knowledge
Nikhil Kandpal, Haikang Deng, Adam Roberts, **Eric Wallace**, Colin Raffel
International Conference on Machine Learning (ICML), 2023.
- [20] InCoder: A Generative Model for Code Infilling and Synthesis
Daniel Fried, Armen Aghajanyan, Jessy Lin, Sida Wang, **Eric Wallace**, Freda Shi, Ruiqi Zhong, Wen-tau Yih, Luke Zettlemoyer, Mike Lewis
International Conference on Learning Representations (ICLR), 2023.
Spotlight Presentation
- [21] Measuring Forgetting of Memorized Training Examples
Matthew Jagielski, Om Thakkar, Florian Tramér, Daphne Ippolito, Katherine Lee, Nicholas Carlini, **Eric Wallace**, Shuang Song, Abhradeep Thakurta, Nicolas Papernot, Chiyuan Zhang
International Conference on Learning Representations (ICLR), 2023.
- [22] Deduplicating Training Data Mitigates Privacy Risks in Language Models
Nikhil Kandpal, **Eric Wallace**, Collin Raffel
International Conference on Machine Learning (ICML), 2022.
- [23] Automated Crossword Solving
Eric Wallace*, Nicholas Tomlin*, Albert Xu*, Kevin Yang*, Eshaan Pathak*, Matt Ginsberg, Dan Klein
Association for Computational Linguistics (ACL), 2022.
First Superhuman Crossword AI
- [24] Analyzing Dynamic Adversarial Training Data in the Limit
Eric Wallace, Adina Williams, Robin Jia, Douwe Kiela
Findings of the Association for Computational Linguistics (ACL Findings), 2022.

- [25] Cutting Down on Prompts and Parameters: Simple Few-Shot Learning with Language Models
Robert L. Logan IV, Ivana Balažević, **Eric Wallace**, Fabio Petroni, Sameer Singh, Sebastian Riedel
ACL Findings 2022; NeurIPS Efficient NLP Workshop.
Best Poster Award
- [26] Calibrate Before Use: Improving Few-shot Performance of Language Models
Tony Z. Zhao*, **Eric Wallace***, Shi Feng, Dan Klein, Sameer Singh
International Conference on Machine Learning (ICML), 2021.
Long Oral Presentation (Top 3%)
- [27] Extracting Training Data from Large Language Models
Nicholas Carlini, Florian Tramèr, **Eric Wallace**, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Úlfar Erlingsson, Alina Oprea, Colin Raffel
USENIX Security Symposium, 2021.
Runner up for PET Award (Test of Time Award)
- [28] Concealed Data Poisoning Attacks on NLP Models
Eric Wallace*, Tony Z. Zhao*, Shi Feng, Sameer Singh
North American Chapter of the Association for Computational Linguistics (NAACL), 2021.
- [29] Detoxifying Language Models Risks Marginalizing Minority Voices
Albert Xu, Eshaan Pathak, **Eric Wallace**, Maarten Sap, Suchin Gururangan, Dan Klein
North American Chapter of the Association for Computational Linguistics (NAACL), 2021.
- [30] Imitation Attacks and Defenses for Black-box Machine Translation Systems
Eric Wallace, Mitchell Stern, Dawn Song
Empirical Methods in Natural Language Processing (EMNLP), 2020.
- [31] Evaluating Models’ Local Decision Boundaries via Contrast Sets
Matt Gardner, Yoav Artzi, ... (other authors hidden) ... **Eric Wallace**, Ally Zhang, Ben Zhou
Findings of the Empirical Methods in Natural Language Processing (EMNLP Findings), 2020.
- [32] AutoPrompt: Eliciting Knowledge from Language Models with Automatically Generated Prompts
Taylor Shin*, Yasaman Razeghi*, Robert L Logan IV*, **Eric Wallace**, Sameer Singh
Empirical Methods in Natural Language Processing (EMNLP), 2020.
- [33] Gradient-based Analysis for NLP Models is Manipulable
Junlin Wang*, Jens Tuyls*, **Eric Wallace**, Sameer Singh
Findings of the Empirical Methods in Natural Language Processing (EMNLP Findings), 2020.
- [34] Train Large, Then Compress: Rethinking Model Size for Efficient Training and Inference of Transformers
Zhuohan Li*, **Eric Wallace***, Sheng Shen*, Kevin Lin*, Kurt Keutzer, Dan Klein, Joseph E. Gonzalez
International Conference on Machine Learning (ICML), 2020.
- [35] Pretrained Transformers Improve Out-of-Distribution Robustness
Dan Hendrycks*, Xiaoyuan Liu*, **Eric Wallace**, Adam Dziedziec, Rishabh Krishnan, Dawn Song
Association for Computational Linguistics (ACL), 2020.
- [36] Universal Adversarial Triggers for Attacking and Analyzing NLP
Eric Wallace, Shi Feng, Nikhil Kandpal, Matt Gardner, Sameer Singh
Empirical Methods in Natural Language Processing (EMNLP), 2019.
- [37] AllenNLP Interpret: A Framework for Explaining Predictions of NLP Models
Eric Wallace, Jens Tuyls, Junlin Wang, Sanjay Subramanian, Matt Gardner, Sameer Singh
Demo at Empirical Methods in Natural Language Processing (EMNLP), 2019.
Best Demo Award
- [38] Do NLP Models Know Numbers? Probing Numeracy in Embeddings
Eric Wallace*, Yizhong Wang*, Sujian Li, Sameer Singh, Matt Gardner
Empirical Methods in Natural Language Processing (EMNLP), 2019.
- [39] Misleading Failures of Partial-input Baselines
Shi Feng, **Eric Wallace**, Jordan Boyd-Graber
Association for Computational Linguistics (ACL), 2019.
- [40] Compositional Questions Do Not Necessitate Multi-hop Reasoning
Sewon Min*, **Eric Wallace***, Sameer Singh, Matt Gardner, Hannaneh Hajishirzi, Luke Zettlemoyer
Association for Computational Linguistics (ACL), 2019.
- [41] Understanding Impacts of High-Order Loss Approximations and Features in Deep Learning Interpretation
Sahil Singla, **Eric Wallace**, Shi Feng, Soheil Feizi.
International Conference on Machine Learning (ICML), 2019.
- [42] Trick Me If You Can: Human-in-the-loop Generation of Adversarial Examples for Question Answering
Eric Wallace, Pedro Rodriguez, Shi Feng, Ikuya Yamada, Jordan Boyd-Graber
Transactions of the Association for Computational Linguistics (TACL), 2019.
- [43] Pathologies of Neural Models Make Interpretations Difficult
Shi Feng, **Eric Wallace**, Alvin Grissom II, Mohit Iyyer, Pedro Rodriguez, Jordan Boyd-Graber
Empirical Methods in Natural Language Processing (EMNLP), 2018.

TEACHING
EXPERIENCE

Courses:

- Co-instructor of Berkeley’s graduate-level NLP (CS 288) with 90 students in Spring 2023. Taught alongside Dan Klein and Kevin Lin. I developed and taught ~10 new lectures on language models and advanced NLP topics (e.g., RLHF, retrieval, vision-language models). I also developed new homeworks, coding assignments, and mentored students.
- Teaching assistant for Berkeley’s CS188: Artificial Intelligence in Summer 2023. Aside from typical TA duties (e.g., leading discussion sections), I co-designed the midterm and final exam.

Tutorials:

- EMNLP, 2020. *Interpreting Predictions of NLP Models*.

Guest Lectures for Courses:

- UC Berkeley 194/294-267, 2024. *Memorization in Large Language Models*
- USC CSCI 699, 2023. *Security & Privacy in NLP*
- UC Berkeley CS294/194-196, 2023. *Intro & Foundations of LLMs*
- Stanford CS 329X, 2023. *Security & Privacy in NLP*
- Washington University in St. Louis CSE 527A, 2022. *Security & Privacy in NLP*
- University of Minnesota CSCI 8980-06, 2022. *Robustness in NLP*
- UC Berkeley CS 288, 2022. *Robustness in NLP*
- ML @ Berkeley, 2022. *Security & Privacy in NLP*
- University of Stuttgart, 2022. *Interpreting Predictions of NLP Models*

MENTORING

Student Research Mentoring

- Alex Wan (2022-2024), UC Berkeley Undergrad. Published [18].
- Carolyn Wang (2023), UC Berkeley Undergrad.
- Arnav Gudibande (2022–2023), UC Berkeley Masters. Published [16]. Now at Perplexity AI.
- Tony Zhao (2020–2021), UC Berkeley Undergrad. Published [26, 28]. Now PhD at Stanford.
- Albert Xu (2020–2021), UC Berkeley Undergrad. Published [23, 29]. Now PhD at USC.
- Eshaan Pathak (2020–2021), UC Berkeley Undergrad. Published [23, 29]. Now at You.com
- Jens Tuyls (2019–2020), UC Irvine Undergrad. Published [33, 37]. Now PhD at Princeton.
- Junlin Wang (2019–2020), UC Irvine Undergrad. Published [33, 37]. Now PhD at Duke.
- Nikhil Kandpal (2019), UMD Undergrad. Published [36]. Now PhD at UNC.

Masters Thesis Advising

- Arnav Gudibande, 2023. *On Imitating Proprietary Language Models*. Chair: Dawn Song.

Other Mentoring

- BAIR Undergraduate Mentoring, 2022–2024.
- Association of Women in EE&CS Office Hours on Grad School, 2023.
- Women in Machine Learning (WiML), 2022–2023. PhD Application Assistance.
- Berkeley Equal Access Assistance Program (EAAA), 2022-2023. PhD Application Assistance.
- Berkeley AI4All, 2022. Instructor

PRESENTATIONS

External Talks

- CVPR Workshop on Adversarial ML for Computer Vision. *Making “GPT-Next” Trustworthy*
- ICLR Workshop on Secure and Trustworthy LLMs. *Making “GPT-Next” Trustworthy*
- ICLR Workshop on Data for Foundation Models. *Making “GPT-Next” Trustworthy*
- IEEE S&P Workshop on Security Architectures for GenAI. *Making “GPT-Next” Trustworthy*
- Simons Institute Workshop on LLMs, 2023. *Memorization in Large Language Models*
- Princeton, 2023. *Memorization in Large Language Models*
- Oracle Labs, 2023. *Memorization in Large Language Models*
- University of Maryland, 2023. *Memorization in Large Language Models*
- University of North Carolina, 2023. *Memorization in Large Language Models*
- USC ISI, 2022. *Emerging Vulnerabilities in Large-scale NLP Models*
- Malicious Life Podcast, 2022. *Hacking Language Models*
- Stanford, 2021. *What Can We Learn from Vulnerabilities of NLP Models?*
- Cornell, 2021. *What Can We Learn from Vulnerabilities of NLP Models?*
- DeepMind, 2021. *What Can We Learn from Vulnerabilities of NLP Models?*
- UT Austin, 2021. *What Can We Learn from Vulnerabilities of NLP Models?*

- CMU, 2021. *What Can We Learn from Vulnerabilities of NLP Models?*

Panels:

- Women in Machine Learning, 2022. *PhD Fellowships Applications*
- ACL Mentoring, 2022. *How to Keep Up with Work in the Field*
- Berkeley AI Hackathon (w/ 1500 participants), 2023. *Future of LLMs—Beyond Hacking*
- USENIX PEPR, 2023. *Privacy Challenges and Opportunities in LLM-Based Chatbots*

ACADEMIC SERVICE

Program Committee Member

- Journals: Computational Linguistics (2023)
- Conferences: ACL (2020, 2021, 2022), ICML (2021, 2023), NeurIPS (2020, 2021), EMNLP (2018, 2019, 2020, 2021, 2022), ACL Rolling Review (2021, 2022), ICLR (2023), NAACL (2021, 2022), COLM (2024)
- Workshops: Distribution Shifts (NeurIPS 2022, ICML 2022, NeurIPS 2023), BlackBox NLP (EMNLP 2022), RobustML Workshop (ICLR 2021), MRQA (EMNLP 2021), NLP for Positive Impact (ACL 2021), SRW (NAACL 2021), DistShift (NeurIPS 2021, NeurIPS 2023)

Area Chair

- Workshops: NextGenAISafety, ICML 2024

Workshop Organizer

- Language Model Memorization (L2M2), ACL 2025
- Future of Decentralization, AI, and Computing Summit, 2023 Berkeley.

Departmental Service

- Berkeley PhD Admissions. 2021–2023
- Berkeley Student Committee for Faculty Hiring. 2023
- Berkeley PhD Visit Days Recruitment. 2021–2024

Academic Grants & Sponsorships

- Led successful award for ~1500 TPUs from Google TRC (\$10M+ USD value)
- Apple Fellowship in AI/ML, 2022–2024

SELECTED MEDIA & PRESS

Extracting Training Data from Diffusion Models [17], MIT Technology Review, TWIML Podcast, Gizmodo, Vice, TechSpot, New Scientist, The Register, Ars Technica, Twitter #1 (3 million views), Twitter #2, Twitter #3,

Automated Crossword Solving [23], Discover, New Scientist, Atlantic, Wired, Slate, BBC, Science Friday, Top of Hacker News, The Register, Berkeley Engineering Magazine, WNPR, Daily Californian, NVIDIA Blog, Neil deGrasse Tyson Podcast, Twitter (1M views)

Extracting Training Data from Large Language Models [27], MIT Technology Review, Wired, Google Blog, BAIR Blog, Nature, Top of Hacker News, Twitter #1, Twitter #2, Twitter #3,