

# Eric Wong

Phone: +1 (339) 223-7159

Email: [exwong@upenn.edu](mailto:exwong@upenn.edu)

Website: <https://www.cis.upenn.edu/~exwong/>

Google Scholar: <https://scholar.google.com/citations?user=pWnTMRkAAAAJ>

*Last updated: 2025/07/13*

## Briefly

How can we make sure that deep learning models are actually doing what we want them to do? My research interests are centered around foundations of reliable machine learning systems: understanding, debugging, and guaranteeing the behavior of data-driven models. I created the first provable defenses that guarantee robustness to adversarial examples and real-world specifications, and currently work on securing modern foundation models. I am also interested in explaining models with provable certificates and scientific applications ranging from surgery to cosmology.

## Education

<b>University of Pennsylvania, Assistant Professor</b>	Philadelphia, PA 2022-current
<b>Massachusetts Institute of Technology, Post-Doctoral Associate</b> <i>Advisor: Aleksander Mądry</i>	Cambridge, MA 2020-2022
<b>Carnegie Mellon University, Ph. D. in Machine Learning</b> Thesis: Provable, structured, and efficient methods for robustness of deep networks to adversarial examples; SCS Dissertation Award — Honorable Mention <i>Advisor: Zico Kolter</i>	Pittsburgh, PA 2015-2020
<b>Carnegie Mellon University, B. S. in Computer Science</b> Double major in Mathematics, minor in Machine Learning	Pittsburgh, PA 2011-2015

## Work experience

2019-2020	<b>Bosch Center for Artificial Intelligence (Renningen, Germany and Pittsburgh, PA)</b> Created a virtual sensor based on neural networks for a fuel injection system in truck engines; formally verified the worst-case error of the system under conservative estimates of physical sensor noise.
2012-2015	<b>CERT Program (Pittsburgh, PA)</b> Migrated secure coding rules from POSIX to C11; analyzed security reports for Java android applications; developed an analysis tool for security vulnerabilities in source code.

## Awards

---

2025	<b>AI2050 Early Career Award</b> <i>Towards Robust Generative AI with Adaptive Risk Evaluations, Schmidt Sciences</i>
2025	<b>NSF Early Career Award</b> <i>CAREER: Certified Explanations for Trustworthy Artificial Intelligence, NSF</i>
2024	<b>Amazon Research Award (AWS AI)</b> <i>Adversarial Manipulation of Prompting Interfaces, Amazon</i>
2023	<b>Area Chair Award (Interpretability and Analysis of Models for NLP)</b> <i>Faithful Chain-of-Thought Reasoning, IJCNLP-AACL Conference</i>
2020	<b>SCS Dissertation Award – Honorable Mention</b> <i>Provable, structured, and efficient methods for robustness of deep networks to adversarial examples, Carnegie Mellon University</i>
2020	<b>Siebel Scholar Fellowship</b> Carnegie Mellon University
2017	<b>Best Defense Paper</b> <i>Provable defenses against adversarial examples via the convex outer adversarial polytope, NeurIPS 2017 ML &amp; Security Workshop</i>
2013	<b>Summer Undergraduate Research Fellowship</b> Carnegie Mellon University

## Publications

---

DMLR 2025	<b>The FIX Benchmark: Extracting Features Interpretable to eXperts</b> Helen Jin, Shreya Havaldar, Chaehyeon Kim, Anton Xue, Weiqiu You, Helen Qu, Marco Gatti, Daniel A. Hashimoto, Bhuvnesh Jain, Amin Madani, Masao Sako, Lyle Ungar, Eric Wong
ICML 2025	<b>Sum-of-Parts Models: Faithful Attributions for Groups of Features</b> Weiqiu You, Helen Qu, Marco Gatti, Bhuvnesh Jain, Eric Wong
ICML 2025	<b>DOLPHIN: A Programmable Framework for Scalable Neurosymbolic Learning</b> Aaditya Naik, Jason Liu, Claire Wang, Amish Sethi, Saikat Dutta, Mayur Naik, Eric Wong
NAACL-Findings 2025	<b>Avoiding Copyright Infringement via Machine Unlearning</b> Guangyao Dou, Zheyuan Liu, Qing Lyu, Kaize Ding, Eric Wong
ICLR 2025	<b>Logicbreaks: A Framework for Understanding Subversion of Rule-based Inference</b> Anton Xue, Avishree Khare, Rajeef Alur, Surbhi Goel, Eric Wong
TMLR 2025	<b>SmoothLLM: Defending Large Language Models Against Jailbreaking Attacks</b> Alexander Robey, Eric Wong, Hamed Hassani, George J. Pappas

SaTML 2025	<b>Jailbreaking Black Box Large Language Models in Twenty Queries</b> Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J. Pappas, Eric Wong
NeurIPS 2024	<b>AR-Pro: Counterfactual Explanations for Anomaly Repair with Formal Properties</b> Xiayan Ji, Anton Xue, Eric Wong, Oleg Sokolsky, Insup Lee
eBioMedicine	<b>Crowd-sourced machine learning prediction of long COVID using data from the National COVID Cohort Collaborative</b> Timothy Bergquist, Johanna Loomba, Emily Pfaff, Fangfang Xia, Zixuan Zhao, Yitan Zhu, Elliot Mitchell, Biplab Bhattacharya, Gaurav Shetty, Tamanna Munia, Grant Delong, Adbul Tariq, Zachary Butzin-Dozier, Yunwen Ji, Haodong Li, Jeremy Coyle, Seraphina Shi, Rachael V. Philips, Andrew Mertens, Romain Pirracchio, Mark van der Laan, John M. Colford Jr., Alan Hubbard, Jifan Gao, Guanhua Chen, Neelay Velingker, Ziyang Li, Yinjun Wu, Adam Stein, Jiani Huang, Zongyu Dai, Qi Long, Mayur Naik, John Holmes, Danielle Mowery, Eric Wong, Ravi Parekh, Emily Getzen, Jake Hightower, Jennifer Blase
NeurIPS 2024	<b>Data-Efficient Learning with Neural Programs</b> Alaia Solko-Breslin, Seewon Choi, Ziyang Li, Neelay Velingker, Rajeev Alur, Mayur Naik, Eric Wong
ICML 2024	<b>Towards Compositionality in Concept Learning</b> Adam Stein, Aaditya Naik, Yinjun Wu, Mayur Naik, Eric Wong
ICML 2024	<b>DISCRET: Synthesizing Faithful Explanations For Treatment Effect Estimation</b> Yinjun Wu, Mayank Keoliya, Kan Chen, Neelay Velingker, Ziyang Li, Emily J Getzen, Qi Long, Mayur Naik, Ravi B Parikh, Eric Wong
NeurIPS 2024	<b>JailbreakBench: An Open Robustness Benchmark for Jailbreaking Large Language Models</b> Patrick Chao, Edoardo Debenedetti, Alexander Robey, Maksym Andriushchenko, Francesco Croce, Vikash Sehwal, Edgar Dobriban, Nicolas Flammarion, George J. Pappas, Florian Tramèr, Hamed Hassani, Eric Wong
ICLR 2024, Tiny Papers (Oral)	<b>Evaluating Groups of Features via Consistency, Contiguity, and Stability</b> Chaehyeon Kim, Weiqiu You, Shreya Havaladar, Eric Wong
ICLR 2024	<b>SalUn: Empowering Machine Unlearning via Gradient-based Weight Saliency in Both Image Classification and Generation</b> Chongyu Fan, Jiancheng Liu, Yihua Zhang, Dennis Wei, Eric Wong, Sijia Liu
CVPR 2024	<b>Initialization Matters for Adversarial Transfer Learning</b> Andong Hua, Jindong Gu, Zhiyu Xue, Nicholas Carlini, Eric Wong, Yao Qin
OOPSLA 2024	<b>TorchQL: A Programming Framework for Integrity Constraints in Machine Learning</b> Aaditya Naik, Adam Stein, Yinjun Wu, Eric Wong, Mayur Naik
EMNLP 2023	<b>Comparing Styles across Languages</b> Shreya Havaladar, Matthew Pressimone, Eric Wong, Lyle Ungar

NeurIPS 2023	<b>Stability Guarantees for Feature Attributions with Multiplicative Smoothing</b> Anton Xue, Rajeev Alur, Eric Wong
ICLR 2023, Tiny Papers	<b>TopEx: Topic-based Explanations for Model Comparison</b> Shreya Havaladar, Adam Stein, Eric Wong, Lyle Ungar
ICML 2023	<b>Do Machine Learning Models Learn Statistical Rules Inferred from Data?</b> Aaditya Naik, Yinjun Wu, Mayur Naik, Eric Wong
DLSP 2023 Keynote	<b>Adversarial Prompting for Black Box Foundation Models</b> Natalie Maus*, Patrick Chao*, Eric Wong, Jacob Gardner
IJCNLP-AAACL, 2023	<b>Faithful Chain-of-Thought Reasoning</b> Qing Lyu*, Shreya Havaladar*, Adam Stein*, Li Zhang, Delip Rao, Eric Wong, Marianna Apidianaki, Chris Callison-Burch
CVPR 2023	<b>A data-based perspective on transfer learning</b> Saachi Jain*, Hadi Salman*, Alaa Khaddaj*, Eric Wong, Sung Min Park, Aleksander Madry
ICLR 2022	<b>Missingness bias in model debugging</b> Saachi Jain*, Hadi Salman*, Pengchuan Zhang, Vibhav Vineet, Sal Vemprala, Aleksander Madry
CVPR 2022	<b>Certified patch robustness via smoothed vision transformers</b> Hadi Salman*, Saachi Jain*, Eric Wong*, Aleksander Madry
OJCS 2022	<b>DeepSplit: Scalable verification of deep neural networks via operator splitting</b> Shaoru Chen*, Eric Wong*, J. Zico Kolter, Mahyar Fazlyab
ICML 2021 (Oral)	<b>Leveraging Sparse Linear Layers for Debuggable Deep Networks</b> Eric Wong*, Shibani Santurkar*, Aleksander Madry
ICLR 2021	<b>Learning perturbation sets for robust machine learning</b> Eric Wong, J. Zico Kolter
ICML 2020	<b>Overfitting in adversarially robust deep learning</b> Leslie Rice*, Eric Wong*, J. Zico Kolter
IEEE IV 2020	<b>Neural network virtual sensors for fuel injection quantities with provable performance specifications</b> Eric Wong, Tim Schneider, Joerg Schmitt, Frank R. Schmidt, J. Zico Kolter
ICLR 2020	<b>Fast is better than free: revisiting adversarial training</b> Eric Wong*, Leslie Rice*, J. Zico Kolter
ICML 2020	<b>Adversarial robustness against the union of multiple perturbation models</b> Pratyush Maini, Eric Wong, J. Zico Kolter
ICML 2019	<b>Wasserstein adversarial examples</b> Eric Wong, Frank R. Schmidt, J. Zico Kolter

NeurIPS 2018	<b>Scaling provable adversarial defenses</b> Eric Wong, Frank R. Schmidt, Jan Hendrik Metzen, J. Zico Kolter
ICML 2018	<b>Provable defenses against adversarial examples via the convex outer adversarial polytope</b> Eric Wong, J. Zico Kolter
ICML 2017	<b>A Semismooth Newton Method for Fast, Generic Convex Programming</b> Alnur Ali*, Eric Wong*, J. Zico Kolter
ICML 2015	<b>An SVD and Derivative Kernel Approach to Learning from Geometric Data</b> Eric Wong, J. Zico Kolter

## Preprints

---

2025	<b>Instruction Following by Boosting Attention of Large Language Models</b> Vitoria Guardieiro, Adam Stein, Avishree Khare, Eric Wong
2025	<b>Benchmarking Misuse Mitigation Against Covert Adversaries</b> Davis Brown, Mahdi Sabbaghi, Luze Sun, Alexander Robey, George J. Pappas, Eric Wong, Hamed Hassani
2025	<b>Probabilistic Stability Guarantees for Feature Attributions</b> Helen Jin, Anton Xue, Weiqiu You, Surbhi Goel, Eric Wong
2025	<b>Adaptively profiling models with task elicitation</b> Davis Brown, Prithvi Balehannina, Helen Jin, Shreya Havaldar, Hamed Hassani, Eric Wong
2025	<b>The Road to Generalizable Neuro-Symbolic Learning Should be Paved with Foundation Models</b> Adam Stein, Aaditya Naik, Neelay Velingker, Eric Wong
2024	<b>Defending Large Language Models against Jailbreak Attacks via Semantic Smoothing</b> Jiabao Ji, Bairu Hou, Alexander Robey, George J. Pappas, Hamed Hassani, Yang Zhang, Eric Wong, Shiyu Chang
2023	<b>Rectifying Group Irregularities in Explanations for Distribution Shift</b> Adam Stein, Yinjun Wu, Eric Wong, Mayur Naik
2023	<b>In-context Example Selection with Influences</b> Tai Nguyen, Eric Wong
2022	<b>When does bias transfer in transfer learning</b> Hadi Salman*, Saachi Jain*, Andrew Ilyas*, Logan Engstrom*, Eric Wong, Aleksander Madry

## Grants

---

2025-06-01 – 2027-05-31	<b>Towards Robust Generative AI with Adaptive Risk Evaluations</b> \$450k, AI2050 Early Career Fellowship, Schmidt Sciences, PI
2025-06-01 – 2030-05-31	<b>CAREER: Certified Explanations for Trustworthy Artificial Intelligence</b> \$675k, NSF, PI
2024-12-01 – 2027-05-31	<b>Harnessing Artificial Intelligence and Language Modeling for Enhancing Innovation and Evaluating Research Claims (HAILMEIER-C)</b> \$5.9M, DARPA, Co-PI with Chris Callison-Burch, Hannaneh Hajishirzi, Andrew Head, Peter Jansen, Chinedum Osuji, Yulia Tsvetkov, Duncan Watts
2025-03-01 – 2027-02-28	<b>TIGER: Trustable Information Generation and Explanation Resilience</b> \$4M, IARPA, Co-PI with Rene Vidal, Chris Callison-Burch, Hamed Hassani, Mark Yatskar, Rama Chellappa, Vishal Patel
2024-05-01 – 2025-04-30	<b>Preventing Complications with Transparent Surgical AI Assistants</b> \$100K, ASSET-IBI, PI with Daniel Hashimoto
2024-07-01 – 2028-06-30	<b>Safe and Explainable AI-enabled Decision Making for Personalized Treatment</b> \$6.85M, ARPA-H, Co-PI with Rajeev Alur, Rajat Deo, Sameed Ahmed M. Khatana, Qi Long, Mayur Naik, Ravi Parikh, Gary Weissman
2023-10-01 – 2027-09-30	<b>SLES: SPECSRL: Specification-guided Perception-enabled Conformal Safe Reinforcement Learning</b> \$1.5M, NSF, Co-PI with Rajeev Alur, Osbert Bastani & Dinesh Jayaraman
2024-05-01 – 2025-04-30	<b>Adversarial Manipulation of Prompting Interfaces</b> \$70K (+\$50K compute), Amazon Research Award, PI
2023-10-01 – 2027-09-30	<b>SHF: Medium: Scallop: A Neurosymbolic Programming Framework for Combining Logic with Deep Learning</b> \$1.2M, NSF, Co-PI with Mayur Naik & Rajeev Alur

## Media

---

2025	<b>Penn lab researches AI cyberbullying capabilities</b> (Penn Today) News article about our cyberbullying research. <a href="https://penntoday.upenn.edu/news/penn-seas-evaluating-large-language-models-cyberbullying-behavior">https://penntoday.upenn.edu/news/penn-seas-evaluating-large-language-models-cyberbullying-behavior</a>
2024	<b>The Llama 3 Herd of Models</b> (Meta) Meta technical report that used our PAIR red-teaming algorithm to develop their Llama3 models <a href="https://arxiv.org/pdf/2407.21783">https://arxiv.org/pdf/2407.21783</a>
2024	<b>Gemini 1.5 Report</b> (Google Deepmind) Google technical report that used our JailBreakBench evaluation to red-team Gemini models <a href="https://arxiv.org/pdf/2403.05530">https://arxiv.org/pdf/2403.05530</a>

- 2024 **Anthropic Sleeper Agents** (Anthropic)  
Google technical report that used our PAIR red-teaming algorithm to test the Claude family of models  
<https://arxiv.org/pdf/2401.05566>
- 2023 **New method reveals how one LLM can be used to jailbreak another** (VentureBeat)  
News article about our PAIR red-teaming algorithm.  
<https://venturebeat.com/ai/new-method-reveals-how-one-llm-can-be-used-to-jailbreak-another/>

## Invitations

---

- 2025 **Dynamic & Stateful Evals of Safety on the Frontier: What can Academics do?**  
Invited Talk, Data in Generative Models, ICML 2025 Workshop
- 2025 **Explanations for Experts via Guarantees and Domain Knowledge: From Attributions to Reasoning**  
Invited talk, Actionable Interpretability, ICML 2025 Workshop
- 2025 **How Jailbreaking 1-Layer Transformers Taught us how to Steer LLMs**  
Invited talk, Methods and Opportunities at Small Scale (MOSS), ICML 2025 Workshop
- 2025 **Explanations for Experts**  
Seminar speaker, New Jersey Institute of Technology, AI & Data Science Stars Seminar Series
- 2025 **What does a Foundation Model (not) Know?**  
Keynote speaker, Towards Knowledgeable Foundation Models @ AAAI 2025 Workshop
- 2024 **Convincing Experts to (not) Trust ML Models**  
Seminar speaker, Cornell AI Seminar
- 2024 **Jailbreaking LLMs: Attack, Defense, and Theory**  
Seminar speaker, University of Maryland
- 2023 **Robustness of Adversarial Attacks for LLM**  
Distinguished speaker, Responsible Machine Learning Summit, UCSB
- 2023 **Adversarial Prompting: Return of the Adversarial Example**  
Keynote speaker, IEEE S&P 2023, 6th Deep Learning Security & Privacy Workshop
- 2023 **From Prompt Engineering to Prompt Science**  
Seminar speaker, Wayne State University
- 2022 **Robustness for the Real World**  
Invited talk, 6th Annual Conference on Information Sciences and Systems (CISS)
- 2022 **Debuggable Deep Networks**  
Invited talk, TrustML Young Scientist Seminar

2021 **Panel Discussion**  
Panelist, ATVA 2021 Workshop on Security and Reliability of Machine Learning

## Patents

---

2024 **Methods, systems, and computer readable media for defending large language models (LLMs) against jailbreaking attacks** (18/907376)

2019 **Method, apparatus and computer program for generating robust automated learning systems and testing trained automated learning systems** (16/173698)

## Teaching Experience

---

2025 **Accelerating Research with Generative AI** - CIS 7000 (UPenn, Instructor)  
New special topics course to prepare students for best research practices augmented with LLMs

2025 **Machine Learning** - CIS 5200 (UPenn, Instructor)  
3rd round teaching CIS 5200

2024 **Mathematics of Machine Learning** - CIS 3333 (UPenn, Instructor)  
Second iteration of the new Mathematics of Machine Learning course under an official course number as a SEAS mathematics elective.

2024 **Machine Learning** - CIS 5200 (UPenn, Instructor)  
2nd round teaching CIS 5200, further consolidation and unification of the course material with the other AI faculty. 150 students.

2023 **Mathematics of Machine Learning** - CIS 3990 (UPenn, Instructor)  
Created a new course that prepares undergraduates for technical research and a graduate level coursework in machine learning. Two students that took the course last semester are now doing ML theory research.

2023 **Machine Learning** - CIS 5200 (UPenn, Instructor)  
Substantially overhauled and updated the Machine Learning course at UPenn to (a) fully autograded assignments using PennGrader, (b) brand new PyTorch-based programming assignments to replace dated Numpy notebooks, (c) expanded to a more balanced set of topics across all of Machine Learning (i.e. Duality/Lagrangian, MCMC, an entire theory module including PAC Learning & VC Theory, other learning paradigms like Online and Active Learning). 128 students.

2022 **Debugging Data & Models** - CIS 7000-005 (UPenn, Instructor)  
Designed new special topics course in the seminar/lecture format on debugging machine learning (7000-005). Taught 25 enrolled students with average instructor/course scores of 3.47/4 and 3.54/4.



2016	<b>Practical Data Science</b> - 15-388/688 (CMU, TA) Designed new assignments, taught recitations, and prepared write-ups for the first iteration of CMU's Practical Data Science course (15-388/688). I was the head TA (out of two TAs) and managed over 300 enrolled students. Received 55 student reviews with an average rating of 4.85/5.
2016-2019	<b>Eberly Center for Teaching Excellence and Educational Innovation</b> - Teaching Seminars (CMU, Participant) Enrolled in teaching seminars at the Eberly Center in CMU to develop personal teaching skills; seminars include "Teaching Inclusively: Leveraging Diversity and Promoting Equity in Your Classroom" and "Helping Students Develop Mastery and Critical Thinking"
2015	<b>Advanced Introduction to Machine Learning</b> - 10-715 (CMU, TA) Taught recitations, held office hours, and created/graded assignments for the second iteration of the Advanced Introduction to Machine Learning course intended for doctoral students in CMU's Machine Learning Department.
2014	<b>Algorithm Design and Analysis</b> - 15-451 (CMU, TA) Taught recitations, conducted oral examinations/office hours, and graded assignments/exams for the computer science department's Algorithm Design and Analysis course (15-451) at CMU.
2014	<b>Pervasive and Mobile Computing Services</b> - 08-766/781 (CMU, TA) Held office hours and graded assignments/exams for the software engineering department's Pervasive and Mobile Computing Services course (08-766/781) at CMU.
2013	<b>Pervasive and Mobile Computing Services</b> - 08-766/781 (CMU, TA) Held office hours and graded assignments/exams for the software engineering department's Pervasive and Mobile Computing Services course (08-766/781) at CMU.
2013	<b>Mobile Development for iOS and Android</b> - 08-723 (CMU, TA) Held office hours and graded assignments/exams for the software engineering department's Mobile Deveopment course (08-723) at CMU.

## Graduate Theses Supervised

---

Fall 2024 – Spring 2025	<b>Prithvi Balehannina</b> (Masters) Masters research project on LLM evaluations
Fall 2024 – Summer 2025	<b>Luze Sun</b> (Masters) Masters research project on LLM security
Fall 2024 – Spring 2026	<b>Cindy Xin</b> (Masters) Masters research project on AI-assisted discovery in cosmology
Spring 2024 – Summer 2025	<b>Rupkatha Hira</b> (Masters) Thesis: Anticipated Summer 2025
Fall 2024 – Spring 2029	<b>Cassandra Goldberg</b> (PhD) Thesis: Anticipated Spring 2029

Fall 2024 – Spring 2029	<b>Davis Brown</b> (PhD) Thesis: Anticipated Spring 2029, co-advised with Hamed Hassani
Fall 2024 – Spring 2029	<b>Vitoria Guardieiro</b> (PhD) Thesis: Anticipated Spring 2029
Spring 2024 – Spring 2026	<b>Adam Stein</b> (PhD) Thesis: Anticipated Spring 2026, co-advised with Mayur Naik
Fall 2023 – Spring 2028	<b>Chaehyeon Kim</b> (PhD) Thesis: Anticipated Spring 2028
Fall 2023 – Spring 2026	<b>Helen Jin</b> (PhD) Thesis: Anticipated Summer 2026
Summer 2023 – Spring 2025	<b>Anton Xue</b> (PhD) Thesis: Anticipated Spring 2025, co-advised with Rajeev Alur
Spring 2023 – Fall 2025	<b>Wei qiu You</b> (PhD) Thesis: Anticipated Spring 2026
Spring 2023 – Spring 2026	<b>Shreya Havaladar</b> (PhD) Thesis: Anticipated Spring 2026, co-advised with Lyle Ungar
Fall 2022 – Spring 2024	<b>Ningyuan Li</b> (Masters) Independent study
Fall 2022 – Spring 2024	<b>Shailesh Sridhar</b> (Masters) Thesis: Controlling for Missingness Bias in Feature Attribution Evaluation
Fall 2022 – Spring 2024	<b>Tai Nguyen</b> (Masters) Thesis: Attribute in-context learning examples with influences

## Undergraduate Projects Supervised

---

Spring 2025 – Fall 2025	<b>Siri Nellutla</b> (Undergraduate) Undergraduate research with VIPER on reasoning over material science documents
Fall 2024 – Spring 2025	<b>Lyuxin (David) Zhang</b> (Undergraduate) Undergraduate research project on in-context learning
Spring 2024 – Summer 2024	<b>Faraz Rahman</b> (Undergraduate) Undergraduate research on understanding features generated during diffusion
Spring 2024 – Spring 2024	<b>Dora Wu</b> (Undergraduate) Thesis: Image Generative Artificial Intelligence: Theory, Applications, and Outlook
Fall 2022 – Spring 2023	<b>Gideon Tesfaye</b> (Undergraduate) Practicum: Using deep learning to compose music

## Penn Service

---

2024 – 2025	<b>IDEAS Search Committee</b> Faculty search committee member
2024 – 2024	<b>PhD Admit Weekend</b> Organizer, Co-Lead with Andrew Head
2023 – 2024	<b>BSE in AI</b> Curriculum Committee
2023 – ongoing	<b>ML+FM Seminar</b> Organizer of a seminar series for researchers in the areas of formal methods and machine learning. Averages 18 attendees weekly.
2023 – 2023	<b>Adhoc Computing Cluster Committee</b> Committee Member
2023 – 2023	<b>PhD Admit Weekend</b> Organizer, Co-Lead with Andrew Head
2022 – ongoing	<b>Locust Cluster</b> Test driver for the SEAS cluster and working with CETS to iron out scalability of the system.

## DEI Service

---

2024-2025	<b>WiML@PennCIS (CIS)</b> Organizer of a new DEI event to help women in CIS form a community. Averages 20 attendees per monthly meeting for women in machine learning at CIS.
2023	<b>WiML Workshop Mentoring (NeurIPS)</b> Volunteered as a mentor for the round table event at the Women in Machine Learning workshop at NeurIPS.
2023	<b>USABE fireside chat (Penn Engineering)</b> Participated in the Faculty Fireside Chat Series, where students may talk to professors and faculty in a small-group setting run by the Underrepresented Student Advisory Board in Engineering (USABE). USABE is an organization that works with SEAS leadership to promote diversity, equity, and inclusion through student advocacy. One of their initiatives includes student-faculty engagement and allowing students to interact in more casual settings with faculty.
2022-2025	<b>Mentorship for Underrepresented Masters/Undergraduates at Penn (CIS)</b> Direct mentorship in research experiences of masters and undergraduate students that are underrepresented in CS (2 woman and 1 Ethiopian)
2021	<b>Graduate Application Assistance Program (MIT)</b> Assisted applicants from under-represented groups with their graduate student applications to MIT's EECS PhD program.
2021-2022	<b>MIT Undergraduate Research Opportunities Program (MIT)</b> Directly supervised an undergraduate for the UROP program at MIT; provided an opportunity for a member of an under-represented group to learn about machine learning and tackle a challenging research project

2020-2022	<b>MEnTorEd Opportunities in Research (METEOR) (MIT)</b> Participated in the METEOR postdoc fellowship selection committee, an effort at CSAIL MIT to increase diversity, equity, and inclusion. Provided confidential technical feedback on candidates based on their application materials.
2019-2020	<b>CMU AI Mentoring Program (CMU)</b> Mentored undergraduate women and minorities in one-on-one meetings to provide career advice and discuss research/graduate school; the mentee is now a PhD student at UC Berkeley.
2019-2020	<b>Teknowledge Mentor (Obama Academy)</b> Taught middle schoolers how to code as part of the Teknowledge outreach program at the Obama Academy; courses were intended to provide early exposure to computer science for under-represented students in low-income neighborhoods of Pittsburgh
2019	<b>Mental Health First Aid Certification (CMU)</b> Underwent training to recognize mental health issues and provide first aid assistance to those in need

## Workshop Organizing

---

2024	<b>3rd New Frontiers in Adversarial Machine Learning</b> Organizer for the 3rd workshop on new directions in adversarial machine learning held at NeurIPS 2024 Website: <a href="https://advml-frontier.github.io/">https://advml-frontier.github.io/</a>
2023	<b>2nd New Frontiers in Adversarial Machine Learning</b> Organizer for the 2nd ICML 2023 workshop on new directions in adversarial machine learning Website: <a href="https://advml-frontier.github.io/">https://advml-frontier.github.io/</a>
2022	<b>Workshop on Adversarial Machine Learning and Beyond</b> Organizer for an AAAI 2022 workshop broadly themed around adversarial machine learning Website: <a href="https://advml-workshop.github.io/aaai2022/">https://advml-workshop.github.io/aaai2022/</a>
2022	<b>New Frontiers in Adversarial Machine Learning</b> Organizer for an ICML 2022 workshop on new directions in adversarial machine learning Website: <a href="https://advml-frontier.github.io/">https://advml-frontier.github.io/</a>
2021	<b>A Blessing in Disguise: The Prospects and Perils of Adversarial Machine Learning</b> Organizer for an ICML 2021 workshop themed around the dangers and benefits of adversarial machine learning Website: <a href="https://advml-workshop.github.io/icml2021/">https://advml-workshop.github.io/icml2021/</a>
2021	<b>Robust and reliable ML in the real world</b> Main organizer for an ICLR 2021 workshop on real world robustness. Website: <a href="https://sites.google.com/connect.hku.hk/robustml-2021/home">https://sites.google.com/connect.hku.hk/robustml-2021/home</a>

## Research Community Service

---

2025-Current	<b>NeurIPS</b> Area Chair
2025	<b>COLM</b> Area Chair
2024-2025	<b>NSF CISE/CNS</b> Panelist
2024-Current	<b>ICML</b> Area Chair
2023	<b>WiML</b> Workshop Mentor
2023	<b>SatML</b> Program Committee
2022	<b>Principles of Distribution Shift Workshop at ICML</b> Program Committee Website: <a href="https://sites.google.com/view/icml-2022-pods">https://sites.google.com/view/icml-2022-pods</a>
2022	<b>AAAI 2023 Doctoral Consortium</b> Program Committee
2022	<b>15th ACM Workshop on Artificial Intelligence and Security</b> Program Committee Website: <a href="https://aisec.cc/">https://aisec.cc/</a>
2021	<b>14th ACM workshop on Artificial Intelligence and Security</b> Program Committee Website: <a href="https://aisec.cc/">https://aisec.cc/</a>
2020	<b>Towards Trustworthy ML: Rethinking Security and Privacy for ML</b> Program Committee Website: <a href="https://trustworthyiclr20.github.io/">https://trustworthyiclr20.github.io/</a>
2020	<b>AAAI</b> Program Committee
2020-2024	<b>ICML, NeurIPS, ICLR</b> Reviewer
2019	<b>Human-Centric Machine Learning Workshop</b> Program Committee Website: <a href="https://sites.google.com/view/hcml-2019">https://sites.google.com/view/hcml-2019</a>
2019	<b>Security and Privacy of Machine Learning Workshop at ICML</b> Program Committee Website: <a href="https://icml2019workshop.github.io/">https://icml2019workshop.github.io/</a>
2019	<b>Adversarial Machine Learning in Real-World Computer Vision Systems at CVPR</b> Technical Program Committee

- 2019      **1st Workshop on Adversarial Learning Methods for Machine Learning and Data Mining at KDD**  
Technical Program Committee  
Website: <https://sites.google.com/view/advml>
- 2019      **Safe Machine Learning Workshop at ICLR**  
Program Committee  
Website: <https://sites.google.com/view/safeml-iclr2019/>