

# Himabindu Lakkaraju

---

**Contact Information** 491 Morgan Hall  
15 Harvard Way  
Boston, MA 02163

Science and Engineering Complex  
150 Western Ave, Suite 6.220  
Boston, MA 02134

E-mail: [hlakkaraju@hbs.edu](mailto:hlakkaraju@hbs.edu); [hlakkaraju@seas.harvard.edu](mailto:hlakkaraju@seas.harvard.edu)  
Webpage: <http://himalakkaraju.github.io>

**Research Interests** Trustworthy Machine Learning (Interpretability, Fairness, Robustness, and Privacy); Large Language Models; Human-AI Interaction; Applications of AI/ML to Decision Making in Healthcare, Law, and Policy.

**Academic & Professional Experience** **Harvard University**  
*Assistant Professor with appointments in the Business School* 01/2020 - Present  
*and the Department of Computer Science*

*Postdoctoral Fellow* 11/2018 - 12/2019

**Simons Institute for the Theory of Computing, UC Berkeley**  
*Visiting Scientist, Summer Cluster on Interpretable Machine Learning* 06/2022 - 08/2022  
*Visiting Graduate Student, Summer Cluster on Algorithmic Fairness* 07/2018 - 08/2018

**Microsoft Research, Redmond**  
*Visiting Researcher* 5/2017 - 6/2017  
*Research Intern* 6/2016 - 9/2016

**University of Chicago**  
*Data Science for Social Good Fellow* 6/2014 - 8/2014

**IBM Research**  
*Research Engineer* 7/2010 - 7/2012

**Advisory Roles** **The Stanford Center for Legal Informatics, Stanford University**  
*Advisory Board Member, Computational Antitrust Project* 01/2020 - Present

**Fiddler AI**  
*Chief AI Research Fellow and Advisor* 06/2021 - 11/2022

**Education** **Stanford University** 9/2012 - 9/2018  
Ph.D. in Computer Science

**Stanford University** 9/2012 - 9/2015  
Master of Science (MS) in Computer Science

**Indian Institute of Science (IISc)** 8/2008 - 7/2010  
Master of Engineering (MEng) in Computer Science & Automation

**Selected Honors & Achievements** **AI2050 Early Career Fellowship** by Schmidt Sciences 2024  
**NSF CAREER Award** 2023  
Named **Kavli Fellow** by the National Academy of Sciences 2023  
**Adobe Data Science Research Award** 2023  
**Best Paper Award, ICML Workshop on Interpretable ML in Healthcare** 2022

<b>Outstanding Paper Award Honorable Mention</b>	2022
<b>NeurIPS Workshop on Trustworthy and Socially Responsible Machine Learning</b>	
<b>JP Morgan Faculty Research Award</b>	2022
Selected as a member of the <b>National AI Advisory Committee</b>	2022
instituted by the US government (could not serve due to citizenship status)	
<b>National Science Foundation (NSF) Amazon Fairness in AI Grant</b>	2021
<b>Google AI for Social Good Research Award</b>	2021
<b>Best Paper Runner Up, ICML Workshop on Algorithmic Recourse</b>	2021
<b>Google Research Award</b>	2020
<b>Amazon Research Award</b>	2020
Co-founded <b>Trustworthy ML Initiative</b> with the goal of enabling easy access to resources on trustworthy ML & to build a community of researchers/practitioners	2020
<b>Hoopes Prize</b> for undergraduate thesis mentoring, Harvard University	2020
Named as one of the <b>35 Innovators Under 35</b> (Global) by MIT Tech Review	2019
Named as an <b>Innovator to Watch</b> by Vanity Fair	2019
Selected for the prestigious <b>Cowles Fellowship</b> by Yale University (declined)	2018
<b>INFORMS Data Mining Best Paper Award</b>	2017
<b>Microsoft Research Dissertation Grant</b>	2017
Named as a <b>Rising Star in Computer Science</b>	2016
<b>Outstanding Reviewer Award</b>	2016
International World Wide Web Conference (WWW)	
<b>Google Anita Borg Fellowship</b> in recognition of research and leadership	2015
<b>Stanford Graduate Fellowship</b> for exceptional academic performance	2013-17
Awarded to top 3% of Stanford Ph.D. students	
<b>Eminence and Excellence Award</b> for outstanding research contributions	2012
IBM Research	
<b>Best Paper Award, SIAM International Conference on Data Mining (SDM)</b>	2011
<b>All India Rank 32</b> (99.82%ile)	2008
Graduate Aptitude Test in Engineering (GATE)	
Entrance examination for IISc & IITs in Computer Science & Engineering	

#### Selected Grants & Fellowships

##### As Faculty

NSF CAREER Award (US\$550,664) – Sole PI	2023 - 2028
AI2050 Early Career Fellowship by Schmidt Futures (US\$300,00) – Sole PI	2023 - 2026
Adobe Data Science Research Award (US\$50,000) – PI	2023 - 2024
D3 Institute at Harvard Grant (US\$600,000) – Sole PI	2022 - 2025
JP Morgan Faculty Research Award (US\$110,000) – Sole PI	2022 - 2024
NSF-Amazon Fairness in AI (FAI) grant (US\$375,000) – co-PI	2021 - 2024
Amazon Faculty Research Award (US\$70,000) – Sole PI	2021 - 2024
Google AI for Social Good Research Award (US\$10,000) – Sole PI	2021 - 2022
Google Research Award (US\$600,000) – PI	2020 - 2024
NSF IIS: Robust Intelligence (RI) Small (US\$450,000) – Harvard PI	2020 - 2023
Bayer Trust in Science Award (US\$100,000) – PI	2020 - 2021

##### As Student

Microsoft Research Dissertation Grant (US\$20,000)	2017
Stanford Graduate Fellowship (tuition + US\$41,700 p.a.)	2013 - 2017
Google Anita Borg Scholarship (US\$10,000)	2015
Facebook Graduate Fellowship Finalist (US\$500)	2013
Indian Institute of Science Graduate Scholarship (tuition + Rs.96,000 p.a.)	2008 - 2010
SAP India Research Grant (Rs.150,000)	2009 - 2010

**Research Articles**      **Total Citations: 7834**      **h-index: 37**      **i10-index: 56**

(\* below indicates equal contribution)

### Book Chapters

- [76] Analyzing Human Decisions and Machine Predictions in Bail Decision Making  
Jon Kleinberg, Himabindu Lakkaraju, Jure Leskovec, Jens Ludwig, Sendhil Mullainathan  
(author names are ordered alphabetically)  
[The Inequality Reader: Contemporary and Foundational Readings in Race, Class, and Gender](#); Third Edition, 2022.

### Articles in Peer-Reviewed Journals

- [75] TalkToModel: Explaining Machine Learning Models with Interactive Natural Language Conversations  
Dylan Slack, Satyapriya Krishna, Himabindu Lakkaraju\*, Sameer Singh\*  
[Nature Machine Intelligence](#) - 2023.  
**Outstanding Paper Award Honorable Mention, NeurIPS Workshop on Trustworthy and Socially Responsible ML, 2022.**
- [74] Evaluating Explainability for Graph Neural Networks  
Chirag Agarwal, Owen Queen, Himabindu Lakkaraju, Marinka Zitnik  
[Nature Scientific Data](#) - 2023.
- [73] When Does Uncertainty Matter?: Understanding the Impact of Predictive Uncertainty in ML Assisted Decision Making  
Sean McGrath, Parth Mehta, Alexandra Zyttek, Isaac Lage, Himabindu Lakkaraju  
[TMLR](#) - Transactions on Machine Learning Research, 2023.  
**Featured in [VentureBeat](#)**
- [72] Human Decisions and Machine Predictions  
Jon Kleinberg, Himabindu Lakkaraju, Jure Leskovec, Jens Ludwig, Sendhil Mullainathan  
[QJE - Quarterly Journal of Economics](#), 2018.  
(author names are ordered alphabetically)  
**Featured in [MIT Technology Review](#), [Harvard Business Review](#), [The New York Times](#), and as Research Spotlight on [National Bureau of Economics front page](#)**
- [71] Mining Digital Footprints to Extract Patterns and Predict Real-Life Outcomes  
Michal Kosinski, Yilun Wang, Himabindu Lakkaraju, Jure Leskovec  
[Psychological Methods](#) - 2016.

### Articles in Peer-Reviewed Conference Proceedings

- [70] In-context Unlearning: Language Models as Few Shot Unlearners  
Martin Pawelczyk, Seth Neel, Himabindu Lakkaraju  
[ICML](#) - International Conference on Machine Learning, 2024.
- [69] Understanding the Effects of Iterative Prompting on Truthfulness  
Satyapriya Krishna, Chirag Agarwal, Himabindu Lakkaraju  
[ICML](#) - International Conference on Machine Learning, 2024.

- [68] Characterizing Data Point Vulnerability as Average-Case Robustness  
Tessa Han, Suraj Srinivas, Himabindu Lakkaraju  
[UAI](#) - International Conference on Uncertainty in Artificial Intelligence, 2024.
- [67] Quantifying Uncertainty in Natural Language Explanations of Language Models  
Sree Harsha Tanneru, Chirag Agarwal, Himabindu Lakkaraju  
[AISTATS](#) - International Conference on Artificial Intelligence and Statistics, 2024.  
**Spotlight Presentation, NeurIPS Workshop on Robustness of Few-shot and Zero-shot Learning in Foundation Models, 2023.**
- [66] Fair Machine Unlearning: Data Removal while Mitigating Disparities  
Alex Oesterling, Jiaqi Ma, Flavio Calmon, Himabindu Lakkaraju  
[AISTATS](#) - International Conference on Artificial Intelligence and Statistics, 2024.
- [65] Investigating the Fairness of Large Language Models for Predictions on Tabular Data  
Yanchen Liu, Srishti Gautam, Jiaqi Ma, Himabindu Lakkaraju  
[NAACL](#) - The North American Chapter of the Association for Computational Linguistics, 2024.
- [64] A Study on the Calibration of In-context Learning  
Hanlin Zhang, Yi-Fan Zhang, Yaodong Yu, Dhruv Madeka, Dean Foster, Eric Xing, Himabindu Lakkaraju, Sham Kakade  
[NAACL](#) - The North American Chapter of the Association for Computational Linguistics, 2024.
- [63] Post hoc Explanations of Language Models can Improve Language Models  
Satyapriya Krishna, Jiaqi Ma, Dylan Slack, Asma Ghandeharioun, Sameer Singh, Himabindu Lakkaraju  
[NeurIPS](#) - Advances in Neural Information Processing Systems, 2023.
- [62] Which Models have Perceptually-Aligned Gradients? An Explanation via Off-Manifold Robustness  
Suraj Srinivas\*, Sebastian Bordt\*, Himabindu Lakkaraju  
[NeurIPS](#) - Advances in Neural Information Processing Systems, 2023.  
**Spotlight Presentation (Top 3%)**
- [61] Verifiable Feature Attributions: A Bridge between Post Hoc Explainability and Inherent Interpretability  
Usha Bhalla\*, Suraj Srinivas\*, Himabindu Lakkaraju  
[NeurIPS](#) - Advances in Neural Information Processing Systems, 2023.
- [60] M4: A Unified XAI Benchmark for Faithfulness Evaluation of Feature Attribution Methods across Metrics, Modalities, and Models  
Xuhong Li, Mengnan Du, Jiamin Chen, Yekun Chai, Himabindu Lakkaraju, Haoyi Xiong  
[NeurIPS](#) - Advances in Neural Information Processing Systems, 2023.
- [59] Towards Bridging the Gaps between the Right to Explanation and the Right to be Forgotten  
Satyapriya Krishna\*, Jiaqi Ma\*, Himabindu Lakkaraju  
[ICML](#) - International Conference on Machine Learning, 2023.
- [58] On the Impact of Actionable Explanations on Social Segregation  
Ruijiang Gao, Himabindu Lakkaraju  
[ICML](#) - International Conference on Machine Learning, 2023.
- [57] On Minimizing the Impact of Dataset Shifts on Actionable Explanations  
Anna Meyer\*, Dan Ley\*, Suraj Srinivas, Himabindu Lakkaraju  
[UAI](#) - Conference on Uncertainty in Artificial Intelligence, 2023.  
**Oral Presentation (Top 5%)**
- [56] Probabilistically Robust Recourse: Navigating the Trade-offs between Costs and Robustness in Algorithmic Recourse  
Martin Pawelczyk, Teresa Datta, Johannes van den Heuvel, Gjergji Kasneci, Himabindu Lakkaraju  
[ICLR](#) - International Conference on Learning Representations, 2023.

- [55] On the Privacy Risks of Algorithmic Recourse  
Martin Pawelczyk, Himabindu Lakkaraju\*, Seth Neel\*  
[AISTATS](#) - International Conference on Artificial Intelligence and Statistics, 2023.
- [54] Which Explanation Should I Choose? A Function Approximation Perspective to Characterizing Post hoc Explanations  
Tessa Han, Suraj Srinivas, Himabindu Lakkaraju  
[NeurIPS](#) - *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.  
**Best Paper Award, ICML Workshop on Interpretable ML in Healthcare, 2022.**
- [53] Flatten the Curve: Efficiently Training Low-Curvature Neural Networks  
Suraj Srinivas, Kyle Matoba, Himabindu Lakkaraju, Francois Fleuret  
[NeurIPS](#) - *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- [52] OpenXAI: Towards a Transparent Evaluation of Model Explanations  
Chirag Agarwal, Satyapriya Krishna, Eshika Saxena, Martin Pawelczyk, Nari Johnson, Isha Puri, Marinka Zitnik, Himabindu Lakkaraju  
[NeurIPS](#) - *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- [51] Data Poisoning Attacks on Off-Policy Evaluation Methods  
Elita Lobo, Harvineet Singh, Marek Petrik, Cynthia Rudin, Himabindu Lakkaraju  
[UAI](#) - *Conference on Uncertainty in Artificial Intelligence*, 2022.  
**Oral Presentation (Top 5%)**
- [50] Exploring Counterfactual Explanations Through the Lens of Adversarial Examples: A Theoretical and Empirical Analysis  
Martin Pawelczyk, Chirag Agarwal, Shalmali Joshi, Sohini Upadhyay, Himabindu Lakkaraju  
[AISTATS](#) - *International Conference on Artificial Intelligence and Statistics*, 2022.
- [49] Probing GNN Explainers: A Rigorous Theoretical and Empirical Analysis of GNN Explanation Methods  
Chirag Agarwal, Marinka Zitnik\*, Himabindu Lakkaraju\*  
[AISTATS](#) - *International Conference on Artificial Intelligence and Statistics*, 2022.
- [48] Fairness via Explanation Quality: Evaluating Disparities in the Quality of Post hoc Explanations  
Jessica Dai, Sohini Upadhyay, Ulrich Aivodji, Stephen Bach, Himabindu Lakkaraju  
[AIES](#) - *AAAI/ACM Conference on AI, Ethics, and Society*, 2022.
- [47] Towards Robust Off-Policy Evaluation via Human Inputs  
Harvineet Singh, Shalmali Joshi, Finale Doshi-Velez, Himabindu Lakkaraju  
[AIES](#) - *AAAI/ACM Conference on AI, Ethics, and Society*, 2022.
- [46] A Human-Centric Perspective on Model Monitoring  
Murtuza N Shergadwala, Himabindu Lakkaraju, Krishnaram Kenthapadi  
[HCOMP](#) - *AAAI Conference on Human Computation and Crowdsourcing*, 2022.
- [45] Towards Robust and Reliable Algorithmic Recourse  
Sohini Upadhyay\*, Shalmali Joshi\*, Himabindu Lakkaraju  
[NeurIPS](#) - *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.  
**Best Paper Runner Up, ICML Workshop on Algorithmic Recourse, 2021.**
- [44] Reliable Post hoc Explanations: Modeling Uncertainty in Explainability  
Dylan Slack, Sophie Hilgard, Sameer Singh, Himabindu Lakkaraju  
[NeurIPS](#) - *Advances in Neural Information Processing Systems*, 2021.
- [43] Counterfactual Explanations Can Be Manipulated  
Dylan Slack, Sophie Hilgard, Himabindu Lakkaraju, Sameer Singh  
[NeurIPS](#) - *Advances in Neural Information Processing Systems*, 2021.
- [42] Learning Models for Algorithmic Recourse  
Alexis Ross, Himabindu Lakkaraju, Osbert Bastani  
[NeurIPS](#) - *Advances in Neural Information Processing Systems*, 2021.

- [41] Towards the Unification and Robustness of Perturbation and Gradient Based Explanations  
Sushant Agarwal, Shahin Jabbari, Chirag Agarwal\*, Sohini Upadhyay\*, Steven Wu, Himabindu Lakkaraju  
*ICML - International Conference on Machine Learning, 2021.*  
Shorter version presented at Foundations of Responsible Computing (FORC), 2022.
- [40] Towards a Unified Framework for Fair and Stable Graph Representation Learning  
Chirag Agarwal, Himabindu Lakkaraju\*, Marinka Zitnik\*  
*UAI - Conference on Uncertainty in Artificial Intelligence, 2021.*  
**Oral Presentation (Top 5%)**
- [39] Does Fair Ranking Improve Minority Outcomes? Understanding the Interplay of Human and Algorithmic Biases in Online Hiring  
Tom Suhr, Sophie Hilgard, Himabindu Lakkaraju  
*AIES - AAAI/ACM Conference on Artificial Intelligence, Ethics, and Society, 2021.*
- [38] Fair influence maximization: A welfare optimization approach  
Aida Rahmattalabi, Shahin Jabbari, Himabindu Lakkaraju, Phebe Vayanos, Eric Rice, Milind Tambe  
*AAAI - AAAI International Conference on Artificial Intelligence, 2021.*
- [37] Beyond Individualized Recourse: Interpretable and Interactive Summaries of Actionable Recourses  
Kaivalya Rawal, Himabindu Lakkaraju  
*NeurIPS - Advances in Neural Information Processing Systems, 2020.*
- [36] Incorporating Interpretable Output Constraints in Bayesian Neural Networks  
Wanqian Yang, Lars Lorch, Moritz Gaule, Himabindu Lakkaraju, Finale Doshi-Velez  
*NeurIPS - Advances in Neural Information Processing Systems, 2020.*  
**Spotlight Presentation (Top 3%)**
- [35] Robust and Stable Black Box Explanations  
Himabindu Lakkaraju, Nino Arsov, Osbert Bastani  
*ICML - International Conference on Machine Learning, 2020*
- [34] How do I fool you?: Manipulating User Trust via Misleading Black Box Explanations  
Himabindu Lakkaraju, Osbert Bastani  
*AIES - AAAI/ACM Conference on Artificial Intelligence, Ethics, and Society, 2020.*  
**Oral Presentation (Top 16.6%)**
- [33] Fooling LIME and SHAP: Adversarial Attacks on Post hoc Explanation Methods  
Dylan Slack, Sophie Hilgard, Emily Jia, Sameer Singh, Himabindu Lakkaraju  
*AIES - AAAI/ACM Conference on Artificial Intelligence, Ethics, and Society, 2020.*  
**Featured in [Harvard Business Review](#) and [deeplearning.ai](#)**  
**Best Paper (Non-Archival) at AAAI Workshop on Safe AI, 2020**  
**Oral Presentation (Top 16.6%)**
- [32] Faithful and Customizable Explanations of Black Box Models  
Himabindu Lakkaraju, Ece Kamar, Rich Caruana, Jure Leskovec  
*AIES - AAAI/ACM Conference on Artificial Intelligence, Ethics, and Society, 2019.*  
**Oral Presentation (Top 10%)**
- [31] The Selective Labels Problem: Evaluating Algorithmic Predictions in the Presence of Unobservables  
Himabindu Lakkaraju, Jon Kleinberg, Jure Leskovec, Jens Ludwig, Sendhil Mullainathan  
*KDD - ACM SIGKDD Conference on Knowledge Discovery and Data Mining, 2017.*  
**Oral Presentation (Top 8.5%)**
- [30] Learning Cost-Effective and Interpretable Treatment Regimes  
Himabindu Lakkaraju, Cynthia Rudin  
*AISTATS - International Conference on Artificial Intelligence and Statistics, 2017.*  
**INFORMS Data Mining Best Paper Award, 2017**



- [29] Identifying Unknown-Unknowns in the Open World: Representations and Policies for Guided Exploration  
Himabindu Lakkaraju, Ece Kamar, Rich Caruana, Eric Horvitz  
[AAAI](#) - *AAAI International Conference on Artificial Intelligence*, 2017.  
**Featured in [Bloomberg Technology](#)**
- [28] Confusions over Time: An Interpretable Bayesian Model for Characterizing Trends in Decision Making  
Himabindu Lakkaraju, Jure Leskovec  
[NIPS](#) - *Advances in Neural Information Processing Systems*, 2016.
- [27] Interpretable Decision Sets: A Joint Framework for Description and Prediction  
Himabindu Lakkaraju, Stephen Bach, Jure Leskovec  
[KDD](#) - *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2016.
- [26] A Machine Learning Framework to Identify Students at Risk of Adverse Academic Outcomes  
Himabindu Lakkaraju, Everaldo Aguiar, Carl Shan, David Miller, Nasir Bhanpuri, Rayid Ghani, Kecia Addison  
[KDD](#) - *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2015.  
**Oral Presentation (Top 8.2%)**
- [25] A Bayesian Framework for Modeling Human Evaluations  
Himabindu Lakkaraju, Jure Leskovec, Jon Kleinberg, Sendhil Mullainathan  
[SDM](#) - *SIAM International Conference on Data Mining*, 2015.  
**Oral Presentation (Top 5%)**
- [24] Who, When, and Why: A Machine Learning Approach to Prioritizing Students at Risk of not Graduating High School on Time  
Everaldo Aguiar, Himabindu Lakkaraju, Nasir Bhanpuri, David Miller, Ben Yuhas, Kecia Addison, Shihching Liu, Marilyn Powell and Rayid Ghani  
[LAK](#) - *Learning Analytics and Knowledge Conference*, 2015.
- [23] What's in a name? Understanding the Interplay between Titles, Content, and Communities in Social Media  
Himabindu Lakkaraju, Julian McAuley, Jure Leskovec  
[ICWSM](#) - *International AAAI Conference on Weblogs and Social Media*, 2013.  
**Featured in [Time](#), [Forbes](#), [Phys.Org](#), [Business Insider](#), [New Scientist](#)**  
**Oral Presentation (Top 3%)**
- [22] Dynamic Multi-Relational Chinese Restaurant Process for Analyzing Influences on Users in Social Media  
Himabindu Lakkaraju, Indrajit Bhattacharya, Chiranjib Bhattacharyya  
[ICDM](#) - *IEEE International Conference on Data Mining*, 2012.  
**Oral Presentation (Top 8.6%)**
- [21] Attention prediction on social media brand pages  
Himabindu Lakkaraju, Jitendra Ajmera  
[CIKM](#) - *ACM Conference on Information and Knowledge Management*, 2011.
- [20] Exploiting Coherence for the Simultaneous Discovery of Latent Facets and associated Sentiments  
Himabindu Lakkaraju, Chiranjib Bhattacharyya, Indrajit Bhattacharya, Srujana Merugu  
[SDM](#) - *SIAM International Conference on Data Mining*, 2011.  
**Best Paper Award**
- [19] TEM: A novel perspective to modeling content on microblogs  
Himabindu Lakkaraju, Hyung-Il-Ahn  
[WWW](#) - *International World Wide Web Conference*, 2011.
- [18] Smart news feeds for social networks using scalable joint latent factor models  
Himabindu Lakkaraju, Angshu Rai, Srujana Merugu  
[WWW](#) - *International World Wide Web Conference*, 2011.

## Selected Preprints, Working Papers, and Workshop Articles

- [17] The Disagreement Problem in Explainable Machine Learning: A Practitioner’s Perspective [\[PDF\]](#) (under review)  
Satyapriya Krishna\*, Tessa Han\*, Alex Gu, Shahin Jabbari, Steven Wu, Himabindu Lakkaraju  
Preliminary version presented at CHI Workshop on Trust and Reliance in Human-AI Teams, 2022; **Featured in Fortune Magazine.**
- [16] Towards Safe Large Language Models for Medicine [\[PDF\]](#) (under review)  
Tessa Han, Aounon Kumar, Chirag Agarwal, Himabindu Lakkaraju
- [15] Are Large Language Models Post Hoc Explainers? [\[PDF\]](#) (under review)  
Nicholas Kroeger, Dan Ley, Satyapriya Krishna, Chirag Agarwal, Himabindu Lakkaraju  
Preliminary version presented at NeurIPS Workshop on XAI in Action: Past, Present, and Future Applications, 2023.
- [14] Certifying LLM Safety against Adversarial Prompting [\[PDF\]](#) (under review)  
Aounon Kumar, Chirag Agarwal, Suraj Srinivas, Aaron Li, Soheil Feizi, Himabindu Lakkaraju
- [13] Accurate, Explainable, and Private Models: Providing Recourse While Minimizing Training Data Leakage [\[PDF\]](#) (under review)  
Catherine Huang, Chelsea Swoopes, Christina Xiao, Jiaqi Ma, Himabindu Lakkaraju  
Preliminary version presented at ICML Workshop on New Frontiers in Adversarial Machine Learning, 2023.
- [12] Analyzing chain-of-thought prompting in Large language models via gradient-based feature Attributions [\[PDF\]](#) (under review)  
Skyler Wu, Eric Shen, Charumathi Badrinath, Jiaqi Ma, Himabindu Lakkaraju  
Preliminary version presented at ICML Workshop on Challenges in Deployable Generative AI, 2023.
- [11] Rethinking Explainability as a Dialogue: A Practitioner’s Perspective [\[PDF\]](#) (under review)  
Himabindu Lakkaraju, Dylan Slack, Yuxin Chen, Chenhao Tan, Sameer Singh  
Preliminary version presented at NeurIPS Workshop on Human-Centered AI, 2022.
- [10] On the Impact of Adversarially Robust Models on Algorithmic Recourse [\[PDF\]](#) (under review)  
Satyapriya Krishna, Chirag Agarwal, Himabindu Lakkaraju  
Preliminary version presented at NeurIPS Workshop on Trustworthy and Socially Responsible ML, 2022.
- [9] Manipulating Large Language Models to Increase Product Visibility [\[PDF\]](#) (working paper)  
Aounon Kumar, Himabindu Lakkaraju
- [8] When Algorithms Explain Themselves: AI Adoption and Accuracy of Experts’ Decisions (working paper)  
Himabindu Lakkaraju, Chiara Farronato
- [7] Can Model Explanations Help Reduce Biases in Real-World Decision Making? (working paper)  
Himabindu Lakkaraju, Paul Hamilton, Sarah Tan
- [6] Operationalizing the Blueprint for an AI Bill of Rights: Understanding and Addressing the Gaps between Research and Policy (working paper)  
Himabindu Lakkaraju, Usha Bhalla, Alex Oesterling, Suresh Venkatasubramanian
- [5] On the Incompatibility Between AI Regulatory Guidelines (working paper)  
Paul Hamilton, Jiaqi Ma, Himabindu Lakkaraju
- [4] An Empirical Study of the Trade-offs between Interpretability and Fairness [\[PDF\]](#)  
Shahin Jabbari, Han-Ching Ou, Himabindu Lakkaraju, Milind Tambe  
ICML Workshop on Human Interpretability in Machine Learning, 2020



- [3] Aspect Specific Sentiment Analysis using Hierarchical Deep Learning [\[PDF\]](#)  
Himabindu Lakkaraju, Richard Socher, Christopher Manning  
NIPS Workshop on Deep Learning and Representation Learning, 2014

## Patents

- [2] Extraction and Grouping of Feature Words  
Chiranjib Bhattacharyya, Himabindu Lakkaraju, Sunil Aravindam, Kaushik Nath  
[US8484228 B2](#)
- [1] Enhancing knowledge bases using rich social media  
Jitendra Ajmera, Shantanu Godbole, Himabindu Lakkaraju, Ashish Verma  
[US20130224714 A1](#)

## Advising & Mentoring

### Current Advisees:

Chirag Agarwal, Postdoctoral Fellow, Harvard University	2020 - Present
Suraj Srinivas, Postdoctoral Fellow, Harvard University	2022 - Present
Aounon Kumar, Postdoctoral Fellow, Harvard University	2023 - Present
Martin Pawelczyk, Postdoctoral Fellow, Harvard University	2023 - Present
Tessa Han, PhD Student, Harvard Medical School	2020 - Present
Satyapriya Krishna, PhD Student, Harvard CS	2021 - Present
Dan Ley, PhD Student, Harvard CS	2022 - Present
Alex Oesterling, PhD Student, Harvard CS	2022 - Present
Usha Bhalla, PhD Student, Harvard CS	2022 - Present
Paul Hamilton, PhD Student, Harvard Business School	2023 - Present
Elita Lobo, PhD Student, UMass Amherst CS	2023 - Present
Sree Harsha Tanneru, Masters Student, Harvard University	2023 - Present
Nikhil Nayak, Masters Student, Harvard University	2023 - Present
Aaron Li, Masters Student, Harvard University	2023 - Present
Yanchen Liu, Masters Student, Harvard University	2023 - Present
Charu Badrinath, Undergrad, Harvard University	2023 - Present
Eric Shen, Undergrad, Harvard University	2023 - Present
Catherine Huang, Undergrad, Harvard University	2023 - Present
Christina Xiao, Undergrad, Harvard University	2023 - Present

### Past Advisees and Interns:

Jiaqi Ma (Postdoc, Harvard University => Assistant Professor, UIUC)  
Dylan Slack (PhD, UC Irvine => Research Scientist, Scale AI)  
Aditya Karan (MS, Harvard University => PhD Student, UIUC CS)  
Kaivalya Rawal (MS, Harvard University => Research Fellow, Oxford University)  
Alexis Ross (Undergraduate, Harvard University => PhD Student, MIT EECS)  
Isha Puri (Undergraduate, Harvard University => PhD Student, MIT EECS)  
Emily Jia (Undergraduate, Harvard University => Data Scientist, Figma)  
Umang Bhatt (Research Intern, Harvard University => Assistant Professor, NYU CDS)  
Ruijiang Gao (Research Intern, Harvard University => Assistant Professor, UT Dallas)  
Harvineet Singh (Research Intern, Harvard University => Postdoc, UCSF/UC Berkeley)  
Jessica Dai (Research Intern, Harvard University => PhD Student, UC Berkeley EECS)  
Tom Suhr (Research Intern, Harvard University => PhD Student, Max Planck Institute)

## Teaching Experience

Instructor, <a href="#">Explainable Artificial Intelligence</a> Department of Computer Science, Harvard University (First ever full-fledged course on this topic)	2019, 2021, 2023
Instructor, <a href="#">Introduction to Data Science and Machine Learning</a> Harvard Business School	2020 - Present
Instructor, <a href="#">A Short Course on Explainable Machine Learning</a> Stanford Center for AI Safety	2022
Instructor, <a href="#">Introduction to ML for Social Scientists</a>	Spring 2020

Harvard Business School & Department of Computer Science		
Instructor, Explainable and Accurate AI for High-Stakes Decision Making	2020 - 2023	
Harvard Online Analytics Program		
Guest Lecture, Explainable ML in the Era of Foundation Models	Spring 2024	
Cornell University: Algorithmic Fairness Course		
Guest Lecture, User Evaluations in Explainable Machine Learning	Spring 2023	
UC Berkeley: <a href="#">Human-Centered AI Course</a>		
Guest Lecture, Explainable ML in the Era of Foundation Models	Spring 2023	
Carnegie Mellon University: <a href="#">Trustworthy AI Course</a>		
Guest Lecture, Evaluating ML Models in the Presence of Unobservables	Spring 2021	
Stanford University: <a href="#">Counterfactuals: The Science of What Ifs?</a>		
Guest Lecture, An Overview of Explainable Machine Learning	Spring 2021	
Harvard University: <a href="#">AI for Social Impact Course</a>		
Guest Lecture, Algorithms for Explainable Machine Learning	Autumn 2020	
Carnegie Mellon University: <a href="#">Advanced Introduction to Machine Learning Course</a>		
Guest Lecture, Explainable Machine Learning in Practice	Autumn 2020	
Carnegie Mellon University: <a href="#">Human-AI Interaction Course</a>		
Guest Lecture, Introduction to Data Science, Stanford Law School	Spring 2016	
Guest Lecture, Algorithms for Submodular Optimization	Winter 2016	
Stanford University: <a href="#">Mining Massive Data Sets Course</a>		
Co-instructor, Introduction to Python Programming	Spring 2015	
Stanford University: <a href="#">Girls Teaching Girls to Code (GTGTC) Initiative</a>		
Teaching Assistant for		
Stanford University: <a href="#">Mining Massive Data Sets Course</a>	Winter 2016	
Stanford University: <a href="#">Social &amp; Information Network Analysis Course</a>	Autumn 2014	
Indian Institute of Science: Machine Learning Course	Autumn 2010	

## Tutorials

<a href="#">Trustworthy Machine Learning in the Era of Foundation Models</a>	ICML, FAccT, KDD 2023
<a href="#">Model Monitoring in Practice: Lessons Learned and Open Challenges</a>	KDD, FAccT 2022
<a href="#">Explainable ML in the Wild: When Not to Trust Your Explanations</a>	FAccT 2021
<a href="#">Explainable ML: Understanding the Limits and Pushing the Boundaries</a> (Invited Tutorial)	CHIL 2021
<a href="#">Explaining Machine Learning Predictions: State-of-the-art, Challenges, and Opportunities</a>	NeurIPS 2020, AAAI 2021

## Invited Talks

Princeton University Workshop on Understanding in Natural and Artificial Minds	2024
Johns Hopkins CS Seminar Series	2024
US Securities and Exchange Commission	2024
MIT Data Science Seminar Series	2024
UPenn Center for Safe, Explainable, and Trustworthy AI Seminar Series	2024
AAAI Workshop on Privacy-Preserving Artificial Intelligence	2024
NSF Workshop on Advanced Automated Systems, Contestability, and the Law	2024
Google, Stanford, and UW Madison Workshop on Securing the Future of GenAI	2023
Yale and Google Joint Workshop on Theory and Practice of Foundation Models	2023
ICML Workshop on Interpretable ML in Healthcare	2023
ICML Workshop on Counterfactuals in Minds and Machines	2023
ICLR Workshop on Trustworthy & Reliable Large-Scale Machine Learning Models	2023
RSS Workshop on Safe Autonomy	2023
Mind and Machine Intelligence Summit, UC Santa Barbara	2023

Cornell University and Weill Cornell Medicine	2023
Kavli Frontiers of Science Symposium	2023
Cohere AI	2023
<b>Keynote</b> at AAAI Workshop on Representation Learning for Responsible Human-Centric AI	2023
<b>Keynote</b> at AAAI Workshop on Deployable AI	2023
INFORMS Annual Meeting	2016 - 2023
NeurIPS Workshop on Women in Machine Learning (WiML)	2022
NeurIPS Workshop on Machine Learning for Health (ML4H)	2022
ICLR Workshop on Privacy, Accountability, Interpretability, Robustness, Reasoning on Structured Data	2022
CVPR Workshop on Explainable AI for Computer Vision	2022
<b>Keynote</b> at WWW Workshop on Explainable AI in Health	2022
ECCV Workshop on Adversarial Robustness in the Real World	2022
Panel Discussion on AI and the Economy, Jointly Organized by U.S. Department of Commerce, NIST, Stanford HAI, and the FinRegLab	2022
Simons Institute (Berkeley) Workshop on Societal Considerations and Applications	2022
Stanford Center for AI Safety Workshop on Explainable AI	2022
Stanford Human-Centered Artificial Intelligence (HAI) Conference	2022
Stanford Digital Econ Seminar	2022
MIT Initiative on the Digital Economy (IDE) Seminar Series	2022
Harvard Data Science Initiative's Annual Conference	2022
Berkman Klein Center, Harvard University	2022
Amazon Alexa Rising Star Speaker Series	2022
University of Southern California	2022
Fireside Chat on Explainability, Fiddler AI	2022
<b>Keynote</b> at ACM CIKM Conference	2021
NIST AI Risk Management Framework Workshop	2021
Pinterest <b>Distinguished Lecture</b>	2021
NeurIPS Workshop on Algorithmic Fairness through the Lens of Causality and Robustness	2021
NeurIPS Workshop on Explainable AI Approaches for Debugging and Diagnosis	2021
NeurIPS Workshop on Human and Machine Decisions	2021
<b>Keynote</b> at ICML Workshop on Interpretable ML in Healthcare	2021
<b>Keynote</b> at KDD Workshop on ML in finance	2021
AI for Good Summit organized by International Telecommunications Union & the United Nations	2021
<b>Keynote</b> at CVPR Workshop on Responsible Computer Vision	2021
<b>Keynote</b> at ICLR Workshop on Responsible AI	2021
<b>Keynote</b> at ASPLOS Workshop on Systems Architecture for Robust, Safe, and Resilient Software	2021
<b>Keynote</b> at MLSys Workshop on Personalized Recommender Systems & Algorithms	2021
University of Cambridge	2021
Neurosym Webinar Series, Jointly Organized by UPenn, MIT, Caltech, and Stanford	2021
Voices of Data Science, UMass Amherst	2021
Max Planck Symposium on Computing and Society	2021
<b>Keynote</b> at CVPR Workshop on Fair, Data-Efficient and Trusted Computer Vision	2020
<b>Keynote</b> at MICCAI Workshop on Interpretability in Medical Imaging	2020
ETH - Center for Law and Economics, Zurich	2020
University of Michigan, Ann Arbor	2019
Harvard CRCS Seminar, Cambridge	2019
AI World Conference & Expo, Cambridge	2019
EmTech MIT Conference, Cambridge	2019
Google DeepMind Annual Summit, Cambridge	2019
Women in Machine Learning Workshop, Boston	2019
ICLR Workshop on Safe Machine Learning, New Orleans	2019
Harvard Data Science Conference, Cambridge	2018
South Park Commons, San Francisco	2018

Computer Science Departmental Seminars at Carnegie Mellon University, UIUC,	2018
Harvard University, Georgia Tech, Yale University, UC San Diego,	
USC, UCLA, UC Irvine, Duke University, Brown University,	
University of Michigan, University of Maryland	
Machine Learning Departmental Seminar at Carnegie Mellon University	2018
Operations Research Departmental Seminars at Columbia University,	2018
Cornell University, Princeton University	
NYU Stern School of Business, New York	2018
MIT Sloan School of Management, Cambridge	2018
Harvard Business School, Boston	2018
UC Berkeley School of Public Health, San Francisco	2018
Microsoft Research, Redmond	2017, 2018
IBM Thomas J. Watson Research Center, New York	2017
Machine Learning Seminar at Duke University, Durham	2017
<b>Keynote</b> at ICML Workshop on Automatic Machine Learning, Sydney, Australia	2017
Stanford Biomedical Data Science Lecture Series, Palo Alto	2017
Stanford Symbolic Systems Coffee Chat Series, Palo Alto	2017
Stanford Data Science Workshop, Palo Alto	2017
Rising Stars Workshop in EECS, Pittsburgh	2016
CodeX Center, Stanford Law School, Palo Alto	2016
KDD Workshop on Data Science for Social Good, New York	2014
University of Chicago Computation Institute, Chicago	2014
Grace Hopper India Chapter, Bangalore, India	2011

**Community Service** **Co-Founder & Chair: [Trustworthy ML Initiative](#)** 2020 - Present  
We launched this initiative to enable easy access to resources on trustworthy ML, to showcase and promote the work of researchers from underrepresented groups, and to build a community of researchers and practitioners working on the topic.

**Panelist and Reviewer:** 2020 - Present  
4 National Science Foundation (NSF) Review Panels,  
Directorate for Computer and Information Science and Engineering (CISE)

**Conference Organization:**  
NeurIPS Conference (Ethics Chair) 2024  
WSDM Conference (Tutorial Chair) 2024  
FAccT Conference (Sponsorship Chair) 2023  
KDD Trustworthy AI Day (Program Chair) 2022  
KDD Deep Learning Day (Program Chair) 2021  
Grace Hopper India Conference (Program Chair) 2011

**Workshop Chair:**  
NeurIPS Workshop on Regulatable Machine Learning 2023  
NeurIPS Workshop on Explainable Artificial Intelligence 2023  
ICML Workshop on New Frontiers in Adversarial Machine Learning 2022  
ICML Workshop on Algorithmic Recourse 2021  
ELLIS Human-Centric Machine Learning Workshop 2021  
Session on Trustworthy Machine Learning at INFORMS 2020  
Session on Fairness in Machine Learning at INFORMS 2019  
ICLR Workshop on Debugging Machine Learning Models 2019  
Workshop for spreading awareness about STEM fields among middle school girls 2016  
Stanford's Girls Teaching Girls To Code (GTGTC) 2015

**Area Chair:**  
ICML - *International Conference on Machine Learning* 2019 - 2024  
NeurIPS - *Advances in Neural Information Processing Systems* 2019 - 2023  
ICLR - *International Conference on Learning Representations* 2020 - 2023  
AISTATS - *International Conference on Artificial Intelligence and Statistics* 2021 - 2023

**Program Committee:**  
AISTATS - *International Conference on Artificial Intelligence and Statistics* 2019 - 2020

FAccT - ACM Conference on Fairness, Accountability, and Transparency	2019 - 2020
AAAI - AAAI International Conference on Artificial Intelligence	2019
ICML - International Conference on Machine Learning	2018
ICLR - International Conference on Learning Representations	2018 - 2019
IJCAI - International Joint Conference on Artificial Intelligence	2018 - 2019
WWW - International World Wide Web Conference	2017 - 2018
NIPS - Advances in Neural Information Processing Systems	2016 - 2017
KDD - ACM SIGKDD Conference on Knowledge Discovery and Data Mining	2015 - 2017
CIKM - ACM Conference on Information and Knowledge Management	2011, 2017
SDM - SIAM International Conference on Data Mining	2015
UAI - Conference on Uncertainty in Artificial Intelligence	2011
AAAI - AAAI conference on Artificial Intelligence	2011

#### **Journal Reviewing and Editing:**

Frontiers in Big Data (Associate Editor)	2021 - 2023
JMLR - Journal of Machine Learning Research	2020 - 2023
MS - Management Science	2021 - 2023
OR - Operations Research	2021 - 2023
TWEB - ACM Transactions on the Web	2017
PLOS ONE - Public Library of Science ONE	2017
TKDD - ACM Transactions on Knowledge Discovery from Data	2016
TKDE - IEEE Transactions on Knowledge and Data Engineering	2015

#### **Other:**

Member, Faculty Hiring Committee, Harvard Business School	2020 - 2023
Member, Ph.D. Student Selection Committee, Harvard University	2020 - 2023
Member, Ph.D. Student Selection Committee, Stanford CS	2016

#### **Selected Media Coverage**

TIME: [Chuck Schumer wants AI to be explainable. It's harder than it sounds](#)  
Science News: [AI chatbots can be tricked into misbehaving. Can scientists stop it?](#)  
Fortune: [Explainable AI & The Disagreement Problem](#)  
Harvard Business Review: [The AI transparency paradox](#)  
MIT Technology Review: [How to upgrade judges with machine learning](#)  
Harvard Business Review: [Solving social problems with machine learning](#)  
The New York Times: [Even Imperfect Algorithms Can Improve the Criminal Justice System](#)  
VentureBeat: [Confidence, uncertainty, and trust in AI affect how humans make decisions](#)  
Wired: [This Agency Wants to Figure Out Exactly How Much You Trust AI](#)  
Bloomberg Technology: [Researchers combat gender and racial bias in AI](#)  
Forbes: [How to craft the perfect Reddit posting](#)  
Time: [How to succeed on Reddit](#)  
Business Insider: [How to execute the perfect Reddit submission](#)  
Phys.org: [Stanford Trio explore success formula for Reddit posts](#)  
International Business Times: [The secret to what makes something go viral](#)  
New Scientist: [Things that make a meme explode](#)  
The Verge: [The math behind successful Reddit submissions](#)  
ACM TechNews: [Stanford trio explore success formula for Reddit posts](#)