

# 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: 2025/01/17*

## 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

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

## Work experience

|           |  |
|-----------|--|
| 2019-2020 | <b>Bosch Center for Artificial Intelligence (Renningen, Germany and Pittsburgh, PA)</b><br>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><br>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><br><i>Towards Robust Generative AI with Adaptive Risk Evaluations, Schmidt Sciences</i>  |
| 2025 | <b>NSF Early Career Award</b><br><i>CAREER: Certified Explanations for Trustworthy Artificial Intelligence, NSF</i>   |
| 2024 | <b>Amazon Research Award (AWS AI)</b><br><i>Adversarial Manipulation of Prompting Interfaces, Amazon</i>  |
| 2023 | <b>Area Chair Award (Interpretability and Analysis of Models for NLP)</b><br><i>Faithful Chain-of-Thought Reasoning, IJCNLP-AAACL Conference</i>  |
| 2020 | <b>SCS Dissertation Award – Honorable Mention</b><br><i>Provable, structured, and efficient methods for robustness of deep networks to adversarial examples, Carnegie Mellon University</i> |
| 2020 | <b>Siebel Scholar Fellowship</b><br>Carnegie Mellon University  |
| 2017 | <b>Best Defense Paper</b><br><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><br>Carnegie Mellon University   |

## Publications

---

|                          |   |
|--------------------------|---|
| NeurIPS 2024             | <b>AR-Pro: Counterfactual Explanations for Anomaly Repair with Formal Properties</b><br>Xiayan Ji, Anton Xue, Eric Wong, Oleg Sokolsky, Insup Lee   |
| AdvML at<br>NeurIPS 2024 | <b>Logicbreaks: A Framework for Understanding Subversion of Rule-based Inference</b><br>Anton Xue, Avishree Khare, Rajeev Alur, Surbhi Goel, Eric Wong  |
| NeurIPS 2024             | <b>Data-Efficient Learning with Neural Programs</b><br>Alaia Solko-Breslin, Seewon Choi, Ziyang Li, Neelay Velingker, Rajeev Alur, Mayur Naik, Eric Wong  |
| ICML 2024                | <b>Towards Compositionality in Concept Learning</b><br>Adam Stein, Aaditya Naik, Yinjun Wu, Mayur Naik, Eric Wong   |
| ICML 2024                | <b>DISCRET: Synthesizing Faithful Explanations For Treatment Effect Estimation</b><br>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><br>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><br>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><br>Chongyu Fan, Jiancheng Liu, Yihua Zhang, Dennis Wei, Eric Wong, Sijia Liu |
| CVPR 2024                     | <b>Initialization Matters for Adversarial Transfer Learning</b><br>Andong Hua, Jindong Gu, Zhiyu Xue, Nicholas Carlini, Eric Wong, Yao Qin  |
| EMNLP 2023                    | <b>Comparing Styles across Languages</b><br>Shreya Havaladar, Matthew Pressimone, Eric Wong, Lyle Ungar   |
| XAIA at NeurIPS 2022          | <b>Sum-of-Parts Models: Faithful Attributions for Groups of Features</b><br>Weiqiu You, Helen Qu, Marco Gatti, Bhuvnesh Jain, Eric Wong   |
| SaTML 2025                    | <b>Jailbreaking Black Box Large Language Models in Twenty Queries</b><br>Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J. Pappas, Eric Wong                                      |
| OOPSLA 2024                   | <b>TorchQL: A Programming Framework for Integrity Constraints in Machine Learning</b><br>Aaditya Naik, Adam Stein, Yijun Wu, Eric Wong, Mayur Naik  |
| NeurIPS 2023                  | <b>Stability Guarantees for Feature Attributions with Multiplicative Smoothing</b><br>Anton Xue, Rameev Alur, Eric Wong   |
| ICLR 2023, Tiny Papers        | <b>TopEx: Topic-based Explanations for Model Comparison</b><br>Shreya Havaladar, Adam Stein, Eric Wong, Lyle Ungar  |
| ICML 2023                     | <b>Do Machine Learning Models Learn Statistical Rules Inferred from Data?</b><br>Aaditya Naik, Yijun Wu, Mayur Naik, Eric Wong  |
| DLSP 2023 Keynote             | <b>Adversarial Prompting for Black Box Foundation Models</b><br>Natalie Maus*, Patrick Chao*, Eric Wong, Jacob Gardner  |
| IJCNLP-AAACL, 2023            | <b>Faithful Chain-of-Thought Reasoning</b><br>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><br>Saachi Jain*, Hadi Salman*, Alaa Khaddaj*, Eric Wong, Sung Min Park, Aleksander Madry   |
| ICLR 2022                     | <b>Missingness bias in model debugging</b><br>Saachi Jain*, Hadi Salman*, Pengchuan Zhang, Vibhav Vineet, Sal Vemprala, Aleksander Madry  |
| CVPR 2022                     | <b>Certified patch robustness via smoothed vision transformers</b><br>Hadi Salman*, Saachi Jain*, Eric Wong*, Aleksander Madry  |
| OJCS 2022                     | <b>DeepSplit: Scalable verification of deep neural networks via operator splitting</b><br>Shaoru Chen*, Eric Wong*, J. Zico Kolter, Mahyar Fazlyab  |

|                     |   |
|---------------------|---|
| ICML 2021<br>(Oral) | <b>Leveraging Sparse Linear Layers for Debuggable Deep Networks</b><br>Eric Wong*, Shibani Santurkar*, Aleksander Madry   |
| ICLR 2021           | <b>Learning perturbation sets for robust machine learning</b><br>Eric Wong, J. Zico Kolter  |
| ICML 2020           | <b>Overfitting in adversarially robust deep learning</b><br>Leslie Rice*, Eric Wong*, J. Zico Kolter  |
| IEEE IV 2020        | <b>Neural network virtual sensors for fuel injection quantities with provable performance specifications</b><br>Eric Wong, Tim Schneider, Joerg Schmitt, Frank R. Schmidt, J. Zico Kolter |
| ICLR 2020           | <b>Fast is better than free: revisiting adversarial training</b><br>Eric Wong*, Leslie Rice*, J. Zico Kolter  |
| ICML 2020           | <b>Adversarial robustness against the union of multiple perturbation models</b><br>Pratyush Maini, Eric Wong, J. Zico Kolter  |
| ICML 2019           | <b>Wasserstein adversarial examples</b><br>Eric Wong, Frank R. Schmidt, J. Zico Kolter  |
| NeurIPS 2018        | <b>Scaling provable adversarial defenses</b><br>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><br>Eric Wong, J. Zico Kolter  |
| ICML 2017           | <b>A Semismooth Newton Method for Fast, Generic Convex Programming</b><br>Alnur Ali*, Eric Wong*, J. Zico Kolter  |
| ICML 2015           | <b>An SVD and Derivative Kernel Approach to Learning from Geometric Data</b><br>Eric Wong, J. Zico Kolter   |

## Preprints

---

|      |   |
|------|---|
| 2024 | <b>The FIX Benchmark: Extracting Features Interpretable to eXperts</b><br>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   |
| 2024 | <b>Crowd-sourced machine learning prediction of long COVID using data from the National COVID Cohort Collaborative</b><br>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 Velinger, 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 |

|      |   |
|------|---|
| 2024 | <b>Avoiding Copyright Infringement via Machine Unlearning</b><br>Guangyao Dou, Zheyuan Liu, Qing Lyu, Kaize Ding, Eric Wong   |
| 2024 | <b>Defending Large Language Models against Jailbreak Attacks via Semantic Smoothing</b><br>Jiabao Ji, Bairu Hou, Alexander Robey, George J. Pappas, Hamed Hassani, Yang Zhang, Eric Wong, Shiyu Chang |
| 2023 | <b>SmoothLLM: Defending Large Language Models Against Jailbreaking Attacks</b><br>Alexander Robey, Eric Wong, Hamed Hassani, George J. Pappas   |
| 2023 | <b>Rectifying Group Irregularities in Explanations for Distribution Shift</b><br>Adam Stein, Yinjun Wu, Eric Wong, Mayur Naik   |
| 2023 | <b>In-context Example Selection with Influences</b><br>Tai Nguyen, Eric Wong  |
| 2022 | <b>When does bias transfer in transfer learning</b><br>Hadi Salman*, Saachi Jain*, Andrew Ilyas*, Logan Engstrom*, Eric Wong, Aleksander Madry  |

## Grants

---

|                            |   |
|----------------------------|---|
| 2025-06-01 –<br>2027-05-31 | <b>Towards Robust Generative AI with Adaptive Risk Evaluations</b><br>\$450k, AI2050 Early Career Fellowship, Schmidt Sciences, PI  |
| 2025-06-01 –<br>2030-05-31 | <b>CAREER: Certified Explanations for Trustworthy Artificial Intelligence</b><br>\$675k, NSF, PI  |
| 2024-12-01 –<br>NaT        | <b>Harnessing Artificial Intelligence and Language Modeling for Enhancing Innovation and Evaluating Research Claims (HAILMEIER-C)</b><br>nan  |
| 2024-05-01 –<br>2025-04-30 | <b>Preventing Complications with Transparent Surgical AI Assistants</b><br>\$100K, ASSET-IBI, PI with Daniel Hashimoto  |
| 2024-07-01 –<br>2028-06-30 | <b>Safe and Explainable AI-enabled Decision Making for Personalized Treatment</b><br>\$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 –<br>2027-09-30 | <b>SLES: SPECSRL: Specification-guided Perception-enabled Conformal Safe Reinforcement Learning</b><br>\$1.5M, NSF, Co-PI with Rajeev Alur, Osbert Bastani & Dinesh Jayaraman                                     |
| 2024-05-01 –<br>2025-04-30 | <b>Adversarial Manipulation of Prompting Interfaces</b><br>\$70K (+\$50K compute), Amazon Research Award, PI  |
| 2023-10-01 –<br>2027-09-30 | <b>SHF: Medium: Scallop: A Neurosymbolic Programming Framework for Combining Logic with Deep Learning</b><br>\$1.2M, NSF, Co-PI with Mayur Naik & Rajeev Alur   |

## Invitations

---

|      |  |
|------|--|
| 2024 | <b>Elements of Trust in Machine Learning</b><br>Seminar speaker, Cornell AI Seminar  |
| 2024 | <b>Elements of Trust in Machine Learning</b><br>Seminar speaker, University of Maryland  |
| 2023 | <b>Robustness of Adversarial Attacks for LLM</b><br>Distinguished speaker, Responsible Machine Learning Summit, UCSB                             |
| 2023 | <b>Adversarial Prompting: Return of the Adversarial Example</b><br>Keynote speaker, IEEE S&P 2023, 6th Deep Learning Security & Privacy Workshop |
| 2023 | <b>From Prompt Engineering to Prompt Science</b><br>Seminar speaker, Wayne State University  |
| 2022 | <b>Robustness for the Real World</b><br>Invited talk, 6th Annual Conference on Information Sciences and Systems (CISS)                           |
| 2022 | <b>Debuggable Deep Networks</b><br>Invited talk, TrustML Young Scientist Seminar   |
| 2021 | <b>Panel Discussion</b><br>Panelist, ATVA 2021 Workshop on Security and Reliability of Machine Learning  |

## Teaching Experience

---

|      |   |
|------|---|
| 2024 | <b>Mathematics of Machine Learning - CIS 3333 (UPenn, Instructor)</b><br>Second iteration of the new Mathematics of Machine Learning course under an official course number as a SEAS mathematics elective.   |
| 2024 | <b>Machine Learning - CIS 5200 (UPenn, Instructor)</b><br>2nd round teaching CIS 5200, further consolidation and unification of the course material with the other AI faculty. 150 students.  |
| 2023 | <b>Mathematics of Machine Learning - CIS 3990 (UPenn, Instructor)</b><br>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 | <b>Machine Learning - CIS 5200 (UPenn, Instructor)</b><br>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 | <b>Debugging Data &amp; Models - CIS 7000-005 (UPenn, Instructor)</b><br>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 - 15-388/688 (CMU, TA)</b><br>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 - Teaching Seminars (CMU, Participant)</b><br>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 - 10-715 (CMU, TA)</b><br>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 - 15-451 (CMU, TA)</b><br>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 - 08-766/781 (CMU, TA)</b><br>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 - 08-766/781 (CMU, TA)</b><br>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 - 08-723 (CMU, TA)</b><br>Held office hours and graded assignments/exams for the software engineering department's Mobile Development course (08-723) at CMU.  |

## Graduate Theses Supervised

---

|                              |  |
|------------------------------|--|
| Fall 2024 –<br>Spring 2029   | <b>Cassandra Goldberg (PhD)</b><br>Thesis: Anticipated Spring 2029                         |
| Fall 2024 –<br>Spring 2029   | <b>Davis Brown (PhD)</b><br>Thesis: Anticipated Spring 2029, co-advised with Hamed Hassani |
| Fall 2024 –<br>Spring 2029   | <b>Vitoria Guardieiro (PhD)</b><br>Thesis: Anticipated Spring 2029                         |
| Spring 2024 –<br>Spring 2026 | <b>Adam Stein (PhD)</b><br>Thesis: Anticipated Spring 2026, co-advised with Mayur Naik     |
| Fall 2023 –<br>Spring 2028   | <b>Chaehyeon Kim (PhD)</b><br>Thesis: Anticipated Spring 2028                              |

|                              |   |
|------------------------------|---|
| Fall 2023 –<br>Spring 2026   | <b>Helen Jin</b> (PhD)<br>Thesis: Anticipated Spring 2026   |
| Summer 2023 –<br>Spring 2025 | <b>Anton Xue</b> (PhD)<br>Thesis: Anticipated Spring 2025, co-advised with Rajeev Alur                          |
| Spring 2023 –<br>Fall 2025   | <b>Wei qiu You</b> (PhD)<br>Thesis: Anticipated Fall 2025   |
| Spring 2023 –<br>Spring 2026 | <b>Shreya Havaladar</b> (PhD)<br>Thesis: Anticipated Spring 2026, co-advised with Lyle Ungar                    |
| Fall 2022 –<br>Spring 2024   | <b>Shailesh Sridhar</b> (Masters)<br>Thesis: Controlling for Missingness Bias in Feature Attribution Evaluation |
| Fall 2022 –<br>Spring 2024   | <b>Tai Nguyen</b> (Masters)<br>Thesis: Attribute in-context learning examples with influences                   |

## Undergraduate Projects Supervised

---

|                              |   |
|------------------------------|---|
| Spring 2024 –<br>Spring 2024 | <b>Dora Wu</b> (Undergraduate)<br>Thesis: Image Generative Artificial Intelligence: Theory, Applications, and Outlook |
| Fall 2022 –<br>Spring 2023   | <b>Gideon Tesfaye</b> (Undergraduate)<br>Practicum: Using deep learning to compose music                              |

## Penn Service

---

|                |  |
|----------------|--|
| 2024 – 2024    | <b>PhD Admit Weekend</b><br>Organizer, Co-Lead with Andrew Head  |
| 2023 – 2024    | <b>BSE in AI</b><br>Curriculum Committee   |
| 2023 – ongoing | <b>ML+FM Seminar</b><br>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><br>Committee Member   |
| 2023 – 2023    | <b>PhD Admit Weekend</b><br>Organizer, Co-Lead with Andrew Head  |
| 2022 – ongoing | <b>Locust Cluster</b><br>Test driver for the SEAS cluster and working with CETS to iron out scalability of the system.                                   |

## DEI Service

---



|           |  |
|-----------|--|
| 2024      | <b>WiML@PennCIS (CIS)</b><br>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><br>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><br>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. |
| 2022-2024 | <b>Mentorship for Underrepresented Masters/Undergraduates at Penn (CIS)</b><br>Direct mentorship in research experiences of masters and undergraduate students that are underrepresented in CS (1 woman and 1 Ethiopian)   |
| 2021      | <b>Graduate Application Assistance Program (MIT)</b><br>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><br>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><br>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><br>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><br>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><br>Underwent training to recognize mental health issues and provide first aid assistance to those in need   |

## Workshop Organizing

---

|      |   |
|------|---|
| 2023 | <b>2nd New Frontiers in Adversarial Machine Learning</b><br>Organizer for the 2nd ICML 2023 workshop on new directions in adversarial machine learning<br>Website: <a href="https://advml-frontier.github.io/">https://advml-frontier.github.io/</a>  |
| 2022 | <b>Workshop on Adversarial Machine Learning and Beyond</b><br>Organizer for an AAAI 2022 workshop broadly themed around adversarial machine learning<br>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><br>Organizer for an ICML 2022 workshop on new directions in adversarial machine learning<br>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><br>Organizer for an ICML 2021 workshop themed around the dangers and benefits of adversarial machine learning<br>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><br>Main organizer for an ICLR 2021 workshop on real world robustness.<br>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

---

|      |   |
|------|---|
| 2024 | <b>NSF SaTC</b><br>Panelist   |
| 2024 | <b>ICML</b><br>Area Chair   |
| 2023 | <b>WiML</b><br>Workshop Mentor  |
| 2023 | <b>SatML</b><br>Program Committee   |
| 2022 | <b>Principles of Distribution Shift Workshop at ICML</b><br>Program Committee<br>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><br>Program Committee   |
| 2022 | <b>15th ACM Workshop on Artificial Intelligence and Security</b><br>Program Committee<br>Website: <a href="https://aisec.cc/">https://aisec.cc/</a>   |
| 2021 | <b>14th ACM workshop on Artificial Intelligence and Security</b><br>Program Committee<br>Website: <a href="https://aisec.cc/">https://aisec.cc/</a>   |

|           |  |
|-----------|--|
| 2020      | <b>Towards Trustworthy ML: Rethinking Security and Privacy for ML</b><br>Program Committee<br>Website: <a href="https://trustworthyiclr20.github.io/">https://trustworthyiclr20.github.io/</a>                                   |
| 2020      | <b>AAAI</b><br>Program Committee   |
| 2020-2024 | <b>ICML, NeurIPS, ICLR</b><br>Reviewer   |
| 2019      | <b>Human-Centric Machine Learning Workshop</b><br>Program Committee<br>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><br>Program Committee<br>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><br>Technical Program Committee   |
| 2019      | <b>1st Workshop on Adversarial Learning Methods for Machine Learning and Data Mining at KDD</b><br>Technical Program Committee<br>Website: <a href="https://sites.google.com/view/advml">https://sites.google.com/view/advml</a> |
| 2019      | <b>Safe Machine Learning Workshop at ICLR</b><br>Program Committee<br>Website: <a href="https://sites.google.com/view/safeml-iclr2019/">https://sites.google.com/view/safeml-iclr2019/</a>                                       |