

Umang Bhatt

CONTACT umang.s.bhatt@gmail.com
INFORMATION Citizenship: USA

Website
Google Scholar

INTERESTS Trustworthy Machine Learning, Responsible AI, Human-Machine Collaboration, AI Policy

ACADEMIC Assistant Professor & Faculty Fellow, **New York University** Oct 2023 – Present
POSITIONS Center for Data Science

Senior Research Associate, **The Alan Turing Institute** Jan 2023 – Present
Project ELSA: European Lighthouse on Secure and Safe AI

Research Fellow, **Harvard University** Jun 2022 – Mar 2023
Center for Research on Computation and Society

EDUCATION **University of Cambridge**, Cambridge, UK
Ph.D. in Engineering Sept 2019 – Nov 2023

Passed with No Corrections

Thesis: [Trustworthy Machine Learning: From Algorithmic Transparency to Decision Support](#)

Committee: [Adrian Weller](#) MBE (Advisor), Neil Lawrence, Eric Horvitz

Affiliations: Machine Learning Group, Computation and Biological Learning Lab

Carnegie Mellon University, Pittsburgh, PA

M.S. in Electrical and Computer Engineering Aug 2017 – May 2019

B.S. in Electrical and Computer Engineering Aug 2015 – May 2019

SELECT **Center for Democracy & Technology Fellowship** 2024 – 2026

FELLOWSHIPS **J.P. Morgan AI PhD Fellowship** (*Declined*) 2022 – 2023

AND AWARDS **The Alan Turing Institute Enrichment Studentship** 2021 – 2022

Mozilla Fellowship 2020 – 2021

Partnership on AI Research Fellowship 2019 – 2020

Leverhulme Center for the Future of Intelligence PhD Scholarship 2019 – 2023

Fully funded by Google DeepMind and the Leverhulme Trust

Best Presentation Award, AAAI Spring Symposium on Interpretable AI for Well-Being 2019

Lovett Family Endowed Scholarship, The Andrew Carnegie Society 2019

Undergraduate Research Presentation Award, CMU 2017

NSF I-Corps Site Award for research commercialization 2017

H. F. McCullough Memorial Scholarship, CMU 2016

JOURNAL AND [1] **Learning Personalized Decision Support Policies**
CONFERENCE *AAAI Conference on Artificial Intelligence (AAAI) 2025 (Oral Presentation)*
PUBLICATIONS **Umang Bhatt***, Valerie Chen*, Katherine Collins, Parameswaran Kamalaruban, Emma Kallina, Adrian Weller, Ameet Talwalkar

[2] **Building Machines that Learn and Think *with* People**
Nature Human Behavior, 2024.
Katherine Collins*, Ilia Sucholutsky*, **Umang Bhatt***, Kartik Chandra*, Lionel Wong*, Mina Lee, Cedegao E. Zhang, Tan Zhi-Xuan, Mark Ho, Vikash Mansinghka, Adrian Weller, Joshua Tenenbaum, Thomas Griffiths

[3] **When Should Algorithms Resign? A Proposal for AI Governance**
IEEE Computer, 2024
Umang Bhatt*, Holli Sargeant*

[4] **Evaluating Language Models for Mathematics through Interactions**
Proceedings of the National Academy of Sciences (PNAS), 2024.
Katherine Collins, Albert Jiang, Simon Frieder, Lionel Wong, Miri Zilka, **Umang Bhatt**, Thomas Lukasiewicz, Yuhuai Wu, Joshua Tenenbaum, William Hart, Timothy Gowers, Wenda Li, Adrian Weller, Mateja Jamnik

- [5] **Large Language Models Must Be Taught to Know What They Don't Know**
Neural Information Processing Systems (NeurIPS), 2024.
 Sanyam Kapoor, Nate Gruver, Manley Roberts, Katherine Collins, Arka Pal, **Umang Bhatt**, Adrian Weller, Samuel Dooley, Micah Goldblum, Andrew Gordon Wilson
- [6] **Comparing Abstraction in Humans and Large Language Models Using Multimodal Serial Reproduction**
46th Annual Conference of the Cognitive Science Society (CogSci) 2024 (Oral Presentation)
 Sreejan Kumar, Raja Marjeh, Byron Zhang, Declan Campbell, Michael Y. Hu, **Umang Bhatt**, Brenden Lake, Thomas L. Griffiths
- [7] **Perspectives on Incorporating Expert Feedback into Model Updates**
Patterns, 2023.
 Valerie Chen*, **Umang Bhatt***, Hoda Heidari, Adrian Weller, Ameet Talwalkar
- [8] **Algorithmic Loafing and Mitigation Strategies in Human-AI Teams**
Computers in Human Behavior: Artificial Humans, 2023.
 Isa Inuwa-Dutse, Alice Toniolo, Adrian Weller, **Umang Bhatt**
- [9] **Selective Concept Models: Permitting Stakeholder Customisation at Test-Time**
AAAI Conference on Human Computation and Crowdsourcing (HCOMP) 2023
 Matthew Barker, Katherine Collins, Krishnamurthy Dvijotham, Adrian Weller, **Umang Bhatt**
- [10] **FeedbackLogs: Recording and Incorporating Stakeholder Feedback into Machine Learning Pipelines**
ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization (EAAMO) 2023
 Matthew Barker, Emma Kallina, Dhananjay Ashok, Katherine Collins, Ashley Casovan, Adrian Weller, Ameet Talwalkar, Valerie Chen, **Umang Bhatt**
- [11] **Human-in-the-Loop Mixup**
Uncertainty in Artificial Intelligence (UAI) 2023 (Oral Presentation)
 Katherine Collins, **Umang Bhatt**, Weiyang Liu, Vihari Piratla, Ilia Sucholutsky, Bradley Love, Adrian Weller
- [12] **On the Informativeness of Supervision Signals**
Uncertainty in Artificial Intelligence (UAI) 2023
 Ilia Sucholutsky, Ruairidh McLennan Battleday, Katherine Collins, Raja Marjeh, Joshua Peterson, Pulkit Singh, **Umang Bhatt**, Nori Jacoby, Adrian Weller, Thomas Griffiths
- [13] **Human Uncertainty in Concept-Based AI Systems**
AAAI/ACM Conference on Artificial Intelligence, Ethics, and Society (AIES) 2023
 Katherine Collins, Matthew Barker, Mateo Espinosa Zarlenga, Naveen Raman, **Umang Bhatt**, Mateja Jamnik, Ilia Sucholutsky, Adrian Weller, Krishnamurthy Dvijotham
- [14] **Iterative Teaching by Data Hallucination**
International Conference on Artificial Intelligence and Statistics (AISTATS) 2023
 Zeju Qiu, Weiyang Liu, Tim Xiao, Zhen Liu, Yucen Luo, **Umang Bhatt**, Adrian Weller, Bernhard Schölkopf
- [15] **Approximating Full Conformal Prediction at Scale via Influence Functions**
AAAI Conference on Artificial Intelligence (AAAI) 2023
 Javier Abad Martinez, **Umang Bhatt**, Adrian Weller, Giovanni Cherubin
- [16] **Towards Robust Metrics for Concept Representation Evaluation**
AAAI Conference on Artificial Intelligence (AAAI) 2023
 Mateo Zarlenga, Pietro Barbiero, Zohreh Shams, D. Kazhdan, **Umang Bhatt**, Adrian Weller, Mateja Jamnik
- [17] **How Transparency Modulates Trust in Artificial Intelligence**
Patterns, 2022.
 John Zerilli, **Umang Bhatt**, Adrian Weller
- [18] **Uncertainty Quantification with Pre-trained Language Models: An Empirical Analysis**
Conference on Empirical Methods in Natural Language Processing (EMNLP) 2022 (Findings)
 Yuxin Xiao, Paul Pu Liang, **Umang Bhatt**, Willie Neiswanger, Ruslan Salakhutdinov, L.P. Morency

- [19] **Eliciting and Learning with Soft Labels from Every Annotator**
AAAI Conference on Human Computation and Crowdsourcing (HCOMP) 2022
 Katherine Collins*, **Umang Bhatt***, Adrian Weller
- [20] **On the Utility of Prediction Sets in Human-AI Teams**
International Joint Conference on Artificial Intelligence (IJCAI) 2022 (Oral Presentation)
 Varun Babbar, **Umang Bhatt**, Adrian Weller
- [21] **Diverse and Amortised Counterfactual Explanations for Uncertainty Estimates**
AAAI Conference on Artificial Intelligence (AAAI) 2022
 Dan Ley, **Umang Bhatt**, Adrian Weller
- [22] **On the Fairness of Causal Algorithmic Recourse**
AAAI Conference on Artificial Intelligence (AAAI) 2022 (Oral Presentation)
 Julius von Kügelgen, Amir Karimi, **Umang Bhatt**, Isabel Valera, Adrian Weller, Bernhard Schölkopf
- [23] **Uncertainty as a Form of Transparency: Measuring and Communicating Uncertainty**
AAAI/ACM Conference on Artificial Intelligence, Ethics, and Society (AIES) 2021
Umang Bhatt, Javier Antorán, Yunfeng Zhang, Q. Vera Liao, Prasanna Sattigeri, Riccardo Fogliato, Gabrielle Melançon, Ranganath Krishnan, Jason Stanley, Omesh Tickoo, Lama Nachman, Rumi Chunara, Madhulika Srikumar, Adrian Weller, Alice Xiang
- [24] **Getting a CLUE: A Method for Explaining Uncertainty Estimates**
International Conference on Learning Representations (ICLR) 2021 (Oral Presentation)
 Javier Antorán, **Umang Bhatt**, Tameem Adel, Adrian Weller, José Miguel Hernández-Lobato
- [25] **FIMAP: Feature Importance by Minimal Adversarial Perturbation**
AAAI Conference on Artificial Intelligence (AAAI) 2021
 Matt Chapman-Rounds, **Umang Bhatt**, Erik Pazos, Marc-Andre Schulz, Kostas Georgatzis
- [26] **Evaluating and Aggregating Feature-based Explanations**
International Joint Conference on Artificial Intelligence (IJCAI) 2020
Umang Bhatt, Adrian Weller, José M.F. Moura
- [27] **Explainable Machine Learning in Deployment**
ACM Conference on Fairness, Accountability, and Transparency (FAccT) 2020
Umang Bhatt, Alice Xiang, Shubham Sharma, Adrian Weller, Ankur Taly, Yunhan Jia, Joydeep Ghosh, Ruchir Puri, José M.F. Moura, Peter Eckersley
- [28] **You Shouldn't Trust Me: Learning Models Which Conceal Unfairness From Multiple Explanation Methods**
European Conference on Artificial Intelligence (ECAI) 2020
 Boty Dimanov, **Umang Bhatt**, Mateja Jamnik, Adrian Weller
- [29] **On Network Science and Mutual Information for Explaining Deep Neural Networks**
IEEE Conference on Acoustics, Speech and Signal Processing (ICASSP) 2020
 Brian Davis*, **Umang Bhatt***, Kartikeya Bhardwaj*, Radu Marculescu, José M.F. Moura
- BOOK CHAPTERS [30] **Trust in Artificial Intelligence: Clinicians are Essential**
Healthcare Information Technology for Cardiovascular Medicine, 2021
Umang Bhatt, Zohreh Shams
- TEACHING EXPERIENCE **Course Instructor, New York University**
- Responsible Data Science (DS-UA 202) Spring 2025
 Co-Taught with Jonathan Colner
 - Responsible Data Science (DS-UA 202) Spring 2024
 Topics: Fairness, Accountability, Transparency, and Privacy
- Teaching Assistant (Supervisor/Grader), University of Cambridge**
- *Inference* (3F8) for J.M. Hernández-Lobato and David Krueger Lent 2023

	<ul style="list-style-type: none"> • <i>Inference</i> (3F8) for Richard Turner and David Krueger • <i>Probabilistic ML</i> (4F13) for Zoubin Ghahramani and J.M. Hernández-Lobato 	Lent 2022 Michaelmas 2020
	Teaching Assistant , <i>Carnegie Mellon University</i>	
	<ul style="list-style-type: none"> • <i>Machine Learning for Engineers - Masters</i> (18-661) for Gauri Joshi • <i>Machine Learning - PhD</i> (10-701) for Ziv Bar-Joseph and Pradeep Ravikumar • <i>Practical Data Science</i> (15-688) for Zico Kolter 	Spring 2019 Fall 2018 Spring 2018
SUPERVISION	Advisor , <i>Alan Turing Institute</i>	
	<ul style="list-style-type: none"> • Mackenzie Jorgensen, Postdoctoral Research Associate • Kendall Brogle, Research Assistant 	Nov 2024 – Present Sept 2024 – Present
	Advisor , <i>New York University</i>	
	<ul style="list-style-type: none"> • Lujain Ibrahim, Visiting PhD Student • Mihir Upadhyay, MS in Data Science • Elaf Almahmoud, MS in Computer Science • Ghina Al Shdaifat, MS in Data Science • Arina Shah, BA in Philosophy • Suha Memon, BS in Computer Science (<i>Next: MS at Penn</i>) 	Nov 2024 – Present Aug 2024 – Present Jun 2024 – Present Apr 2024 – Present Feb 2024 – Present Jan 2024 – June 2024
	Thesis Co-Supervisor/Mentor , <i>University of Cambridge</i>	
	<ul style="list-style-type: none"> • Matthew Barker, MEng in Information Engineering (<i>Next: Trustwise</i>) • Vivek Palaniappan, MEng in Information Engineering (<i>Next: Citadel</i>) • Katherine Collins, MPhil in Machine Learning (<i>Next: PhD at Cambridge</i>) Departmental Distinction Awardee • Varun Babbar, MEng in Information Engineering (<i>Next: PhD at Duke</i>) 2022 Engineering Division F Thesis Prize Winner • Javier Abad Martinez, Research Assistant (<i>Next: PhD at ETH Zurich</i>) • Dan Ley, MEng in Information Engineering (<i>Next: PhD at Harvard</i>) 2021 Engineering Division F Thesis Prize Winner 	May 2022 – Jun 2023 May 2022 – Jun 2023 Nov 2021 – Sept 2022 May 2021 – Jun 2022 Nov 2021 – Feb 2022 May 2020 – Aug 2021
SELECT INVITED TALKS	<ul style="list-style-type: none"> • TU Dortmund, <i>RC Trustworthy Data Science Seminar</i> • Google, <i>Responsible AI Talk Series</i> • Finnish Center for Artificial Intelligence, <i>ELLIS HCML Workshop</i> • KTH Royal Institute of Technology, <i>RPL Summer School</i> • US Department of Defense, <i>CDAO Seminar</i> • KAUST, <i>Rising Stars in AI Symposium</i> • Instituto dos Advogados de São Paulo - IASP, <i>Seminar</i> • House of Lords Select Committee on AI in Weapon Systems • Indian Institute of Technology Bombay, <i>C-MInDS Seminar</i> • University of Chicago, <i>Chicago Human+AI Lab</i> • Massachusetts Institute of Technology, <i>AI Ethics and Policy Group</i> • Birkbeck, University of London, <i>Workshop on Human Behavioral Aspects of XAI</i> • DeepMind, <i>Ethics Research Team</i> • University of Bristol, <i>REPHRAIN Masterclass</i> • University of Manchester, <i>Advances in Data Science and AI Conference</i> • Leverhulme Centre for the Future of Intelligence, <i>Seminar Series</i> • Ada Lovelace Institute, <i>Brown Bag Lunch Seminar</i> • Von Hügel Institute, <i>Research Salon on Artificial Intelligence</i> • Pennsylvania State University, <i>Young Achievers Symposium</i> • UK Ministry of Defense, <i>AI Safety Workshop</i> • Vanguard Health, <i>MedTech Insight Podcast</i> • Technical University of Denmark, <i>Trustworthiness and Interpretability in ML Seminar</i> • Harvard SEAS, <i>Guest Lecture for COMPSCI 282BR: Explainability in ML</i> • Imperial College, <i>Explainable AI Seminar</i> • Robust and Responsible AI Developers (<i>Keynote</i>) • ICML Workshop on Extending Explainable AI (<i>Keynote</i>) 	Jul 2024 Jul 2024 Jun 2024 Jun 2024 Feb 2024 Feb 2024 Dec 2023 Jun 2023 Jan 2023 Nov 2022 Oct 2022 Sept 2022 July 2022 June 2022 June 2022 June 2022 June 2022 June 2022 Apr 2022 Mar 2022 June 2021 Apr 2021 Apr 2021 Feb 2021 July 2020 July 2020

	<ul style="list-style-type: none"> • Mozilla All-Hands Meeting (<i>Keynote</i>) • QuantumBlack (McKinsey), <i>AI Seminar</i> • Fiddler Labs, <i>Explainable AI Seminar</i> 	June 2020 May 2020 May 2019
PROFESSIONAL SERVICE	Finance Chair , ACM Conference on Fairness, Accountability, and Transparency (FAccT) Associate Chair , International Conference on Machine Learning (ICML) Co-Organizer <ul style="list-style-type: none"> • DALI Meeting: Communication of Uncertainty Workshop • Deep Learning Indaba: Responsible AI Practical Workshop • ELSA Workshop on Generative AI and Creative Arts • Deep Learning Indaba: Responsible AI Practical Workshop • ICML Workshop on Human-Machine Collaboration and Teaming • ELLIS Workshop on Human-Centric Machine Learning • ICML Workshop on Human Interpretability in Machine Learning • IBM + Partnership on AI Workshop on Explainable AI 	2024 2024 Apr 2025 Sept 2024 Mar 2024 Sept 2023 July 2022 May 2021 July 2020 Feb 2020
	Program Committee (Reviewer) <ul style="list-style-type: none"> • 2025: AISTATS, FAccT (Area Chair), ICML, UAI • 2024: AIES, FAccT, NeurIPS • 2023: FAccT, ICLR, ICML, IJCAI, NeurIPS, UIST, WebConf • 2022: AAAI, AISTATS, FAccT, ICLR, ICML, NeurIPS, UAI • 2021: AAAI, AISTATS, FAccT, ICAIF, ICLR, ICML, KDD, NeurIPS, UAI • 2020: ICAIF, NeurIPS 	
	Reviewer: ACM Computing Surveys, ACM Transactions on Computer-Human Interaction (TOCHI), ACM Transactions on Interactive Intelligent Systems (TiiS), Artificial Intelligence Journal (AIJ), Harvard Data Science Review (HDSR), IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI), Journal of Artificial Intelligence Research (JAIR), Journal of Machine Learning Research (JMLR), Transactions on Machine Learning Research (TMLR), Management Science	
PROFESSIONAL EXPERIENCE	Scientific Advisor, Multiple Startups <ul style="list-style-type: none"> • Starlab (Barcelona, Spain), Reflective Works (New York, NY) Advisor, Slalom , New York, NY <ul style="list-style-type: none"> • Responsible AI Client Advisory Board Advisory Panel, Bourne-Epsom Protocol , London, UK <ul style="list-style-type: none"> • Expert guidance on the technical aspects of AI in education Advisor, OECD , Paris, FR <ul style="list-style-type: none"> • Expert Group on AI Risk & Accountability Advisory Chief Scientist, Trustwise , London, UK <ul style="list-style-type: none"> • Building safety guardrails for large language models Subject Matter Expert, Accenture , New York, NY <ul style="list-style-type: none"> • Advising on Responsible AI strategy Advisor, Responsible AI Institute , Austin, TX <ul style="list-style-type: none"> • Building a scalable certification program for AI systems Advisor, Credo AI , Palo Alto, CA <ul style="list-style-type: none"> • Scoped an AI governance and auditing platform; backed by AI Fund Research Assistant, Carnegie Mellon University , Pittsburgh, PA <ul style="list-style-type: none"> • Advisors: José M.F. Moura (ECE), Pradeep Ravikumar (MLD), and Zico Kolter (CSD) Student Fellow, .406 Ventures , Boston, MA <ul style="list-style-type: none"> • Sourced startups and performed first-round due diligence on ventures Intern, Microsoft , Redmond, WA <ul style="list-style-type: none"> • Project: explainable conversational agents for technical hardware documentation Co-Founder, Perceptense , Pittsburgh, PA <ul style="list-style-type: none"> • Built products to harvest vehicular telematics data; pipeline acquired by Honda Motors 	Nov 2024 – Present Nov 2023 – Present Jul 2023 – Present Apr 2023 – Present Jan 2023 – Present Jan 2023 – Apr 2024 Oct 2021 – Jan 2023 Jun 2020 – Jun 2021 Jan 2017 – Sept 2019 July 2017 – Jun 2019 May 2018 – Aug 2018 Jan 2017 – May 2018