

# A. STEVIE BERGMAN

Public Interest Technologist & Scientist

✉ [brown.edu](mailto:asteviebergman@brown.edu)

🌐 [asteviebergman.com](http://asteviebergman.com)

📍 New York City

🌐 [/asteviebergman](https://www.linkedin.com/in/asteviebergman)

## SUMMARY

Public interest scientist and multidisciplinary research expert with years of experience in the tech industry, government, and a commitment to human rights and public service. Experienced research and team lead, as well as people manager. Areas of expertise include AI ethics, governance and tech policy, democratic and participatory AI, implications of generative AI, algorithmic fairness, mis/disinformation, AI evaluation, representation, multilingual AI, and data annotation, labor, and curation practices. Proficient educator, with diverse interdisciplinary experience in areas such as legal work and international aid, physics research, and science & tech communication.

## CURRENT POSITIONS

### Head of Applied Systems

Center for AI Standards and Innovation (CAISI)

11/2024 – Present

National Institute of Standards and Technology  
Department of Commerce, US Government

Created, built, and manage a team of researchers focused on measuring AI systems in application, post-deployment, and in adoption. This role is on the CAISI leadership team, and both a research and people manager position.

### Visiting Assistant Professor

Center for Technological Responsibility, Reimagination, and Redesign

10/2024 – Present

Fall 2025: *Forces of influence in AI governance* (DATA 2030) graduate course

Data Science Institute, Brown University

## PRIOR ROLES

4/2022 – 10/2024

### Senior Research Scientist, Ethics Research Team

Google DeepMind

Rigorous multidisciplinary research employing mixed methods on complex, open-ended problems across a host of interdisciplinary areas including: international technology policy and ethical AI development, democratic and participatory AI, AI evaluation, representation and fairness in AI, the impact of emerging technologies on non-Western, at-risk or in-conflict communities, societal impacts of generative AI (eg., mis- and disinformation), and the implications of design choices (eg., data curation and AI annotation with fair labor practices). Role involved,

- Leading teams of colleagues (research scientists, engineers, etc) on complex, multi-stream research projects with several dependent deliverables and high-level buy-in,
- Expert consultation to legal and policy teams,
- Internal partnership with product and research teams across Google,
- External collaborations and engagements, such as close collaboration with the National Institute of Standards and Technology, providing feedback on the AI Risk Management Framework (2022-24), [fire-side chats](#), dialogues, workshops, and panels on [participation in](#) and governance of AI systems, [representation in AI](#), as well as the many [risks of disinformation from generative AI](#). See [Public Engagements](#).

1/2020 – 4/2022

### Research Scientist, AI Fairness, Responsible AI

Meta

Practical application of AI ethics principles to Meta products with a keen focus on countries and communities outside the West, and at-risk for offline harm. Flagship research focused on Arabic dialects, hateful and violent speech, annotation/labor practices (see: *Towards Responsible Natural Language Annotation for Varieties of Arabic*), and the dynamics technology design choices foment in the region. Additional research on African languages and Portuguese. Further work in this role included:

- Development of novel, responsible research methodologies to address open-ended questions and recommend high-impact mitigation strategies.
- Utilization of interdisciplinary research methods including qualitative & quantitative data analysis techniques, systematic root cause analysis, robust, ethical, and inclusive data collection & representative sampling, and surveys.
- Production of many internal research reports on these investigations & their impacts.
- Technical subject-matter expert on the Responsible AI Education team, producing and delivering clear, high-impact education projects, driving adoption across the company.
- Creating and co-leading weekly AI Ethics Foundational Reading & Discussion group.

Spring & Fall 2021

### Lecturer, Ethics of AI Graduate Seminar

Princeton University

Seminar speaker on the practical application of AI ethics with a focus on methods for accountability & the complications involved with measurements for benchmarking. After taking this seminar in 2019, was invited back as a lead to teach on the topic.

Spring 2019

### Creator, Producer, and Host of AI & Human Rights Podcast Miniseries

Princeton University

Educational podcast commissioned by Princeton's Center for Information & Technology Policy, on the intersection of digital technologies and human rights, framed by a CITP/UN conference in April 2019. Required research into the ethical implications of emerging technologies and a grasp of technology policy and its implications for society. Available on [Apple](#), [Spotify](#), & [Soundcloud](#).

7/2015 – 10/2019	<b>Creator, Co-Host, Co-Producer of These Vibes Are Too Cosmic</b> Weekly, prime time science & music radio show featuring science news, local events, and live, longform interviews with experts from across the sciences. Full list of 100+ interviews available for streaming on <a href="http://thesevibesaretoocosmic.com">thesevibesaretoocosmic.com</a> .	<b>WPRB Princeton</b>
10/2013 – 7/2014	<b>Fulbright Fellow, Indonesia</b> Independent research in theoretical high energy physics (complex scalar field dynamics with BPS domain walls), with university collaborator.	<b>Institut Teknologi Bandung</b>
Summer 2013	<b>Research Assistant, Large Hadron Collider at CERN</b> Testing & commissioning the inner B-Layer of the ATLAS Pixel Detector in Geneva, Switzerland.	<b>Columbia University</b>
2/2011 – 4/2013	<b>US Peace Corps, Science Education Volunteer</b> Mathematics, physics, and computer teacher at Lacor Secondary School and Warr Girls. Creator & Co-director of <a href="#">GirlTech Uganda</a> and Secondary Education PCV trainer.	<b>Uganda</b>
8/2009 – 12/2010	<b>Vaccine Litigation Paralegal</b> Assisted attorneys in understanding the science involved in their work, as well as research, documentation, and argument for cases at the US Court of Federal Claims.	<b>US Department of Justice</b>
<b>Other Projects</b>		
2020	Flatbush Mutual Aid, New York City, NY	
2017	Princeton Citizen Scientists, Founding member and Day of Action coordinator, Princeton University	
2016 – 2019	Prison Teaching Initiative, Weekly Tutor, Garden State Prison, NJ	
2018 – 2019	Physics Dept. Climate and Inclusion Committee, Princeton University	

## PUBLICATIONS IN TECH & SOCIETY

---

AIES 2024	<i>Gaps in the Safety Evaluation of Generative AI</i> Maribeth Rauh, Nahema Marchal, Arianna Manzini, Lisa Anne Hendricks, Ramona Comanescu, Canfer Akbulut, Tom Stepleton, Juan Mateos-Garcia, <b>Stevie Bergman</b> , Jackie Kay, Conor Griffin, Ben Bariach, Iason Gabriel, Verena Rieser, William Isaac, Laura Weidinger <a href="http://aaai.org/index.php/AIES/.../31717">aaai.org/index.php/AIES/.../31717</a> This paper is a sub-publication of <i>Sociotechnical Safety Evaluation of Generative AI Systems</i> , listed below.
EMNLP 2024	<i>STAR: SocioTechnical Approach to Red Teaming Language Models</i> Laura Weidinger, John Mellor, Bernat Guillén Pegueroles, Nahema Marchal, Ravin Kumar, Kristian Lum, Canfer Akbulut, Mark Diaz, <b>Stevie Bergman</b> , Mikel Rodriguez, Verena Rieser, William Isaac <a href="http://aclanthology.org/2024.emnlp-main.1200">aclanthology.org/2024.emnlp-main.1200</a>
Nature Scientific Reports 2024	<i>STELA: A Community Centred Approach to Norm Elicitation for Agent Alignment</i> <b>Stevie Bergman</b> , Nahema Marchal, John Mellor, Shakir Mohamed, Iason Gabriel, William Isaac <a href="http://nature.com/articles/s41598-024-56648-4">nature.com/articles/s41598-024-56648-4</a>
Released Apr 2024 Future submissions for publication	<i>Ethics of Advanced Assistants: Access and Opportunity</i> (Chapter) <b>A. Stevie Bergman</b> , Renee Shelby, Iason Gabriel Full paper: <a href="#">the-ethics-of-advanced-ai-assistants-2024-i.pdf</a> Blog post: <a href="http://blog/the-ethics-of-advanced-ai-assistants">blog/the-ethics-of-advanced-ai-assistants</a>
CHI 2024	<i>The Illusion of Artificial Inclusion</i> William Agnew, <b>A. Stevie Bergman</b> , Jennifer Chien, Mark Diaz, Seliem El-Sayed, Jaylen Pittman, Shakir Mohamed, Kevin R. McKee <a href="http://arxiv.org/pdf/2401.08572.pdf">arxiv.org/pdf/2401.08572.pdf</a> This paper was a <i>CHI 2024 Editors' Choice</i> .
arXiv 2023 Two submissions for publication	<i>Sociotechnical Safety Evaluation of Generative AI Systems</i> Laura Weidinger, Maribeth Rauh, Nahema Marchal, Arianna Manzini, Lisa Anne Hendricks, Juan Mateos-Garcia, <b>Stevie Bergman</b> , Jackie Kay, Conor Griffin, Ben Bariach, Iason Gabriel, Verena Rieser and William Isaac Full paper: <a href="http://arxiv.org/pdf/2310.11986.pdf">arxiv.org/pdf/2310.11986.pdf</a> Blog post: <a href="http://blog/evaluating-social-and-ethical-risks-from-generative-ai">blog/evaluating-social-and-ethical-risks-from-generative-ai</a>
ACM FAccT 2023	<i>Representation in AI Evaluation</i> <b>A. Stevie Bergman</b> , Lisa Anne Hendricks, Maribeth Rauh, Boxi Wu, William Agnew, Markus Kunesch, Isabella Duan, Iason Gabriel, William Isaac <a href="http://dl.acm.org/doi/abs/10.1145/3593013.3594019">dl.acm.org/doi/abs/10.1145/3593013.3594019</a>
ACL 2022	<i>Towards Responsible Natural Language Annotation for Varieties of Arabic</i> <b>A. Stevie Bergman</b> & Mona Diab <a href="http://aclanthology.org/2022.findings-acl.31">aclanthology.org/2022.findings-acl.31</a>
ACM FAccT 2022	<i>Adaptive Sampling Strategies to Construct Equitable Training Datasets</i> William Cai, Ro Encarnacion, Bobbie Chern, Sam Corbett-Davies, Miranda Bogen, <b>Stevie Bergman</b> , Sharad Goel <a href="http://dl.acm.org/doi/abs/10.1145/3531146.3533203">dl.acm.org/doi/abs/10.1145/3531146.3533203</a>

SIGDIAL 2022	<i>Guiding the Release of Safer E2E Conversational AI through Value Sensitive Design</i> <b>A. Stevie Bergman</b> , Gavin Abercrombie, Shannon Spruit, Dirk Hovy, Emily Dinan, Y-Lan Boureau, Verena Rieser <a href="https://iris.unibocconi.it/bitstream/11565/4053244/1/2022.sigdial-1.4.pdf">iris.unibocconi.it/bitstream/11565/4053244/1/2022.sigdial-1.4.pdf</a>
ACL 2022	<i>SafetyKit: First Aid for Measuring Safety in Open-domain Conversational Systems</i> Emily Dinan, Gavin Abercrombie, <b>A. Stevie Bergman</b> , Shannon Spruit, Dirk Hovy, Y-Lan Boureau, Verena Rieser <a href="https://iris.unibocconi.it/bitstream/11565/4053224/1/2022.acl-long.284.pdf">iris.unibocconi.it/bitstream/11565/4053224/1/2022.acl-long.284.pdf</a>
arXiv 2021	<i>Fairness On The Ground: Applying Algorithmic Fairness Approaches to Production Systems</i> Chloè Bakalar, Renata Barreto, <b>Stevie Bergman</b> , Miranda Bogen, Bobbie Chern, Sam Corbett-Davies, Melissa Hall, Isabel Kloumann, Michelle Lam, Joaquin Quiñero Candela, Manish Raghavan, Joshua Simons, Jonathan Tannen, Edmund Tong, Kate Vredenburg, Jiejing Zhao <a href="https://arxiv.org/pdf/2103.06172v2.pdf">arxiv.org/pdf/2103.06172v2.pdf</a>

Full list of journal articles available on [Google Scholar](https://scholar.google.com/), including publications in physics.

## SELECT PUBLIC ENGAGEMENTS

---

February 2025 Guest Lecturer	<b>Forces of Influence in AI</b> Guest lecture in the Philosophy of Computing course at Barnard University on the "the forces of influence in AI" – an ecosystem-view of the many different actors influencing AI governance across the globe.	<b>Philosophy of Computing, Barnard University</b>
May 2024 Keynote Speaker	<b>Third Workshop on Safety for E2E Conversational AI</b> Workshop focusing on the contextual and culturally dependent risks with end-to-end open-domain dialogue agents, such as not being able to respond reliably in safety critical situations. Keynote will be on the topic of culturally situated and sociotechnical evaluations of these models, as well as effective representation and community inclusion.	<b>LREC-COLING 2024 Conference</b>
April 2024 Workshop	<b>Operationalizing the NIST AI Risk Management Framework: A Diversity of Human Factors</b> Invitation-only workshop bringing together experts to build on NIST's AI Risk Management Framework to work towards developing operational, reproducible, and scalable evaluation tools and standards for mapping and assessing critical aspects of human involvement throughout the lifecycle of AI systems.	<b>NIST &amp; Northeastern University Ethics Institute</b>
November 2023 Presentation	<b>AI Governance, Evaluations, and Representation</b> Invited presentation to CDT non-resident fellows on Google DeepMind policy planning, and sociotechnical research with policy implications.	<b>Center for Democracy &amp; Technology</b>
October 2023 Fireside Chat	<b>"Risk Management (or Measurement), for whom, why and how?"</b> Fireside chat with Zachary Lipton, facilitated by Nicol Turner Lee of the Brookings Institution, at the <a href="#">NIST Workshop on Operationalizing the Measure Function of the AI Risk Management Framework</a> . Organized by Northwestern Center for Advancing Safety of Machine Intelligence (CASMI).	<b>NIST &amp; CASMI</b>
October 2023 Workshop	<b>Operationalizing the Measure Function of the NIST AI RMF</b> Two day, invite-only workshop to discuss the complexities of putting the sociotechnical principles outlined in the AI Risk Management Framework into practice, organized by Northwestern Center for Advancing Safety of Machine Intelligence (CASMI).	<b>NIST &amp; CASMI</b>
July 2023 Panel	<b>Potential Benefits and Risks of Generative AI</b> Panelist alongside representatives from McKinsey and Meta, fielding audience questions on the growing risks (and potential benefits) of generative AI – particularly mis/disinformation risks and threats to democracy. Full recording can be found <a href="#">here</a> .	<b>Cambridge University, Disinformation Summit</b>
2023 (ongoing) Dialogue	<b>Public Technology Leadership Collaborative (PTLC)</b> A collaborative group of civil servants, industry technologists, and academics, exchange insights for responsible development and governance of AI technologies.	<b>Data &amp; Society</b>
June 2023 Dialogue	<b>Discussions between AI ethics and AI safety researchers</b> Two part dialogue alongside the 2023 FAccT conference, between sociotechnical AI ethics researchers (who tend to engage with current risks of AI) and AI safety researchers (who tend to be far-long-term engaged).	<b>ACM FAccT 2023</b>
June 2023 Presentation	<b>Representation in AI Evaluations</b> Presentation on the paper <i>Representation in AI Evaluation</i> at the Fairness Accountability and Transparency conference in Chicago, 2023.	<b>ACM FAccT 2023</b>
April 2023 Lecture	<b>Experience AI lesson on Data Bias</b> Conversational lecture on bias in data for the Experience AI lessons for schoolchildren, in collaboration with Raspberry Pi. (Not yet public.)	<b>Raspberry Pi &amp; DeepMind</b>

April 2023 Conference & Workshop	<b>Envisioning equitable representation in AI</b> At this conference, given extended time to present an engaging workshop walking participants through the complexities and pitfalls in thinking through effective representation in AI design, and the consequences of different design choices. This workshop presented and extended FAcCT 2023 publication on <i>Representation in AI Evaluation</i> .	<b>Cambridge, Many Worlds of AI</b>
October 2022 Workshop & Presentation	<b>Representation in AI Annotation</b> Invite-only work-in-progress workshop at which I presented research on considerations, complexities, and implications of representation in AI data annotation.	<b>REAL ML</b>
October 2022 Panel	<b>Building the NIST AI Risk Management Framework (Workshop #3)</b> Panelist discussing participatory & sociotechnical methods for integration into the (at the time) upcoming AI Risk Management Framework from the National Institute of Standards and Technology. Recording can be found <a href="#">here</a> .	<b>NIST</b>
July 2022 Dialogue	<b>Dialogue on an Institution for Assessing Digital Harms</b> Conversation across organizations and disciplines to explore the prospects for a trusted, neutral global institution that could develop agreed definitions and measurements of digital harms.	<b>New America Foundation</b>
October 2022 Workshop	<b>Assessing Social Impacts of General Purpose AI Systems</b> Workshop bringing together academic and industry experts to brainstorm the requirements for effective sociotechnical evaluations of a generative AI systems.	<b>Hugging Face</b>
2019 (ongoing) Dialogue	<b>TechSoc Discussion Group</b> Monthly reading and discussion group on technology and society topics, organized by the Center for Information Technology Policy.	<b>Princeton CITP</b>
2020 (ongoing) Dialogue	<b>Privacy Research Group</b> Weekly discussion group on information law and technology policy topics – typically works in progress from group members – organized by New York University’s Information Law Institute.	<b>NYU ILI</b>

---

## EDUCATION

10/2019	<b>Doctorate in Physics</b> Observational cosmology & experimental physics. NSF Graduate Research Fellowship (2015-18), Joseph Henry Merit Prize (2014) <b>Thesis description:</b> Measuring the statistical polarization in the cosmic microwave background radiation from the early Universe among instrumental, environmental, and astrophysical noise, via constructing the most sensitive radio telescopes ever built, the SPIDER instrument, & flying it on a balloon above Antarctica. <b>Work included:</b> Close teamwork on an experiment with many potential single-point failures, data analysis and experimental research methods.	<b>Princeton University</b>
5/2009	<b>BA Physics, Minor in Astrophysics</b> cum laude, Highest Honors, Phi Beta Kappa, Sigma Xi, Waterman Prize for Outstanding Senior in Physics Year abroad – Oxford University	<b>Smith College</b>

---

## LANGUAGES

Python, SQL, C/C++, some French, some Acholi and Alur, learning Levantine Arabic

---

## REFERENCES

Available upon request.