# Umang Bhatt

| | |
|---|---|
| CONTACT INFORMATION — usb20@cam.ac.uk / umangsbhatt.github.io | Citizenship: USA |

CONTACT
INFORMATION

usb20@cam.ac.uk
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

| | |
|---|---|
| Enrichment Student, **The Alan Turing Institute** | Oct 2021 – Present |
| Student Fellow, **Leverhulme Center for the Future of Intelligence** | Sept 2019 – 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*
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] **On Network Science and Mutual Information for Explaining Deep Neural Networks**<br>*IEEE Conference on Acoustics, Speech and Signal Processing (ICASSP) 2020*<br>Brian Davis\*, ***Umang Bhatt***\*, Kartikeya Bhardwaj\*, Radu Marculescu, José M.F. Moura |
| | [10] **Building Human-Machine Trust via Interpretability**<br>*AAAI International Conference on Artificial Intelligence (AAAI) 2019* (Student Abstract)<br>***Umang Bhatt***, Pradeep Ravikumar, José M.F. Moura |
| | [11] **A Robot's Expressive Language Affects Human Strategy and Perceptions in a Competitive Game**<br>*IEEE Conference on Robot and Human Interactive Communication (ROMAN) 2019*<br>Aaron Roth, Samantha Reig, ***Umang Bhatt***, Johnathan Schulgach, Tamara Amin, Afsaneh Doryab, Fei Fang, Manuela Veloso |
| JOURNAL ARTICLES | [12] **How Transparency Modulates Trust in Artificial Intelligence**<br>*Patterns.* Cell Press 2022<br>John Zerilli, ***Umang Bhatt***, Adrian Weller |
| SELECT WORKSHOP PUBLICATIONS | [13] **Do Concept Bottleneck Models Learn As Intended?**<br>*ICLR Workshop on Responsible AI 2021*<br>Andrei Margeloiu\*, Matt Ashman\*, ***Umang Bhatt***\*, Yanzhi Chen, Mateja Jamnik, Adrian Weller |
| | [14] **Counterfactual Accuracies for Alternative Models**<br>*ICLR Workshop on Machine Learning in Real Life 2020*<br>***Umang Bhatt***, Adrian Weller, Muhammad Bilal Zafar, Krishna Gummadi |
| BOOK CHAPTERS | [15] **Challenges in Deploying Explainable Machine Learning**<br>*xxAI – Beyond explainable Artificial Intelligence.* Springer 2022<br>***Umang Bhatt***, Alice Xiang, Shubham Sharma, Joydeep Ghosh, Ruchir Puri, José M.F. Moura, Peter Eckersley, Adrian Weller |
| | [16] **Trust in Artificial Intelligence: Clinicians are Essential**<br>*Healthcare Information Technology for Cardiovascular Medicine.* Springer 2021<br>***Umang Bhatt***, Zohreh Shams |

| | | |
|---|---|---|
| SELECT FELLOWSHIPS AND AWARDS | **J.P. Morgan AI PhD Fellowship** | 2022 − 2023 |
| | **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 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 |

| | | |
|---|---|---|
| GRANTS | Co-Investigator, "Social Explainability (SOXAI) for Trustworthy AI" | 2021−2022 |
| | £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) | |

| | | |
|---|---|---|
| SELECT INVITED TALKS | *Beyond Feature Importance: Explaining Uncertainty Estimates* | |
| | • Huawei Strategy & Technology Workshop | Oct 2021 |
| | • International Conference on Learning Representations (ICLR) | May 2021 |
| | • Technical University of Denmark | Apr 2021 |
| | • Harvard SEAS | Apr 2021 |
| | *Challenges in Deploying Explainable Machine Learning* | |
| | • Imperial College Explainable AI Seminar | Feb 2021 |

| | | |
|---|---|---|
| | • Robust and Responsible AI (Rsqrd) Developers Meetup | July 2020 |
| | • Keynote: ICML Workshop on Extending Explainable AI | July 2020 |
| | • Keynote: Mozilla All-Hands Meeting | June 2020 |
| | • QuantumBlack (McKinsey) AI Seminar | May 2020 |
| | • ACM Conference on Fairness, Accountability, and Transparency (FAccT) | Jan 2020 |

*Aggregating Feature-based Model Explanations*

| | |
|---|---|
| • International Joint Conference on Artificial Intelligence (IJCAI) | July 2020 |
| • Fiddler Labs | May 2019 |
| • AAAI Spring Symposium on Interpretable AI for Well-being | Mar 2019 |
| • Element AI | Dec 2018 |

<table>
<tr><td>TEACHING<br>EXPERIENCE</td><td colspan="2"><b>Thesis Co-Supervisor/Mentor</b>, University of Cambridge</td></tr>
<tr><td></td><td>• Katherine Collins, MPhil in Machine Learning and Machine Intelligence</td><td>Nov 2021 – Present</td></tr>
<tr><td></td><td>• Varun Babbar, MEng in Information Engineering</td><td>May 2021 – Present</td></tr>
<tr><td></td><td>• Javier Abad Martinez, Research Assistantship</td><td>Nov 2021 – Feb 2022</td></tr>
<tr><td></td><td>• Dan Ley, MEng in Information Engineering <i>(Next: JP Morgan)</i></td><td>May 2020 – Aug 2021</td></tr>
<tr><td></td><td>• Charlie Rogers Smith, Research Assistantship <i>(Next: Future of Humanity Institute)</i></td><td>June 2021 – Aug 2021</td></tr>
</table>

**Teaching Assistant (Supervisor/Grader)**, University of Cambridge

| | |
|---|---|
| • *Inference* (3F8) for Richard Turner and David Krueger | Lent 2022 |
| • *Probabilistic ML* (4F13) for Zoubin Gharamani and J.M. Hernández-Lobato | Michaelmas 2020 |

**Teaching Assistant**, Carnegie Mellon University

| | |
|---|---|
| • *Machine Learning for Engineers - Masters* (18-661) for Gauri Joshi | Spring 2019 |
| • *Machine Learning - PhD* (10-701) for Ziv Bar-Joseph and Pradeep Ravikumar | Fall 2018 |
| • *Practical Data Science* (15-688) for Zico Kolter | Spring 2018 |
| • *Principles of Imperative Computation* (15-122) for Illiano Cervesato | Fall 2017 |
| • *Principles of Computing* (15-110) for Margret Reid-Miller | Spring 2017 |

<table>
<tr><td>PROFESSIONAL<br>SERVICE</td><td colspan="2"><b>Co-Organizer</b></td></tr>
<tr><td></td><td>• ELLIS Workshop on Human-Centric Machine Learning</td><td>May 2021</td></tr>
<tr><td></td><td>• ICML Workshop on Human Interpretability in Machine Learning</td><td>July 2020</td></tr>
<tr><td></td><td>• IBM + Partnership on AI Workshop on Explainable AI</td><td>Feb 2020</td></tr>
</table>

**Program Committee**

| | |
|---|---|
| • UAI – Conference on Uncertainty in Artificial Intelligence | 2021, 2022 |
| • ICML – International Conference on Machine Learning | 2021, 2022 |
| • FAccT – ACM Conference on Fairness, Accountability, and Transparency | 2021, 2022 |
| • AAAI – International Conference on Artificial Intelligence | 2021, 2022 |
| • AISTATS – Conference on Artificial Intelligence and Statistics | 2021, 2022 |
| • ICLR – International Conference on Learning Representations | 2021, 2022 |
| • NeurIPS – Conference on Neural Information Processing Systems | 2020, 2021 |
| • ICAIF – ACM Conference on Artificial Intelligence and Finance | 2020, 2021 |
| • KDD – ACM Conference on Knowledge Discovery and Data Mining | 2021 |

**Reviewer**

- Journal of Artificial Intelligence Research (JAIR)
- ACM Transactions on Computer-Human Interaction (TOCHI)
- Artificial Intelligence Journal (AIJ)
- ACM Transactions on Interactive Intelligent Systems (TiiS)

<table>
<tr><td>PROFESSIONAL<br>EXPERIENCE</td><td>Advisor, <b>Responsible AI Institute</b>, Austin, TX</td><td>Oct 2021 – Present</td></tr>
<tr><td></td><td>• Building tools and certification programs for Responsible AI</td><td></td></tr>
<tr><td></td><td>Advisor, <b>Credo AI</b>, Palo Alto, CA</td><td>June 2020 – June 2021</td></tr>
<tr><td></td><td>• Scoped an AI governance and auditing platform; backed by AI Fund</td><td></td></tr>
<tr><td></td><td>Student Fellow, <b>.406 Ventures</b>, Boston, MA</td><td>July 2018 – June 2020</td></tr>
<tr><td></td><td>• Sourced startups and performed first-round due diligence on ventures</td><td></td></tr>
<tr><td></td><td>Program Management Intern, <b>Microsoft</b>, Redmond, WA</td><td>May 2018 – Aug 2018</td></tr>
<tr><td></td><td>• Project: explainable conversational agents for technical hardware documentation</td><td></td></tr>
<tr><td></td><td>Co-Founder, <b>Percepsense</b>, Pittsburgh, PA</td><td>Jan 2017 – May 2018</td></tr>
<tr><td></td><td>• Built products to harvest vehicular telematics data; pipeline acquired by <b><i>Honda Motors</i></b></td><td></td></tr>
</table>