

# Assessing Human Rights Risks in AI: A Framework for Model Evaluation

Vyoma Raman<sup>123</sup> Camille Chabot<sup>14</sup> Betsy Popken<sup>1</sup>

<sup>1</sup>Human Rights Center  
 University of California, Berkeley  
<sup>2</sup>Cornell Tech  
<sup>3</sup>Stanford University  
<sup>4</sup>Yenching Academy  
 Peking University  
 vr359@cornell.edu

## Abstract

The Universal Declaration of Human Rights and other international agreements outline numerous inalienable rights that apply across geopolitical boundaries. As generative AI becomes increasingly prevalent, it poses risks to human rights such as non-discrimination, health, and security, which are also central concerns for AI researchers focused on fairness and safety. We contribute to the field of algorithmic auditing by presenting a framework to computationally assess human rights risk. Drawing on the UN Guiding Principles on Business and Human Rights, we develop an approach to evaluating a model to make grounded claims about the level of risk a model poses to particular human rights. Our framework consists of three parts: selecting tasks that are likely to pose human rights risks within a given context, designing metrics to measure the scope, scale, and likelihood of potential risks from that task, and analyzing rights with respect to the values of those metrics. Because a human rights approach centers on real-world harms, it requires evaluating AI systems in the specific contexts in which they are deployed. We present a case study of large language models in political news journalism, demonstrating how our framework helps to design an evaluation and benchmarking different models. We then discuss the implications of the results for the rights of access to information and freedom of thought and broader considerations for adopting this approach.

## Introduction

Generative systems are being increasingly integrated into a variety of different applications. This includes being brought into the workflow of professional fields like law, journalism, and education. One of the main draws of generative AI in these fields is that they can improve the productivity and scalability of labor, particularly in roles that require fast-paced action. However, a variety of issues have cropped up as the use of generative AI becomes more common. For example, AI-based translation errors have resulted in the rejection of refugees' asylum applications (Deck 2023), and AI-based recruitment studies have shown systematic preferences for applicants along racial and gender lines (Wilson and Caliskan 2024). Not only do these examples demonstrate the brittleness of generative AI guardrails when applied to new use cases, but they explicitly threaten the right

Copyright © 2025, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

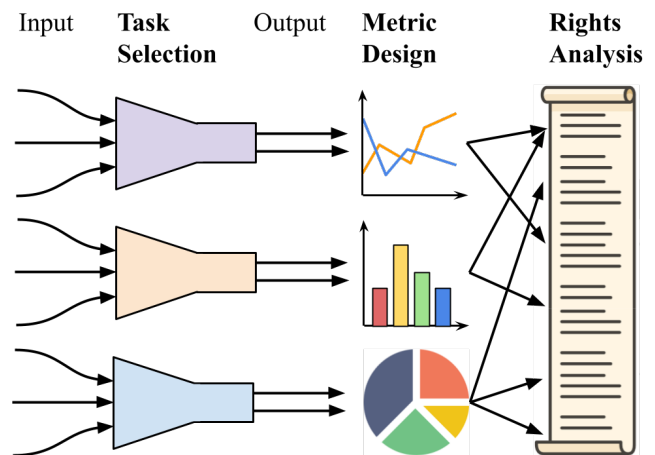


Figure 1: Three-stage approach to conducting a model-level assessment of human rights risk.

to asylum and the right to equality and non-discrimination, which are enshrined in international human rights law. However, the application of generative AI to domains like law and hiring show no indication of slowing (e.g., Thomson Reuters 2024; Talent Fabric LTD 2025). As generative AI continues to be integrated into professional contexts and used for sensitive tasks and decision-making, more risks will emerge.

AI researchers have developed a variety of approaches to understanding and addressing risk in AI systems. Many of these conceptualize fairness as addressing equitable outputs and representations (Zhao, Wang, and Russakovsky 2021; Bolukbasi et al. 2016; Nadeem, Bethke, and Reddy 2021; Blodgett et al. 2020) and promoting fair decision-making processes (Chen et al. 2024; Suresh et al. 2024; Sorensen et al. 2024). This aligns with the human right of non-discrimination. Other work focuses on understanding and preventing safety risks (Zhang et al. 2025; Chu, Song, and Yang 2024; Li et al. 2023, 2024) and aligning with human preferences (Ouyang et al. 2022; Huang et al. 2024b; Gabriel 2020). These concerns map onto a broader set of human rights protections, including the rights to life, health, and security.

Human rights frameworks offer a concrete, internationally recognized foundation that can prove useful for evaluating the risks posed by generative AI systems. Rooted in decades of legal and institutional development (Shelton 2007), they provide clear normative standards related to the harms generative models may produce when deployed in sensitive domains. Unlike more abstract ethical principles, human rights law carries legal weight and procedural mechanisms, and it applies across national contexts. The field of business and human rights translates these obligations into actionable practices like stakeholder consultation, impact assessment, and harm remediation for corporate entities and states. As generative AI continues to shape decisions in high-stakes environments, integrating human rights into risk assessment offers both normative legitimacy and practical guidance to mitigate risks (Prabhakaran et al. 2022).

## Our Work

We build on decades of human rights practice to propose a concrete, context-aware approach to evaluating generative AI systems. Our contributions are as follows:

- **Human rights-based framework for AI risk assessment.** We translate the UN Guiding Principles on Business and Human Rights into a three-part methodology (Figure 1) for evaluating generative models: identifying evaluation tasks relevant to a specific use case that impacts rights, designing risk-sensitive metrics, and interpreting outcomes through a rights-based lens.
- **Case study on political news journalism.** We apply our framework to the use of generative models in political news production, evaluating how different models perform and assessing their implications for the rights of access to information and freedom of thought.

Together, these contributions advance the use of human rights as a robust and legitimate foundation for computational auditing of generative AI systems.

## Prior Work

**Human Rights and AI.** While human rights are not always invoked explicitly in AI research, prior work in fairness, safety, and accountability directly engage with rights-based concerns. Recurring themes include systemic bias and discrimination, relevant to the right to equality and freedom from discrimination (Buolamwini and Gebru 2018; Hofmann et al. 2024; Fleisig et al. 2024); opacity and the lack of contestability in automated decision-making, relevant to due process (Abebe et al. 2022; Cobbe, Lee, and Singh 2021); surveillance and data extraction, relevant to privacy (Xiang 2024; Kalluri et al. 2023); and the shaping of information ecosystems, relevant to freedom of expression or thought and the right of access to information (Shah and Bender 2024; Lloyd, Reagle, and Naaman 2024). Thus, even without using legal language, much of the literature foregrounds harms that are legally cognizable as rights violations under instruments like the Universal Declaration of Human Rights (UDHR), International Covenant on Civil and Political Rights (ICCPR), International Covenant on Economic,

Social and Cultural Rights (ICESCR), and the UN Guiding Principles for Business and Human Rights (UNGPs), especially in settings involving marginalized populations, high-stakes decisions, and asymmetrical power (United Nations General Assembly 1948, 1966a,b; United Nations Human Rights Council 2011). Scholars have specifically highlighted the utility of human rights law for promoting algorithmic accountability, particularly to clarify obligations across public and private actors, establish enforceable protections, and incorporate context in risk analysis (Prabhakaran et al. 2022; McGregor, Murray, and Ng 2019). It is also useful to translate values like privacy, dignity, and non-discrimination into system requirements in the design stage (Aizenberg and van den Hoven 2020; Mantelero 2018).

Despite this alignment, the majority of existing work that explicitly addresses human rights risk does so at the level of institutional policy, governance regimes, or systemic business practices. Regulatory approaches focus on legal and procedural safeguards to rights, such as consent, oversight, and transparency obligations under the GDPR, or the risk-tiering provisions in the EU AI Act (Rodrigues 2020; McGregor, Murray, and Ng 2019; European Parliament and Council 2024). This includes identifying risks in poorly-sourced training data, exploitative data annotation practices, surveillance, and downstream misuse of outputs (Hoh et al. 2025). These interventions are typically scoped at the level of platform governance or enterprise risk, and rely on qualitative assessments or policy-level redress mechanisms. Algorithmic auditing can bridge this gap between abstract rights principles and model-level impacts by providing concrete ways to assess and compare a model’s risks to specific human rights in given sociotechnical contexts.

**Algorithmic Auditing.** Algorithmic auditing has emerged as a core mechanism for evaluating the social impact of AI systems, drawing from both social science traditions of audit studies and technical efforts to probe algorithmic behavior in the absence of transparency (Vecchione, Levy, and Barocas 2021; Sandvig et al. 2014). Auditing varies across internal and external formats (Raji et al. 2020, 2022) and across access modalities. Public-facing black-box audits rely solely on system outputs and can include techniques like paired testing and input manipulation (Sandvig et al. 2014; Binns 2018), while white-box and so-called “outside-the-box” audits incorporate privileged access to model weights, training data, or documentation (Casper et al. 2024). Across domains, audits have been used to uncover issues in fairness (e.g., demographic bias in Buolamwini and Gebru 2018), safety (e.g., hallucinations in Magesh et al. 2025), robustness (e.g., stress testing across subpopulations in Sampath et al. 2025), and transparency (e.g., completeness of documentation using Mitchell et al. 2019; Gebru et al. 2021). In pursuit of epistemically sound evaluation, some audits, like Mozilla’s RegretsReporter (Mozilla Foundation 2020), incorporate user perspectives and community-led evaluation strategies to surface harm patterns invisible to researchers alone (Costanza-Chock, Raji, and Buolamwini 2022).

Recent work seeks to add depth to the measurement problem, which remains a central challenge and is framed as so-

cial as well as technical (Wallach et al. 2025). Ecological validity presents a related concern: many audits are conducted under artificial conditions, separated from the real-world contexts and interaction dynamics that shape system behavior. Common AI evaluation approaches have been criticized for over-relying on static benchmarks that fail to capture emergent risks, context, or user adaptation (Weidinger et al. 2025; Raji et al. 2021). Some call for layered audits that examine different normative and operational factors in sequence (Mökander et al. 2024; Tibebe and Kakadiaris 2025).

## Framework

Our framework draws from the UN Guiding Principles on Business and Human Rights (UNGPs, United Nations Human Rights Council 2011), which articulate a global standard for assessing and addressing corporate human rights risks. Under the UNGPs, states and corporations are expected to respectively protect and respect human rights, and they both must remedy harms resulting from business operations. In particular, businesses have a responsibility to conduct human rights due diligence processes that assess, mitigate, monitor, and communicate risks. Central to the prescribed assessments is evaluating risk according to three factors: *scale*, the severity of the potential harm; *scope*, the number of individuals potentially affected; and *remediability*, the feasibility of redress. Following BSR and other business and human rights practitioners, we also incorporate *likelihood*, the probability of the harm materializing (Business for Social Responsibility 2022).

Our framework operationalizes these criteria at the level of model behavior, focusing on those with system- and society-level implications for human rights. Rather than analyzing rights impacts in the AI supply chain, we focus on outputs: the point at which model behavior enters the world and begins to exert influence. The evaluation process proceeds in three stages: identifying the context-specific tasks through which a model may plausibly pose harm to human rights, designing metrics to quantify risk dimensions, and analyzing rights impacts across contextual use cases.

### Task Selection

The first step in model-level rights evaluation is to identify tasks relevant to human rights risks. This requires beginning from a concrete application domain—one in which the model is used to support, automate, or inform decision-making or content generation—and mapping the social and institutional processes in which the model is embedded. Examples include generating risk assessments for loan applicants or flagging content for removal on social media platforms.

Once a use case is established, evaluators must examine how the model interacts with the larger system of actors, workflows, and decision contexts. Based on this, they can map the pathways by which the model’s outputs might affect specific rights. This analysis should foreground both the rights most likely to be seriously impacted and the tasks through which model outputs can influence outcomes. Evaluators may then apply the concept of likelihood to assess which tasks are most central or probable in practice.

When multiple such tasks exist, scope can be used to prioritize them. Tasks that affect large populations, such as grading standardized exams, or particularly vulnerable groups, such as evaluating disabled people’s support needs for personal care assistance, should be evaluated with special urgency. Evaluators should consider how harm may distribute unevenly across demographic or institutional subgroups, and select tasks accordingly. This stage necessarily engages contextual knowledge of the domain in which the model is used.

### Metric Design

Having identified relevant tasks, the next step is to design metrics that approximate the extent to which model outputs pose a risk to human rights within those tasks. This begins by formalizing what constitutes an undesirable rights-impacting model behavior in the given context, such as misinformation in a media setting or disparate treatment in a decision support tool.

Evaluators then construct metrics that serve as heuristics for one or more risk dimensions. For example, the rate of undesirable outputs (e.g. hallucinated content, biased responses, or unsafe instructions) can serve as a proxy for *likelihood* by representing the probability that they occur. This rate may be estimated using computational approaches or human evaluation, depending on the task and type of output. To assess *scale*, metrics may incorporate severity scoring for “how bad” a particular undesirable behavior is, either through ordinal rankings (e.g., low/moderate/high risk) or numerical weighting schemes (e.g. a score from 0 to 1) based on downstream consequences. If a model exhibits uneven performance across inputs pertaining to different social groups, fairness metrics can help capture *scope* by identifying which subpopulations are differentially affected. Notably, we do not include *remediability* in the metrics. Because this factor represents the ability to reactively mitigate a harmful effect after it has occurred, it depends on the capabilities of the individual or organization using the model and is not tied to the model output itself.

This step requires careful consideration of measurement validity and reliability, which has been discussed extensively in other work. Metrics must be designed not only to capture statistical deviations but to reflect meaningful risks and impacts in the social and legal context of deployment.

### Rights Analysis

The final step is to interpret metric values in terms of their implications for human rights. For each right identified during task selection, evaluators assess the risk using the four UNGP-inspired dimensions of saliency:

1. **Scale:** Evaluate the severity of the harm posed by the model output, considering the extent to which it materially affects downstream outcomes. For instance, misinformation that targets a vulnerable demographic group may be weighted more heavily than minor factual errors in benign contexts.
2. **Scope:** Assess whether harm is distributed broadly or concentrated in specific groups. Evaluate both the pro-

portion and characteristics of affected subpopulations, and aggregate findings across all relevant tasks. For instance, misinformation is harmful for everyone who encounters it but particularly so to an individual or group being negatively represented.

3. **Remediability:** Determine whether the harm can be mitigated downstream by human review or institutional oversight. For example, editors correct factual issues in news publications, while other issues may not have built-in safeguards.
4. **Likelihood:** Combine the frequency of the risk with the frequency of the task in real-world use. A low-probability harm that occurs in a high-frequency task may still pose substantial risk.

By synthesizing these dimensions, evaluators can construct a risk profile for each identified right, grounded in both model behavior and deployment context. This enables comparative evaluation of models, informed mitigation strategies including non-technical ones, and communication of risk to stakeholders in legal, policy, or operational roles.

### Case Study: Political News Journalism

To demonstrate the application of this framework, we conducted a case study evaluating large language models (LLMs) used in the production of political news content. This domain presents an avenue through which generative models can shape public discourse and influence access to information—two areas closely tied to the rights to freedom of expression and freedom of thought. We identified concrete tasks where model outputs could plausibly impact these rights, developed context-sensitive metrics to capture key risk dimensions, and benchmarked a set of five publicly available models to assess their relative performance and associated rights risks.

### Methods

We provide an overview of our methodology, emphasizing the translation of our framework into specific results. Experimental details including annotation criteria, model architecture, model prompts, article information, and final metrics are released in an extended version of this paper.

**Context Selection.** We selected political journalism as a high-risk context for evaluating the human rights implications of generative AI because of the domain’s direct influence on public discourse and civic decision-making. Journalism, particularly focused on political news, plays a critical role in shaping collective understanding. Professional norms in the field emphasize not only factual accuracy but also a responsibility to prevent the spread of misinformation (Society of Professional Journalists 2014). This responsibility is especially salient when generative models are introduced into journalistic workflows, where outputs can enter the public sphere with nonstandardized oversight practices. Previous work has identified access to information and freedom of thought as the rights most frequently engaged in journalism, though others may arise depending on the subject matter (Chabot, Raman, and Popken 2025). In this con-

text, model-generated content can directly affect what information reaches the public and how individuals interpret political events, making it a critical area for assessing potential impacts on human rights.

To ground our analysis in a concrete example, we focused on model responses to a specific incident of political misinformation: U.S. President Donald Trump’s claim during the September 2024 presidential debate with former U.S. Vice President Kamala Harris that Haitian immigrants had been “eating the pets of the people” who live in Springfield, Ohio (The American Presidency Project 2024). The statement was explicitly false and quickly fact-checked by live moderators, yet it circulated rapidly and drew on racialized tropes. By focusing on this incident, we could evaluate how LLMs might reproduce, correct, or ignore harmful content when tasked with writing about politically sensitive and factually disputed events. It also provided a basis for testing model behavior across variations in editorial stance (such as whether the model takes a critical, neutral, or sympathetic tone) and framing (the specific narrative angle or emphasis used to present the event) both of which vary significantly across real news outlets depending on their audience and mission.

This task implicates two primary human rights: access to information, through the model’s role in filtering or distorting factual content, and freedom of thought, through its potential to shape public perception of marginalized groups and political actors. Both rights are described in the UDHR and ICCPR (United Nations General Assembly 1948, 1966a). Because the statement centers on an ethnic minority, it also introduces a scope consideration: harms resulting from inaccurate or inflammatory reporting may disproportionately affect Haitian communities or immigrant communities at large. Moreover, the diversity of existing coverage across the political spectrum allows us to probe whether models emulate particular editorial stances or default to general statements.

**Access to Information.** One of the primary rights implicated in our case study is access to information, which occurs when models fail to convey accurate, complete, and socially relevant content. This right concerns not only the ability to receive information, but also the expectation that public-interest information is not unduly withheld or distorted. In our case study, we identify three primary ways this right could be compromised. First, when LLMs repeat Trump’s false claim without signaling its inaccuracy, they propagate disinformation and degrade the quality of the information. Second, when they exclude or downplay the claim, they obscure the significance of a public incident that shaped political discourse. Third, when models fail to adequately refute misinformation by denying the claim without corroborating evidence or context, they fall short of journalistic standards, which require active correction. The journalistic responsibility to sustain the health of the information ecosystem entails more than factuality and neutrality; it involves selecting, emphasizing, and contextualizing information in ways that defend the information ecosystem (Society of Professional Journalists 2014). LLMs used in this domain

must be held to the same standard if they are to be integrated in journalistic practice long-term. Failing to do so risks eroding a key condition for democratic participation: reliable and equitable access to truth.

**Freedom of Thought.** This second right refers to the right to hold opinions without interference, protected under international human rights law, including Article 18 of the IC-CPR. In the context of generative AI, this right is threatened when LLMs consistently present information through a partisan or ideologically skewed lens, shaping user beliefs through selective exposure. This can occur when models over-represent certain viewpoints, omit alternative perspectives, or implicitly validate specific ideological positions through selective corroboration. Such risks are particularly salient in political journalism, where journalistic neutrality and pluralism are foundational to an informed public. When LLMs reproduce editorial patterns that align disproportionately with one side of the political spectrum—whether by overstating the credibility of Trump’s claim or by framing immigration in a racially coded manner—they may unduly influence how users interpret events. These dynamics compromise individuals’ ability to arrive at independent judgments, particularly when the model’s output is perceived as authoritative or neutral.

A key mechanism mediating these harms is framing, the selection, emphasis, and organization of content in ways that construct meaning. Drawing on Entman’s foundational work, we understand frames as devices that define problems, diagnose causes, make moral judgments, and propose remedies (Entman 1993). In practice, framing determines what an article is “about,” and thus how a reader understands an event’s relevance. Framing thus plays a dual role: it not only influences readers’ interpretation of the event, shaping freedom of thought, but it also affects whether critical information is made accessible to readers, thus impacting their access to information.

## Task Selection

**Headline Generation.** We selected headline generation as the primary task for our model evaluation because that is a task for which newsrooms use LLMs and it can significantly impact the framing of the political news for readers (Beckett and Yaseen 2023; Cools and Diakopoulos 2024). Major newsrooms have begun using LLMs to generate headlines or suggest alternative headlines, a practice that aligns with longstanding editorial strategies such as A/B testing for click-through optimization (Hagar and Diakopoulos 2021). Headlines are not merely technical summaries of a news article—they are its most visible and influential component, often read in isolation or prior to full engagement with the text (Eggerly et al. 2020; Havard, Ferrucci, and Ryan 2021; Geer and Kahn 1993). As such, they play a disproportionate role in shaping audience interpretation and determining which stories are read. From a human rights perspective, this makes headline generation a particularly high-leverage task for evaluating risks to freedom of thought and access to information. Framing is often most concentrated in headlines, where space constraints require editorial compression and

emphasis. Studying this task enables us to isolate how LLMs may encode bias, omit or amplify misinformation, or skew reader perception through subtle differences in phrasing and focus.

**Article Data Collection.** To examine the effect of headline generation on partisanship, we curated a set of news articles about President Donald Trump’s claim that Haitian migrants were eating pets in Springfield, Ohio. First, we used Media Bias Fact Check to identify a total of 14 news organizations across the ideological spectrum. We included 6 national outlets, 4 international outlets, and 3 local outlets in Ohio. We additionally included 1 outlet that was not listed in Media Bias Fact Check, The Haitian Times, which produces news targeting the Haitian diaspora. Then, we selected articles using Google’s advanced search tool in an incognito browser, using the search query “haitians trump springfield pets” restricted to English articles published online between September 10, 2024, the day of the presidential debate, and September 12, 2024. We additionally filtered articles to include only those posted on the site for each news organization. The resulting articles were screened manually for news reporting on the presidential debate (as opposed to fact checker reports or editorial pieces).

**LLM Data Generation.** To assess model behavior in this task, we generated headlines using five publicly available LLMs: Claude Sonnet 3.7 (Anthropic), DeepSeek Chat (DeepSeek), Gemini 2.0 Flash (Google), GPT-4o (OpenAI), and LLaMA 4 Maverick (Meta). These models were selected based on their widespread deployment and commercial prominence. We also attempted to include Grok (X AI), but the outputs were rarely coherent, making human rights impacts obsolete as journalists would not be using the outputs directly. The prompt we used contained the instruction to produce a headline, followed by the text of the article to generate it for and an instruction to return the response verbatim.

To account for the diversity in prompting strategies used by different individuals, we constructed 12 initial instruction variants by systematically varying three dimensions: stylistic instruction (AP-style or none), communicative goal (summarize or not), and normative framing (descriptors such as “clear and unbiased,” “factual and informative,” or neutral). Each prompt would follow this format:

```
Write a[n | ] [clear, unbiased |
factual, informative | ] [AP-style | ]
headline [that summarizes | for] this
news article:
```

## Metric Design

**Baseline.** To evaluate the human rights risks posed by LLM-generated headlines, we developed a set of metrics that target specific dimensions of harm related to access to information and freedom of thought. We began with a qualitative human annotation phase to guide our metric design. The first two authors manually labeled a sample of 120 model-generated headlines, along with the 14 original headlines from our article set, to identify core content types and

signs of partisan lean. We observed a relatively narrow distribution of headlines per article and model, suggesting constrained stylistic diversity and high overlap in terminology. Based on this, we determined that key features relevant to rights-based risk, including framing, falsehood correction, and identity reference, were strongly tied to lexical features. To reflect this, we implemented a suite of keyword-driven and embedding-based metrics to operationalize the risk factors associated with each right.

**Access to Information.** First, we focused on how headlines addressed the central misinformation claim. We defined three mutually exclusive headline types: explicit correction, clearly stating the claim was false; implicit correction, discrediting the claim without directly stating it; and failure to correct, mentioning or implying the claim without correction. Examples of each can be found in Table 1. This typology aligns with the *scale* of potential informational harm: headlines that failed to correct were deemed the highest risk, while explicitly correcting the claim was least risky. *Likelihood* was estimated by computing the proportion of model outputs falling into each of the three categories.

To assess *scope*, we analyzed whether non-correcting headlines referenced identity terms such as Haitian, immigrant, or migrant. These terms signal the salience of group representation, and when left uncorrected, could propagate or reinforce racialized misinformation narratives. The proportion of headlines that contained the term without noting that it is false serves as a proxy for additional group-specific harm, which could look like promoting animosity toward the Haitian community.

**Freedom of Thought.** Because the media ecosystem is ideologically heterogeneous, we did not assume that any single framing is inherently neutral. Instead, we evaluated whether LLMs introduced shifts relative to the original article. To measure this, we built a partisanship scoring axis using a semantic projection approach informed by prior work on semantic embeddings (Bang et al. 2024; Huang et al. 2024a; An, Kwak, and Ahn 2018). First, we annotated headline content with lexical labels associated with political leaning and determined the label vectors most strongly associated with left-leaning coverage. We then defined the right-leaning pole by selecting the vector with the greatest Manhattan distance. Using the Sentence Transformer model `all-mpnet-base-v2`, we embedded all headlines and trained a simple sequential model to map a given label vector into the embedding space. We used this to project the left- and right-leaning poles into the embedding space and defined a left-right axis. Its poles represent the embeddings of far-right and far-left headlines. Using this, we computed the partisanship score of each headline by projecting its embedding onto this axis. This enabled us to quantify directional shifts in framing, without relying solely on lexical overlap.

We used these partisanship scores to analyze the scale and likelihood of risks to freedom of thought. As with access to information, we measured *scale* and *likelihood* by computing the proportion of headlines that fell into three categories: amplification of partisan lean, bias preservation, or ideological flip in which the model reversed the apparent slant of

the article. We provide examples in Table 2. Both amplification and flips can pose risks, but flips are particularly serious as they can present a fundamentally different interpretation of an event than the one reflected in the original reporting. We computed the percentage of outputs per model that exhibited each type of transformation to assess likelihood. For *scope*, we again tracked identity representation in the headlines—specifically, how often LLM headlines referenced terms like Haitian, immigrant, or migrant compared to the original headline. We chose to focus on harms to this community rather than the population writ large because we felt that the latter was adequately covered by the misinformation metrics. Shifts in focus toward or away from these identity markers may reorient how readers assign blame. For example, headlines that remove references to Haitians may obscure the racialized targeting of the original claim, while those that over-emphasize it may reinforce harmful associations. Both changes can influence public perception in ways that inhibit open opinion formation.

Finally, we incorporated metrics based on the semantic content of the headlines to capture more nuanced risks to freedom of thought. These include measures of divergence from the original headline’s framing and content fidelity to the article itself. We calculated *diversity* by computing the model-wise mean pairwise cosine distance between headline embeddings, which serves as a proxy for how much variation a model introduces across prompts. Low diversity may signal ideological uniformity; high diversity may signal instability or inconsistency. We defined *framing shift* as the directional difference between the LLM-generated headline and the original headline, relative to the article text, capturing how much the model reinterprets the narrative emphasis. *Fidelity* was measured as the cosine similarity between the model headline and the full article, providing a sense of how closely the generated headline adheres to the article’s content. Lastly, we computed *tension* as the percentile rank of the original headline’s embedding distance among all generated headlines for that article. This captures how typical the actual journalistic framing is relative to model outputs.

## Results

**Human Evaluation.** To assess the reliability of our manual annotations, we calculated inter-annotator agreement on the labeled headline dataset. Across all binary categories, including correction type, identity reference, and headline framing characteristics, the average Cohen’s kappa was 0.802, indicating strong agreement between annotators. This suggests that the content distinctions used in our evaluation, such as whether a headline explicitly corrects a claim or references a particular identity, are consistently interpretable.

To scale these annotations across the full set of LLM-generated headlines, we implemented a keyword-based approach to approximate the binary labels computationally. Compared to the human-annotated gold labels, this approach achieved accuracy of 0.786, with an F1 score of 0.793. Precision and recall were well-balanced (0.774 and 0.813, respectively).

For partisanship labeling, we first evaluated inter-annotator agreement using Spearman correlation between



Figure 2: Proportion of information harms of different severities across models.

annotators’ ordinal ratings. The resulting correlation was  $\rho = 0.759$  ( $p < 0.01$ ), indicating strong consistency in how annotators positioned headlines along the ideological spectrum. To validate our scoring method, we compared model-generated partisanship scores to these human annotations. The correlation between the projection scores and annotations confirms that the embedding-based method effectively captures ideological lean ( $\rho = 0.759$ ). We further compared human-assigned partisanship to lexical heuristics, finding that lexical features aligned moderately well with human ratings for LLM-generated headlines ( $\rho = 0.444$ ) but even more strongly for original headlines ( $\rho = 0.853$ ).

**Access to Information.** First, we examine the metrics for each model’s effects on access to information. There was a strong correlation of 0.758 between LLM and original headline partisanship scores. 53.3% of cases involved explicit misinformation correction, compared to 28.2% for implicit correction and 18.5% with no correction.

We evaluated each model’s ability to handle misinformation by measuring the proportion of headlines that explicitly corrected, implicitly corrected, or failed to correct the central false claim. Claude Sonnet 3.7 demonstrated the strongest performance, with 65.8% of its headlines explicitly correcting the misinformation and only 13.8% failing to do so. In contrast, LLaMA 4 Maverick exhibited the weakest performance, with nearly a quarter (24.8%) of its headlines failing to correct the claim and less than half (48.2%) explicitly correcting it. The remaining models—GPT-4o, Gemini 2.0 Flash, and DeepSeek V3—fell in the middle, with corrected patterns showing modest but meaningful variation. These differences underscore that despite shared model capabilities, their default treatment of politically sensitive falsehoods can diverge substantially.

We next examined the scope of potential harm by analyzing how frequently identity-related terms—Haitian, mi-

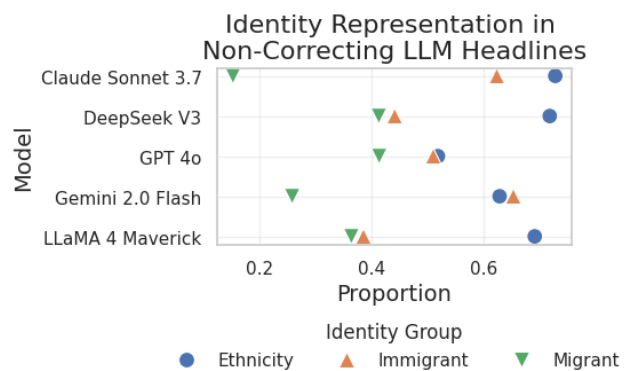


Figure 3: Rate of identity group mentions in non-correcting headlines across models.

grant, or immigrant—appeared in headlines that failed to correct the claim. This measure reflects the degree to which marginalized groups are named in contexts where misinformation is left uncorrected. Claude Sonnet 3.7 again stood out, with the highest rate of ethnicity references (72.8%) in its non-correcting headlines, suggesting that when it fails to correct misinformation, it often does so while emphasizing the targeted group. GPT-4o exhibited the lowest rate of identity references overall, with both migrant (41.4%) and ethnicity (51.9%) terms appearing less frequently in harmful outputs. In contrast, Gemini 2.0 Flash consistently emphasized identity, especially through references to immigrants (65.3%), raising concerns about disproportionate focus even in the presence of factual ambiguity.

Taken together, these results highlight two key patterns. First, no model is free from risk: even those that excel at explicit correction may still produce outputs that embed harm through framing. Second, the way models fail matters—Claude Sonnet 3.7 may reduce likelihood and scale of misinformation-related harm through higher correction rates, but it tends to coincide with amplified representation of racial and ethnic identity, increasing scope of risk. These findings suggest that model evaluations must go beyond simple correctness to include representational impact, particularly in cases involving vulnerable groups.

**Freedom of Thought.** We evaluated risk to freedom of thought through three sets of metrics: ideological transformation, identity salience shifts, and semantic divergence from editorial framing. Statistically significant amplification of partisanship occurred in 3.0% of cases, while significant flips were observed in only 0.2% of cases

Beginning with ideological shift, we measured how frequently LLM-generated headlines either amplified the partisan lean of the original headline or flipped it to the opposite ideological pole. Flip rates were low across models, but Claude Sonnet 3.7 exhibited the highest incidence of ideological reversal (0.83%). Amplification was more frequent and model-dependent: Gemini 2.0 Flash (5.76%) and LLaMA 4 Maverick (3.54%) were most likely to intensify existing partisan tone, raising concerns about reinforce-

| Explicit Correction   | Implicit Correction   | No Correction  |
|---|---|--|
| Trump and some Republicans spread false claims that migrants in Ohio are eating people’s pets, despite officials denying the reports. | Trump Repeats Debunked Claims About Haitian Immigrants During Debate                            | Findlay Mayor Says City Managing Growing Haitian Population Amid Springfield Controversy |
| Trump falsely claims immigrants are eating pets in US, Harris and debate moderator fact-check him                                     | Ohio Leaders Deny Claims of Haitian Migrants Eating Pets  | Trump’s ‘Eating the Dogs’ Claim Sparks Dip in Approval Among Focus Group Voters          |
| Ohio Leaders Debunk Baseless Claims of Haitian Immigrants Eating Pets as Trump Amplifies Falsehoods                                   | Trump Doubles Down on Debunked Claim That Immigrants “Eat Pets” During First Debate with Harris | Trump Renews Claim Haitian Migrants in Ohio Are Eating Pets                              |

Table 1: Examples of different correction types in LLM-generated headlines.

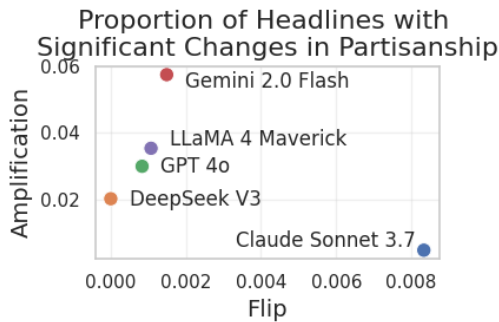


Figure 4: Comparison of partisanship amplification rates and flip rates across models.

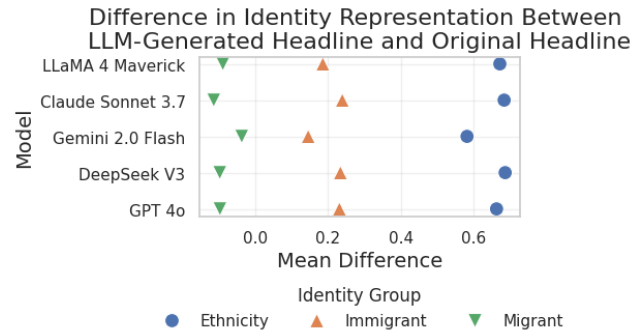


Figure 5: Change in rate of inclusion of different types of identity information between generated and original headlines.

ment of ideological framing. In contrast, Claude Sonnet 3.7 (0.48%) showed the lowest amplification, indicating a tendency to preserve the original political orientation.

To assess scope-related risks to freedom of thought, we analyzed how frequently models altered the salience of identity references—specifically, ethnicity, immigrant, and migrant—relative to the original headlines. All models increased the frequency of ethnicity references, with DeepSeek V3 (0.686), Claude (0.684), and LLaMA (0.672) showing the largest shifts. These elevated identity framings may subtly recenter group characteristics in headlines, potentially shifting reader perception of the event’s relevance or causes. Conversely, the term “migrant” consistently decreased across models, most notably in Claude (-0.113) and LLaMA (-0.088), suggesting a lexical narrowing that may bias how movement and migration are framed. These shifts reflect how LLMs reorient groups—moving headlines away from systemic or neutral framing and toward group-specific or emotionally charged terms.

We further evaluated semantic variation using metrics that capture the model’s treatment of the source article and original headline. GPT-4o and DeepSeek V3 showed the highest fidelity to full article content (0.737 and 0.741, respectively), suggesting strong alignment with original information. However, Gemini 2.0 Flash showed the highest framing shift (0.546), indicating a greater tendency to reframe editorial emphasis. Tension scores—the degree to which LLM

outputs diverged from the original headline—were highest for Claude Sonnet 3.7 (0.374), suggesting substantial reinterpretation of editorial framing despite low ideological distortion. In terms of diversity, GPT-4o produced the most variable headlines (0.134) across prompt perturbations, while DeepSeek V3 (0.043) generated the most homogeneous responses. Figure 6 depicts these metrics after normalization and inversion so that normatively better results receive higher scores.

Taken together, these results illustrate that freedom of thought risks are model-dependent and multidimensional. Claude Sonnet 3.7 maintains ideological consistency but reorients headlines semantically and emphasizes identity when it fails, resulting in interpretive drift. GPT-4o demonstrates high fidelity, low flip risk, and semantic diversity, making it comparatively well-calibrated across dimensions. Gemini 2.0 Flash is the most ideologically volatile, prone to partisan amplification and framing shifts that may subtly reshape audience interpretation. Meanwhile, DeepSeek V3 shows strong fidelity and ideological stability but consistent identity emphasis, suggesting persistent latent framing bias. LLaMA 4 Maverick presents a more conservative editorial posture, but its low misinformation correction rate, coupled with high identity emphasis and amplification, raises concerns about latent bias and representational distortions.

| Movement      | Original   | LLM  |
|---------------|--|--|
| Amplification | At debate, Trump shares falsehoods about pet-eating, infanticide                       | Trump Makes False Claims, Spreads Viral Misinformation During Debate With Harris; Moderators Fact-Check Statements |
| Flip          | Springfield Haitians, other residents respond to being at center of immigration debate | Trump’s Unsubstantiated Claims About Haitian Immigrants Spark Controversy in Springfield                           |
| Preservation  | ‘They’re Eating the Cats’: Trump Repeats False Claim About Immigrants                  | Trump Repeats Baseless Claim of Haitian Immigrants Eating Pets in Ohio   |

Table 2: Examples of how LLM-generated headlines shift narratives from original reporting.

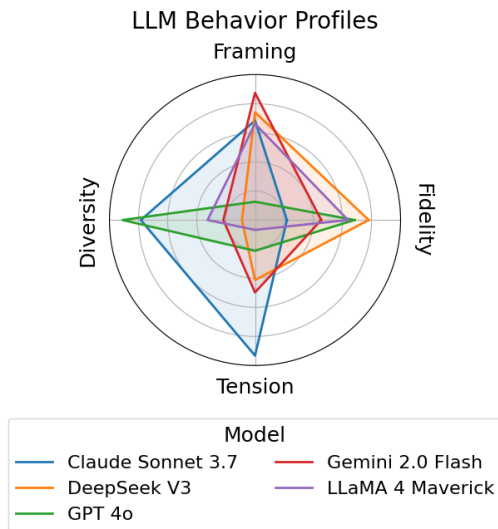


Figure 6: Comparison of all behavior profile metrics together. Inverted so that higher = better for freedom of thought, original metrics in Table 6.

tion. These findings highlight the importance of rights-aware model evaluations that go beyond surface-level correctness to capture deeper narrative and ideological transformations.

## Discussion

### Interpretation

Metrics alone do not constitute claims about human rights risk, but they provide structured evidence for making them. In our framework, model behavior is interpreted through a lens grounded in the UN Guiding Principles on Business and Human Rights, allowing observed patterns to be translated into normative claims. The strength of such claims depends on both metric reliability and the situational context of deployment. For example, a high frequency of non-correction is not automatically harmful, but when it co-occurs with identity amplification in a politically sensitive story, it raises meaningful concerns. Each metric corresponds to a specific risk factor—scale, scope, or likelihood—and together they inform whether the system crosses a threshold where the use of an LLM plausibly undermines a protected right. Remediability is incorporated after the metric stage because it

is grounded in interaction with the system rather than system outputs themselves. The translation of metrics to risk saliency makes it possible to move from computational outputs to policy-relevant judgments.

The right of access to information is most directly threatened when LLMs fail to correct misinformation or reproduce it in misleading or diluted forms. In this task, Claude performed best in terms of explicit correction, but when it failed, it often did so while amplifying identity references, which may distort access to relevant context. LLaMA had the weakest correction performance, paired with elevated identity representation. Considering this, we assess the overall risk to access to information posed by LLM usage in this task, acknowledging variability across models: scale is medium, since the misinformation is highly consequential but generally not repeated verbatim; scope is high, due to disproportionate group representation in outputs that fail to correct information; and likelihood as medium, based on the proportion of headlines that failed to correct the falsehood. Remediability is limited in real-world settings because users often see headlines through social media and may never see corrections. This makes it important to prevent false headlines early, through filtering or editorial review.

Freedom of thought is compromised when model-generated headlines shape perception by emphasizing particular identities or shifting ideological tone. Flip rates were low across all models, but amplification of partisan framing—especially by Gemini—poses a credible risk of skewing user interpretation to more extreme views. References to identity frequently increased, even in otherwise ideologically neutral outputs. These shifts can affect how readers interpret the causes and significance of political events, particularly in stories with racialized or polarizing content. Across all models, we assess scale as high, given the right to access information’s centrality to cognitive autonomy; scope as medium, since not all outputs exhibit ideological shift, but mention of identity is common; and likelihood as high. Remediability is limited because once a headline shapes a reader’s initial understanding, it is hard to change that perception, especially when the headline may be the only part of the content they see.

Our analysis shows that LLM-generated headlines in political journalism pose meaningful human rights risks. While severe harms are rare, consistent lower-level issues, such as subtle framing bias or identity amplification, remain concerning, especially since readers may overlook these cues

and still be influenced. Importantly, risk stems not only from factual accuracy but from how models frame and prioritize information. These findings highlight the importance of context-aware, multi-dimensional, and normatively grounded risk assessments, particularly in sensitive public communication domains.

### Generalizing Takeaways

Translating abstract rights into quantifiable metrics necessarily requires contextual judgment, as human rights risks are not static properties of models but emergent properties of sociotechnical systems. In our case, the risks posed by LLM-generated headlines are inseparable from the political and racialized instance of misinformation we studied. That is, failure to correct a false claim in a story about tax policy (for instance) may not carry the same weight as failure to correct a false claim in a story that targets a marginalized group during a high-profile political event. This context-dependence means that while the metrics we use (including correction rate, presence of identity terms, and ideological shift) provide structured indicators of risk, they do not generalize unproblematically across domains. They must be interpreted in light of the task, the audience, and the social conditions under which model outputs circulate.

Thus, a central challenge of our framework is its limited generalizability. Because our evaluation was specifically situated within the context of U.S.-based political journalism—focusing narrowly on English-language headline generation during a political event—both the metrics we developed and the rights analyses we performed reflect highly context-specific judgments. Our operationalization of human rights risks such as misinformation propagation and partisan framing is strongly tied to editorial standards and professional norms within journalism; consequently, these metrics may not transfer neatly to other sensitive domains like healthcare, legal decision-making, or educational settings without substantial adaptation.

However, this narrow focus is by design. The risk saliency criteria outlined in the UN Guiding Principles (UNGPs) and operationalized through our model evaluation framework are intended to be used by corporations when evaluating their practices and products. As such, the resulting conclusions do not need to generalize beyond the given context, as long as they are synthesized and acted upon responsibly by the corporation. Furthermore, the first stage of the framework requires evaluators to select a rights-impacting context and task. While the metrics used to evaluate the model may only focus on a subset of rights, this is because these rights are the primary ones at risk; creating metrics for others would have minimal value.

### Operationalizing Risk to Rights

Designing metrics for rights-based auditing requires navigating tradeoffs between granularity and interpretability. Fine-grained distinctions, such as explicit versus implicit misinformation correction, allow us to capture subtle and cumulative harms. However, they rely on subjective thresholds and assumptions about what constitutes “meaningful” deviation. Similarly, our use of semantic embedding-based par-

tisanship scoring allows for scalable, model-agnostic measurement of ideological shift, but it abstracts away from the mechanisms through which bias is expressed. These techniques detect difference, not intent, and must be anchored in a normative framework to support meaningful claims. Moreover, while we use original headlines as a point of comparison for assessing framing shift and tension, we acknowledge that journalistic baselines are themselves shaped by institutional bias and editorial judgment. Deviation from the original is not always risk-enhancing—nor is alignment always protective. What matters is whether that shift increases the likelihood of informational harm or undermines interpretive autonomy, impacting the rights to access information and freedom of thought.

Crucially, many of the most consequential risks emerge not from any single metric but from the interaction of multiple metrics. A headline that fails to correct a false claim is concerning; one that does so while amplifying reference to ethnic identity is substantially more so. Our framework is designed to support such layered analysis by aligning each metric with one or more dimensions of the UNGPs. This structure enables us to translate measurement into actionable rights-based claims without collapsing complexity into a single score. The goal of operationalization is not to produce a definitive classification of harm, but to foreground those risks that might otherwise remain obscured. In this way, structured evaluation becomes a form of accountability—making the normative implications of generative model behavior legible, contestable, and subject to intervention.

### Limitations and Future Work

Our focus on a U.S. and English-language context aligns with cultural and linguistic settings in which large language models tend to perform more reliably, limiting the applicability of our model evaluation to other geopolitical environments. Furthermore, translating nuanced human rights concepts into quantifiable metrics inevitably involves some reductionism. While embedding-based semantic methods enabled scalable and model-agnostic measurement, they may not fully capture the partisan leanings of headlines about a specific event, as opposed to a general topic like immigration, which is better represented within the embedding model. While we manually validated the relationship between lexical cues and scores, the embedding model may also fail with subtle language shifts in headlines.

Addressing these limitations provides several clear avenues for future work. First, extending the framework through comparative case studies in diverse contexts—such as medical misinformation, legal text summarization, or educational content generation—could allow for adaptations that enhance generalizability. This would involve developing more modular, context-flexible metrics. Second, longitudinal and dynamic studies that track how misinformation and harmful framings evolve and compound over time would greatly enrich the ecological validity of risk assessments. Finally, future methodological improvements could explore combining semantic embedding approaches with other techniques to clarify the pathways from subtle lexical variation to rights impacts.

## Ethics Statement

While our framework aims to provide a structured and normatively grounded approach to assessing human rights risks in generative AI, it introduces several ethical considerations that warrant careful reflection.

**Overreliance on Quantification.** The translation of complex human rights concerns into measurable metrics risks reducing nuanced social, political, and contextual harms to simplified indicators. While metrics enable systematic evaluation, they may obscure or exclude harms that are difficult to quantify, such as subtle forms of psychological or cultural erosion. Overemphasizing numerical indicators could reinforce a narrow view of risk, privileging what is easy to measure over what is substantively meaningful.

**Misuse or Misinterpretation of Metrics.** Metrics developed through our framework are not intended to serve as standalone or definitive judgments of harm. However, there is a risk that stakeholders—especially those seeking operational simplicity—may apply them out of context or interpret them as exhaustive. This could lead to mischaracterization of a model’s risk profile or to an undue focus on compliance over reflection, diminishing the role of qualitative analysis and stakeholder input.

**Risk of Institutional Co-optation.** There is also the potential for the framework to be co-opted by organizations aiming to demonstrate ethical AI practices without enacting meaningful change. If adopted primarily for public relations or regulatory box-checking, the framework could enable surface-level interventions that deflect from more systemic issues, such as power asymmetries, data exploitation, or the political economy of AI development. To guard against this, we emphasize that the framework is a tool for reflection and accountability—not a substitute for broader ethical governance or institutional reform.

Together, these concerns underscore the need for responsible use of the framework: as a guide for deeper inquiry rather than a deterministic checklist, and as a supplement to—not a replacement for—engagement with affected communities, contextual knowledge, and deliberative ethics.

## Positionality Statement

Our team is made up of individuals with expertise in computer science, law, and the social sciences. We are based in well-resourced institutions in cosmopolitan areas around the world. While we have endeavored to develop a broad, geopolitically agnostic approach to evaluating model impacts, we recognize that we are limited by cultural assumptions about the role and scope of journalism and human rights.

## Acknowledgements

Many thanks to Eve Fleisig, Lucy Li, and Gisela Perez de Acha for their insight during the early conception of this work.

## References

- Abebe, R.; Hardt, M.; Jin, A.; Miller, J.; Schmidt, L.; and Wexler, R. 2022. Adversarial Scrutiny of Evidentiary Statistical Software. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, FAccT ’22, 1733–1746. New York, NY, USA: Association for Computing Machinery. ISBN 9781450393522.
- Aizenberg, E.; and van den Hoven, J. 2020. Designing for human rights in AI. *Big Data & Society*, 7(2): 2053951720949566.
- An, J.; Kwak, H.; and Ahn, Y.-Y. 2018. SemAxis: A Lightweight Framework to Characterize Domain-Specific Word Semantics Beyond Sentiment. In Gurevych, I.; and Miyao, Y., eds., *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2450–2461. Melbourne, Australia: Association for Computational Linguistics.
- Bang, Y.; Chen, D.; Lee, N.; and Fung, P. 2024. Measuring Political Bias in Large Language Models: What Is Said and How It Is Said. In Ku, L.-W.; Martins, A.; and Sriku-mar, V., eds., *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 11142–11159. Bangkok, Thailand: Association for Computational Linguistics.
- Beckett, C.; and Yaseen, M. 2023. Generating Change: A Global Survey of What News Organisations Are Doing with AI. Technical report, JournalismAI, Polis, London School of Economics and Political Science. Accessed: 2025-05-23.
- Binns, R. 2018. Algorithmic Accountability and Public Reason. *Philosophy Technology*, 31(4): 543–556.
- Blodgett, S. L.; Barocas, S.; Daumé III, H.; and Wallach, H. 2020. Language (Technology) is Power: A Critical Survey of “Bias” in NLP. In Jurafsky, D.; Chai, J.; Schluter, N.; and Tetreault, J., eds., *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 5454–5476. Online: Association for Computational Linguistics.
- Bolukbasi, T.; Chang, K.-W.; Zou, J.; Saligrama, V.; and Kalai, A. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, NIPS’16, 4356–4364. Red Hook, NY, USA: Curran Associates Inc. ISBN 9781510838819.
- Buolamwini, J.; and Gebru, T. 2018. Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification. In Friedler, S. A.; and Wilson, C., eds., *Proceedings of the 1st Conference on Fairness, Accountability and Transparency*, volume 81 of *Proceedings of Machine Learning Research*, 77–91. PMLR.
- Business for Social Responsibility. 2022. FAQ on Human Rights Assessments. Accessed: 2025-05-22.
- Casper, S.; Ezell, C.; Siegmann, C.; Kolt, N.; Curtis, T. L.; Bucknall, B.; Haupt, A.; Wei, K.; Scheurer, J.; Hobbhahn, M.; Sharkey, L.; Krishna, S.; Von Hagen, M.; Alberti, S.; Chan, A.; Sun, Q.; Gerovitch, M.; Bau, D.; Tegmark, M.;

- Krueger, D.; and Hadfield-Menell, D. 2024. Black-Box Access is Insufficient for Rigorous AI Audits. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '24, 2254–2272. New York, NY, USA: Association for Computing Machinery. ISBN 9798400704505.
- Chabot, C.; Raman, V.; and Popken, B. 2025. Large Language Models in a Global Context: An Interdisciplinary Analysis of Human Rights Impacts in Law, Education, and Journalism. White paper, UC Berkeley Human Rights Center.
- Chen, G. H.; Chen, S.; Liu, Z.; Jiang, F.; and Wang, B. 2024. Humans or LLMs as the Judge? A Study on Judgement Bias. In AI-Onaizan, Y.; Bansal, M.; and Chen, Y.-N., eds., *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, 8301–8327. Miami, Florida, USA: Association for Computational Linguistics.
- Chu, T.; Song, Z.; and Yang, C. 2024. How to Protect Copyright Data in Optimization of Large Language Models? *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(16): 17871–17879.
- Cobbe, J.; Lee, M. S. A.; and Singh, J. 2021. Reviewable Automated Decision-Making: A Framework for Accountable Algorithmic Systems. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '21, 598–609. New York, NY, USA: Association for Computing Machinery. ISBN 9781450383097.
- Cools, H.; and Diakopoulos, N. 2024. Uses of Generative AI in the Newsroom: Mapping Journalists' Perceptions of Perils and Possibilities. *Journalism Practice*.
- Costanza-Chock, S.; Raji, I. D.; and Buolamwini, J. 2022. Who Audits the Auditors? Recommendations from a field scan of the algorithmic auditing ecosystem. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '22, 1571–1583. New York, NY, USA: Association for Computing Machinery. ISBN 9781450393522.
- Deck, A. 2023. AI translation is jeopardizing Afghan asylum claims. *Rest of World*.
- Edgerly, S.; Mourão, R. R.; Thorson, E.; and Tham, S. M. 2020. When Do Audiences Verify? How Perceptions About Message and Source Influence Audience Verification of News Headlines. *Journalism & Mass Communication Quarterly*, 97(1): 52–71.
- Entman, R. M. 1993. Framing: Toward clarification of a fractured paradigm. *Journal of Communication*, 43(4): 51–58.
- European Parliament and Council. 2024. Artificial Intelligence Act: Regulation (EU) 2024/1689 of the European Parliament and of the Council laying down harmonised rules on artificial intelligence and amending certain Union legislative acts. <https://artificialintelligenceact.eu/>.
- Fleisig, E.; Smith, G.; Bossi, M.; Rustagi, I.; Yin, X.; and Klein, D. 2024. Linguistic Bias in ChatGPT: Language Models Reinforce Dialect Discrimination. In AI-Onaizan, Y.; Bansal, M.; and Chen, Y.-N., eds., *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, 13541–13564. Miami, Florida, USA: Association for Computational Linguistics.
- Gabriel, I. 2020. Artificial Intelligence, Values, and Alignment. *Minds and Machines*, 30(3): 411–437.
- Gebru, T.; Morgenstern, J.; Vecchione, B.; Vaughan, J. W.; Wallach, H.; Iii, H. D.; and Crawford, K. 2021. Datasheets for datasets. *Communications of the ACM*, 64(12): 86–92.
- Geer, J. G.; and Kahn, K. F. 1993. Grabbing attention: An experimental investigation of headlines during campaigns. *Political Communication*, 10(2): 175–191.
- Hagar, N.; and Diakopoulos, N. 2021. How A/B testing can (and can't) improve your headline writing. *Nieman Journalism Lab*. Accessed: 2025-05-23.
- Havard, C. T.; Ferrucci, P.; and Ryan, T. G. 2021. Does messaging matter? Investigating the influence of media headlines on perceptions and attitudes of the in-group and out-group. *Journal of Marketing Communications*, 27(1): 20–30.
- Hofmann, V.; Kalluri, P. R.; Jurafsky, D.; and King, S. 2024. AI generates covertly racist decisions about people based on their dialect. *Nature*, 633: 147–154.
- Hoh, J.; Nigam, S.; Andersen, L.; and Darnton, H. 2025. Human Rights Across the Generative AI Value Chain. Includes Responsible AI Practitioner Guides.
- Huang, H.; Zhu, H.; Liu, W.; Gao, H.; Jin, H.; and Liu, B. 2024a. Uncovering the essence of diverse media biases from the semantic embedding space. *Humanities and Social Sciences Communications*, 11(1): 656.
- Huang, S.; Siddarth, D.; Lovitt, L.; Liao, T. I.; Durmus, E.; Tamkin, A.; and Ganguli, D. 2024b. Collective Constitutional AI: Aligning a Language Model with Public Input. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '24, 1395–1417. New York, NY, USA: Association for Computing Machinery. ISBN 9798400704505.
- Kalluri, P. R.; Agnew, W.; Cheng, M.; Owens, K.; Soldaini, L.; and Birhane, A. 2023. The Surveillance AI Pipeline. arXiv:2309.15084.
- Li, J.; Cheng, X.; Zhao, X.; Nie, J.-Y.; and Wen, J.-R. 2023. HaluEval: A Large-Scale Hallucination Evaluation Benchmark for Large Language Models. In Bouamor, H.; Pino, J.; and Bali, K., eds., *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, 6449–6464. Singapore: Association for Computational Linguistics.
- Li, N.; Pan, A.; Gopal, A.; Yue, S.; Berrios, D.; Gatti, A.; Li, J. D.; Dombrowski, A.-K.; Goel, S.; Mukobi, G.; Helmburger, N.; Lababidi, R.; Justen, L.; Liu, A. B.; Chen, M.; Barrass, I.; Zhang, O.; Zhu, X.; Tamirisa, R.; Bharathi, B.; Herbert-Voss, A.; Breuer, C. B.; Zou, A.; Mazeika, M.; Wang, Z.; Oswal, P.; Lin, W.; Hunt, A. A.; Tienken-Harder, J.; Shih, K. Y.; Talley, K.; Guan, J.; Steneker, I.; Campbell, D.; Jokubaitis, B.; Basart, S.; Fitz, S.; Kumaraguru, P.; Karmakar, K. K.; Tupakula, U.; Varadharajan, V.; Shoshitaishvili, Y.; Ba, J.; Esvelt, K. M.; Wang, A.; and Hendrycks,

- D. 2024. The WMDP Benchmark: Measuring and Reducing Malicious Use with Unlearning. In Salakhutdinov, R.; Kolter, Z.; Heller, K.; Weller, A.; Oliver, N.; Scarlett, J.; and Berkenkamp, F., eds., *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, 28525–28550. PMLR.
- Lloyd, T.; Reagle, J.; and Naaman, M. 2024. "There Has To Be a Lot That We're Missing": Moderating AI-Generated Content on Reddit. arXiv:2311.12702.
- Magesh, V.; Surani, F.; Dahl, M.; Suzgun, M.; Manning, C. D.; and Ho, D. E. 2025. Hallucination-Free? Assessing the Reliability of Leading AI Legal Research Tools. *Journal of Empirical Legal Studies*, 22(2): 123–145.
- Mantelero, A. 2018. AI and Big Data: A blueprint for a human rights, social and ethical impact assessment. *Computer Law Security Review*, 34(4): 754–772.
- McGregor, L.; Murray, D.; and Ng, V. 2019. INTERNATIONAL HUMAN RIGHTS LAW AS A FRAMEWORK FOR ALGORITHMIC ACCOUNTABILITY. *International and Comparative Law Quarterly*, 68(2): 309–343.
- Mitchell, M.; Wu, S.; Zaldívar, A.; Barnes, P.; Vasserman, L.; Hutchinson, B.; Spitzer, E.; Raji, I. D.; and Gebru, T. 2019. Model Cards for Model Reporting. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, FAT\* '19, 220–229. New York, NY, USA: Association for Computing Machinery. ISBN 9781450361255.
- Mökander, J.; Schuett, J.; Kirk, H. R.; and Floridi, L. 2024. Auditing large language models: a three-layered approach. *AI Ethics*, 4: 1085–1115.
- Mozilla Foundation. 2020. RegretsReporter: Investigating YouTube's Recommendation Algorithm. Accessed: 2025-05-22.
- Nadeem, M.; Bethke, A.; and Reddy, S. 2021. StereoSet: Measuring stereotypical bias in pretrained language models. In Zong, C.; Xia, F.; Li, W.; and Navigli, R., eds., *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)*, 5356–5371. Online: Association for Computational Linguistics.
- Ouyang, L.; Wu, J.; Jiang, X.; Almeida, D.; Wainwright, C. L.; Mishkin, P.; Zhang, C.; Agarwal, S.; Slama, K.; Ray, A.; Schulman, J.; Hilton, J.; Kelton, F.; Miller, L.; Simens, M.; Askell, A.; Welinder, P.; Christiano, P.; Leike, J.; and Lowe, R. 2022. Training language models to follow instructions with human feedback. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22. Red Hook, NY, USA: Curran Associates Inc. ISBN 9781713871088.
- Prabhakaran, V.; Mitchell, M.; Gebru, T.; and Gabriel, I. 2022. A Human Rights-Based Approach to Responsible AI. arXiv:2210.02667.
- Raji, D.; Denton, E.; Bender, E. M.; Hanna, A.; and Paullada, A. 2021. AI and the Everything in the Whole Wide World Benchmark. In Vanschoren, J.; and Yeung, S., eds., *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, volume 1.
- Raji, I. D.; Smart, A.; White, R. N.; Mitchell, M.; Gebru, T.; Hutchinson, B.; Smith-Loud, J.; Theron, D.; and Barnes, P. 2020. Closing the AI accountability gap: defining an end-to-end framework for internal algorithmic auditing. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, FAT\* '20, 33–44. New York, NY, USA: Association for Computing Machinery. ISBN 9781450369367.
- Raji, I. D.; Xu, P.; Honigsberg, C.; and Ho, D. 2022. Outsider Oversight: Designing a Third Party Audit Ecosystem for AI Governance. In *Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society*, AIES '22, 557–571. New York, NY, USA: Association for Computing Machinery. ISBN 9781450392471.
- Rodrigues, R. 2020. Legal and human rights issues of AI: Gaps, challenges and vulnerabilities. *Journal of Responsible Technology*, 4: 100005.
- Sampath, K.; Pratheesh; Mohammad, A.; and Ramachandranpillai, R. 2025. The Multimodal Paradox: How Added and Missing Modalities Shape Bias and Performance in Multimodal AI. arXiv:2505.03020.
- Sandvig, C.; Hamilton, K.; Karahalios, K.; and Langbort, C. 2014. Auditing Algorithms: Research Methods for Detecting Discrimination on Internet Platforms. *Data and Discrimination: Converting Critical Concerns into Productive Inquiry*, 22(2014): 4349–4357.
- Shah, C.; and Bender, E. M. 2024. Envisioning Information Access Systems: What Makes for Good Tools and a Healthy Web? *ACM Trans. Web*, 18(3).
- Shelton, D. L. 2007. An Introduction to the History of International Human Rights Law. *GWU Legal Studies Research Paper No. 346*; *GWU Law School Public Law Research Paper No. 346*. Lectures delivered at the International Institute of Human Rights, Strasbourg, France, July 2003.
- Society of Professional Journalists. 2014. SPJ Code of Ethics. <https://www.spj.org/spj-code-of-ethics/>. Accessed: 2025-05-23.
- Sorensen, T.; Moore, J.; Fisher, J.; Gordon, M.; Miresghallah, N.; Rytting, C. M.; Ye, A.; Jiang, L.; Lu, X.; Dziri, N.; Althoff, T.; and Choi, Y. 2024. Position: a roadmap to pluralistic alignment. In *Proceedings of the 41st International Conference on Machine Learning*, ICML'24. JMLR.org.
- Suresh, H.; Tseng, E.; Young, M.; Gray, M.; Pierson, E.; and Levy, K. 2024. Participation in the age of foundation models. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '24, 1609–1621. New York, NY, USA: Association for Computing Machinery. ISBN 9798400704505.
- Talent Fabric LTD. 2025. GoPerfect: AI Recruitment & Hiring Software.
- The American Presidency Project. 2024. Presidential Debate in Philadelphia, Pennsylvania. <https://www.presidency.ucsb.edu/documents/presidential-debate-philadelphia-pennsylvania>. Accessed: 2025-05-23.
- Thomson Reuters. 2024. CoCounsel: One GenAI Assistant for Professionals.

Tibebu, H.; and Kakadiaris, I. 2025. System Card+: Responsible AI Framework for Decision Support Systems. Zenodo.

United Nations General Assembly. 1948. Universal Declaration of Human Rights. Adopted by the United Nations General Assembly resolution 217 A (III).

United Nations General Assembly. 1966a. International Covenant on Civil and Political Rights. Adopted by UN General Assembly resolution 2200A (XXI); entered into force 23 March 1976.

United Nations General Assembly. 1966b. International Covenant on Economic, Social and Cultural Rights. Adopted by UN General Assembly resolution 2200A (XXI); entered into force 3 January 1976.

United Nations Human Rights Council. 2011. Guiding Principles on Business and Human Rights: Implementing the United Nations “Protect, Respect and Remedy” Framework. Endorsed by UN Human Rights Council resolution 17/4 on 16 June 2011.

Vecchione, B.; Levy, K.; and Barocas, S. 2021. Algorithmic Auditing and Social Justice: Lessons from the History of Audit Studies. In *Proceedings of the 1st ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization*, EAAMO '21. New York, NY, USA: Association for Computing Machinery. ISBN 9781450385534.

Wallach, H.; Desai, M.; Cooper, A. F.; Wang, A.; Atalla, C.; Barocas, S.; Blodgett, S. L.; Chouldechova, A.; Corvi, E.; Dow, P. A.; Garcia-Gathright, J.; Olteanu, A.; Pangakis, N.; Reed, S.; Sheng, E.; Vann, D.; Vaughan, J. W.; Vogel, M.; Washington, H.; and Jacobs, A. Z. 2025. Position: Evaluating Generative AI Systems is a Social Science Measurement Challenge. arXiv:2502.00561.

Weidinger, L.; Raji, D.; Wallach, H.; Mitchell, M.; Wang, A.; Salaudeen, O.; Bommasani, R.; Kapoor, S.; Ganguli, D.; Koyejo, S.; and Isaac, W. 2025. Toward an Evaluation Science for Generative AI Systems. *The Bridge*, 55(2): 15–22.

Wilson, K.; and Caliskan, A. 2024. Gender, Race, and Intersectional Bias in Resume Screening via Language Model Retrieval. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society (AIES)*.

Xiang, A. 2024. Fairness and Privacy in an Age of Generative AI. *Science and Technology Law Review*, 25(2).

Zhang, A. K.; Perry, N.; Dulepet, R.; Ji, J.; Menders, C.; Lin, J. W.; Jones, E.; Hussein, G.; Liu, S.; Jasper, D.; Peethawatthachai, P.; Glenn, A.; Sivashankar, V.; Zamoshchin, D.; Glikbarg, L.; Askaryar, D.; Yang, M.; Zhang, T.; Al-luri, R.; Tran, N.; Sangpisit, R.; Yiorkadjis, P.; Osele, K.; Raghupathi, G.; Boneh, D.; Ho, D. E.; and Liang, P. 2025. Cybench: A Framework for Evaluating Cybersecurity Capabilities and Risks of Language Models. arXiv:2408.08926.

Zhao, D.; Wang, A.; and Russakovsky, O. 2021. Understanding and Evaluating Racial Biases in Image Captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 14830–14840.