# **Eric Wong**

Phone: +1 (339) 223-7159 Email: exwong@upenn.edu

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

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

Last updated: 2025/02/17

#### **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**

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

#### Work experience

2019-2020	Bosch Center for Artifical Intelligence (Renningen, Germany and Pittsburgh, PA) 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 conversative 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 toolfor security vulnerabilities in source code.

	1	
<b>△ TA7</b>	ard	S
$\Delta vv$	aru	כו

2025	AI2050 Early Career Award Towards Robust Generative AI with Adaptive Risk Evaluations, Schmidt Sciences
2025	NSF Early Career Award CAREER: Certified Explanations for Trustworthy Artificial Intelligence, NSF
2024	Amazon Research Award (AWS AI) Adversarial Manipulation of Prompting Interfaces, Amazon
2023	Area Chair Award (Interpretability and Analysis of Models for NLP) Faithful Chain-of-Thought Reasoning, IJCNLP-AACL Conference
2020	SCS Dissertation Award – Honorable Mention  Provable, structured, and efficient methods for robustness of deep networks to adversarial examples, Carnegie Mellon University
2020	Siebel Scholar Fellowship Carnegie Mellon University
2017	<b>Best Defense Paper</b> Provable defenses against adversarial examples via the convex outer adversarial polytope, NeurIPS 2017 ML & Security Workshop
2013	Summer Undergraduate Research Fellowship Carnegie Mellon University
Publications	
NeurIPS 2024	<b>AR-Pro: Counterfactual Explanations for Anomaly Repair with Formal Properties</b> Xiayan Ji, Anton Xue, Eric Wong, Oleg Sokolsky, Insup Lee
ICLR 2025	<b>Logicbreaks: A Framework for Understanding Subversion of Rule-based Inference</b> Anton Xue, Avishree Khare, Rajeev Alur, Surbhi Goel, Eric Wong
eBioMedicine	Crowd-sourced machine learning prediction of long COVID using data from the National COVID Cohort Collaborative 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
NAACL- Findings 2025	Avoiding Copyright Infringement via Machine Unlearning Guangyao Dou, Zheyuan Liu, Qing Lyu, Kaize Ding, Eric Wong
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 Towards Compositionality in Concept Learning** Adam Stein, Aaditya Naik, Yinjun Wu, Mayur Naik, Eric Wong **ICML 2024** DISCRET: Synthesizing Faithful Explanations For Treatment Effect Estimation Yinjun Wu, Mayank Keoliya, Kan Chen, Neelay Velingker, Ziyang Li, Emily J Getzen, Qi Long, Mayur Naik, Ravi B Parikh, Eric Wong NeurIPS 2024 JailbreakBench: An Open Robustness Benchmark for Jailbreaking Large Language Models Patrick Chao, Edoardo Debenedetti, Alexander Robey, Maksym Andriushchenko, Francesco Croce, Vikash Sehwag, Edgar Dobriban, Nicolas Flammarion, George J. Pappas, Florian Tramer, Hamed Hassani, Eric Wong ICLR 2024, Tiny Evaluating Groups of Features via Consistency, Contiguity, and Stability Chaehyeon Kim, Weiqiu You, Shreya Havaldar, Eric Wong Papers (Oral) ICLR 2024 SalUn: Empowering Machine Unlearning via Gradient-based Weight Saliency in **Both Image Classification and Generation** Chongyu Fan, Jiancheng Liu, Yihua Zhang, Dennis Wei, Eric Wong, Sijia Liu **CVPR 2024** Initialization Matters for Adversarial Transfer Learning Andong Hua, Jindong Gu, Zhiyu Xue, Nicholas Carlini, Eric Wong, Yao Qin **EMNLP 2023 Comparing Styles across Languages** Shreya Havaldar, Matthew Pressimone, Eric Wong, Lyle Ungar SaTML 2025 Jailbreaking Black Box Large Language Models in Twenty Queries Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J. Pappas, Eric Wong OOPSLA 2024 TorchQL: A Programming Framework for Integrity Constraints in Machine Learn-Aaditya Naik, Adam Stein, Yinjun Wu, Eric Wong, Mayur Naik NeurIPS 2023 Stability Guarantees for Feature Attributions with Multiplicative Smoothing Anton Xue, Rajeev Alur, Eric Wong TopEx: Topic-based Explanations for Model Comparison ICLR 2023, Tiny **Papers** Shreya Havaldar, Adam Stein, Eric Wong, Lyle Ungar Do Machine Learning Models Learn Statistical Rules Inferred from Data? ICML 2023 Aaditya Naik, Yinjun Wu, Mayur Naik, Eric Wong **DLSP 2023** Adversarial Prompting for Black Box Foundation Models Keynote Natalie Maus\*, Patrick Chao\*, Eric Wong, Jacob Gardner IJCNLP-AACL, Faithful Chain-of-Thought Reasoning 2023 Qing Lyu\*, Shreya Havaldar\*, Adam Stein\*, Li Zhang, Delip Rao, Eric Wong, Marianna Apidianaki, Chris Callison-Burch **CVPR 2023** A data-based perspective on transfer learning Saachi Jain\*, Hadi Salman\*, Alaa Khaddaj\*, Eric Wong, Sung Min Park, Aleksander Madry

2024	The FIX Benchmark: Extracting Features Interpretable to eXperts Helen Jin, Shreya Havaldar, Chaehyeon Kim, Anton Xue, Weiqiu You, Helen Qu, Marco Gatti, Daniel A. Hashimoto, Bhuvnesh Jain, Amin Madani, Masao Sako, Lyle
Preprints	
ICML 2015	<b>An SVD and Derivative Kernel Approach to Learning from Geometric Data</b> Eric Wong, J. Zico Kolter
ICML 2017	A Semismooth Newton Method for Fast, Generic Convex Programming Alnur Ali*, Eric Wong*, J. Zico Kolter
ICML 2018	Provable defenses against adversarial examples via the convex outer adversarial polytope Eric Wong, J. Zico Kolter
NeurIPS 2018	Scaling provable adversarial defenses Eric Wong, Frank R. Schmidt, Jan Hendrik Metzen, J. Zico Kolter
ICML 2019	Wasserstein adversarial examples Eric Wong, Frank R. Schmidt, J. Zico Kolter
ICML 2020	Adversarial robustness against the union of multiple perturbation models Pratyush Maini, Eric Wong, J. Zico Kolter
ICLR 2020	Fast is better than free: revisiting adversarial training Eric Wong*, Leslie Rice*, J. Zico Kolter
IEEE IV 2020	Neural network virtual sensors for fuel injection quantities with provable performance specifications Eric Wong, Tim Schneider, Joerg Schmitt, Frank R. Schmidt, J. Zico Kolter
ICML 2020	Overfitting in adversarially robust deep learning Leslie Rice*, Eric Wong*, J. Zico Kolter
ICLR 2021	<b>Learning perturbation sets for robust machine learning</b> Eric Wong, J. Zico Kolter
ICML 2021 (Oral)	Leveraging Sparse Linear Layers for Debuggable Deep Networks Eric Wong*, Shibani Santurkar*, 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
CVPR 2022	Certified patch robustness via smoothed vision transformers Hadi Salman*, Saachi Jain*, Eric Wong*, Aleksander Madry
ICLR 2022	Missingness bias in model debugging Saachi Jain*, Hadi Salman*, Pengchuan Zhang, Vibhav Vineet, Sal Vemprala, Aleksander Madry

2024	Defending Large Language Models against Jailbreak Attacks via Semantic Smoothing Jiabao Ji, Bairu Hou, Alexander Robey, George J. Pappas, Hamed Hassani, Yang Zhang, Eric Wong, Shiyu Chang
2023	Sum-of-Parts Models: Faithful Attributions for Groups of Features Weiqiu You, Helen Qu, Marco Gatti, Bhuvnesh Jain, Eric Wong
2023	SmoothLLM: Defending Large Language Models Against Jailbreaking Attacks Alexander Robey, Eric Wong, Hamed Hassani, George J. Pappas
2023	Rectifying Group Irregularities in Explanations for Distribution Shift Adam Stein, Yinjun Wu, Eric Wong, Mayur Naik
2023	In-context Example Selection with Influences Tai Nguyen, Eric Wong
2022	When does bias transfer in transfer learning Hadi Salman*, Saachi Jain*, Andrew Ilyas*, Logan Engstrom*, Eric Wong, Aleksander Madry
Grants	
2025-06-01 – 2027-05-31	Towards Robust Generative AI with Adaptive Risk Evaluations \$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	Harnessing Artificial Intelligence and Language Modeling for Enhancing Innovation and Evaluating Research Claims (HAILMEIER-C) \$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	TIGER: Trustable Information Generation and Explanation Resilience \$4M, IARPA, Co-PI with Rene Vidal, Chris Callison-Burch, Hamed Hassani, Mark Yatskar, Rama Chellappa, Vishal Patel
2024-05-01 – 2025-04-30	Preventing Complications with Transparent Surgical AI Assistants \$100K, ASSET-IBI, PI with Daniel Hashimoto
2024-07-01 – 2028-06-30	Safe and Explainable AI-enabled Decision Making for Personalized Treatment \$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	SLES: SPECSRL: Specification-guided Perception-enabled Conformal Safe Reinforcement Learning \$1.5M, NSF, Co-PI with Rajeev Alur, Osbert Bastani & Dinesh Jayaraman
2024-05-01 -	Adversarial Manipulation of Prompting Interfaces

2023-10-01 -	SHF: Medium: Scallop: A Neurosymbolic Programming Framework for Combin-
2027-09-30	ing Logic with Deep Learning
	\$1.2M, NSF, Co-PI with Mayur Naik & Rajeev Alur

## Invitations

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

## **Teaching Experience**

2025	Machine Learning - CIS 5200 (UPenn, Instructor) nan
2024	<b>Mathematics of Machine Learning</b> - 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 Practical Data Science - 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 **Eberly Center for Teaching Excellence and Educational Innovation** - Teaching Seminary (CML) Participants

inars (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 Advanced Introduction to Machine Learning - 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 Algorithm Design and Analysis - 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 Pervasive and Mobile Computing Services - 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 **Pervasive and Mobile Computing Services** - 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 Mobile Development for iOS and Android - 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 – **Cassandra Goldberg** (PhD) Spring 2029 Thesis: Anticipated Spring 2029

Fall 2024 – Spring 2029	Davis Brown (PhD) Thesis: Anticipated Spring 2029, co-advised with Hamed Hassani
Fall 2024 – Spring 2029	Vitoria Guardieiro (PhD) Thesis: Anticipated Spring 2029
Spring 2024 – Spring 2026	Adam Stein (PhD) Thesis: Anticipated Spring 2026, co-advised with Mayur Naik
Fall 2023 – Spring 2028	Chaehyeon Kim (PhD) Thesis: Anticipated Spring 2028
Fall 2023 – Spring 2026	Helen Jin (PhD) Thesis: Anticipated Spring 2026
Summer 2023 – Spring 2025	Anton Xue (PhD) Thesis: Anticipated Spring 2025, co-advised with Rajeev Alur
Spring 2023 – Fall 2025	Weiqiu You (PhD) Thesis: Anticipated Fall 2025
Spring 2023 – Spring 2026	Shreya Havaldar (PhD) Thesis: Anticipated Spring 2026, co-advised with Lyle Ungar
Fall 2022 – Spring 2024	Shailesh Sridhar (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 2024 – Spring 2024	<b>Dora Wu</b> (Undergraduate) Thesis: Image Generative Artificial Intelligence: Theory, Applications, and Outlook
Fall 2022 – Spring 2023	Gideon Tesfaye (Undergraduate) Practicum: Using deep learning to compose music

## Penn Service

2024 – 2025	IDEAS Search Committee Faculty search committee member
2024 – 2024	PhD Admit Weekend Organizer, Co-Lead with Andrew Head
2023 – 2024	BSE in AI Curriculum Committee
2023 – ongoing	ML+FM Seminar Organizer of a seminar series for researchers in the areas of formal methods and machine learning. Averages 18 attendees weekly.

2023 – 2023 Adhoc Computing Cluster Committee

Commitee Member

2023 – 2023 PhD Admit Weekend

Organizer, Co-Lead with Andrew Head

2022 – ongoing Locust Cluster

Test driver for the SEAS cluster and working with CETS to iron out scalability of the

system.

#### **DEI Service**

2024 WiML@PennCIS (CIS)

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 WiML Workshop Mentoring (NeurIPS)

Volunteered as a mentor for the round table event at the Women in Machine Learning

workshop at NeurIPS.

2023 **USABE fireside chat** (Penn Engineering)

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 inter-

act in more casual settings with faculty.

2022-2024 Mentorship for Underrepresented Masters/Undergraduates at Penn (CIS)

Direct mentorship in research experiences of masters and undergraduate students

that are underrepresented in CS (1 woman and 1 Ethiopean)

2021 Graduate Application Assistance Program (MIT)

Assisted applicants from under-represented groups with their graduate student ap-

plications to MIT's EECS PhD program.

2021-2022 MIT Undergraduate Research Opportunities Program (MIT)

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 MEnTorEd Opportunities in Research (METEOR) (MIT)

Participated in the METEOR postdoc fellowship selection committee, an effort at CSAIL MIT to increase diversity, equity, and inclusion. Provided confidential tech-

nical feedback on candidates based on their application materials.

2019-2020 CMU AI Mentoring Program (CMU)

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 **Tecknowledge Mentor** (Obama Academy)

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 Mental Health First Aid Certification (CMU)

Underwent training to recognize mental health issues and provide first aid assistance to those in need

### **Workshop Organizing**

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

### **Research Community Service**

2024	NSF SaTC Panelist
2024	ICML Area Chair
2023	WiML Workshop Mentor

Website: https://sites.google.com/connect.hku.hk/robustml-2021/home

2023 SatML **Program Committee** 2022 Principles of Distribution Shift Workshop at ICML **Program Committee** Website: https://sites.google.com/view/icml-2022-pods 2022 **AAAI 2023 Doctoral Consortium Program Committee** 2022 15th ACM Workshop on Artificial Intelligence and Security **Program Committee** Website: https://aisec.cc/ 2021 14th ACM workshop on Artificial Intelligence and Security **Program Committee** Website: https://aisec.cc/ 2020 Towards Trustworthy ML: Rethinking Security and Privacy for ML **Program Committee** Website: https://trustworthyiclr20.github.io/ 2020 **AAAI Program Committee** 2020-2024 ICML, NeurIPS, ICLR Reviewer 2019 **Human-Centric Machine Learning Workshop Program Committee** Website: https://sites.google.com/view/hcml-2019 2019 Security and Privacy of Machine Learning Workshop at ICML **Program Committee** Website: https://icml2019workshop.github.io/ 2019 Adversarial Machine Learning in Real-World Computer Vision Systems at CVPR **Technical Program Commitee** 2019 1st Workshop on Adversarial Learning Methods for Machine Learning and Data Mining at KDD **Technical Program Commitee** Website: https://sites.google.com/view/advml

Website: https://sites.google.com/view/safeml-iclr2019/

Safe Machine Learning Workshop at ICLR

**Program Committee** 

2019