

# Umang Bhatt

---

CONTACT [usb20@cam.ac.uk](mailto:usb20@cam.ac.uk)  
INFORMATION [umangsbhatt.github.io](https://umangsbhatt.github.io)

Citizenship: USA

EDUCATION **University of Cambridge**, Cambridge, UK  
Ph.D. in Engineering (Machine Learning) Sept 2019 – Present  
Advisor: [Adrian Weller](#)  
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  
Advisor: [José M.F. Moura](#)  
B.S. in Electrical and Computer Engineering Aug 2015 – May 2019

POSITIONS Summer Fellow, **Harvard Center for Research on Computation and Society** From June 2022  
Hosts: Milind Tambe (SEAS) and Himabindu Lakkaraju (HBS)  
Enrichment Student, **The Alan Turing Institute** Oct 2021 – Present  
Fellow, **Mozilla Foundation** Oct 2020 – Dec 2021  
Research Fellow, **Partnership on AI** June 2019 – Sept 2020  
Research Assistant, **Carnegie Mellon University** Jan 2017 – Sept 2019  
Mentors: Pradeep Ravikumar (MLD), Zico Kolter (CSD), Fei Fang (ISR), and Radu Marculescu (ECE)

PEER-REVIEWED CONFERENCE PUBLICATIONS

- [1] **Diverse and Amortised Counterfactual Explanations for Uncertainty Estimates**  
*AAAI International Conference on Artificial Intelligence (AAAI) 2022*  
Dan Ley, **Umang Bhatt**, Adrian Weller
- [2] **On the Fairness of Causal Algorithmic Recourse**  
*AAAI International Conference on Artificial Intelligence (AAAI) 2022 (Oral)*  
Julius von Kügelgen, Amir Karimi, **Umang Bhatt**, Isabel Valera, Adrian Weller, Bernhard Schölkopf
- [3] **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
- [4] **Getting a CLUE: A Method for Explaining Uncertainty Estimates**  
*International Conference on Learning Representations (ICLR) 2021 (Oral)*  
Javier Antorán, **Umang Bhatt**, Tameem Adel, Adrian Weller, José Miguel Hernández-Lobato
- [5] **FIMAP: Feature Importance by Minimal Adversarial Perturbation**  
*AAAI International Conference on Artificial Intelligence (AAAI) 2021*  
Matt Chapman-Rounds, **Umang Bhatt**, Erik Pazos, Marc-Andre Schulz, Kostas Georgatzis
- [6] **Evaluating and Aggregating Feature-based Explanations**  
*International Joint Conference on Artificial Intelligence (IJCAI) 2020*  
**Umang Bhatt**, Adrian Weller, José M.F. Moura
- [7] **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
- [8] **You Shouldn't Trust Me: Learning Models Which Conceal Unfairness From Multiple Explanation Methods**  
*European Conference on Artificial Intelligence (ECAI) 2020*  
Botty Dimanov, **Umang Bhatt**, Mateja Jamnik, Adrian Weller

	[9] <b>On Network Science and Mutual Information for Explaining Deep Neural Networks</b> <i>IEEE Conference on Acoustics, Speech and Signal Processing (ICASSP) 2020</i> Brian Davis*, <b>Umang Bhatt*</b> , Kartikeya Bhardwaj*, Radu Marculescu, José M.F. Moura	
	[10] <b>Building Human-Machine Trust via Interpretability</b> <i>AAAI International Conference on Artificial Intelligence (AAAI) 2019 (Student Abstract)</i> <b>Umang Bhatt</b> , Pradeep Ravikumar, José M.F. Moura	
	[11] <b>A Robot's Expressive Language Affects Human Strategy and Perceptions in a Competitive Game</b> <i>IEEE Conference on Robot and Human Interactive Communication (ROMAN) 2019</i> Aaron Roth, Samantha Reig, <b>Umang Bhatt</b> , Johnathan Schulgach, Tamara Amin, Afsaneh Doryab, Fei Fang, Manuela Veloso	
JOURNAL PAPERS	[12] <b>How Transparency Modulates Trust in Artificial Intelligence</b> <i>Patterns</i> . Cell Press 2022 John Zerilli, <b>Umang Bhatt</b> , Adrian Weller	
SELECT WORKSHOP PUBLICATIONS	[13] <b>Do Concept Bottleneck Models Learn As Intended?</b> <i>ICLR Workshop on Responsible AI 2021</i> Andrei Margeloiu*, Matt Ashman*, <b>Umang Bhatt*</b> , Yanzhi Chen, Mateja Jamnik, Adrian Weller	
	[14] <b>Counterfactual Accuracies for Alternative Models</b> <i>ICLR Workshop on Machine Learning in Real Life 2020</i> <b>Umang Bhatt</b> , Adrian Weller, Muhammad Bilal Zafar, Krishna Gummadi	
BOOK CHAPTERS	[15] <b>Challenges in Deploying Explainable Machine Learning</b> <i>xxAI – Beyond explainable Artificial Intelligence</i> . Springer 2022 <b>Umang Bhatt</b> , Alice Xiang, Shubham Sharma, Joydeep Ghosh, Ruchir Puri, José M.F. Moura, Peter Eckersley, Adrian Weller	
	[16] <b>Trust in Artificial Intelligence: Clinicians are Essential</b> <i>Healthcare Information Technology for Cardiovascular Medicine</i> . Springer 2021 <b>Umang Bhatt</b> , Zohreh Shams	
SELECT FELLOWSHIPS AND AWARDS	J.P. Morgan AI PhD Fellowship  The Alan Turing Institute Enrichment Studentship  Mozilla Fellowship  Partnership on AI Research Fellowship Leverhulme Center for the Future of Intelligence PhD Scholarship Fully funded by DeepMind and the Leverhulme Trust Best Presentation Award, AAAI Spring Symposium on Interpretable AI for Well-Being Lovett Family Endowed Scholarship, The Andrew Carnegie Society Undergraduate Research Presentation Award, CMU NSF I-Corps Site Award for research commercialization H. F. McCullough Memorial Scholarship, CMU	2022 – 2023 2021 – 2022 2020 – 2021 2019 – 2020 2019 – 2023  2019 2019 2017 2017 2016
GRANTS	Co-Investigator, “Social Explainability (SOXAI) for Trustworthy AI” £107,000 from EPSRC via Research Centre on Privacy, Harm Reduction, and Adversarial Influence PI: Frens Kroeger (Coventry); Other CIs: James Hancock (Stanford) and Beate Grawemeyer (Coventry)	2021–2022
SELECT INVITED TALKS	<ul style="list-style-type: none"> <li>• UK Ministry of Defense, <i>AI Safety Workshop</i></li> <li>• Huawei, <i>Strategy &amp; Technology Workshop</i></li> <li>• Technical University of Denmark, <i>Trustworthiness and Interpretability in ML Seminar</i></li> <li>• Harvard SEAS (<i>Guest Lecture</i>)</li> <li>• Imperial College, <i>Explainable AI Seminar</i></li> <li>• Cambridge Observatory for Human-Machine Collaboration (<i>Keynote</i>)</li> <li>• Robust and Responsible AI Developers (<i>Keynote</i>)</li> </ul>	Mar 2022 Oct 2021 Apr 2021 Apr 2021 Feb 2021 Sept 2020 July 2020

	<ul style="list-style-type: none"> <li>• ICML Workshop on Extending Explainable AI (<i>Keynote</i>)</li> <li>• Mozilla All-Hands Meeting (<i>Keynote</i>)</li> <li>• QuantumBlack (McKinsey), <i>AI Seminar</i></li> <li>• Fiddler Labs, <i>Explainable AI Seminar</i></li> <li>• AILA, <i>Symposium on Creating a Fair and Ethical Future</i></li> <li>• AAAI, <i>Spring Symposium on AI and Society</i></li> <li>• University of Chicago, <i>Data Science for Social Good Conference</i></li> </ul>	<p>July 2020</p> <p>June 2020</p> <p>May 2020</p> <p>May 2019</p> <p>Oct 2018</p> <p>Mar 2018</p> <p>Oct 2017</p>
TEACHING EXPERIENCE	<p><b>Thesis Co-Supervisor/Mentor</b>, University of Cambridge</p> <ul style="list-style-type: none"> <li>• Katherine Collins, MPhil in Machine Learning and Machine Intelligence</li> <li>• Varun Babbar, MEng in Information Engineering</li> <li>• Javier Abad, Research Assistantship (<i>Next: PhD at ETH Zurich</i>)</li> <li>• Dan Ley, MEng in Information Engineering (<i>Next: PhD at Harvard</i>)</li> <li>• Charlie Rogers Smith, Research Assistantship (<i>Next: Future of Humanity Institute</i>)</li> </ul> <p><b>Teaching Assistant (Supervisor/Grader)</b>, University of Cambridge</p> <ul style="list-style-type: none"> <li>• <i>Inference</i> (3F8) for Richard Turner and David Krueger</li> <li>• <i>Probabilistic ML</i> (4F13) for Zoubin Ghahramani and J.M. Hernández-Lobato</li> </ul> <p><b>Teaching Assistant</b>, Carnegie Mellon University</p> <ul style="list-style-type: none"> <li>• <i>Machine Learning for Engineers - Masters</i> (18-661) for Gauri Joshi</li> <li>• <i>Machine Learning - PhD</i> (10-701) for Ziv Bar-Joseph and Pradeep Ravikumar</li> <li>• <i>Practical Data Science</i> (15-688) for Zico Kolter</li> <li>• <i>Principles of Imperative Computation</i> (15-122) for Illiano Cervesato</li> <li>• <i>Principles of Computing</i> (15-110) for Margret Reid-Miller</li> </ul>	<p>Nov 2021 – Present</p> <p>May 2021 – Present</p> <p>Nov 2021 – Feb 2022</p> <p>May 2020 – Aug 2021</p> <p>June 2021 – Aug 2021</p> <p>Lent 2022</p> <p>Michaelmas 2020</p> <p>Spring 2019</p> <p>Fall 2018</p> <p>Spring 2018</p> <p>Fall 2017</p> <p>Spring 2017</p>
PROFESSIONAL SERVICE	<p><b>Co-Organizer</b></p> <ul style="list-style-type: none"> <li>• ICML Workshop on Human-Machine Collaboration and Teaming</li> <li>• ELLIS Workshop on Human-Centric Machine Learning</li> <li>• ICML Workshop on Human Interpretability in Machine Learning</li> <li>• IBM + Partnership on AI Workshop on Explainable AI</li> </ul> <p><b>Program Committee</b></p> <ul style="list-style-type: none"> <li>• UAI – Conference on Uncertainty in Artificial Intelligence</li> <li>• ICML – International Conference on Machine Learning</li> <li>• FAccT – ACM Conference on Fairness, Accountability, and Transparency</li> <li>• AAAI – International Conference on Artificial Intelligence</li> <li>• AISTATS – Conference on Artificial Intelligence and Statistics</li> <li>• ICLR – International Conference on Learning Representations</li> <li>• NeurIPS – Conference on Neural Information Processing Systems</li> <li>• ICAIF – ACM Conference on Artificial Intelligence and Finance</li> <li>• KDD – ACM Conference on Knowledge Discovery and Data Mining</li> </ul> <p><b>Reviewer</b></p> <ul style="list-style-type: none"> <li>• Journal of Artificial Intelligence Research (JAIR)</li> <li>• ACM Transactions on Computer-Human Interaction (TOCHI)</li> <li>• Artificial Intelligence Journal (AIJ)</li> <li>• ACM Transactions on Interactive Intelligent Systems (TiiS)</li> </ul>	<p>July 2022</p> <p>May 2021</p> <p>July 2020</p> <p>Feb 2020</p> <p>2021, 2022</p> <p>2021, 2022</p> <p>2021, 2022</p> <p>2021, 2022</p> <p>2021, 2022</p> <p>2021, 2022</p> <p>2020, 2021</p> <p>2020, 2021</p> <p>2021</p>
PROFESSIONAL EXPERIENCE	<p>Advisor, <b>Responsible AI Institute</b>, Austin, TX</p> <ul style="list-style-type: none"> <li>• Building tools and certification programs for Responsible AI</li> </ul> <p>Advisor, <b>Credo AI</b>, Palo Alto, CA</p> <ul style="list-style-type: none"> <li>• Scoped an AI governance and auditing platform; backed by AI Fund</li> </ul> <p>Student Fellow, <b>.406 Ventures</b>, Boston, MA</p> <ul style="list-style-type: none"> <li>• Sourced startups and performed first-round due diligence on ventures</li> </ul> <p>Program Management Intern, <b>Microsoft</b>, Redmond, WA</p> <ul style="list-style-type: none"> <li>• Project: explainable conversational agents for technical hardware documentation</li> </ul> <p>Co-Founder, <b>Perceptsense</b>, Pittsburgh, PA</p> <ul style="list-style-type: none"> <li>• Built products to harvest vehicular telematics data; pipeline acquired by <b>Honda Motors</b></li> </ul>	<p>Oct 2021 – Present</p> <p>June 2020 – June 2021</p> <p>July 2018 – June 2020</p> <p>May 2018 – Aug 2018</p> <p>Jan 2017 – May 2018</p>