



THE UNIVERSITY *of* EDINBURGH
informatics



THE UNIVERSITY *of* EDINBURGH
Edinburgh Futures Institute

**Centre for
Technomoral Futures**

On the Futility of Evaluation

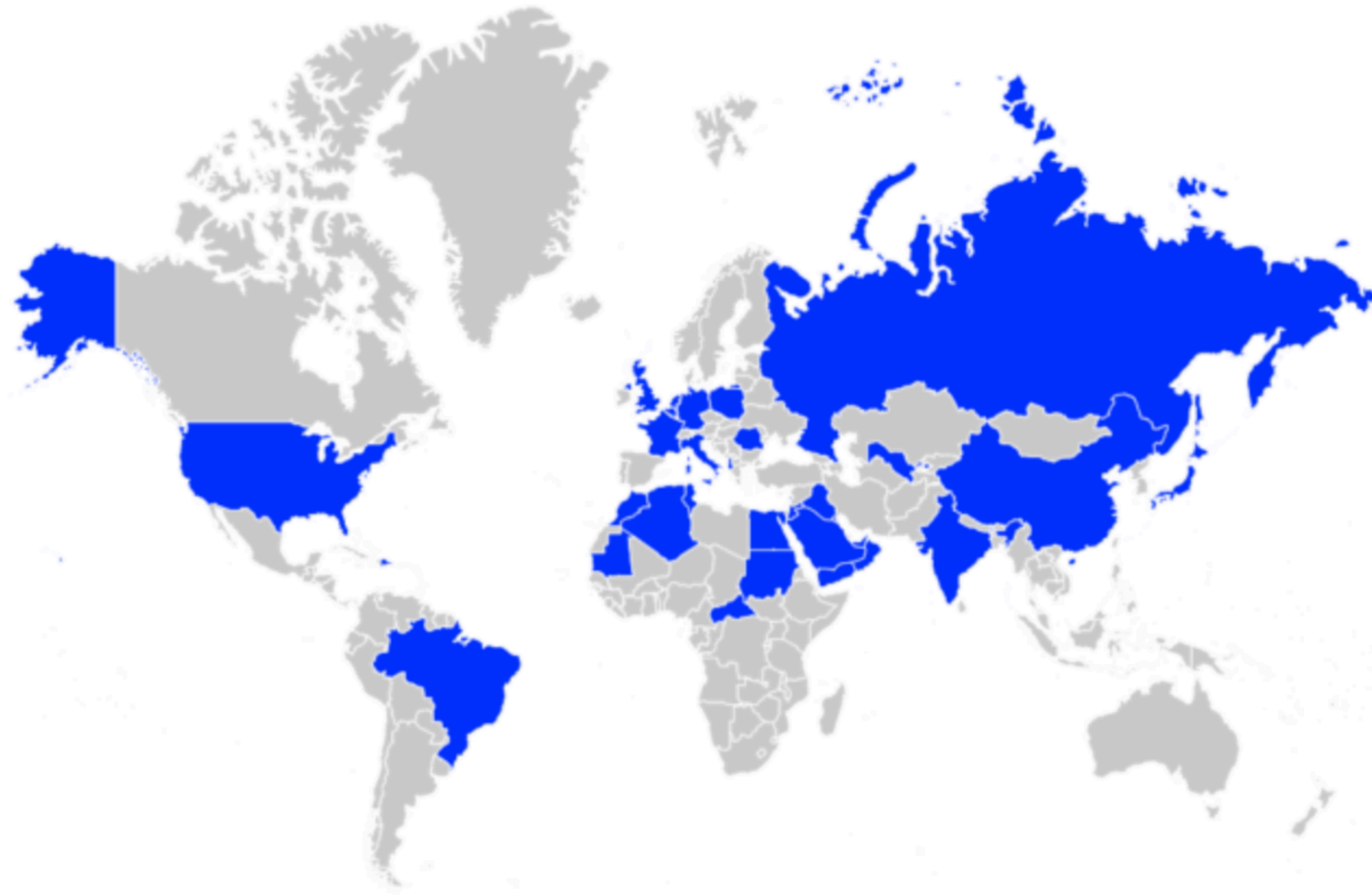
17-Nov-2025

Zeeraak Talat
z@zeeraak.org | [@zeeraaktalat](https://www.instagram.com/zeeraaktalat)
www: zeeraak.org

SHADES: Towards a Multilingual Assessment of Stereotypes in Large Language Models

May 12, 2025

On the Futility of Evals



SHADES: Towards a Multilingual Assessment of Stereotypes in Large Language Models

**Margaret Mitchell¹, Hamdan Al-Ali², Giuseppe Attanasio³, Ioana Baldini⁴,
Miruna Clinciu^{5,6}, Jordan Clive⁷, Pieter Delobelle^{8,45}, Manan Dey⁹,
Kaustubh Dhole¹⁰, Timm Dill¹¹, Amirbek Djanibekov², Tair Djanibekov¹²,
Jad Doughman², Ritam Dutt¹³, Jessica Zosa Forde¹⁴, Jay Gala²,
Avijit Ghosh¹, Sil Hamilton¹⁵, Carolin Holtermann¹¹, Jerry Huang^{16,17},
Lucie-Aimée Kaffee¹, Janavi Kasera¹⁸, Tanmay Laud^{19,20}, Anne Lauscher¹¹,
Roberto Luis López²¹, Jonibek Mansurov², Maraim Masoud²², Sagnik Mukherjee²³,
Nurdaulet Mukhituly², Nikita Nangia²⁴, Shangrui Nie²⁵, Anaelia Ovalle²⁶, Giada Pistilli¹,
Esther Ploeger²⁷, Jeremy Qin^{16,17,28}, Dragomir Radev²⁹, Vipul Raheja³⁰, Beatrice Savoldi³¹,
Shanya Sharma³², Xudong Shen³³, Karolina Stańczak^{16,34}, Arjun Subramonian²⁶,
Kaiser Sun³⁵, Eliza Szczechla³⁶, Tiago Timponi Torrent^{37,38}, Deepak Tunuguntla³⁹,
Emilio Villa-Cueva², Marcelo Viridiano⁴⁰, Oskar van der Wal⁴¹, Adina Yakefu¹,
Kayo Yin⁴², Mike Zhang²⁷, Sydney Zink⁴³, Aurélie Névéal⁴⁴, Zeerak Talat⁶**

**You Reap What You Sow:
On the Challenges of Bias Evaluation Under Multilingual Settings**

**Zeerak Talat¹, Aurélie Névéol², Stella Biderman^{3,4}, Miruna Clinciu^{5,6,7}, Manan Dey⁸,
Shayne Longpre⁹, Alexandra Sasha Luccioni¹⁰, Maraim Masoud¹¹, Margaret Mitchell¹⁰,
Dragomir Radev¹², Shanya Sharma¹³, Arjun Subramonian^{14,15}, Jaesung Tae^{10,12},
Samson Tan^{16,17}, Deepak Tunuguntla¹⁸, Oskar van der Wal¹⁹**

¹Digital Democracies Institute, Simon Fraser University ²Université Paris-Saclay, CNRS, LISN
³Booz Allen Hamilton ⁴EleutherAI ⁵Edinburgh Centre for Robotics ⁶Heriot-Watt University
⁷University of Edinburgh ⁸SAP ⁹MIT ¹⁰Hugging Face ¹¹Adapt Centre, Trinity College Dublin
¹²Yale University ¹³Walmart Labs, India ¹⁴University of California, Los Angeles ¹⁵Queer-in-AI
¹⁶AWS AI Research and Education ¹⁷National University of Singapore ¹⁸Independent
Researcher ¹⁹University of Amsterdam

**Evaluating the Social Impact of Generative AI Systems
in Systems and Society**

Irene Solaiman^{1*}	Zeerak Talat^{2*}	William Agnew³	Lama Ahmad⁴
Dylan Baker⁵	Su Lin Blodgett⁶	Canyu Chen⁷	Hal Daumé III⁸
Jesse Dodge⁹	Isabella Duan¹⁰	Felix Friedrich^{11,12}	Avijit Ghosh¹
Usman Gohar¹³	Sara Hooker¹⁴	Yacine Jernite¹	Ria Kalluri¹⁵
Alberto Lusoli¹⁶	Alina Leidinger¹⁷	Michelle Lin^{18,19}	Xiuzhu Lin¹¹
Sasha Luccioni¹	Jennifer Mickel²¹	Margaret Mitchell¹	Jessica Newman²²
Anaelia Ovalle²²	Marie-Therese Png²³	Shubham Singh²⁴	Andrew Strait²⁵
	Lukas Struppek^{11,26}	Arjun Subramonian²²	

¹Hugging Face, ²Mohamed Bin Zayed University of Artificial Intelligence, ³Carnegie Mellon University, ⁴OpenAI, ⁵DAIR, ⁶Microsoft Research, ⁷Illinois Institute of Technology, ⁸University of Maryland, ⁹Allen Institute for AI, ¹⁰University of Chicago, ¹¹TU Darmstadt, ¹²hessian.AI, ¹³Iowa State University, ¹⁴Cohere for AI, ¹⁵Stanford University, ¹⁶Simon Fraser University, ¹⁷University of Amsterdam, ¹⁸Mila - Quebec AI Institute, ¹⁹McGill University, ²¹University of Texas at Austin, ²¹University of California, Berkeley, ²²University of California, Los Angeles, ²³Oxford University, ²⁴University of Illinois Chicago, ²⁵Ada Lovelace Institute, ²⁶DFKI

**You Reap What You Sow:
On the Challenges of Bias Evaluation Under Multilingual Settings**

**Zeerak Talat¹, Aurélie Névéol², Stella Biderman^{3,4}, Miruna Clinciu^{5,6,7}, Manan Dey⁸,
Shayne Longpre⁹, Alexandra Sasha Luccioni¹⁰, Maraim Masoud¹¹, Margaret Mitchell¹⁰,
Dragomir Radev¹², Shanya Sharma¹³, Arjun Subramonian^{14,15}, Jaesung Tae^{10,12},
Samson Tan^{16,17}, Deepak Tunuguntla¹⁸, Oskar van der Wal¹⁹**

¹Digital Democracies Institute, Simon Fraser University ²Université Paris-Saclay, CNRS, LISN
³Booz Allen Hamilton ⁴EleutherAI ⁵Edinburgh Centre for Robotics ⁶Heriot-Watt University
⁷University of Edinburgh ⁸SAP ⁹MIT ¹⁰Hugging Face ¹¹Adapt Centre, Trinity College Dublin
¹²Yale University ¹³Walmart Labs, India ¹⁴University of California, Los Angeles ¹⁵Queer-in-AI
¹⁶AWS AI Research and Education ¹⁷National University of Singapore ¹⁸Independent
Researcher ¹⁹University of Amsterdam

**Evaluating the Social Impact of Generative AI Systems
in Systems and Society**

Irene Solaiman^{1*}	Zeerak Talat^{2*}	William Agnew³	Lama Ahmad⁴
Dylan Baker⁵	Su Lin Blodgett⁶	Canyu Chen⁷	Hal Daumé III⁸
Jesse Dodge⁹	Isabella Duan¹⁰	Felix Friedrich^{11,12}	Avijit Ghosh¹
Usman Gohar¹³	Sara Hooker¹⁴	Yacine Jernite¹	Ria Kalluri¹⁵
Alberto Lusoli¹⁶	Alina Leidinger¹⁷	Michelle Lin^{18,19}	Xiuzhu Lin¹¹
Sasha Luccioni¹	Jennifer Mickel²¹	Margaret Mitchell¹	Jessica Newman²²
Anaelia Ovalle²²	Marie-Therese Png²³	Shubham Singh²⁴	Andrew Strait²⁵
	Lukas Struppek^{11,26}	Arjun Subramonian²²	

¹Hugging Face, ²Mohamed Bin Zayed University of Artificial Intelligence, ³Carnegie Mellon University, ⁴OpenAI, ⁵DAIR, ⁶Microsoft Research, ⁷Illinois Institute of Technology, ⁸University of Maryland, ⁹Allen Institute for AI, ¹⁰University of Chicago, ¹¹TU Darmstadt, ¹²hessian.AI, ¹³Iowa State University, ¹⁴Cohere for AI, ¹⁵Stanford University, ¹⁶Simon Fraser University, ¹⁷University of Amsterdam, ¹⁸Mila - Quebec AI Institute, ¹⁹McGill University, ²¹University of Texas at Austin, ²¹University of California, Berkeley, ²²University of California, Los Angeles, ²³Oxford University, ²⁴University of Illinois Chicago, ²⁵Ada Lovelace Institute, ²⁶DFKI



RESEARCH COMMUNITY

EvalEval Coalition

We are a research community developing scientifically grounded **research outputs** and robust **deployment infrastructure** for **broader impact evaluations**.

WHO EVALUATES AI’S SOCIAL IMPACTS? MAPPING COVERAGE AND GAPS IN FIRST AND THIRD PARTY EVALUATIONS

Anka Reuel^{1,*}, Avijit Ghosh^{2,*}, Jenny Chim^{3,*}



Andrew Tran^{4,◇}, Yanan Long^{5,◇}, Jennifer Mickel^{6,◇}, Usman Gohar^{7,◇},
Srishti Yadav^{8,◇}, Pawan Sasanka Ammanamanchi^{4,◇}

Mowafak Allaham⁹, Hossein A. Rahmani¹⁰, Mubashara Akhtar¹¹, Felix Friedrich¹², Robert Scholz¹³,
Michael Alexander Riegler¹⁴, Jan Batzner^{15,16}, Eliya Habba¹⁷, Arushi Saxena¹⁸, Anastassia Kornilova¹⁹,
Kevin Wei²⁰, Prajna Soni²¹, Yohan Mathew⁴, Kevin Klyman¹, Jeba Sania²⁰, Subramanyam Sahoo²²,
Olivia Beyer Bruvik¹, Pouya Sadeghi²³, Sujata Goswami²⁴, Angelina Wang²⁵,

Yacine Jernite^{2,†}, Zeerak Talat^{26,†}, Stella Biderman^{6,†},
Mykel Kochenderfer^{1,†}, Sanmi Koyejo^{1,†}, Irene Solaiman^{2,†}

ABSTRACT

Foundation models are increasingly central to high-stakes AI systems, and governance frameworks now depend on evaluations to assess their risks and capabilities. Although general capability evaluations are widespread, social impact assessments covering bias, fairness, privacy, environmental costs, and labor practices remain uneven across the AI ecosystem. To characterize this landscape, we conduct the first comprehensive analysis of both first-party and third-party social impact evaluation reporting across a wide range of model developers. Our study examines 186 first-party release reports and 183 post-release evaluation sources, and complements this quantitative analysis with interviews of model developers. We find a clear division of evaluation labor: first-party reporting is sparse, often superficial, and has declined over time in key areas such as environmental impact and bias, while third-party evaluators including academic researchers, nonprofits, and independent organizations provide broader and more rigorous coverage of bias, harmful content, and performance disparities. However, this complementarity has limits. Only model developers can authoritatively report on data provenance, content moderation labor, financial costs, and training infrastructure, yet interviews reveal that these disclosures are often deprioritized unless tied to product adoption or regulatory compliance. Our findings indicate that current evaluation practices leave major gaps in assessing AI’s societal impacts, highlighting the urgent need for policies that promote developer transparency, strengthen independent evaluation ecosystems, and create shared infrastructure to aggregate and compare third-party evaluations in a consistent and accessible way.

 [Dataset](#)  [Code](#)