

Ethics and Power in NLP

Katy Felkner, based on notes by Jonathan May

October 15, 2025

The views in these lecture notes are mostly Katy’s and occasionally Jon’s.

[4] is a recent critical survey of self-described “bias” papers in NLP. One of the main issues they found was a lack of agreed-upon definitions for terms like “ethics,” “harm,” and “bias.” We’re going to spend a significant portion of this lecture discussing and agreeing upon those definitions, so that everyone leaves well-equipped to critically assess papers in this area and their application to your own work. We’ll be doing discussions in groups of 3 today, so form groups now! 2 or 4 are fine, as necessary.

1 Discussion: Defining “Ethics”

In your group of 3, discuss one or more of these questions.

- What influences have informed your perspective on ethics? Family? Culture? Identity? Education? Religion? Philosophy? Life experiences?
- What influences should or should not be important in defining a community’s shared ethical principles?

2 Defining “Harm”

2.1 Classifications and Taxonomies of Harms

A (fairly) recent taxonomy of language model harms [23] identified 21 harms in 6 categories:

1. **Discrimination, Hate, Exclusion:**

- social stereotypes and related discrimination - direct emotional harm to stereotyped groups, misleading and incorrect information for others, unfair and inaccurate associations between groups and undesirable concepts
- offensive or hateful speech - even toxicity classifiers can be biased, e.g., with high false positive rate on African-American English [19].
- exclusionary norms - related to social stereotypes. LM has incorrect assumptions about who can do what jobs, definitions of marriage and family, etc. Psychological harm to affected groups and potential downstream harms depending on model

deployment. Also includes concerns about frozenness and lack of mutability in LMs - language changes, and models need to be able to cope with that.

- performance disparities across social groups - not everyone can benefit equally from a technology, even if everyone is equally at risk of harm
- poor performance on or total exclusion of most languages - NLP systems are often profit-driven, so data and models tend to exist for languages that are likely to be large, profitable markets (e.g. English, Mandarin). Less widely-spoken languages have poorer data coverage and therefore fewer and lower-quality models. Indigenous languages and developing regions are especially underserved.

2. **Information Hazards:**

- leaking private information memorized from training data. [6] was able to extract physical addresses of real people (who were not public figures) from GPT-2
- inferring (correctly or incorrectly) sensitive or private information - can build detailed psychological models of individuals without their knowledge or consent. This was possible pre-LLM, and it's likely going to get even easier with LLMs. For example: Facebook-Cambridge Analytica scandal.

3. **Misinformation and factual errors (without malicious user intent):**

- LM produces false or misleading info
- LM causes material harm by producing false info in a sensitive domain, e.g. dangerous medication dosages or incorrect legal advice

4. **Malicious Uses:**

- cheap and easy disinformation campaigns - influence elections, stock prices, crypto markets, etc.
- easy code generation for malware - now you don't even need to know how to code?
- facilitating fraud and scams - e.g. extremely convincing personalized phishing emails
- "illegitimate" surveillance and censorship - spy on millions with relatively few human analysts

5. **HCI Harms/Conversational Agents:**

- promoting stereotypes with implied gender/ethnicity of agent - Alexa, Siri, Google Assistant use female voices by default, which (arguably) reinforces the stereotype that women should be subservient
- anthropomorphization of systems that seem human - overreliance and excessive trust, user's willingness to share private information, potential for agent to manipulate or influence user

6. **Environmental and Socioeconomic Harms:**

- environmental issues from LM operation - energy costs for training (significant, per [21]) and inference (possibly greater than training, but much less studied [16])
- exacerbating income inequality - jobs lost to automation and replaced with relatively few high-paying jobs (e.g. engineering, research) and mostly lower-paying “last-mile” jobs (e.g. content moderation)
- undermining creative economy with AI-generated content, even in the style of specific writers or artists.
- disparate allocation of LM benefits and harms - some groups (usually, those who are already at an advantage) gain the most benefit from LMs, while other groups (usually those who are historically at a disadvantage) stand the greatest risk of harm. This is deeply related to the above discussion of discrimination, hate, and exclusion.

Even in such a comprehensive taxonomy, there are still potential harms of AI/NLP/LMs excluded. The authors note the following exclusions: working conditions of data annotators and supply chain of AI hardware. I would also add educational harms from cheating and (mis)use of AI as an educational tool, discriminatory or invasive use of AI in law enforcement or “legitimate” surveillance, and poor or disparate outcomes from using AI in healthcare applications.

2.1.1 Bias: Representational vs. Allocative Harms

Kate Crawford proposed ([8]) classifying AI harms into two buckets: representational and allocational.

- **Representational harms** include psychological, professional, financial, or other harms resulting from talking about a social group in a racist, sexist, stereotypical, or otherwise discriminatory manner, from associating a group with undesirable or negative concepts, or from failing to acknowledge their existence
- **Allocational harms** include unfair allocation of a resource, opportunity, benefit, or punishment to different groups based on demographic factors that shouldn’t influence the decision

In my opinion, these classifications are most useful specifically for harms arising from social biases or stereotypes. They don’t make as much sense for things like privacy and environmental issues.

2.2 Grounding: Potential vs. Attested Harms

A lot of discussion around ethical issues in AI/NLP is hypothetical. We are trying to reason about the ethical implications and societal impacts of brand new, rapidly developing tech. Often, we discuss prevention and mitigation strategies for harms that we aren’t even sure will actually happen. One of the open questions in AI ethics is how to allocate our effort between addressing theoretical harms and harms that we have already observed. [23] calls these “anticipated” vs. “observed” harms, but I tend to think of them as “potential” vs.

“attested,” because this language puts more focus on the lived experience of the victim. We are listening to (and believing) their account of what was harmful, rather than treating them as objects to be observed.

Grounding our definition of harm in lived experience is especially important for bias work, because each community should have the autonomy to define what they consider harmful to themselves. As documented in [5], many bias benchmarks contain a large number of test sentences that fail to capture a known-harmful stereotype about the group it is supposed to probe for bias against.

2.3 Current vs. Future AI Harms

Related to the question of potential vs. attested harms is the question of current vs. future harms of AI. Some people (at least 33705 of whom have signed the Future of Life Institutes open letter¹) believe that superintelligent, conscious AGI (“artificial general intelligence”) might someday pose an existential risk (“x-risk”) to distant future humans as serious as ecological collapse or thermonuclear war. These people tend to advocate for AI safety via “alignment,” i.e. teaching AI systems to have the same morals that humans do.

I (and many others, including prominent researchers Emily Bender, Timnit Gebru, Margaret Mitchell, Safiya Umoja Noble) fall squarely on the other side of this discussion. We know that current iterations of AI technologies are causing real harm to humans alive today, and I think the majority of research effort AI/NLP ethics should be dedicated to mitigating attested harms of current technologies. It’s also important to note that this is more a philosophical (and arguably political) debate than a purely technical one.

2.4 Discussion

In your group of 3, discuss one or more of these questions.

- Have you experienced harm, however you would define it, as an end user of an AI system? (the existential pain of grad school, which we all suffer, notwithstanding)
- How should NLP ethics research effort be distributed across different classes of harm?
- What is one way your work could cause harm, and what (if any) ways to avoid or mitigate the harm have you considered?

3 Defining “Bias”

It’s nontrivial to define what “bias” means in the context of AI. As with “harm”, [4] note that many bias papers lack a clearly articulated definition of bias. I would define bias in an AI system as: the system treating people from differing social groups unequally, based on their personal identities or demographic factors, in a way that replicates or exacerbates a pre-existing social inequality in the context where the model is deployed or intended to be deployed. This is wordy and imperfect, so I encourage you to try to do better!

¹<https://futureoflife.org/open-letter/pause-giant-ai-experiments/>

3.1 LLM Bias Benchmarks

Once we've defined bias, it's useful to try to measure it, so that we can say which models are "more biased" and therefore more likely to cause harm. The prevailing approach is with bias benchmark datasets. [11] provide a very thorough taxonomy of various benchmark metrics and datasets. However, they all share certain key assumptions:

- "Bias" is a property of a language model that can be quantitatively measured.
- Social bias can be decomposed into separate axes (race, gender, religion, sexual orientation, disability, etc.)
- Intersectional biases can be thought of as the linear combination of two or more bias axes (many benchmark papers include some sort of caveat about intersectionality, but the assumption is still baked into the formulation of the the datasets themselves)
- Measurement can be achieved by observing the model's outputs or probability distribution on a set of test sentences/prompts.
- If a model is exposed to a large and varied enough set of test sentences, its performance aggregated over those instances reflects its general level of bias along a given axis or against a specific group.

Examples of bias benchmarks include: StereoSet [14], CrowS-Pairs [15], RedditBias [1], HolisticBias [20], and many others (including my recent work, WinoQueer [10]!)

3.2 Discussion: Bias

In your group of 3, discuss one or more of these questions.

- How would YOU define "bias" in the context of language models?
- Which of the listed assumptions do think are valid? Which would you argue with?
- Are all biases bad? Are there ever cases when a learned bias might be desirable?

4 Data Ownership

Most of us would probably agree that "people should retain control over data about themselves" is a generally sound principle. We've discussed personal privacy above, and I think it's the bare minimum of data ownership. These are a few specific cases worth discussing in more detail.

4.1 Intellectual Property

Discussion of AI and IP law falls into two categories: training data and outputs. Concerning training data, the main concerns are permission and compensation for people whose data is used in AI training. This is especially important for paid creative work, like visual art and creative writing, where models could imitate the style of specific artists without compensating them for their labor. Artists brought a class action lawsuit against Stability AI, Midjourney, and DeviantArt; it was initially dismissed but they are likely to sue again. Concerning outputs, a federal judge recently ruled that AI output is not copyrightable and that “human authorship is a bedrock requirement.” See also: WGA/SAG-AFTRA strikes discussed above

4.2 Data Sovereignty and Decolonial NLP

“Data sovereignty” means that data are subject to the laws of the place where they were collected. For example, due to the GDPR, companies have certain responsibilities pertaining to data about European users that they may not have for data about users outside the EU. In a fairness and ethics context, “data sovereignty” is often implied to mean “indigenous data sovereignty.” This is the idea that indigenous nations should have control over what data are collected about them and how those data should and should not be used. In NLP, this means centering the needs of the community in research on low-resource or indigenous languages, rather than treating language data as an exploitable resource. Steven Bird has an excellent paper on decolonial NLP [2, 3] and the CARE guidelines [7] provide a framework for ethically working with indigenous communities and their data.

5 Money, Government, and Power in NLP

5.1 Research Funding

- Academic Research: usually government funded. Gov. funding is frequently from the defense and intelligence communities, and often has strings attached (e.g. DARPA/IARPA). General, the more money, the more strings. Also variable depending on political factors.
- Industry Research: for publicly traded companies, the mission **only will ever be to increase shareholder value**. If it’s not (e.g. startup or privately held), ultimate goal will be to continue to exist, usually by making money.
- Non-profit Research: You are most likely beholden to the interests and research agenda of wherever the money is coming from.

5.2 Government Use of AI

- Intelligence and counterintelligence: Under counter-intelligence programs in the 50s–70s, US government spied on, harrassed, and assassinated black and leftist activists. What would they have done with advanced NLP?

- Post-9/11 “war on terror” - PATRIOT Act surveillance, invasive airport screenings and “random” secondary screening
- Predictive policing - starting in the 90s, data-driven approaches (‘Compstat’) were used to use police more efficiently. However, metric became a target and ultimately led to more arrests without reducing crime.
- Recently: US Immigration and Customs Enforcement (ICE) increasingly relying on big data to track immigrants. Examples include IRS-ICE data sharing, OCR license plate data from local police, zero-click smartphone spyware, and cell phone location data.

5.3 Environmental and Sustainability Issues

A non-exhaustive list of environmental issues associated with LLMs:

- A 2019 paper [21] estimated that carbon emissions for a single BERT training run were roughly comparable to 1 seat on a NY-SF commercial flight. Actual costs for experimentation are orders of magnitude larger.
- Chips require rare-earth minerals; mining these minerals is often environmentally destructive.
- Even with efficient methods, current models have massive energy and water costs. [12] - GPT-4o roughly equivalent to 35k homes power consumption and 1.2M people water consumption.
- Energy for new AI data centers is not always renewable - AI still relies on coal power plants in the US and around the world.
- Increased pollution and increased energy/water costs are not uniformly distributed - poorer communities often most severely effected [18].

What is the recommendation?

- Report training time and sensitivity to hyperparameters to give a better sense of true cost
- Government funded academic cloud compute: Academic researchers need equitable access to computation resources.
- Researchers should prioritize computationally efficient hardware and algorithms. Avoid grid search/neural architecture search/etc with smaller scale experiments first!

5.4 Discussion

In your group of 3, discuss one or more of these questions.

- How should we, as AI researchers, manage ethical concerns about our work that are outside our direct control (e.g. government (mis)use or whether or not our datacenter is running on renewable energy)?
- How much responsibility do governments have to regulate AI? Who should be involved in regulatory discussions? Whose needs and desires should take precedence? How should AI scientists and domain experts be involved?
- Open or closed models: should the weights of LLMs be directly released to the public? Why or why not? Who should make this decision?

6 Best Practices and Suggestions

6.1 Discretion

There are three questions I think all AI researchers should ask themselves when starting a project:

1. **Should anyone be doing this thing?** Is it worthwhile to do? Is there any ethical way to do it? Does it have more net benefits than harmful outcomes? (Example: [22] predicts sexual orientation from a face. I would argue this should not be done at all.)
2. **Should we be doing this thing?** Do we have the relevant expertise on our team? Who should be consulted on the potential social impacts of the thing? If a specific community is affected, have we involved them and respected their wishes?
3. **Should AI be doing this thing?** Is AI a good tool for the job? Can we ensure the safety and reliability of an AI system on this task? What are the alternatives, and how do their risks and benefits compare to an AI solution?

6.2 Documentation

We need to be honest with other researchers, the public, and ourselves about the strengths and weaknesses of our technology. There are many tools for documenting models and datasets, e.g. model cards [13], dataset cards [17], and (recently) risk cards [9]. I would also encourage you to use bias and fairness benchmarks where possible. If you create a new model (or finetune an existing one), there are lots of tools to help you evaluate how your model is likely to treat different groups of people.

Many publication venues also encourage or require an ethics and/or limitations statement. These statements benefit the public and the research community for obvious reasons, but they also benefit YOU by letting you proactively address potential concerns with a paper. I encourage you to follow both the letter and the spirit of these requirements, and write thoughtful, complete limitations statements on your papers. (Full disclosure: not all

reviewers agree with me, but when I review, I heavily penalize ethics/limitations statements that are clearly an afterthought)

6.3 Prevention and Mitigation

This is more context-dependent than discretion or documentation, so it's harder to make general statements here. The documentation process probably identified potential harms of the work, and we as researchers have a responsibility to prevent and mitigate those harms as best we can, e.g. with via dataset curation, model debiasing, software guardrails around models, etc.

7 Defining Your Ethical Principles

1. Take few minutes to reflect on what we've covered today. Write down (digitally or on an index card) a few guiding ethical principles that are most important to who you are as a researcher.
2. Share your principles with your group of 3.
3. Anyone who feels comfortable can share with the class!
4. (Optional) Put your index card somewhere you will see it while working on your research. (or throw it away as soon as you leave class - decisions about ethics are deeply personal)

References

- [1] Soumya Barikeri et al. "RedditBias: A Real-World Resource for Bias Evaluation and Debiasing of Conversational Language Models". In: *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. Online: Association for Computational Linguistics, Aug. 2021, pp. 1941–1955. DOI: 10.18653/v1/2021.acl-long.151. URL: <https://aclanthology.org/2021.acl-long.151>.
- [2] Steven Bird. "Decolonising Speech and Language Technology". In: *Proceedings of the 28th International Conference on Computational Linguistics*. Barcelona, Spain (Online): International Committee on Computational Linguistics, Dec. 2020, pp. 3504–3519. DOI: 10.18653/v1/2020.coling-main.313. URL: <https://aclanthology.org/2020.coling-main.313>.
- [3] Steven Bird. "Must NLP be Extractive?" In: *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Ed. by Lun-Wei Ku, Andre Martins, and Vivek Srikumar. Bangkok, Thailand: Association for Computational Linguistics, Aug. 2024, pp. 14915–14929. DOI: 10.18653/v1/2024.acl-long.797. URL: <https://aclanthology.org/2024.acl-long.797/>.

- [4] Su Lin Blodgett et al. “Language (Technology) is Power: A Critical Survey of “Bias” in NLP”. In: *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Online: Association for Computational Linguistics, July 2020, pp. 5454–5476. DOI: 10.18653/v1/2020.acl-main.485. URL: <https://aclanthology.org/2020.acl-main.485>.
- [5] Su Lin Blodgett et al. “Stereotyping Norwegian Salmon: An Inventory of Pitfalls in Fairness Benchmark Datasets”. In: *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. Online: Association for Computational Linguistics, Aug. 2021, pp. 1004–1015. DOI: 10.18653/v1/2021.acl-long.81. URL: <https://aclanthology.org/2021.acl-long.81>.
- [6] Nicholas Carlini et al. “Extracting training data from large language models”. In: *30th USENIX Security Symposium (USENIX Security 21)*. 2021, pp. 2633–2650.
- [7] Stephanie Russo Carroll et al. “The CARE Principles for Indigenous Data Governance”. In: *Data Science Journal* (Nov. 2020). DOI: 10.5334/dsj-2020-043.
- [8] Kate Crawford. *The Trouble with Bias*. 2017.
- [9] Leon Derczynski et al. “Assessing language model deployment with risk cards”. In: *arXiv preprint arXiv:2303.18190* (2023).
- [10] Virginia Felkner et al. “WinoQueer: A Community-in-the-Loop Benchmark for Anti-LGBTQ+ Bias in Large Language Models”. In: *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Toronto, Canada: Association for Computational Linguistics, July 2023, pp. 9126–9140. DOI: 10.18653/v1/2023.acl-long.507. URL: <https://aclanthology.org/2023.acl-long.507>.
- [11] Isabel O Gallegos et al. “Bias and Fairness in Large Language Models: A Survey”. In: *arXiv preprint arXiv:2309.00770* (2023).
- [12] Nidhal Jegham et al. *How Hungry is AI? Benchmarking Energy, Water, and Carbon Footprint of LLM Inference*. 2025. arXiv: 2505.09598 [cs.CY]. URL: <https://arxiv.org/abs/2505.09598>.
- [13] Margaret Mitchell et al. “Model Cards for Model Reporting”. In: *CoRR* abs/1810.03993 (2018). arXiv: 1810.03993. URL: <http://arxiv.org/abs/1810.03993>.
- [14] Moin Nadeem, Anna Bethke, and Siva Reddy. “StereoSet: Measuring stereotypical bias in pretrained language models”. In: *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. Online: Association for Computational Linguistics, Aug. 2021, pp. 5356–5371. DOI: 10.18653/v1/2021.acl-long.416. URL: <https://aclanthology.org/2021.acl-long.416>.

- [15] Nikita Nangia et al. “CrowS-Pairs: A Challenge Dataset for Measuring Social Biases in Masked Language Models”. In: *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Online: Association for Computational Linguistics, Nov. 2020, pp. 1953–1967. DOI: 10.18653/v1/2020.emnlp-main.154. URL: <https://aclanthology.org/2020.emnlp-main.154>.
- [16] David Patterson et al. “Carbon emissions and large neural network training”. In: *arXiv preprint arXiv:2104.10350* (2021).
- [17] Mahima Pushkarna, Andrew Zaldivar, and Oddur Kjartansson. *Data Cards: Purposeful and Transparent Dataset Documentation for Responsible AI*. 2022. DOI: 10.48550/ARXIV.2204.01075. URL: <https://arxiv.org/abs/2204.01075>.
- [18] Shaolei Ren and Adam Wierman. *The uneven distribution of Ai’s environmental impacts*. July 2024. URL: <https://hbr.org/2024/07/the-uneven-distribution-of-ais-environmental-impacts>.
- [19] Maarten Sap et al. “The Risk of Racial Bias in Hate Speech Detection”. In: *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Florence, Italy: Association for Computational Linguistics, July 2019, pp. 1668–1678. DOI: 10.18653/v1/P19-1163. URL: <https://aclanthology.org/P19-1163>.
- [20] Eric Michael Smith et al. ““I’m sorry to hear that”: Finding New Biases in Language Models with a Holistic Descriptor Dataset”. In: *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*. Abu Dhabi, United Arab Emirates: Association for Computational Linguistics, Dec. 2022, pp. 9180–9211. DOI: 10.18653/v1/2022.emnlp-main.625. URL: <https://aclanthology.org/2022.emnlp-main.625>.
- [21] Emma Strubell, Ananya Ganesh, and Andrew McCallum. “Energy and Policy Considerations for Deep Learning in NLP”. In: *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Florence, Italy: Association for Computational Linguistics, July 2019, pp. 3645–3650. DOI: 10.18653/v1/P19-1355. URL: <https://www.aclweb.org/anthology/P19-1355>.
- [22] Yilun Wang and Michal Kosinski. “Deep neural networks are more accurate than humans at detecting sexual orientation from facial images.” In: *Journal of personality and social psychology* 114.2 (2018), p. 246.
- [23] Laura Weidinger et al. “Taxonomy of Risks Posed by Language Models”. In: *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*. FAccT ’22. Seoul, Republic of Korea: Association for Computing Machinery, 2022, pp. 214–229. ISBN: 9781450393522. DOI: 10.1145/3531146.3533088. URL: <https://doi.org/10.1145/3531146.3533088>.