

Umang Bhatt

CONTACT umangbhatt@nyu.edu
 INFORMATION Citizenship: USA

[Website](#)
[Google Scholar](#)

ACADEMIC Assistant Professor and Faculty Fellow, **New York University** Oct 2023 – Present
 POSITIONS Center for Data Science
 Research Associate, **The Alan Turing Institute** Jan 2023 – Present
 Project ELSA: European Lighthouse on Secure and Safe AI
 Research Fellow, **Harvard University** June 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] **When Should Algorithms Resign? A Proposal for AI Governance**
 CONFERENCE *IEEE Computer* (Forthcoming), 2024
 PUBLICATIONS **Umang Bhatt***, Holli Sargeant*
 [2] **Building Machines that Learn and Think *with* People**
Nature Human Behavior (Forthcoming), 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 Tenen-
 baum, Thomas Griffiths
 [3] **Evaluating Language Models for Mathematics through Interactions**
Proceedings of the National Academy of Sciences, 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
 [4] **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

- [5] **Perspectives on Incorporating Expert Feedback into Model Updates**
Patterns, 2023.
Valerie Chen*, **Umang Bhatt***, Hoda Heidari, Adrian Weller, Ameet Talwalkar
- [6] **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**
- [7] **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**
- [8] **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**
- [9] **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
- [10] **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
- [11] **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
- [12] **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
- [13] **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
- [14] **How Transparency Modulates Trust in Artificial Intelligence**
Patterns, 2022.
John Zerilli, **Umang Bhatt**, Adrian Weller
- [15] **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
- [16] **Eliciting and Learning with Soft Labels from Every Annotator**
AAAI Conference on Human Computation and Crowdsourcing (HCOMP) 2022
Katherine Collins*, **Umang Bhatt***, Adrian Weller
- [17] **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
- [18] **Diverse and Amortised Counterfactual Explanations for Uncertainty Estimates**
AAAI Conference on Artificial Intelligence (AAAI) 2022
Dan Ley, **Umang Bhatt**, Adrian Weller
- [19] **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

	<p>[20] Uncertainty as a Form of Transparency: Measuring and Communicating Uncertainty <i>AAAI/ACM Conference on Artificial Intelligence, Ethics, and Society (AIES) 2021</i> 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</p> <p>[21] Getting a CLUE: A Method for Explaining Uncertainty Estimates <i>International Conference on Learning Representations (ICLR) 2021 (Oral Presentation)</i> Javier Antorán, Umang Bhatt, Tameem Adel, Adrian Weller, José Miguel Hernández-Lobato</p> <p>[22] FIMAP: Feature Importance by Minimal Adversarial Perturbation <i>AAAI Conference on Artificial Intelligence (AAAI) 2021</i> Matt Chapman-Rounds, Umang Bhatt, Erik Pazos, Marc-Andre Schulz, Kostas Georgatzis</p> <p>[23] Evaluating and Aggregating Feature-based Explanations <i>International Joint Conference on Artificial Intelligence (IJCAI) 2020</i> Umang Bhatt, Adrian Weller, José M.F. Moura</p> <p>[24] Explainable Machine Learning in Deployment <i>ACM Conference on Fairness, Accountability, and Transparency (FAccT) 2020</i> Umang Bhatt, Alice Xiang, Shubham Sharma, Adrian Weller, Ankur Taly, Yunhan Jia, Joydeep Ghosh, Ruchir Puri, José M.F. Moura, Peter Eckersley</p> <p>[25] You Shouldn't Trust Me: Learning Models Which Conceal Unfairness From Multiple Explanation Methods <i>European Conference on Artificial Intelligence (ECAI) 2020</i> Botty Dimanov, Umang Bhatt, Mateja Jamnik, Adrian Weller</p> <p>[26] On Network Science and Mutual Information for Explaining Deep Neural Networks <i>IEEE Conference on Acoustics, Speech and Signal Processing (ICASSP) 2020</i> Brian Davis*, Umang Bhatt*, Kartikeya Bhardwaj*, Radu Marculescu, José M.F. Moura</p>
BOOK CHAPTERS	<p>[27] Trust in Artificial Intelligence: Clinicians are Essential <i>Healthcare Information Technology for Cardiovascular Medicine</i>, 2021 Umang Bhatt, Zohreh Shams</p>
TEACHING EXPERIENCE	<p>Course Instructor, New York University</p> <ul style="list-style-type: none"> Responsible Data Science (DS-UA 202) Spring 2024 Topics: Fairness, Accountability, Transparency, and Privacy; 158 undergraduates enrolled <p>Teaching Assistant (Supervisor/Grader), University of Cambridge</p> <ul style="list-style-type: none"> <i>Inference</i> (3F8) for J.M. Hernández-Lobato and David Krueger Lent 2023 <i>Inference</i> (3F8) for Richard Turner and David Krueger Lent 2022 <i>Probabilistic ML</i> (4F13) for Zoubin Ghahramani and J.M. Hernández-Lobato Michaelmas 2020 <p>Teaching Assistant, Carnegie Mellon University</p> <ul style="list-style-type: none"> <i>Machine Learning for Engineers - Masters</i> (18-661) for Gauri Joshi Spring 2019 <i>Machine Learning - PhD</i> (10-701) for Ziv Bar-Joseph and Pradeep Ravikumar Fall 2018 <i>Practical Data Science</i> (15-688) for Zico Kolter Spring 2018 <i>Principles of Imperative Computation</i> (15-122) for Illiano Cervesato Fall 2017 <i>Principles of Computing</i> (15-110) for Margret Reid-Miller Spring 2017
STUDENT SUPERVISION	<p>Advisor, New York University</p> <ul style="list-style-type: none"> Elaf Almahmoud, MS in Computer Science June 2024 – Present Ghina Al Shdaifat, MS in Data Science Apr 2024 – Present Anna Mitrofanova, BA in Business and Data Science Feb 2024 – Present Arina Shah, BA in Philosophy and Data Science Feb 2024 – Present Dean's Undergraduate Research Fund Grantee Kendall Brogle, BS in Data Science (<i>Next: Responsible AI Institute</i>) Feb 2024 – Sept 2024 Suha Memon, BS in Computer Science (<i>Next: MS at Penn</i>) Jan 2024 – June 2024

Thesis Co-Supervisor/Mentor, University of Cambridge		
	• Matthew Barker, MEng in Information Engineering (<i>Next: Trustwise</i>)	May 2022 – June 2023
	• Vivek Palaniappan, MEng in Information Engineering (<i>Next: Citadel</i>)	May 2022 – June 2023
	• Katherine Collins, MPhil in Machine Learning (<i>Next: PhD at Cambridge</i>) Departmental Distinction Awardee	Nov 2021 – Sept 2022
	• Varun Babbar, MEng in Information Engineering (<i>Next: PhD at Duke</i>) 2022 Engineering Division F Thesis Prize Winner	May 2021 – June 2022
	• Javier Abad Martinez, Research Assistant (<i>Next: PhD at ETH Zurich</i>)	Nov 2021 – Feb 2022
	• Dan Ley, MEng in Information Engineering (<i>Next: PhD at Harvard</i>) 2021 Engineering Division F Thesis Prize Winner	May 2020 – Aug 2021
SELECT INVITED TALKS	• TU Dortmund, <i>RC Trustworthy Data Science Seminar</i>	Jul 2024
	• Google, <i>Responsible AI Talk Series</i>	Jul 2024
	• Finnish Center for Artificial Intelligence, <i>ELLIS HCML Workshop</i>	Jun 2024
	• KTH Royal Institute of Technology, <i>RPL Summer School</i>	Jun 2024
	• US Department of Defense, <i>CDAO Seminar</i>	Feb 2024
	• KAUST, <i>Rising Stars in AI Symposium</i>	Feb 2024
	• Instituto dos Advogados de São Paulo - IASP, <i>Seminar</i>	Dec 2023
	• House of Lords Select Committee on AI in Weapon Systems	Jun 2023
	• Indian Institute of Technology Bombay, <i>C-MInDS Seminar</i>	Jan 2023
	• University of Chicago, <i>Chicago Human+AI Lab</i>	Nov 2022
	• Massachusetts Institute of Technology, <i>AI Ethics and Policy Group</i>	Oct 2022
	• Birkbeck, University of London, <i>Workshop on Human Behavioral Aspects of XAI</i>	Sept 2022
	• Harvard Business School, <i>Digital, Data, and Design (D³) Institute</i>	Sept 2022
	• DeepMind, <i>Ethics Research Team</i>	July 2022
	• University of Bristol, <i>REPHRAIN Masterclass</i>	June 2022
	• University of Cambridge, <i>Guest Lecture for MSt in AI Ethics and Society</i>	June 2022
	• University of Manchester, <i>Advances in Data Science and AI Conference</i>	June 2022
	• Leverhulme Centre for the Future of Intelligence, <i>Seminar Series</i>	June 2022
	• Ada Lovelace Institute, <i>Brown Bag Lunch Seminar</i>	June 2022
	• Von Hügel Institute, <i>Research Salon on Artificial Intelligence</i>	June 2022
	• Pennsylvania State University, <i>Young Achievers Symposium</i>	Apr 2022
	• UK Ministry of Defense, <i>AI Safety Workshop</i>	Mar 2022
	• Vanguard Health, <i>MedTech Insight Podcast</i>	June 2021
	• Technical University of Denmark, <i>Trustworthiness and Interpretability in ML Seminar</i>	Apr 2021
	• Harvard SEAS, <i>Guest Lecture for COMPSCI 282BR: Explainability in ML</i>	Apr 2021
	• Imperial College, <i>Explainable AI Seminar</i>	Feb 2021
	• Cambridge Observatory for Human-Machine Collaboration	Sept 2020
	• Robust and Responsible AI Developers (<i>Keynote</i>)	July 2020
	• ICML Workshop on Extending Explainable AI (<i>Keynote</i>)	July 2020
	• Mozilla All-Hands Meeting (<i>Keynote</i>)	June 2020
	• QuantumBlack (McKinsey), <i>AI Seminar</i>	May 2020
	• Fiddler Labs, <i>Explainable AI Seminar</i>	May 2019
PROFESSIONAL SERVICE	Finance Chair , ACM Conference on Fairness, Accountability, and Transparency (FAccT)	2024
	Associate Chair , International Conference on Machine Learning (ICML)	2024
	Co-Organizer	
	• Deep Learning Indaba: Responsible AI Practical Workshop	Sept 2024
	• ELSA Workshop on Generative AI and Creative Arts	Mar 2024
	• Deep Learning Indaba: Responsible AI Practical Workshop	Sept 2023
	• ICML Workshop on Human-Machine Collaboration and Teaming	July 2022
	• ELLIS Workshop on Human-Centric Machine Learning	May 2021
	• ICML Workshop on Human Interpretability in Machine Learning	July 2020
	• IBM + Partnership on AI Workshop on Explainable AI	Feb 2020

Program Committee (Reviewer)

- 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 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

OTHER EXPERIENCE	Advisory Panel, Bourne-Epsom Protocol , London, UK	Jul 2023 – Present
	• Expert guidance on the technical aspects of AI in education	
	Advisor, OECD , Paris, FR	Apr 2023 – Present
	• Expert Group on AI Risk & Accountability	
	Advisory Chief Scientist, Trustwise , London, UK	Jan 2023 – Present
	• Building safety guardrails for large language models	
	Subject Matter Expert, Accenture , New York, NY	Jan 2023 – Apr 2024
	• Advising on Responsible AI strategy	
	Advisor, Responsible AI Institute , Austin, TX	Oct 2021 – Jan 2023
	• Building a scalable certification program for AI systems	
	Advisor, Credo AI , Palo Alto, CA	June 2020 – June 2021
	• Scoped an AI governance and auditing platform; backed by AI Fund	
	Research Assistant, Carnegie Mellon University , Pittsburgh, PA	Jan 2017 – Sept 2019
	• Advisors: José M.F. Moura (ECE), Pradeep Ravikumar (MLD), and Zico Kolter (CSD)	
	Student Fellow, .406 Ventures , Boston, MA	July 2017 – June 2019
	• Sourced startups and performed first-round due diligence on ventures	
	Intern, Microsoft , Redmond, WA	May 2018 – Aug 2018
	• Project: explainable conversational agents for technical hardware documentation	
	Co-Founder, Perceptsense , Pittsburgh, PA	Jan 2017 – May 2018
	• Built products to harvest vehicular telematics data; pipeline acquired by Honda Motors	
	Intern, Groupon , Chicago, IL	May 2017 – Aug 2017
	• Project: personalized and targeted pricing algorithms to drive purchase frequency	
	Intern, InquisitHealth , New York, NY	June 2016 – Sept 2016
	• Project: speaker diarization and adverse event detection in doctor-patient conversations	