

THE PROGRAM GUIDE FOR THE

Eighth AAAI/ACM Conference on

ARTIFICIAL INTELLIGENCE, ETHICS, & SOCIETY

OCTOBER 20-22, 2025 | MADRID, SPAIN

Follow AIES 2025 on X! @AIESConf

PROGRAM CONTENTS

Acknowledgements	3
Welcome	4
Program Overview	7
Detailed Program	8
Monday	8
Tuesday	16
Wednesday	24
Sponsors	34

ACKNOWLEDGEMENTS

The Association for the Advancement of Artificial Intelligence acknowledges and thanks the following individuals for their generous contributions of time and energy to the successful creation and planning of the Eighth Annual AAAI / ACM Conference on Artificial Intelligence, Ethics, and Society.

CONFERENCE CO-CHAIRS

Kush Varshney IBM Research
Theodore Lechterman IE University

CONFERENCE PROGRAM CO-CHAIRS

Emanuelle Burton University of Illinois Nicholas Mattei Tulane University Andrés Páez Universidad de los Andes

STUDENT PROGRAM CHAIRS

Francesca Palmiotto *IE University* **Tracie Farrell** *The Open University*

TRAVEL AWARD CHAIR

Wenbin Zhang Florida International University

The AIES-25 Program Committee welcomes you to Madrid!



WELCOME FROM THE CONFERENCE CO-CHAIRS

AIES seeks to create and sustain a multidisciplinary community of scholars who think deeply about the impact of AI systems on humans, societies, and the world more generally. The astounding increase in submissions this year shows us that the ethical and social implications of AI are now no longer a niche area of research, but in fact a mainstream academic and industry concerns. Although this might suggest that the implications of AI are widening and deepening, we are encouraged by the growing attention from the scientific community and the rigor with which scholars from various fields investigate these issues. That multidisciplinary rigor is on full display in this year's program, which underwent the most competitive selection procedure to date. The program chairs and program committee did an amazing job!

We are excited to host AIES in Madrid this year, the first time the conference has come to either Spain or continental Europe, which offer rich and unique approaches to AI research, development, regulation, and critical reflection. We look forward to welcoming all of you to IE University, one of the region's most dynamic hubs for education and research, and hope you enjoy seeing the incredible work that is taking place in the field and the other opportunities on offer in Spain's capital city. We are grateful to our sponsors for their generous support, which enables us to keep registration fees low and support the AIES student program. We also thank the AAAI staff for doing much of the heavy lifting in organizing the conference. We hope that the conversations at AIES continue to drive the work we need to do to ensure that the path forward is a good one despite the many challenges that lie ahead.

Theodore Lechterman *IE University* · **Kush R. Varshney** *IBM Research*



WELCOME FROM THE

CONFERENCE PROGRAM CO-CHAIRS

Welcome to the 8th annual AAAI/ACM Conference on Artificial Intelligence, Ethics, and Society in Madrid! We are grateful for the support from our hosts at IE University, and the innumerable hours put in by Chesley Grove, Meredith Ellison, the rest of the AAAI staff, and the general chairs in arranging the event so that we could focus on bringing you a high-quality program.

When AIES first convened in 2018, AI ethics was something of a niche topic, both in the public imagination and in the field of AI itself. That first meeting was a small gathering in 2018 co-located with the main AAAI conference and included 57 papers. Seven years later, the technical capacities of AI have expanded exponentially (and continue to do so) and AI has been integrated into more and more aspects of our shared lives, AIES has grown to a selective, full-fledged, independent conference where a transdisciplinary collection of researchers gather to learn from each other, and the need for our work here has grown ever more pressing.

This year, after desk-rejecting a handful of submissions, we sent 748 papers out for review and accepted 238 of them; hands down, both are record numbers. We can see, from the table below, just how much attention is converging on cross-disciplinary AI topics have grown across the areas of computer science, philosophy, social sciences, and law during the conference lifetime.

We are *extremely* grateful to the 42 Senior Program Committee Members, the 648 Program Committee Members, and program committee members from all over the world who volunteered their time to provide peer reviews. The papers will be presented through a combination of oral sessions and poster sessions. One form of presentation is not more prestigious than the other; the session assignments were made to have a coherent and engaging program within the constraints of three days.

Year	Submitted	Accepted	PC Size
2025	748	238	690
2024	468	150	219
2023	237	68	139
2022	235	79	163
2021	278	106	93
2020	211	71	80
2019	220	70	88
2018	162	61	44

Continued on next page

We are pleased to have two distinguished keynote speakers: Miriam Fernandez and Emma Ruttkamp-Bloem. With the keynotes we have two panels, the first Panel on Policy and Governance: Beyond the Brussels Effect: Competing Visions of Al Governance with members Manuel Muñiz, Maria Eriksson, David Leslie, and Urs Gasser. The second, Pedagogy Panel: How (and to Whom) Do We Teach Al Ethics? with Emanuelle Burton, Shannon Vallor, Julienne LaChance, and Samer Hassan.

We hope you find the conference generative and rewarding, and that you walk away with new friends and new ideas.

Emanuelle Burton *University of Illinois, Chicago* · **Nicholas Mattei** *Tulane University* **Andrés Páez** *Universidad de los Andes*



WELCOME FROM THE STUDENT PROGRAM CHAIRS

Student research is a glimpse into the future — where one can see the direction of travel and the new lines of inquiry that are being opened by exceptional students from all over the world. Every year, the AIES Student Program aims to facilitate student participation by offering financial support, connection with our wider community of world class researchers, and 1:2:1 mentorship. We invite you to welcome students into our interdisciplinary community by engaging with student work during the poster sessions and learning more about the exciting new perspectives they are bringing to the subject of Artificial Intelligence, Ethics and Society.

We appreciate those senior scholars who are serving as mentors to the students during the conference, and also wish to thank NSF, SIGAI, and AAAI for generous financial support.

Tracie Farrell Open University, UK · Francesca Palmiotto IE University, Madrid

PROGRAM OVERVIEW

	MONDAY, OCTOBER 20	TUESDAY, OCTOBER 21	WEDNESDAY, OCTOBER 22
9:00 AM	Designation		
9:15 AM	Registration Keynote 1: Miriam Fernandez	Keynote 1:	Discussion Session 4:
9:30 AM 9:45 AM		Stereotypes, Fairness, and Oppression in LLM-based Social Interactions	
10:00 AM 10:15 AM 10:30 AM 10:45 AM 11:00 AM	Discussion Session 1: Shifting Power in Al Ecosystems	Paper Session 3: Navigating Differences and Vulnerability	Poster Session 3 & Reception
11:15 AM 11:30 AM	Refreshment Break	Refreshment Break	
11:45 AM 12:00 PM 12:15 PM 12:30 PM 12:45 PM	Paper Session 1: Redescribing AI	Pedagogy Roundtable	Paper Session 6: Integrating AI into the Workplace
1:00 PM 1:15 PM 1:30 PM 1:45 PM	Policy Panel	Paper Session 4: Dataset Creation and Knowledge Production	Keynote 2: Emma Ruttkamp-Bloem
2:00 PM 2:15 PM 2:30 PM 2:45 PM 3:00 PM	Lunch	Lunch	Lunch
3:15 PM 3:30 PM 3:45 PM 4:00 PM 4:15 PM	Discussion Session 2: Evaluating LLMs in the Context of Patient Autonomy and Human Rights	Discussion Session 3: Alignment with Whom? AI, Lay Judgment, and Expertise	Poster Session 4 & Reception
4:30 PM 4:45 PM	Refreshment Break	Refreshment Break	
5:00 PM 5:15 PM 5:30 PM 5:45 PM	Paper Session 2: Mitigating Bias	Paper Session 5: Al and the Relational Self	Paper Session 7: Risk Management
6:00 PM 6:15 PM 6:30 PM 6:45 PM 7:00 PM 7:15 PM	Poster Session 1 & Reception	Poster Session 2 & Reception	Closing Remarks
7:30 PM 7:45 PM 8:00 PM			

AIES-25 DETAILED PROGRAM

9:00-9:30	Registration
9:30-10:00	Opening Remarks & Paper Awards // Auditorium
10:00-11:15	Discussion Session 1: Shifting Power in AI Ecosystems // Auditorium Discussion Chair: Tracie Farrell
98	The Silicon Sovereignty Paradox: Navigating Fluid State-Corporate Power Dynamics in the Age of Al Maanya Singh; Anka Reuel
438	Al Policy for Whom? Reclaiming Governance from Capitalist Capture Petter Ericson; Rachele Carli; Jason Tucker; Virginia Dignum
611	\$100,000 or the robot gets it! Tech Workers' Resistance Guide: Tech Worker Actions, History, Risks, Impacts, and the Case for a Radical Flank Mohamed Abdalla
126	Algorithmic Fairness Beyond Legally Protected Groups and When Group Labels Are Unknown Abdoul Jalil Djiberou Mahamadou, Judy Wawira Gichoya and Artem A. Trotsyuk
11:15-11:45	Refreshment Break
11:45-1:00	Paper Session 2: Redescribing AI // Auditorium
225	The Stories We Govern By: AI, Risk, and the Power of Imaginaries Ninell Oldenburg; Gleb Papyshev
534	Why (not) use AI? Analyzing People's Reasoning and Conditions for AI Acceptability Jimin Mun; Wei Bin Au Yeong; Wesley Hanwen Deng; Jana Schaich Borg; Maarten Sap
644	Responsible AI Practices: Histories, Definitions, Barriers and Future Directions Lorenn P. Ruster
799	Measuring What Matters: Connecting AI Ethics Evaluations to System Attributes, Hazards, and Harms Shalaleh Rismani; Renee Shelby; Leah Davis; Negar Rostamzadeh; A Jung Moon

1:00–2:00 Headline Panel on Policy and Governance // Auditorium

Beyond the Brussels Effect: Competing Visions of AI Governance

Summary: While the EU's AI Act has established a new regulatory template, competing governance models are emerging in both theory and practice — from market-driven self-regulation to state-led oversight, and from expert-driven risk frameworks to participatory and rights-based approaches. This panel examines how these divergent visions reflect deeper questions about democratic legitimacy and the role of expertise in AI policy, and what the implications of this governance diversity might be for effective AI oversight and other governance challenges.



Panel Chair: Ted Lechterman, Ph.D.

UNESCO Chair in AI Ethics & Governance, IE University

Bio: Theodore "Ted" Lechterman holds the UNESCO Chair in AI Ethics & Governance at IE University in Spain, where he is Assistant Professor of Philosophy. His research spans political philosophy and applied ethics, focusing on artificial intelligence and democratic ideals, corporate responsibility in technology, and the ethics of private efforts to solve public problems. His scholarship appears in outlets such as *Oxford Handbook of AI Governance, Stanford Encyclopedia of Philosophy, Journal of Business Ethics, Organization Studies*, and *History of Political Thought*.

Lechterman is the author of *The Tyranny of Generosity: Why Philanthropy Corrupts Our Politics and How We Can Fix It* (Oxford University Press, 2022), which received an honorable mention for the ECPR Political Theory Prize. His second book, *Recoding Democracy: Al and the Fight for Democracy's Future*, is under contract with Polity. His dynamic lecturing and dedication to student development have earned multiple teaching awards, while his innovations in Al ethics education have been recognized by the *Financial Times*, the QS Reimagine Education Awards, and the Government of Spain.

He serves as a member of UNESCO AI Ethics Experts Without Borders, officer of the ECPR Political Theory Standing Group, managing director of Compass Ethics, and co-chair of the 8th AAAI/ ACM Conference on AI, Ethics, and Society. He regularly contributes to public debates and advises organizations on navigating ethical frontiers in business, technology, and governance.

Lechterman holds degrees from Harvard (A.B.) and Princeton (M.A., Ph.D.) and completed postdoctoral fellowships at Stanford, Goethe University Frankfurt, the Hertie School, and the University of Oxford, where he was an inaugural fellow at the Institute for Ethics in Al.



Maria Eriksson PhD.

Research Fellow, European Centre for Algorithmic Transparency; European Commission, Joint Research Centre

Bio: Maria Eriksson is a Research Fellow at the Joint Research Centre of the European Commission, based at the European Centre for Algorithmic Transparency (ECAT) in Seville, Spain. She is also an Affiliated Researcher at the Department of Arts and Cultural Sciences/Digital Cultures at Lund University, Sweden. Her research is located at the intersection of media studies, social anthropology, science and technology studies, and policy work. She has published extensively

on the sociotechnical impacts of AI technologies, as well as the role of algorithmic systems within the

cultural and creative industries. Currently, she contributes with scientific expertise to the European Commission's implementation of the EU AI Act and Digital Services Act.

Maria will speak in her personal capacity as a researcher working at the Joint Research Centre, meaning, opinions and viewpoints should not, in any way, be interpreted as the official position of the European Commission.



Professor David Leslie

Director of Ethics & Responsible Innovation Research, The Alan Turing Institute Professor of Ethics, Technology & Society, Queen Mary University of London

Bio: David Leslie is the Director of Ethics and Responsible Innovation Research at The Alan Turing Institute and Professor of Ethics, Technology and Society at Queen Mary University of London. He previously taught at Princeton's University Center for Human Values, Yale's programme in Ethics, Politics and Economics and at Harvard's Committee on Degrees in Social Studies, where he received over a dozen teaching awards including the 2014 Stanley Hoffman Prize for Teaching Excellence. David

is the author of the UK Government's official guidance on the responsible design and implementation of AI systems in the public sector, Understanding artificial intelligence ethics and safety (2019) and a principal co-author of Explaining decisions made with AI (2020), a co-badged guidance on AI explainability published by the UK's Information Commissioner's Office and The Alan Turing Institute. After serving as an elected member of the Bureau of the Council of Europe's (CoE) Ad Hoc Committee on Artificial Intelligence (CAHAI) (2021–2022), he was appointed, in 2022, as Specialist Advisor to the CoE's Committee on AI where he has led the writing of its Human Rights, Democracy and the Rule of Law Impact Assessment for AI (2024), which accompanies its AI Convention. He also serves on UNESCO's High-Level Expert Group steering the implementation of its Recommendation on the Ethics of Artificial Intelligence.



Urs Gasser Professor and Dean, Technical University of Munich

Bio: Urs Gasser is a Professor of Public Policy at the Technical University of Munich (TUM), where he is also Dean of the TUM School of Social Sciences and Technology and Rector of the Munich School of Politics and Public Policy. Prior to joining TUM, he was the Executive Director of the Berkman Klein Center at Harvard University. Gasser has advised governments worldwide on technology policy, including serving on Angela Merkel's Digital Council and as the current Chair of Thailand's International Policy Advisory Panel on AI. He is also the co-author, with

Viktor Mayer-Schönberger, of the book "Guardrails: Guiding Human Decisions in the Age of AI" (2024) and, with John Palfrey, of the "Advanced Introduction to Law and Digital Technologies" (2025).

2:00-3:15 Lunch Break (on own)

3:15-4:30	Discussion Session 2: Evaluating LLMs in the Context of Patient Autonomy and Human Rights // Auditorium Discussion Chair: Kristel Clayville
70	Privacy in Image Datasets: A Case Study on Pregnancy Ultrasounds Rawisara Lohanimit; Yankun Wu; Amelia Katirai; Yuta Nakashima; Noa Garcia
487	Assessing Human Rights Risks in AI: A Framework for Model Evaluation Vyoma Raman; Camille Chabot; Betsy Popken
497	Write on Paper, Wrong in Practice: Why LLMs Still Struggle with writing clinical notes Kristina L. Kupferschmidt; Kieran O'Doherty; Joshua A. Skorburg
513	Principles and Policy Recommendations for Comprehensive Genetic Data Governance Vivek Ramanan; Ria Vinod; Cole Williams; Sohini Ramachandran; Suresh Venkatasubramanian
4:30-5:00	Refreshment Break
5:00-6:15	Paper Session 2: Mitigating Bias // Auditorium
135	Fairness of Automatic Speech Recognition: Looking Through a Philosophical Lens Anna Seo Gyeong Choi; Hoon Choi
306	Biased Al Outputs Can Impact Humans' Implicit Bias: A Case Study of the Impact of Gender-Biased Text-to-Image Generators Mattea Sim; Natalie Grace Brigham; Tadayoshi Kohno; Tessa E. S. Charlesworth; Aylin Caliskan
542	Collective Agency in Art-making: Towards Community-centric Design of Text-to-Image (T2I) AI Tools Abdullah Hasan Safir; Noshin Tahsin; Pratyasha Saha; Dipannita Nandi; Zulkarin Jahangir; Cecily Morrison; Syed Ishtiaque Ahmed; Nusrat Jahan Mim
833	Exposing Al Bias by Crowdsourcing: Democratizing Critique of Large Language Models Hangzhi Guo; Pranav Narayanan Venkit; Eunchae Jang; Mukund Srinath; Wenbo Zhang; Bonam Mingole; Vipul Gupta; Kush R. Varshney; S. Shyam Sundar; Amulya Yadav
6:15-8:00	Poster Session 1 & Reception // Level 24
	Posters
14	From Categorical to Contextual: Interpreting High-Risk Classification for Profiling-Based AI Recommender Systems in the EU AI Act <i>Luca Nannini</i>
33	Fairness and Sparsity within Rashomon sets: Enumeration-Free Exploration and Characterization Lucas Langlade, Julien Ferry, Gabriel Laberge and Thibaut Vidal
37	Towards Experience-Centered AI: A Framework for Integrating Lived Experience in Design and Development Sanjana Gautam, Mohit Chandra, Ankolika De, Tatiana Chakravorti, Girik Malik and Munmun De Choudhury

74	Dataset-to-Dataset Evaluation Before (and Without) Sharing Data Keren Fuentes, Mimee Xu and Irene Y. Chen
103	ValuesRAG: Enhancing Cultural Alignment Through Retrieval-Augmented Contextual Learning Wonduk Seo, Zonghao Yuan and Yi Bu
113	Labeling in Their Shoes: Improving Text Annotation with Cognitive Empathy Priming Sung Hyun Kwon, Jessica Clark, Il-Horn Hann and Jui Ramaprasad
116	S-DAT: A Multilingual, GenAI-Driven Framework for Automated Divergent Thinking Assessment Jennifer Haase, Paul H. P. Hanel and Sebastian Pokutta
119	Localizing Persona Representations in LLMs Celia Cintas, Miriam Rateike, Erik Miehling, Elizabeth Daly and Skyler Speakman
128	Should LLMs be WEIRD? Exploring WEIRDness and Human Rights in Large Language Models <i>Ke Zhou, Marios Constantinides and Daniele Quercia</i>
130	Emotional Plausibility vs. Emotional Truth: Designing Against Affective Misinformation in Conversational Al Maalvika Bhat and Duri Long
132	Improving LLM Group Fairness on Tabular Data via In-Context Learning Valeriia Cherepanova, Chia-Jung Lee, Nil-Jana Akpinar, Riccardo Fogliato, Martin Bertran Lopez, Michael Kearns and James Zou
137	Understanding Endogenous Data Drift in Adaptive Models with Recourse-Seeking Users Bo-Yi Liu, Zhi-Xuan Liu, Kuan Lun Chen, Shih-Yu Tsai, Jie Gao and Hao-Tsung Yang
154	Explaining the Reputational Risks of Al-Mediated Communication: Messages labeled as Al-assisted are viewed as less diagnostic of the sender's moral character Pranav Khadpe, Kimi Wenzel, George Loewenstein and Geoff Kaufman
178	Importance of User Control in Data-Centric Steering for Healthcare Experts Aditya Bhattacharya, Simone Stumpf and Katrien Verbert
187	Sound Check: Auditing Recent Audio Dataset Practices William Agnew, Julia Barnett, Annie Chu, Rachel Hong, Michael Feffer, Robin Netzorg, Harry Jiang, Ezra Awumey and Sauvik Das
190	Highlight All the Phrases: Enhancing LLM Transparency through Visual Factuality Indicators Hyo Jin Do, Rachel Ostrand, Werner Geyer, Keerthiram Murugesan, Dennis Wei and Justin Weisz
191	Reward-on-the-Line: A Novel Offline Reinforcement Learning Method for Building Legal Conversational Agents Xubo Lin, Mingze Wang, Grace Hui Yang and Daniel Chen
192	Accountability Capture: How Record-Keeping to Support Al Transparency and Accountability (Re)shapes Algorithmic Oversight Shreya Chappidi, Jennifer Cobbe, Chris Norval, Anjali Mazumder and Jatinder Singh
196	The Term 'Agent' Has Been Diluted Beyond Utility and Requires Redefinition Brinnae Bent

198	An Audit and Analysis of LLM-Assisted Health Misinformation Jailbreaks Against LLMs Ayana Hussain, Patrick Zhao and Nicholas Vincent
200	Value Creation and Value Capture in AI: A Triple Helix Model Antoni Lorente
210	The Disparate Effects of Partial Information in Bayesian Strategic Learning Srikanth Avasarala, Serena Wang and Juba Ziani
211	Disaggregated Health Data in LLMs: Evaluating Data Equity in the Context of Asian American Representation Uvini Balasuriya Mudiyanselage, Bharat Jayprakash, Kookjin Lee and K. Hazel Kwon
216	FairPOT: Balancing AUC Performance and Fairness with Proportional Optimal Transport Pengxi Liu, Yi Shen, Matthew M. Engelhard, Benjamin A. Goldstein, Michael J. Pencina, Nicoleta J. Economou–Zavlanos and Michael M. Zavlanos
219	Beyond "Fairness": Rethinking the use of Algorithmic Predictions in Criminal Justice <i>Tashmia Sabera</i>
220	Effort-aware Fairness: Incorporating a Philosophy-informed, Human-centered Notion of Effort into Algorithmic Fairness Metrics Tin Trung Nguyen, Jiannan Xu, Zora Che, Phuong-Anh Nguyen-Le, Rushil Dandamudi, Donald Braman, Furong Huang, Hal Daumé and Zubin Jelveh
226	Perceived Risks and Benefits of Disclosing ADHD to AI-based Educational Technologies: Semi-structured Interviews Oriane Pierrès, Alireza Darvishy and Markus Christen
228	Experimental evidence that AI-managed workers tolerate lower pay without demotivation Mengchen Dong, Levin Brinkmann, Omar Sherif, Shihan Wang, Xinyu Zhang, Jean-François Bonnefon and Iyad Rahwan
231	Al Governance in the Context of the EU Al Act Byeong-Je Kim, Seunghoo Jeong, Bong-Kyung Cho and Ji-Bum Chung
233	A Human-Centered Approach to Identifying Promises, Risks, & Challenges of Text-to-Image Generative AI in Radiology Katelyn Morrison, Arpit Mathur, Aidan Bradshaw, Tom Wartmann, Steven Lundi, Afrooz Zandifar, Weichang Dai, Kayhan Batmanghelich, Motahhare Eslami and Adam Perer
243	Centring the margins: mapping AI systems as systems of power Garfield Benjamin
249	Beyond Technocratic XAI: The Who, What & How in Explanation Design Ruchira Dhar, Stephanie Brandl, Ninell Oldenburg and Anders Søgaard
253	Legal Affiliates' Views on Algorithmic Decision Making Styliani Kleanthous and Maria Kasinidou
258	A Comprehensive Evaluation of the Sensitivity of Density-Ratio Estimation Based Fairness Measurement in Regression Abdalwahab Almajed, Maryam Tabar and Peyman Najafirad

260	Toward A Causal Framework for Modeling Perception Jose M. Alvarez and Salvatore Ruggieri
262	Can We Trust AI Benchmarks? An Interdisciplinary Review of Current Issues in AI Evaluation Maria Eriksson, Erasmo Purificato, Arman Noroozian, João Vinagre, Guillaume Chaslot, Emilia Gomez and David Fernandez-Llorca
267	Decentralising LLM Alignment: A Case for Context, Pluralism, and Participation Oriane Peter and Kate Devlin
272	"Just a strange pic": Evaluating 'safety' in GenAl Image safety annotation tasks from diverse annotators' perspectives Ding Wang, Mark Díaz, Charvi Rastogi, Aida Davani Davani, Vinodkumar Prabhakaran, Pushkar Mishra, Roma Patel, Alicia Parrish, Zoe Ashwood, Michela Paganini, Tian Huey Teh, Verena Rieser and Lora Aroyo
273	Bridging the Communication Gap: Evaluating AI Labeling Practices for Trustworthy AI Development Raphael Fischer, Magdalena Wischnewski, Alexander van der Staay, Katharina Poitz, Christian Janiesch and Thomas Liebig
278	Who Foots the Bill?: State-Backed No-Fault Compensation for Experimental Harms in Artificial Intelligence Regulatory Sandboxes Ha-Chi Tran
280	LLM-based Simulations of Human Behavior in Psychological Research Santiago Flórez Sánchez
287	Unravelling Responsible AI: An Umbrella Review Gisela Reyes-Cruz, Elvira Perez Vallejos, Pepita Barnard, Eike Schneiders, Marisela Tachiquin, Dominic Price, Damian Eke, Liz Dowthwaite, Aislinn Gomez Bergin, Virginia Portillo and Joel Fischer
291	Towards Interactive Evaluations for Interaction Harms in Human-Al Systems Lujain Ibrahim, Saffron Huang, Lama Ahmad, Umang Bhatt and Markus Anderljung
509	Incident Analysis for AI Agents Carson Ezell, Xavier Roberts-Gaal and Alan Chan
569	RelAltionship Building: Analyzing Recruitment Strategies for Participatory Al Eugene Kim, Vaibhav Balloli, Berelian Karimian, Elizabeth Bondi-Kelly and Benjamin Fish
759	Moral Agents Unlike Us Jen Semler
843	Known Unknowns and Unknown Unknowns: Designing a Scalable Adverse Event Reporting System for Al Lindsey A. Gailmard, Drew Spence, Christie Lawrence and Daniel E. Ho

	Student Posters
13	Leveraging AI and ESG to Combat Corporate Corruption: An Integrated Framework e Oludolapo Makinde
59	Beyond Templates: Understanding and Addressing Human-AI Interaction Harms through Practitioner Assumptions Julia De Miguel Velázquez
7	Detecting, Classifying, and Mitigating Undesired Behaviors in Pretrained Large Language Models Jan Batzner
53	Supporting Marginalized Learners with GenAl Hamayoon Behmanush
8	Machine Legibility and Epistemic Governance in Malaysia's Smart Cities: A Postcolonial Analysis of Algorithmic Knowledge, Identity, and State Power Hesam Nourooz Pour
25	AI Ethics in Cyborg Anthropology: Examining AI-Driven Job Displacement Among Women and Marginalized Groups Tessina Grant Moloney
42	Towards Responsible AI Governance in the Brazilian Judiciary Bruno Fonseca
44	Between Code and Creed: Islamic Ethical In-Betweenness on AI in Indonesia Daphne Wong-A-Foe
56	A Survey of Large Language Model Use and its Technical Limitations in Military Systems through a Decolonial Lens Sonia Fereidooni
72	Resistance to Change as a Diagnostic Insight: An Interdisciplinary Examination of Stakeholder Opposition to AI in Primary Care Teresa Sides
73	Deploying AI in Uncertain Environments: A Technical Limitation or a Human Characteristic? Marc Elliott
695	Learning to Unlearn, Failing to Forget? Assessing Machine Unlearning Through Ethics and Epistemology Igra Aslam, Donal Khosrowi and Rahul Nagshi

AIES-25 DETAILED PROGRAM

TUESDAY, OCTOBER 21

9:00-10:00 Keynote 1: Miriam Fernandez // Auditorium

Responsible AI and the Urgent Challenge of Technology-Facilitated Gender-Based Violence

Abstract: Technology is now embedded in nearly every aspect of our lives, shaping how we work, communicate, and connect. Yet, while online spaces have become central to social interaction, they remain profoundly unsafe for many. Women and girls, in particular, continue to face disproportionate risks from technology-facilitated violence, including harassment, stalking, deepfake pornography, and algorithmic discrimination.

In this talk, I explore how current technologies, from social media platforms and generative AI tools to IoT devices, perpetuate and amplify gender-based harms. I will discuss some of the initiatives and research efforts aimed at tackling these issues, as well as the technical, social, and legal challenges that continue to hinder meaningful protection for women and girls worldwide.



Bio: Miriam Fernandez is a Professor of Responsible Artificial Intelligence at the Knowledge Media Institute (KMi), Open University (OU), UK. Her research agenda revolves around advancing Responsible AI, ensuring that technological innovation aligns with ethical principles and societal values. Her pioneering work spans diverse domains, from algorithmic transparency and fairness to the societal implications of AI deployment. By integrating AI techniques with a human-centred approach, she fosters solutions that prioritise social responsibility, transparency, and inclusivity. With a portfolio of more than 100 scientific articles in some of the best conferences and journals in her field, and having won numerous external grants supporting her research, Professor Fernandez

has significantly influenced the discourse in the field of technology development and its impact on society. Her commitment to education is demonstrated through her leadership of OUAnalyse, a strategic initiative leveraging machine-learning methods for the early identification of students at risk. This technology, currently supporting the Open University's 200K student body, has been highly awarded for its transformative impact on student outcomes. Professor Fernandez is also Equality and Diversity Champion for both KMi and the OU, where she leads the Responsible AI stream of the Center for Protecting Women Online, a flagship initiative that plays a critical role in mitigating the harmful effects of technology on women and girls worldwide.

10:00-11:15	Paper Session 3: Navigating Differences and Vulnerability // Auditorium
159	Toward A Taxonomy of Algorithmic Harms for Disability: A Systematic Review Lining Wang; Vaishnav Kameswaran; Hernisa Kacorri
347	Documenting Patterns of Exoticism of Marginalized Populations within Text-to-Image Generators Sourojit Ghosh; Sanjana Gautam; Pranav Narayanan Venkit; Avijit Ghosh
721	Disability Across Cultures: A Human-Centered Audit of Ableism in Western and Indic LLMs Mahika Phutane; Aditya Vashistha

Govern With, Not For: Understanding the Stuttering Community's Preferences and

Goals for Speech AI Data Governance in the US and China

Jingjin Li; Peiyao Liu; Rebecca Lietz; Ningjing Tang; Norman Makoto Su; Shaomei Wu

11:15–11:45 Refreshment Break

11:45-12:45 Pedagogy Panel // Auditorium

How (and to Whom) Do We Teach AI Ethics?

Summary: The teaching of AI ethics—whether in a classroom or in a corporate setting— is challenging for several reasons. Even when students or colleagues approach the task with good will, many practical hurdles remain: how do we translate the modes of knowledge and inquiry that are necessary for this type of work to students (and colleagues) more accustomed to the high-information field of computer science? How do we balance our own moral intuitions and convictions against the work of thinking through the particulars of messy and challenging situations? How do we ensure that marginalized voices are not only included, but granted equal credence? This round table will encompass both fundamental questions of teaching philosophy and practical, concrete strategies that emerge from those philosophies.



Panel Chair: Emanuelle Burton, PhD.

Bio: Emanuelle Burton is senior lecturer in the department of computer science at the University of Illinois Chicago, where she teaches courses in ethics. She is coauthor of Computing and Technology Ethics: Engaging Through Science Fiction, published in 2023 by MIT Press, and solo author of several articles on ethics in fantasy literature for children. She holds as PhD in religion and literature from the University of Chicago Divinity School.



Shannon Vallor

Bio: Prof. Shannon Vallor is the Baillie Gifford Chair in the Ethics of Data and Artificial Intelligence at the Edinburgh Futures Institute (EFI) at the University of Edinburgh, where she is also appointed in Philosophy. She directs EFI's Centre for Technomoral Futures, and is co-Director of the UKRI's BRAID (Bridging Responsible AI Divides) programme. Professor Vallor's research explores how AI and data science reshape human character and capabilities. From 2018–2020 she served as a Visiting Research Scientist and AI Ethicist at Google, and she is a standing member of Stanford University's One Hundred Year Study of Artificial Intelligence (AI100). She

is the author of Technology and the Virtues: A Philosophical Guide to a Future Worth Wanting (Oxford University Press, 2016) and The AI Mirror: How to Reclaim Our Humanity in an Age of Machine Thinking (Oxford University Press, 2024).



Julienne LaChance

Bio: Julienne is a Senior Research Scientist at SonyAI, Sony's organization dedicated to state-of-the-art AI "moonshot" projects. She is a manager on the AI Ethics team, and currently leads explorations at the intersection of security engineering and generative AI. She holds five STEM degrees, having completed a PhD/postdoc at Princeton University, and has prior industry experience as a software/firmware engineer. In addition to genAI/security, Julienne's projects span (1) large-scale, human-centric visual data collection, (2) AI deployment on a global scale, particularly in underserved regions of the world, and (3) educational initiatives in AI Ethics.



Samer Hassan

Bio: Samer Hassan (he/they) is an activist and interdisciplinary researcher, Faculty Associate at Harvard University's Berkman Klein Center for Internet & Society, and Associate Professor at Universidad Complutense de Madrid (Spain). Their research bridges social and computer sciences, focusing on online communities and democratic governance in digital platforms. They led "P2P Models", a €1.5M EU ERC project exploring decentralized autonomous organizations as free/open-source infrastructures for governance experimentation. A trained grassroots facilitator, Hassan uses participatory methods in CS classrooms to engage students in tech ethics and agile practices.

12:45-2:00	Paper Session 4: Dataset Creation and Knowledge Production // Auditorium
168	From Sandbox to Laboratory: Refractive and Ethical Human / AI Knowledge Co-Production Through Game-ful Research Practices <i>Lea Stöter</i>
261	Disciplinary Practices in the Generation of Text Synthetic Data: A Critical Discourse Analysis Adriana Alvarado Garcia; Nishanshi Atulkumar Shukla; Muneeza Azmat; Marisol Wong-Villacres
576	From Big Data to Valued Data: A Dataset Value Taxonomy for AI-Native Empirical Research Scott Seidenberger; Anindya Maiti
739	Methodological Considerations for Centering Data Workers' Epistemic Authority in Al Research Milagros Miceli; Adio-Adet Dinika; Krystal Kauffman; Camilla Salim Wagner; Laurenz Sachenbacher; Alex Hanna; Timnit Gebru

2:00-3:15 **Lunch Break (on own)**

3:15-4:30 Discussion Session 3: Alignment with Whom? AI, Lay Judgment, and Expertise // Auditorium Discussion Chair: Sanmay Das

Diverse Human Value Alignment for Large Language Models via Ethical Reasoning Jiahao Wang; Songkai Xue; Jinghui Li; Xiaozhen Wang

11

451	Aggregation Problems in Machine Ethics and Al Alignment Kevin Baum; Marija Slavkovik
564	Against Al Jurisprudence: Large Language Models and the False Promises of Empirical Judging Dasha Pruss; Jessie Allen
661	Do Your Guardrails Even Guard? Method for Evaluating Effectiveness of Guardrails in Aligning LLM Outputs with Expert User Expectations Anindya Das Antar; Xun Huan; Nikola Banovic
4:30-5:00	Refreshment Break
5:00-6:15	Paper Session 5: Al and the Relational Self // Auditorium
57	Toward an Ethic of Synthetic Relationality: Identity, Intimacy, and Risk in AI-Mediated Roleplay Environments
422	Santhosh G S; Akshay Govind S; Gokul S Krishnan; Balaraman Ravindran; Sriraam Natarajan Anthropomorphism as Social Affordance: Charting the Co-Animation of Chatbots into Social "Agents" Takuya Maeda; Luke Stark
801	The Intercepted Self: How Generative AI Challenges the Dynamics of the Relational Self Sandrine R. Schiller; Camilo Miguel Signorelli; Filippos Stamatiou
826	The Heterogeneous Effects of AI Companionship: An Empirical Model of Chatbot Usage and Loneliness and a Typology of User Archetypes Auren R. Liu; Pat Pataranutaporn; Pattie Maes
6:15-8:00	Poster Session 2 & Reception // Level 24
	Posters
10	A case for data valuation transparency via DValCards Keziah Naggita and Julienne LaChance
15	When Less Regulation Means More Complexity: The EU AI Liability Directive Withdrawal and its Impact on European Technological Competitiveness <i>Luca Nannini</i>
309	Too Focused on Accuracy to Notice the Fallout: Towards Socially Responsible Fake News Detection Esma Aïmeur, Gilles Brassard and Dorsaf Sallami
312	Steerable Pluralism: Pluralistic Alignment via Few-Shot Comparative Regression Jadie Adams, Brian Hu, Emily Veenhuis, David Joy, Bharadwaj Ravichandran, Aaron Bray, Anthony Hoogs and Arslan Basharat
318	What's Individual about Individual Fairness? Shai Ben-David, Pascale Gourdeau, Tosca Lechner and Ruth Urner

320	Scenarios in Computing Research: A Systematic Review of the Use of Scenario Methods for Exploring the Future of Computing Technologies in Society Julia Barnett, Kimon Kieslich, Jasmine Sinchai and Nicholas Diakopoulos
327	Rethinking Optimization: A Systems-Based Approach to Social Externalities Pegah Nokhiz, Aravinda Kanchana Ruwanpathirana and Helen Nissenbaum
328	Complete Categorization of Instinct-Exploiting Data-Explanations and their Generation with Large-Language Models Taro Higuchi and Einoshin Suzuki
329	Dead Zone of Accountability: Why Social Claims in Machine Learning Research Should Be Articulated and Defended <i>Tianqi Kou, Dana Calacci and Cindy Lin</i>
331	Exporting Autonomy, Importing Dependency: The Geopolitical Work of "Sovereign AI" <i>Heesoo Jang</i>
332	Exploring "Just Noticeable" Group Fairness in Rankings Mallak Alkhathlan, Hilson Shrestha, Lane Harrison and Elke Rundensteiner
335	Same Stereotypes, Different Term? Understanding the "Global South" in AI Ethics Evani Radiya-Dixit and Angèle Christin
336	Making Teams and Influencing Agents: Efficiently Coordinating Decision Trees for Interpretable Multi-Agent Reinforcement Learning Rex Chen, Stephanie Milani, Zhicheng Zhang, Norman Sadeh and Fei Fang
343	Informing AI Risk Assessment with News Media: Analyzing National and Political Variation in the Coverage of AI Risks Mowafak Allaham, Kimon Kieslich and Nicholas Diakopoulos
344	Burying The Lead: Adjusting Goals to Manage Functional Limitations of AI Tools in Healthcare Jacqueline Kernahan, Richard Bartels, Mark de Reuver, Daniel Oberski and Roel Dobbe
346	Bias is a Math Problem, AI Bias is a Technical Problem: 10-year Literature Review of AI/LLM Bias Research Reveals Narrow [Gender-Centric] Conceptions of `Bias', and Academia-Industry Gap Sourojit Ghosh and Kyra Wilson
348	Reframing Al-for-Good: Radical Questioning in Al for Human Trafficking Interventions Pratheeksha Nair, Gabriel Lefebvre, Maryam Molamohammadi, Sophia Garrel and Reihaneh Rabbany
349	Toward AI Matching Policies in Homeless Services: A Qualitative Study with Policymakers Caroline M. Johnston, Olga Koumoundouros, Angel Hsing-Chi Hwang, Laura Onasch-Vera, Eric Rice and Phebe Vayanos
354	GRAILS - A Framework for Embedding Ethical Safeguards in Software Applications for Responsible AI <i>Apurva Kulkarni and Chandrashekar Ramanathan</i>
361	Systemizing Multiplicity: The Curious Case of Arbitrariness in Machine Learning Prakhar Ganesh, Afaf Taik and Golnoosh Farnadi

362	TechOps: Technical Documentation Templates for the AI Act Laura Lucaj, Alex Loosley, Håkan Jonsson, Urs Gasser and Patrick van der Smagt
365	Will AI Take My Job? Evolving Perceptions of Automation and Labor Risk in Latin America Andrea Cremaschi, Dae-Jin Lee and Manuele Leonelli
374	Normative Moral Pluralism for AI: A Framework for Deliberation in Complex Moral Contexts David Doron Yaacov
377	Unintended Impacts of Automation for Integration? Simulating Integration Outcomes of Algorithm-Based Refugee Allocation in Germany Jakob Kappenberger, Clara Strasser Ceballos, Frederic Gerdon, Daria Szafran, Florian Rupp, Kai Eckert, Heiner Stuckenschmidt, Ruben Bach, Frauke Kreuter and Christoph Kern
381	The Socratic Dialogue as a Method for Virtue Ethics in Al: A Case Study Maarten Wilders, Íñigo Martínez de Rituerto de Troya and Roel Dobbe
386	Beauty and the Bias: Exploring the Impact of Attractiveness on Multimodal Large Language Models Aditya Gulati, Moreno D'Incà, Nicu Sebe, Bruno Lepri and Nuria Oliver
393	Comparing Human and LLM Ethical Analyses: A Case Study in Computational Social Science Research Spencer Hey, Julie Walsh and Eni Mustafaraj
395	Beyond Proxy Variables: Extending Refugee Allocation Algorithms for Equitable Predictions Clara Strasser Ceballos, Marcus Novotny and Christoph Kern
396	Trust Formation in Healthcare AI: An Exploration of Older Adults' Perspectives Önder Celik, Marlene Kulla and Justyna Stypinska
401	Why Do Decision Makers (Not) Use AI? A Cross-Domain Analysis of Factors Impacting AI Adoption Rebecca Yu, Valerie Chen, Ameet Talwalkar and Hoda Heidari
402	NomicLaw: Emergent Trust and Strategic Argumentation in LLMs during Collaborative Law-Making Asutosh Hota and Jussi Jokinen
405	Aligning AI Systems with Human Values: A Method for Identifying and Specifying Values Liv Ziegfield, Esther Kox, Ivana Akrum and Marlijin Heijnen
409	(Don't) Tell Me What To Do: A Retrospective Study on the Use of a Data Ethics Framework Sophia Worth, Georgia Panagiotidou and Elena Simperl
416	Toward Responsible ASR for African American English Speakers: A Scoping Review of Bias and Equity in Speech Technology Jay Cunningham, Jainaba Jawara, Jeffrey Basoah, Adinawa Adjagbodjou, Kowe Kadoma and Aaleyah Lewis
417	Advancing NLP Data Equity: Practitioner Responsibility and Accountability in NLP Data Practices Jay Cunningham, Kevin Shao, Rock Yuren Pang and Nat Mengist
418	On the Misalignment Between Legal Notions and Statistical Metrics of Intersectional Fairness Deborah Dormah Kanubala and Isabel Valera

423	Fairness-Aware Post-Processing in Supervised Classification: L1/L2 Norm and Optimal Swapping Methods Flore Vancompernolle Vromman, Sylvain Courtain, Pierre Leleux and Marco Saerens
428	The Case for "Thick Evaluations" of Cultural Representation in AI Rida Qadri, Mark Díaz, Ding Wang and Michael Madaio
432	The AI Power Disparity Index: Toward a Compound Measure of AI Actors' Power to Shape the AI Ecosystem Rachel Kim, Blaine Kuehnert, Seth Lazar, Ranjit Singh and Hoda Heidari
433	Adaptive Accountability in Networked Multi-Agent Systems Saad Alqithami
434	Hide or Highlight: Understanding the Impact of Factuality Expression on User Trust Hyo Jin Do and Werner Geyer
4 39	Accountability Framework for Healthcare AI Systems: Towards Joint Accountability in Decision Making Prachi Bagave, Marcus Westberg, Marijn Janssen and Aaron Yi Ding
442	A Taxonomy of Questions for Critical Reflection in Machine-Assisted Decision-Making Simon Fischer, Hanna Schraffenberger, Serge Thill and Pim Haselager
444	Who Owns The Robot?: Four Ethical and Socio-technical Questions about Wellbeing Robots in the Rea World through Community Engagement Minja Axelsson, Jiaee Cheong, Rune Nyrup and Hatice Gunes
446	A Moral Agency Framework for Legitimate Integration of AI in Bureaucracies Chris Schmitz and Joanna Bryson
457	Documenting Deployment with Fabric: A Repository of Real-World AI Governance Mackenzie Jorgensen, Kendall Brogle, Katherine Collins, Lujain Ibrahim, Arina Shah, Petra Ivanovic, Noah Broestl, Gabriel Piles, Paul Dongha, Hatim Abdulhussein, Adrian Weller, Jillian Powers and Umang Bhatt
458	Embodied AI at the Margins: Postcolonial Ethics for Intelligent Robotic Systems Atmadeep Ghoshal, Martim Brandao, Ruba Abu-Salma and Sanjay Modgil
460	Needle in a Patched Haystack: Evaluating Saliency Maps for Vision LLMs Bastien Zimmermann and Matthieu Boussard
474	Who pays the RENT? Implications of Spatial Inequality for Prediction-Based Allocation Policies <i>Tasfia Mashiat, Patrick Fowler and Sanmay Das</i>
613	Street-Level AI: Are Large Language Models Ready for Real-World Judgements? Gaurab Pokharel, Shafkat Farabi, Patrick J. Fowler, Sanmay Das

	Student Posters
37	Auditing and Validating Fairness and Ethics in Machine Learning Systems Disa Sariola
43	Lay-stakeholder centric sociotechnical mechanisms for addressing the impacts of generative Al <i>Julia Barnett</i>
49	Advancing Fairness in Generative AI through Intrinsic and Extrinsic Bias Evaluation and Mitigation <i>Mina Arzaghi</i>
75	How Can Large Language Models Be More Reliable? Yael Moros Daval
23	Governing AI Proactively: Cooperative Models of Anticipation and Accountability Jared Katzman
9	Human Centered AI for Research Ethics and Transparency Tatiana Chakravorti
21	The Dilemmas of Moral Autonomy in the Transfer of Human Consciousness to Al Denis Chiriac
58	Critical explorations on the socio-ethical implications of Creative AI for artists Anna-Kaisa Kaila
70	Binary Optics: Colonial Classification of Trans Lives in Algorithmic Technologies Christoffer Koch Andersen
14	Platform Vernaculars: How Al Image Generators Create New Forms of Visual Bias Jinu K Varghese and Padma Rani
20	Fairer Datasets for Advancing Responsible AI Systems Siddharth Jaiswal

AIES-25 DETAILED PROGRAM

9:00-10:15	Discussion Session 4: Stereotypes, Fairness, and Oppression in LLM-based Social Interactions // Auditorium Discussion Chair: Andrés Páez
101	AI, Normality, and Oppressive Things Ting-an Lin; Linus Ta-Lun Huang
175	Interactional Fairness in LLM Multi-Agent Systems: An Evaluation Framework Ruta Binkyte
195	Social Misattributions in Conversations with Large Language Models Andrea Ferrario; Alberto Termine; Alessandro Facchini
213	SESGO: Spanish Evaluation of Stereotypical Generative Outputs Melissa Robles; Catalina Bernal; Denniss Raigoso; Mateo Dulce Rubio
10:15-11:45	Poster Session 3 & Refreshment Break // Floor -4
	Posters
97	Aligning AI with Public Values: Deliberation and Decision-Making for Governing Multimodal LLMs in Political Video Analysis
	Tanusree Sharma, Yujin Potter, Zachary Kilhoffer, Yun Huang, Dawn Song and Yang Wang
302	How Deep Is Representational Bias in LLMs? The Cases of Caste and Religion Agrima Seth, Monojit Choudhary, Sunayana Sitaram, Kentaro Toyama, Aditya Vashistha and Kalika Bali
316	Social Scientists on the Role of AI in Research
	Tatiana Chakravorti, Xinyu Wang, Pranav Narayanan Venkit, Sai Koneru, Kevin Munger and Sarah Rajtmajer
308	Toward Valid Measurement Of (Un)fairness For Generative AI: A Proposal For Systematization Through The Lens Of Fair Equality of Chances Kimberly Le Truong, Annette Zimmermann and Hoda Heidari
345	RAI Advocacy: Communicative Strategies for Advancing Responsible AI in Large Technology Companies Jordan Duran, Samir Passi and Mihaela Vorvoreanu
476	Operationalizing critical data approaches in fraud detection in Latin America Ana Paula Moritz and Alayna Kennedy
477	Concept Creep in Safe Artificial Intelligence Laura Fearnley and Ibrahim Habli

481	A case for specialisation in non-human entities El-Mahdi El-Mhamdi, Lê-Nguyên Hoang and Mariame Tighanimine
484	GermanPartiesQA: Benchmarking Commercial Large Language Models and AI Companions for Political Alignment and Sycophancy Jan Batzner, Volker Stocker, Stefan Schmid and Gjergji Kasneci
490	Exclusive Flux: A Review of Flux's Generation of LGBTQ+ Couples Lynn Vonderhaar, Kayla Taylor, Jennifer Wojton and Omar Ochoa
491	Disclosure and Evaluation as Fairness Interventions for General-Purpose AI Vyoma Raman, Judy Hanwen Shen, Andy K. Zhang, Lindsey Gailmard, Rishi Bommasani, Angelina Wang and Daniel Ho
500	Responsible AI in Media Organizations: Four Case Studies Implementing Ethical Tools in Practice Maaike Harbers, Sophie Horsman, Pascal Wiggers, Huib Aldewereld, Marcio Fuckner, Coert Van Gemeren, Oumaima Hajri, Floor Schukking and Nathalie Stembert
502	Ten Insights From Other Domains That Inform Responsible AI Frameworks Dunstan Allison-Hope, Patrick Gage Kelley, Reena Jana, Angela McKay and Allison Woodruff
505	Mapping Moral Reasoning in LLMs: A Multi-Dimensional Analysis of Safety Principle Conflicts Sachit Mahajan
510	SycEval: Evaluating LLM Sycophancy Aaron Fanous, Jacob Goldberg, Ank Agrawal, Joanna Lin, Anson Zhou, Sonnet Xu, Vasiliki Bikia, Roxana Daneshjou and Sanmi Koyejo
514	An Adaptive Responsible AI Governance Framework for Decentralized Organizations Kiana Jafari Meimandi, Anka Reuel, Gabriela Aranguiz-Dias, Hatim Rahama, Ala-Eddine Ayadi, Xavier Boullier, Jérémy Verdo, Louis Montanie and Mykel Kochenderfer
519	Participatory AI and the EU AI Act Chiara Ullstein, Simon Jarvers, Michel Hohendanner, Orestis Papakyriakopoulos and Jens Grossklags
527	A Mathematical Philosophy of Explanations in Mechanistic Interpretability Kola Ayonrinde and Louis Jaburi
530	Aligning Agent Policies with Preferences: Human-Centered Interpretable Reinforcement Learning Stephanie Milani, Zhicheng Zhang, Nicholay Topin, Lirong Xia and Fei Fang
535	Al Managing Agent-Based Healthcare Processes Simon Grange, Safa Alameri, Pearl Rwauya and Rami Bahsoon
537	Re-imagining Virtual Communities: Ethical Guidelines for Studying Black Twitter Christina Chance and Kai-Wei Chang
541	Evaluating Goal Drift in Language Model Agents Rauno Arike, Elizabeth Donoway, Henning Bartsch and Marius Hobbhahn
548	A Critical Look at a Critical Care Dataset: MIMIC-IV's Construction, Contents, & Consequences <i>Pinar Barlas</i>

551	Into the Void: Understanding Online Health Information in Low-Web Data Languages Hellina Hailu Nigatu, Nuredin Ali Abdelkadir, Fiker Tewelde, Stevie Chancellor and Daricia Wilkinson
554	"Stealing is wrong and too hard to pull off successfully": Inclusion of Normative versus Capacity Information in Robot Command Rejection Alyssa Hanson, Gordon Briggs, Ruchen Wen, Yifei Zhu and Tom Williams
557	Al Agents and the Law Mark Riedl and Deven Desai
558	A Principled Approach for Data Bias Mitigation Bruno Scarone, Alfredo Viola, Renée J. Miller and Ricardo Baeza-Yates
566	Fairness in Federated Learning : Fairness for whom? Afaf Taik, Khaoula Chehbouni and Golnoosh Farnadi
588	Can Improved Data Representation Support AI for Health Equity? A Visual Approach Steven Vethman, Jildau Bouwman and Cor Veenman
590	User Privacy and Large Language Models: An Analysis of Frontier Developers' Privacy Policies Jennifer King, Kevin Klyman, Emily Capstick, Tiffany Saade and Victoria Hsieh
594	Modeling Strategic Risk in School Choice: A Case for Transparent Design Mayesha Tasnim, Paul Verhagen, Tobias Blanke, Erman Acar and Sennay Ghebreab
596	Making Sense of AI Ethics and Governance Investments Marianna Ganapini, Francesca Rossi, Brian Gohering, Nich Berente and Marialen Bevilacqua
602	Al-Powered Detection of Inappropriate Language in Medical School Curricula Chiman Salavati, Shannon Song, Scott A. Hale, Roberto E. Montenegro, Shiri Dori-Hacohen and Fabricio Murai
606	You Don't Need Robust Machine Learning to Manage Adversarial Attack Risks Edward Raff, Michel Benaroch and Andrew Farris
609	Interactive AI and Human Behavior: Challenges and Pathways for AI Governance Yulu Pi, Cagatay Turkay and Daniel Bogiatzis-Gibbons
627	Do Al Companies Make Good on Voluntary Commitments to the White House? Jennifer Wang, Kayla Huang, Kevin Klyman and Rishi Bommasani
649	When Algorithms Fail: The Case for Moral Repair Pak-Hang Wong and Gernot Rieder
652	Bridging Research Gaps Between Academic Research and Legal Investigations of Algorithmic Discrimination Colleen Chien, Anna Zink and Irene Chen
660	Model Misalignment and Language Change: Traces of Al-Associated Language in Unscripted Spoken English Bryce Anderson, Riley Galpin and Tom S Juzek

664	Demographic-Agnostic Fairness without Harm Zhongteng Cai, Mohammad Mahdi Khalili and Xueru Zhang
667	VISION: Robust and Interpretable Code Vulnerability Detection Leveraging Counterfactual Augmentation David Egea, Barproda Halder and Sanghamitra Dutta
673	Towards automating deliberation? The idea of deliberative democracy embedded in Google's Habermas Machine Nicolás Palomo Hernández
682	Sacred or Synthetic? Evaluating LLM Reliability and Abstention for Religious Questions Farah Atif, Nursultan Askarbekuly, Kareem Darwish and Monojit Choudhury
686	An Investigation into Black and Brown Communities' Engagement with Data & Technology Ebtesam Al Haque, Gabriella Thompson, Angela D.R. Smith and Brittany Johnson
699	Investigating Political and Demographic Associations in Large Language Models Through Moral Foundations Theory Nicole Smith-Vaniz, Harper Lyon, Lorraine Steigner, Ben Armstrong and Nicholas Mattei
717	Al Supply Chains: An Emerging Ecosystem of Al Actors, Products, and Services Aspen Hopkins, Sarah Cen, Isabella Struckman, Andrew Ilyas, Luis Videgaray and Aleksander Madr
734	What Comes After Harm? Mapping Reparative Actions in AI through Justice Frameworks Sijia Xiao, Haodi Zou, Alice Qian Zhang, Deepak Kumar, Hong Shen, Jason Hong and Motahhare Eslami
735	Machine Learning and Public Health: Identifying and Mitigating Algorithmic Bias through a Systematic Review Sara Altamirano, Arjan Vreeken and Sennay Ghebreab
750	A Multidimensional Approach to Ethical AI Auditing Sónia Teixeira, Atia Cortés, Dilhan Thilakarathne, Gianmarco Gori, Marco Minici, Monowar Bhuyan, Nina Khairova, Tosin Adewumi, Devvjiit Bhuyan, Jack O'Keefe, Carmela Comito, João Gama and Virginia Dignum
	Student Posters
34	On Forecasting Lags in AI Risk Evaluation Paolo Bova
61	Fairness in Social Media Platforms: Modeling Behavior and Designing Interventions Salima Jaoua
64	From Model Multiplicity to Prompt Multiplicity: Emerging Arbitrariness Concerns in the Age of Generative Al Prakhar Ganesh
46	Beyond Automation: Understanding Fairness, Ethics, and Human Discretion in AI-driven Societal Decisions Gaurab Pokharel

71	The Ethics of Surveillance AI: Framing Data as a Socio-collective Good in Mitigating Data Colonialism Abiola Azeez
60	Regulatory Policies on Ethics Evaluations for Large-Scale AI Systems Neha Gupta
35	Building Preference Aware Trustworthy AI Through Fairness and Interpretability Jingyu Hu
36	Human Models for Planning Behavioral Interventions with Reinforcement Learning Eura Nofshin
41	The Design and Analysis of Algorithmic Vaccine Allocation Frameworks Jeffrey Keithley
39	Voice AI and Hermeneutical Injustice at the Border Suvradip Maitra

11:45–1:00 Paper Session 6: Integrating AI into the Workplace // Auditorium 164 AI Self-preferencing in Algorithmic Hiring: Empirical Evidence and Insights *Jiannan Xu; Gujie Li; Jane Yi Jiang* 296 No Thoughts Just AI: Biased LLM Hiring Recommendations Alter Human Decision Making and Limit Human Autonomy *Kyra Wilson; Mattea Sim; Anna-Maria Gueorguieva; Aylin Caliskan* 355 How LLM Counselors Violate Ethical Standards in Mental Health Practice: A Practitioner-Informed Framework *Zainab Iftikhar; Amy Xiao; Sean Ransom; Jeff Huang; Harini Suresh* 787 The 'Wild West' of Medicine: Exploring the emergence of 'grassroots' AI governance in radiology

1:00-2:00 Keynote 2: Emma Ruttkamp-Bloem // Auditorium

Bhargavi Ganesh; Daniel S. Schiff; Stuart Anderson

The Future of AI Ethics

Abstract: All ethics has evolved through various phases over the past decade, and currently is at risk of simply disappearing off the radar in debates on 'responsible Al'. I offer an interpretation of All ethics to turn this situation around.

Until about 2020, AI ethics was portrayed as a principle-driven ethics, which resulted in reams of documents containing abstract principles and little or no information on how to action these principles. The realisation around 2021, that, given their abstract nature, most of these documents were basically impactless, led to a kind of 'technical turn' in AI ethics, where the call was to move from the 'what' (principles) to the 'how' of AI ethics. This turn culminated, among other things, in the change in responsible AI narratives from ethics to safety, and from harm to risk. This is worrying for reasons to do with the so-called 'AI power paradox'.

I argue that in fact, AI practitioners and policy-makers should be focused on the 'why' of AI ethics, which implies interpreting AI ethics as a dynamic reasoning system focused on building resilience against

potential harm from engagement with AI technology, and therefore, ultimately focused on determining the conditions for establishing societies that can thrive in their engagement with AI technology, rather than being focused on regulating the technology.



Bio: Emma Ruttkamp-Bloem is a philosopher of science and technology, an AI ethics policy advisor, and a machine ethics researcher. Currently, she is the Head of the Department of Philosophy, University of Pretoria, and leads the AI ethics group at the South African Centre for AI Research (CAIR). She is the current Chairperson of the UNESCO World Commission on the Ethics of Scientific Knowledge and Technology (COMEST). Emma joined the WEF's Global Future Council on Autonomous Systems in 2025. She is a member of the Global Academic Network, Centre for AI and Digital Policy, Washington DC, and has worked on AI governance projects with the African Union Development Agency and the African Commission on Human and People's Rights. She

was a member of the UN Secretary General's Al Advisory Body 2023-2024. Emma led the UNESCO Ad Hoc Expert Group that prepared the draft of the 2021 UNESCO Recommendation on the Ethics of AI and contributed to development of its implementation instruments. Emma continues to work with UNESCO as a member of UNESCO' AI Ethics without Borders and Women4EthicalAI initiatives. She is a member of various international AI ethics advisory boards ranging from academia (e.g., the Wallenberg AI, Autonomous Systems and Software Programme Human Sciences), the inter-governmental sector (e.g., as expert advisory board member for the Global Commission on Responsible Artificial Intelligence in the Military Domain), to the private sector (e.g., SAP SE). She is an associate editor for the Journal of Science and Engineering Ethics, and a member of the editorial board of the Journal of AI Law and Regulation and the Cambridge Forum on AI: Law and Governance.

Lunch Break (on own) 2:00-3:15

3:15-4:45 Poster Session 4 & Refreshment Break // Floor -4

Posters

80 The Politics of AI Systems are Inextricable from Their Supply Chains: Public Values Versus the Digital Political Economy

Ben Gansky

99 A Framework for Situating Innovations, Opportunities, and Challenges in Advancing

Vertical Systems with Large AI Models

Gaurav Verma, Jiawei Zhou, Mohit Chandra, Srijan Kumar, Munmun De Choudhury

212 Think Outside the Data: Colonial Biases and Systemic Issues in Automated Moderation Pipelines for

Low-Resource Languages

Farhana Shahid, Mona Elswah and Aditya Vashistha

246 Enhancing Image Comprehension: The Impact of AI-Generated Explanations on Perception of Altered

and Synthetic Media

Saquib Ahmed, Tejo Gayathri Busireddy and Sanorita Dey

461	From explaining to diagnosing: a justice-oriented framework of explainable AI for bias detection <i>Miriam Fahimi, Laura State and Atoosa Kasirzadeh</i>
736	Ethical Classification of Non-Coding Contributions in Open-Source Projects via Large Language Models Sergio Cobos and Javier Luis Canovas Izquierdo
738	Mimetic AI Systems: Understanding and Regulating the Use of Generative Models for Impersonation Norman Bukingolts
747	Automating Data Governance with Generative Al Linus W. Dietz, Arif Wider and Simon Harrer
755	The AI Model Risk Catalog: What Developers and Researchers Miss About Real-World AI Harms <i>Pooja S. B. Rao, Sanja Šcepanovic, Dinesh Babu Jayagopi, Mauro Cherubini and Daniele Quercia</i>
756	Responsible AI Governance in the Public Sector: Explaining Contextual Dynamics through a Realist Synthesis Review Ana Gagua, Haiko van der Voort, Nihit Goyal and Alexander Verbraeck
760	PETLP: A Privacy-by-Design Pipeline for Social Media Data in AI Research Nick Oh, Giorgos Vrakas, Siân Brooke, Sasha Morineire and Toju Duke
762	A Longitudinal Randomized Control Study of Companion Chatbot Use: Anthropomorphism and Its Mediating Role on Social Impacts Rose E. Guingrich and Michael S.A. Graziano
763	Addressing Bias in LLMs: Strategies and Application to Fair AI-based Recruitment Alejandro Peña Almansa, Julian Fierrez Aguilar, Aythami Morales Moreno, Gonzalo Mancera, Miguel Lopez-Duran and Ruben Tolosana
764	Al Safety and Security Enable Innovation in Emerging Economies Joanna Wiaterek and Jared Perlo
767	Adoption of Explainable Natural Language Processing: Perspectives from Industry and Academia on Practices and Challenges Mahdi Dhaini, Tobias Müller, Roksoliana Rabets and Gjergji Kasneci
768	Empirical Analysis of Privacy-Fairness-Accuracy Trade-offs in Federated Learning: A Step Towards Responsible Al Dawood Wasif, Dian Chen, Sindhuja Madabushi, Nithin Alluru, Terrence Moore and Jin-Hee Cho
769	When Explainability Meets Privacy: An Investigation at the Intersection of Post-hoc Explainability and Differential Privacy in the Context of Natural Language Processing Mahdi Dhaini, Stephen Meisenbacher, Ege Erdogan, Florian Matthes and Gjergji Kasneci
775	The Transparency Dilemma: An Experiment on How AI Disclosures Affect Credibility Perceptions and Engagement Across Topics Sophie Morosoli, Emma van der Goot, Valeria Resendez, Claes de Vreese and Natali Helberger
777	What Are Chatbots' Stereotypes About? A Data-Driven Analysis of Large Language Models' Content Associations with Social Categories Gandalf Nicolas and Aylin Caliskan

795	Matters of Explanation: Rethinking Explainability with Tangible, Embodied, Material Interactions Goda Klumbyte and Claude Draude
796	Epistemic Destabilization: AI-Driven Knowledge Generation and the Collapse of Validation Systems Bhavneet Singh
798	"One of Silicon Valley's Most Divisive Topics": How the Media Discusses Openness in Al Tamara Paris, Jin L.C. Guo and Ajung Moon
799	Measuring What Matters: Connecting Al Ethics Evaluations to System Attributes, Hazards, and Harms Shalaleh Rismani, Renee Shelby, Leah Davis, Negar Rostamzadeh and Ajung Moon
800	"Accessibility people, you go work on that thing of yours over there": Addressing Disability Inclusion in AI Product Organizations Sanika Moharana, Cynthia Bennett, Erin Buehler, Michael Madaio, Vinita Tibdewal and Shaun Kane
805	Bias Amplification in Stable Diffusion's Representation of Stigma Through Skin Tones and Their Homogeneity Kyra Wilson, Sourojit Ghosh and Aylin Caliskan
809	Whose Personae? Synthetic Persona Experiments in LLM Research and Pathways to Transparency Jan Batzner, Volker Stocker, Bingjun Tang, Anusha Natarajan, Qinhao Chen, Stefan Schmid and Gjergji Kasneci
813	Co-producing Al: Toward an Augmented, Participatory Lifecycle Rashid Mushkani, Hugo Berard, Toumadher Ammar, Cassandre Chatonnier and Shin Koseki
815	IndiCASA: A Dataset and Bias Evaluation Framework for LLMs Using Contrastive Embedding Similarity in the Indian Context Santhosh G S, Akshay Govind S, Gokul S Krishnan, Balaraman Ravindran and Sriraam Natarajan
817	Explanation Difference: Bridging Procedural and Distributional Fairness Joe Germino, Yuying Zhao, Tyler Derr, Nuno Moniz and Nitesh V. Chawla
823	No Such Thing as Free Brain Time: For a Pigouvian Tax on Attention Capture Hamza Belgroun, Franck Michel and Fabien Gandon
828	Models and Algorithms for Balancing Efficiency and Equity in Vaccine Allocation Jeffrey Keithley, Madeline Bonner and Sriram V. Pemmaraju
834	Emergent AI Surveillance: Overlearned Person Re-Identification and Its Mitigation in Law Enforcement Context An Thi Nguyen, Radina Stoykova and Eric Arazo
844	Stop the Non-consensual Use of Intimate Images in Research Princessa Cintaqia, Arshia Arya, Elissa M. Redmiles, Deepak Kumar, Allison McDonald and Lucy Qin
863	When the Past Misleads: Rethinking Training Data Expansion Under Temporal Distribution Shifts Chengyuan Yao, Yunxuan Tang, Christopher Brooks, Rene Kizilcec and Renzhe Yu
870	Al and the Social Contract Chee Hae Chung and Daniel Schiff

876	Designing Stakeholder-Based Pedagogy for AI Ethics Education: Insights from a Multi-Institutional Case Study
	Yifan Zhang, Harini Suresh and Julia Netter
879	Reflective Agency: Ethical and Empirical Framework for AI-Mediated Self-Reflection Systems Minsol Kim, Wendy Wang, Rosalind Picard, Nathan Barczi, Jennifer Long and Pattie Maes
882	Laissez-Faire Harms: Algorithmic Biases in Generative Language Models Evan Shieh, Faye Marie Vassel, Cassidy Sugimoto and Thema Monroe-White
884	Understanding Privacy Norms Around LLM-Based Chatbots: A Contextual Integrity Perspective Sarah Tran, Hongfan Lu, Isaac Slaughter, Bernease Herman, Aayushi Dangol, Yue Fu, Lufei Chen, Biniyam Gebreyohannes, Bill Howe, Alexis Hiniker, Nicholas Weber and Robert Wolfe
886	Do students rely on Al? Analysis of student-ChatGPT conversations from a field study Jiayu Zheng, Lingxin Hao, Kelun Lu, Ashi Garg, Mike Reese, Melo-Jean Yap, I-Jeng Wang, Xingyun Wu, Wenrui Huang, Jenna Hoffman, Ariane Kelly, My Le, Ryan Zhang, Yanyu Lin, Muhammad Faayez and Anqi Liu
887	PRAC3 (Privacy, Reputation, Accountability, Consent, Credit, Compensation): Voice Actors in AI Data-economy Tanusree Sharma, Yihao Zhou and Visar Berisha
896	AI-OCI: A Novel Framework for Assessing AI's Workforce Impact Using LLMs Frederick Awuah-Gyasi and Trilce Estrada
898	Agents without Agency: Anthropological and Sociological Lessons for Contemporary Al Research and Policy Greta Timaite and Michael Castelle
900	Toward Needs-Conscious Design: Co-Designing a Human-Centered Framework for Al-Mediated Communication Robert Wolfe, Aayushi Dangol, Jaewon Kim and Alexis Hiniker
901	From Efficiency to Equity: Measuring Fairness in Preference Learning Shreeyash Gowaikar, Hugo Berard, Rashid Mushkani and Shin Koseki
906	Responsible AI in the OSS: Reconciling Innovation with Risk Assessment and Disclosure Mahasweta Chakraborti, Bert Prestoza, Nicholas Vincent, Vladimir Filkov and Seth Frey
909	Invisible Filters: Cultural Bias in Hiring Evaluations Using Large Language Models Pooja S. B. Rao, Laxminarayen Nagarajan Venkatesan, Mauro Cherubini and Dinesh Babu Jayagopi
924	Making Bodies: Assumptions in the Design and Validation of Motion Capture Technology Mona Sloane, Abigail Jacobs and Emanuel Moss
	Student Posters
10	On the Computational, Informational, and Physical Foundations for AI Safety Robin Young
26	Predictability in Autonomous Driving Systems Felix Marti-Perez

76	Precedent-Based Professional Role Ethics for AI Decision Analysis Cristopher Rauch
45	Safe and Explainable Machine Learning for High-Impact Decisions Khadija Zanna
51	Algorithmic Recourse Kshitij Kayastha
30	Bridging Liability Gaps in the Age of AI: The Case for No-Fault Compensation Schemes <i>Chi Ha Tran</i>
33	To Regulate Or To Be Regulated? – How to Protect the 'Freedom of Mind' Against Emotional AI Surveillance Öznur Uguz
65	Improvising Accountability: The Everyday Governance Work of Responsible AI in the Public Sector Ana Gagua
66	FATE-Compliant ML Architecture with Blockchain-Verifiable Auditing: A Governance Framework for Ethical Compliance in FinTech Samah Kareem
17	Al Contextual Framework: A Zoning Approach to Ethical Al Deployment Yao Xie
4:45-6:00	Paper Session 7: Risk Management // Auditorium
71	Bubble, Bubble, Al's Rumble: Why Global Financial Regulatory Incident Re-porting is Our Shield Against Systemic Stumbles Anchal Gupta; Gleb Papyshev; James T. Kwok
183	A Closer Look at the Existing Risks of Generative AI: Mapping the Who, What, and How of Real-World Incidents Megan Li; Wendy Bickersteth; Ningjing Tang; Lorrie Cranor; Jason Hong; Hong Shen; Hoda Heidari
314	The Model Hears You: Audio Language Model Deployments Should Consider the Principle of Least Privilege Luxi He; Xiangyu Qi; Michel Liao; Inyoung Cheong; Prateek Mittal; Danqi Chen; Peter Henderson
641	When in Doubt, Cascade: Towards Building Efficient and Capable Guardrails Manish Nagireddy; Inkit Padhi; Soumya Ghosh; Prasanna Sattigeri
6:00-6:30	Closing Remarks // Auditorium

AIES-25 SPONSORS

AIES 2025 would like to thank the generous sponsors who allowed us to support Ph.D. students, invited speakers, social events, and to reduce the registration fee.

PLATINUM



GOLD



SILVER



BRONZE





JPMorganChase

ORGANIZED BY





IN PARTNERSHIP WITH



GENEROUS FINANCIAL SUPPORT FOR THE STUDENT PROGRAM PROVIDED BY



