

AI Governance and Accountability: An Analysis of Anthropic's Claude

Aman Priyanshu

*Privacy Engineering,
School of Computer Science,
Carnegie Mellon University
Email: apriyans@andrew.cmu.edu*

Yash Maurya

*Privacy Engineering,
School of Computer Science,
Carnegie Mellon University
Email: ymaurya@andrew.cmu.edu*

Zuofei Hong

*Privacy Engineering,
School of Computer Science,
Carnegie Mellon University
Email: zuofeih@andrew.cmu.edu*

Abstract—As AI systems become increasingly prevalent and impactful, the need for effective AI governance and accountability measures is paramount. This paper examines the AI governance landscape, focusing on Anthropic's Claude, a foundational AI model. We analyze Claude through the lens of the NIST AI Risk Management Framework and the EU AI Act, identifying potential threats and proposing mitigation strategies. The paper highlights the importance of transparency, rigorous benchmarking, and comprehensive data handling processes in ensuring the responsible development and deployment of AI systems. We conclude by discussing the social impact of AI governance and the ethical considerations surrounding AI accountability.

1. Introduction

Artificial Intelligence (AI) has become an integral part of modern society, pervading diverse domains from complex computational tasks and generating reports to mass communication, hiring decisions, and marketing efforts. As AI systems continue to grow in sophistication and influence, their impact expands across numerous spheres, shaping decision-making processes, information dissemination, and human interactions on an unprecedented scale. In this new era of AI, foundation models have assumed a significant role. Models, such as Anthropic's Claude, a large language model (LLM) capable of understanding and generating human-like text, exhibit unique potential for quick, effective, and scalable communication efforts.



Figure 1. Anthropic's Claude is one of the most popular large language model chatbots available to the everyday consumer. This paper presents a study of its practices and conduct through the lens of AI governance.

The customer reach of these LLMs, like Claude, has drastically increased over the years, and their influence on individuals' lives is expected to continue growing. Anthropic has announced several partnerships with prominent companies such as Scale [1], Zoom [2], BCG [3], AWS [4], Accenture [4], SKT Telecom [5], and Keif Studio [6], further amplifying the impact of their AI systems on people's lives, often without their knowledge of interacting with an AI system.

These LLMs are crucial as they underpin many AI systems, influencing outcomes and decision-making processes in areas that directly affect individuals and societies. Due to their unprecedented potential for impact, these foundation models must be evaluated for the risks and challenges they may pose to society.

These challenges motivate the need for AI governance - the processes, policies, and practices aimed at ensuring the responsible development, deployment,



Figure 2. Rapidly growing customer visits on their Claude’s web interface

and use of AI systems. The importance of AI governance lies in its ability to ensure the responsible development and deployment of AI systems, safeguarding against potential harms and unintended consequences. Accountability is a key aspect of AI governance, as it helps establish trust and ensures that AI systems are designed and used in an ethical and transparent manner. Frameworks such as the NIST AI Risk Management Framework and the EU AI Act provide guidelines and standards for assessing, categorizing, and managing AI risks. These frameworks enable stakeholders to develop appropriate governance measures by offering structured approaches to identify, analyze, and mitigate potential threats associated with AI systems.

In this paper, we analyze Anthropic’s Claude through the lens of these frameworks, identifying potential threats and proposing mitigation strategies. We also focus on their Constitutional AI paradigm. By examining Claude as not only an AI product but a foundational model, we aim to provide insights that can inform the broader AI governance discourse and contribute to the responsible advancement of AI technologies. The key objectives of this paper are:

- 1) To analyze Anthropic’s Claude through established AI governance frameworks like:
 - NIST
 - EU AI Act
- 2) To identify potential threats and risks posed by Claude

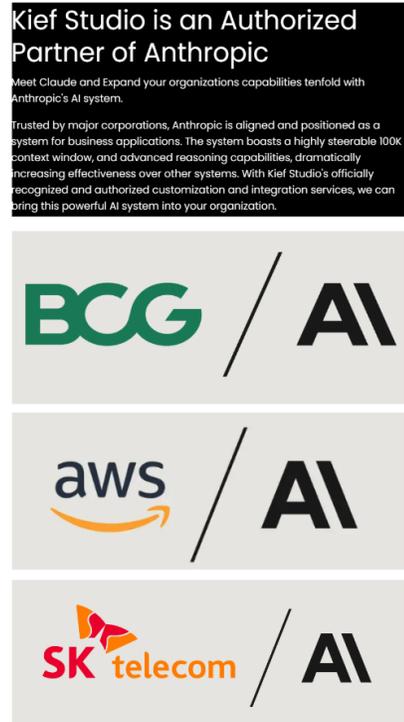


Figure 3. Some of Anthropic’s Partnerships

- 3) To propose mitigation strategies for these risks
- 4) To examine Anthropic’s Constitutional AI paradigm
- 5) To provide insights for broader AI governance discourse

2. Organization of this Report

This report is organized into several key sections to provide a comprehensive analysis of AI governance and accountability, with a focus on Anthropic’s Claude model. The **introduction** sets the stage by highlighting the growing importance of AI governance as AI systems become increasingly prevalent and influential in various domains. It emphasizes the role of foundation models, such as Claude, in shaping decision-making processes and the need for effective governance measures to ensure responsible AI development and deployment.

The **literature review** section explores the current state of AI governance, discussing various frameworks

and guidelines, such as the NIST AI Risk Management Framework and the EU AI Act. It also dives into recent literature on AI governance themes, knowledge gaps, and future agendas, as well as the challenges arising from the growing number of AI ethics documents produced by corporations, governments, and NGOs.

The **preliminaries** section provides a brief overview of key concepts, including artificial intelligence, large language models, and Anthropic’s Claude. It also introduces Constitutional AI, a framework employed by Anthropic to align the model’s outputs with predefined ethical principles and values.

The **threat analysis** section forms the main crux of the report, identifying and discussing potential threats and issues associated with Claude. This section focuses on specific risks, such as the lack of transparency in privacy policies, potential for hallucinations and biases in outputs, concerns about third-party data usage, and the implications of Constitutional AI. The analysis is conducted through the lens of the **NIST AI Risk Management Framework**, examining aspects of governance, risk mapping, and impact characterization. Additionally, the **EU AI Act** is used to categorize the identified risks based on their severity and potential consequences.

Building upon the threat analysis, the report proposes **mitigation strategies and resolution approaches** to address the identified risks. These strategies include enhancing transparency in privacy policies, establishing rigorous benchmarks for hallucination and bias, and developing comprehensive remediation processes for data deletion and model unlearning. The discussion section explores the broader implications of these mitigation strategies for the AI governance landscape and the social impact of AI systems.

The **conclusion** summarizes the key findings and highlights the importance of ongoing collaboration, adaptation, and learning in the evolution of AI governance. It emphasizes the need for aligning AI systems with ethical principles and societal values to foster public trust and support the responsible advancement of AI technologies.

Finally, the report acknowledges its **limitations and discusses ethical considerations** in the development and deployment of AI systems. It stresses the importance of prioritizing ethical principles throughout the AI lifecycle and engaging in ongoing research and stakeholder collaboration to address the ethical implications of AI and develop robust governance frameworks.

We organize this report in this manner, to provide a comprehensive and structured analysis of AI governance and accountability, focusing on the specific risks associated with Anthropic’s Claude model (plus its Constitutional AI efforts) and propose actionable mitigation strategies to ensure responsible AI development and deployment.

3. Literature Review

AI governance has gained significant attention in recent years, with various frameworks and guidelines proposed to address the risks and challenges associated with AI systems. The NIST AI Risk Management Framework [7] provides a comprehensive approach to identifying, assessing, and managing AI risks, emphasizing the importance of governance, risk mapping, and impact characterization. Similarly, the EU AI Act [8] categorizes AI systems based on their risk levels, imposing specific requirements and obligations for high-risk systems.

Beyond regulations and risk frameworks, recent literature has further explored the themes, knowledge gaps, and future agendas in AI governance [9]. Key themes identified include technology, stakeholders and context, regulation, and processes. However, knowledge gaps remain, such as limited understanding of AI governance implementation, lack of attention to context, uncertain effectiveness of ethical principles and regulation, and insufficient operationalization of processes [10]. The growing number of AI ethics documents produced by corporations, governments, and NGOs since 2016 raises important considerations [10], [11]. Challenges may arise from the relative homogeneity of the documents’ creators, and the var-

ied impacts and success factors of these documents on the AI governance landscape warrant examination [11]. Translating AI ethical principles into practicable governance processes is crucial, and a concise AI governance definition can help identify the constituent parts of this complex problem [12].

Previous works on AI auditing and self-governance have highlighted the need for transparency, accountability, and continuous monitoring of AI systems. For example, Raji et al. [13] propose a framework for closing the AI accountability gap, emphasizing the importance of external audits and stakeholder engagement. Additionally, the development of privacy regulations, such as the General Data Protection Regulation (GDPR) [14], [15], has demonstrated the importance of proactive measures and the need for ongoing adaptation to address emerging risks.

As AI governance continues to evolve, it is essential to learn from the successes and challenges of privacy regulations and apply these lessons to the development of AI accountability measures. We study Claude, one of the most popular AI models, through the lens of prior literature and recommended frameworks, as it has the capacity for large-scale harm if not studied ethically [16], [17].

4. Preliminaries

Artificial Intelligence (AI) refers to the development of computer systems that can perform tasks that typically require human intelligence, such as visual perception, speech recognition, decision-making, and language translation. AI encompasses various sub-fields, including machine learning, natural language processing, and computer vision.

Large Language Models (LLMs) are a type of AI model that have gained significant attention in recent years. LLMs are trained on vast amounts of text data, enabling them to generate human-like text, answer questions, and perform various language-related tasks. These models, such as OpenAI's GPT series [18] and Google's BERT [19], have demonstrated remarkable

capabilities and have been applied to a wide range of applications.

Anthropic's Claude is a foundational AI model that aims to push the boundaries of AI capabilities while prioritizing safety and ethical considerations. Claude is designed to be a multi-purpose AI assistant, capable of engaging in open-ended conversations, answering questions, and assisting with various tasks. One of the key features of Claude is its grounding in Constitutional AI [20], [21], a framework that aims to ensure the model's outputs align with predefined ethical principles and values.

4.1. Constitutional AI

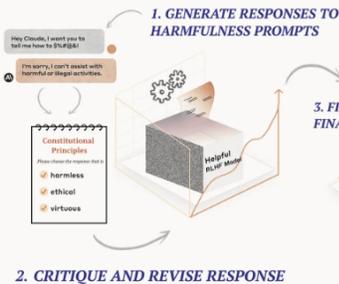
Anthropic's Constitutional AI incorporates a set of principles, or a "constitution," to guide the model's behavior during the training process (Figure 4) [20]. The constitution is used in two phases: first, the model is trained to critique and revise its own responses based on the principles; second, the model undergoes reinforcement learning using AI-generated feedback derived from the principles, rather than human feedback. This approach has been shown to produce models that are both more helpful and more harmless compared to models trained solely with human feedback [21].

The principles in Anthropic's constitution are drawn from various sources, including the UN Declaration of Human Rights [22], trust and safety best practices inspired from Apple's Terms of Service, DeepMind's Sparrow principles [23], and values that encourage the consideration of Non-Western perspectives. These principles cover a wide range of topics, from the protection of human rights and the promotion of equality to the avoidance of harmful, deceptive, or offensive content.

In October 2023, Anthropic partnered with the Collective Intelligence Project to run a public input process involving approximately 1,000 Americans to draft a constitution for an AI system [24]. The raw data from the survey is presented in the report by the Polis Center [25]. The resulting publicly sourced constitution [26] had a moderate degree of overlap

1. Supervised Learning (SL) Stage

Revises harmful AI responses through iterative self-critique and fine-tuning.



2. Reinforcement Learning (RL) Stage

Uses AI evaluations of responses according to constitutional principles to generate preference data for harmlessness and uses it to train a new model via Reinforcement Learning from AI Feedback.

3. FINE-TUNE WITH SL ON THE FINAL REVISED RESPONSES

4. AI GENERATES DATASET OF PREFERENCES FOR HARMLESSNESS



Which is better according to constitutional principles?

4. TRAIN PREFERENCE MODEL

5. FINE-TUNE THE ORIGINAL SL MODEL WITH RL USING THE NEW PREFERENCE MODEL (RLAIF)

Figure 4. Anthropic's Constitutional AI training process

with Anthropic's in-house constitution (roughly 50% overlap in concepts and values). However, the public constitution focused more on objectivity, impartiality, and accessibility, and tended to promote desired behavior rather than avoid undesired behavior.

Anthropic trained two models using Constitutional AI: one with the publicly sourced constitution and another with their in-house constitution. The models performed similarly on language understanding and math tasks, and were perceived as equally helpful and harmless by human evaluators. However, the model trained with the public constitution showed lower bias scores (Figure: 5) across nine social dimensions, particularly in the areas of disability status and physical appearance.

5. Threat Analysis

5.1. Identified Issues

Through our analysis of Anthropic's Claude, we have identified several potential threats and issues that warrant attention. One significant concern is the lack of transparency in Anthropic's privacy policies, particularly regarding the collection and use of personal data

for model training [28]. WeThe company's policies fail to provide clear and accessible information about data handling practices, making it difficult for users to make informed decisions about their data. Anthropic automatically collects browser information, mobile network, IP address (including information about the location of the device derived from your IP address), and identifiers (including device or advertising identifiers, probabilistic identifiers, and other unique personal or online identifiers). The inadequate transparency about personal data usage in training, employing complex terminology and lacking transparency in its trust and safety review criteria, raises concerns about data security and privacy.

Another issue is the potential for hallucinations in Claude's outputs, which can lead users to believe inaccurate or misleading information. While Anthropic claims to have reduced hallucination rates compared to competitors, the lack of open-source benchmarks and validation hinders the ability to independently verify these claims. Anthropic has not released their benchmark dataset, preventing open-source comparisons. Furthermore, Anthropic's claim that Constitutional AI will employ AI itself to train out harmful model outputs is questionable, as prior research shows significant stereotype propagation in such cases [29].

AI Bias scores from the Bias Benchmark for QA

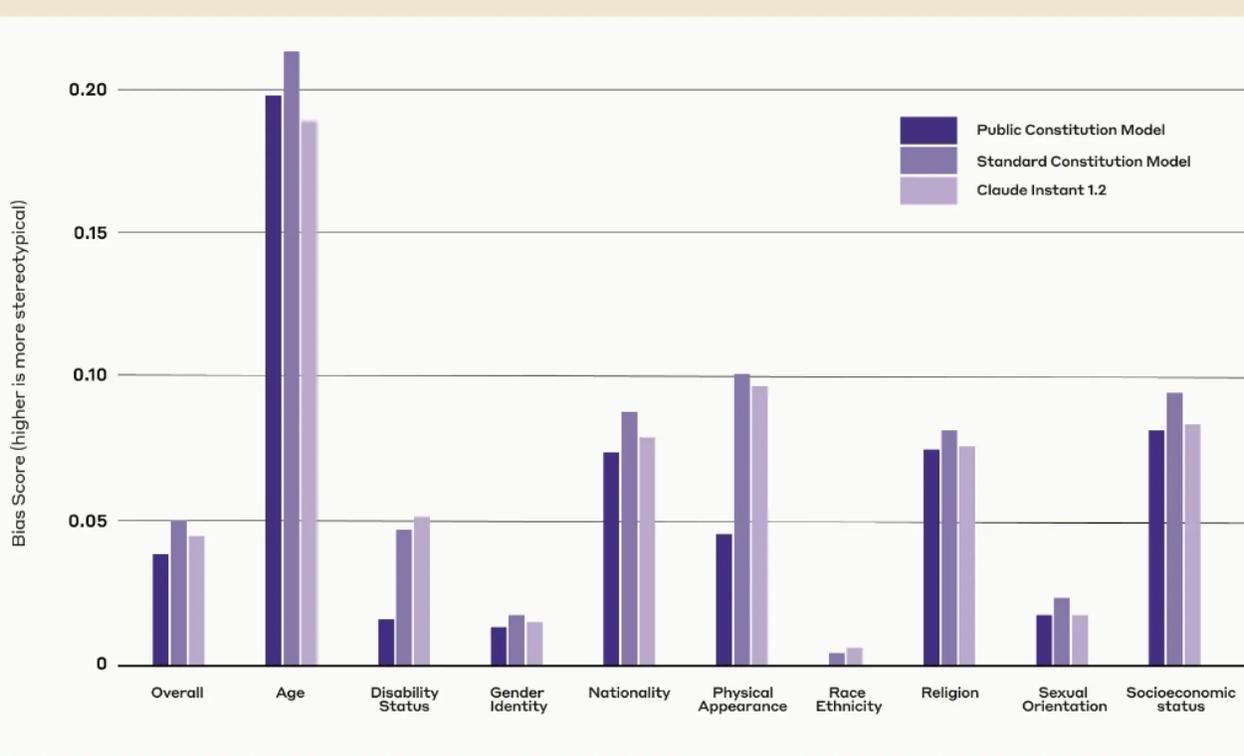


Figure 5. "BBQ [27] bias scores. Higher scores indicate more negative stereotype bias (lower is better). We used the same methods, code, and controls from our previously published work. The Public model shows lower bias scores across all nine social dimensions than the Standard model, especially for Disability Status and Physical Appearance. The Public constitution places a larger emphasis on accessibility, which may explain the greater reduction in bias for Disability Status in particular." [24]

Anthropic's partnerships with tech giants such as Google and Amazon raise concerns about third-party data usage and its implications for user privacy and data security. The reliance on partner policies and the lack of clear accountability mechanisms create uncertainties about data handling practices and potential risks. Anthropic claims AI Trust and Safety commitment but partners with companies that have their own data requirements. There is insufficient disclosure on Amazon and other partners training with Anthropic data. The policies lack a defined accountability structure, emphasizing responsibility without clear accountability mechanisms.

Potential biases and unequal benefits are another area of concern. Claude's bias benchmark, specific to Q&A since 2022, lacks updates and may be outdated

with stronger progress on red-teaming these past two years [27]. Anthropic fails to disclose training data, potentially giving certain groups predisposed advantages. Biased AI can lead to unequal outcomes, particularly when implemented in government agencies like DHS and USCIS as shown in Figure 8, posing a high risk of discrimination [30].

Anthropic's limited engagement with relevant AI actors is another area of concern. While the company has worked with certain organizations to implement AI risk management frameworks and called for funding towards AI safety research, their engagement appears limited compared to their competitors.

Lastly, the insufficient context understanding and impact characterization across various domains raises concerns about the effectiveness of Anthropic's ap-

Foundation Model Transparency Index Scores by Major Dimensions of Transparency, 2023

Source: 2023 Foundation Model Transparency Index

	Meta Llama 2	BigScience BLOOMZ	OpenAI GPT-4	stability.ai Stable Diffusion 2	Google PaLM 2	ANTHROPIC Claude 2	cohere Command	AI21labs Jurassic-2	Inflection Inflection-1	amazon Titan Text	Average
Data	40%	60%	20%	40%	20%	0%	20%	0%	0%	0%	20%
Labor	29%	86%	14%	14%	0%	29%	0%	0%	0%	0%	17%
Compute	57%	14%	14%	57%	14%	0%	14%	0%	0%	0%	17%
Methods	75%	100%	50%	100%	75%	75%	0%	0%	0%	0%	48%
Model Basics	100%	100%	50%	83%	67%	67%	50%	33%	50%	33%	63%
Model Access	100%	100%	67%	100%	33%	33%	67%	33%	0%	33%	57%
Capabilities	60%	80%	100%	40%	80%	80%	60%	60%	40%	20%	62%
Risks	57%	0%	57%	14%	29%	29%	29%	29%	0%	0%	24%
Mitigations	60%	0%	60%	0%	40%	40%	20%	0%	20%	20%	26%
Distribution	71%	71%	57%	71%	71%	57%	57%	43%	43%	43%	59%
Usage Policy	40%	20%	80%	40%	60%	60%	40%	20%	60%	20%	44%
Feedback	33%	33%	33%	33%	33%	33%	33%	33%	33%	0%	30%
Impact	14%	14%	14%	14%	14%	0%	14%	14%	14%	0%	11%
Average	57%	52%	47%	47%	41%	39%	31%	20%	20%	13%	

Scores for the 10 foundation model providers broken down by 13 subdomains, each of which have three or more indicators. Analysis at the level of major subdomains reveals actionable insight into what types of transparency or opacity lead to the above findings.

Figure 6. In depth review of Claude’s feature for Foundation Model Transparency as presented in Stanford’s Foundation Model Transparency Index [28].

Foundation Model Transparency Index Total Scores, 2023

Source: 2023 Foundation Model Transparency Index

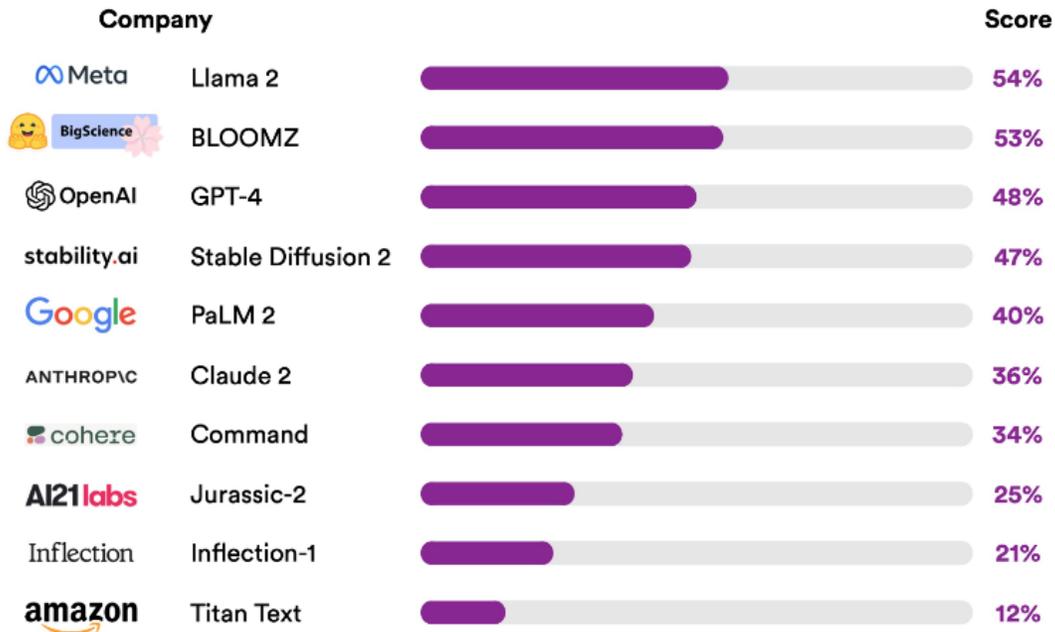


Figure 7. Results of Stanford’s Foundation Model Transparency Index, places Claude really low in comparison [28].

Homeland Security is testing AI to help with immigration, trafficking investigations, and disaster relief



/ DHS is rolling out a pilot program in partnership with OpenAI, Anthropic, and Meta.

Figure 8. Department of Homeland Security working with Anthropic and other AI organizations on Pilot Programs

proach. Although the company documents and discloses their motivations and priority towards AI safety, the lack of comprehensive context understanding and impact characterization underscores the need for a more thorough approach to AI governance.

These identified issues highlight the necessity for increased transparency, accountability, and proactive measures to address potential risks associated with Anthropic's Claude. The lack of clear data usage policies, validation against open-source benchmarks, and insufficient engagement with relevant AI actors emphasizes the importance of a more comprehensive approach to AI governance. As Anthropic continues to develop and deploy its AI systems, it is crucial to address these concerns to ensure responsible and ethical AI practices.

5.1.1. Constitutional AI. Anthropic's Constitutional AI approach, which aims to instill fixed ethical values across all cultures, raises significant concerns. By enforcing a universal set of principles, it risks suppressing diverse perspectives, oversimplifying complex societal dynamics, and favoring certain moral frameworks while marginalizing others. The static nature of this "constitution" may struggle to adapt to evolving norms and address the nuances of translating abstract ethics into algorithms, potentially leading to unintended discriminatory consequences. Furthermore, the lack of transparency and clear public accountability mechanisms, combined with the rigidity in navigating ethical dilemmas involving conflicting principles, undermines its ability to provide nuanced ethical guidance. While well-intentioned, the Constitutional AI model's one-size-fits-all approach may inadvertently perpetuate biases encoded into its fixed framework,

highlighting the need for a more dynamic, inclusive, and contextually aware ethical paradigm for responsible AI development and deployment across diverse moral landscapes.

6. NIST Framework Analysis

When analyzed through the lens of the NIST AI Risk Management Framework [7], the identified threats and issues in Anthropic's Claude can be mapped to various aspects of the framework. In terms of governance, Anthropic has defined its own AI Safety Levels and provides default opt-out options for data usage in model training. However, the company's policies lack clear accountability mechanisms, making it difficult to ensure responsible AI development and deployment.

The risk mapping and impact characterization aspects of the NIST framework reveal that Anthropic fails to appropriately disclose its objectives of AI Trust and Safety, leaving users uncertain about the risks and benefits associated with third-party software and data. While Anthropic uses a 2022 Q&A benchmark for social bias exploration and provides model access for red-teaming and safety research, its safety-centric claims lack the proactive approach demonstrated by competitors like OpenAI.

6.1. NIST "Govern" (Governance Analysis)

- 1) **Policies, processes, procedures, & practices:**
 - a) Defined own AI Safety Levels and discloses their current models' risks
 - b) Default opt-out for data usage in model training.
 - c) Insufficient disclosure regarding the use of personal data in model training, employing complex terminology and lacking transparency in its trust and safety review criteria.
- 2) **Accountability structure:** The policies lack a defined accountability structure, emphasizing responsibility without clear accountability.

Despite Anthropic’s strong recommendations in response to the NTIA’s call [31], their policies fail to specify clear accountability mechanisms.

- 3) **3rd Party Considerations:** Despite Anthropic’s repeated emphasis on Trust and Safety and data protection, they often defer to their partners’ policies, leaving users to decipher whether data usage is permitted. They do have a Acceptable Use Policy for API usage.
- 4) **Cultural considerations & communicated AI risks:** Presented their system prompt publicly focusing on transparency. Promoted for larger funding towards AI Safety Research. Provide special access to researchers seeking to red-team/alignment check their models.
- 5) **Engagement with relevant AI Actors:** Working with NIST to implement their AI Risk Management Framework. They, also called for \$15 Million funding for NIST’s Trustworthy and Responsible AI Resource Center. Announced partnership with Google and Amazon to build for AI Safety.

6.2. NIST ”MAP”

- 1) **Context is established and understood:** Yes, Anthropic documents and discloses their motivations and priority towards AI Safety. This can be seen through their partnerships, compliance, release of own AI Safety Levels, and also bias and multilingual performance benchmarks across their models
- 2) **Categorization of the AI system is performed:** Anthropic releases ASL stage for each of their models, as presented in Figure 9. They specify tasks for biases and benchmark models for those.
- 3) **Risks and benefits are mapped to third-party software and data:** They fail to appropriately disclose or align their objectives of AI Trust and Safety with those of their partners. For their API users, they do have a Acceptable Use Policy.

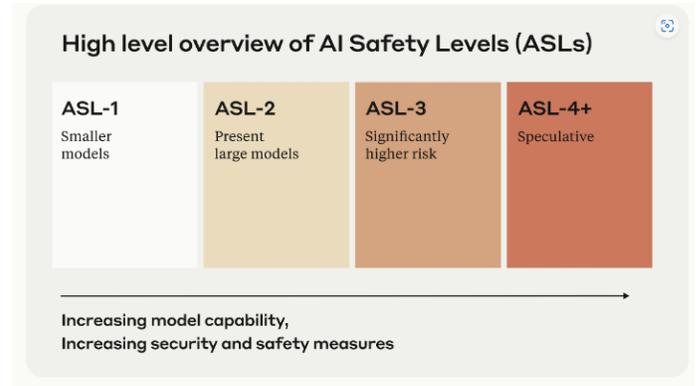


Figure 9. High level overview of AI Safety Levels defined by Anthropic [32]

- 4) **Impacts to individuals, groups, communities, organizations, and society are characterized:** Anthropic uses a 2022 Q&A benchmark for social bias exploration and provides model access for red-teaming and safety research. However, their safety-centric claims lack OpenAI’s proactive approach, which includes a curated red-teaming network actively probing for open-ended and subtle biases.

6.3. NIST ”Manage”

We take a look at the NIST AI Risk Management Framework’s Manage function which explores Anthropic’s responsibility to prioritize and respond to documented risks, plan and implement strategies to maximize benefits and minimize negative impacts, manage risks from third-party entities, and regularly monitor and document responses to identified risks [7], [33]. By addressing these aspects, Anthropic can enhance its capacity to manage Claude’s risks and ensure responsible AI development and deployment. The following points highlight the key aspects of the Manage function that Anthropic should address:

- 1) **AI risks based on impact assessments and other analytical output from the Map and Measure functions are prioritized, responded to, and managed:** Anthropic needs to prioritize and respond to the documented

risks based on their potential impact, likelihood, and available resources. This includes determining whether Claude achieves its intended purpose and stated objectives, considering the risks associated with harmful usage, automation, hallucinations, biases, and weak transparency in data usage policies. Anthropic should develop, plan, and document responses to the most significant risks, which may include mitigating, transferring, sharing, avoiding, or accepting them.

- 2) **Strategies to maximize benefits and minimize negative impacts are planned, prepared, implemented, and documented, and informed by stakeholder input:** Anthropic should plan, prepare, implement, and document strategies to maximize the benefits and minimize the negative impacts of Claude, informed by stakeholder input. This involves considering the resources required to manage risks, along with viable alternative systems, approaches, or methods, and the related reduction in severity of impact or likelihood of each potential action. Mechanisms should be in place and applied to sustain the value of Claude and to supersede, disengage, or deactivate the system if it demonstrates performance or outcomes inconsistent with its intended use.
- 3) **Risks from third-party entities are managed:** Anthropic must manage risks from third-party entities, such as Google and Amazon, by regularly monitoring and applying risk controls. This is particularly important given the concerns raised about third-party data usage and its implications for user privacy and data security.
- 4) **Responses to identified and measured risks are documented and monitored regularly:** Anthropic should document and regularly monitor responses to identified and measured risks. This includes implementing post-deployment system monitoring plans, capturing and evaluating user and stakeholder feed-

back, establishing mechanisms for appeal and override, decommissioning, incident response, and change management. Measurable continuous improvement activities should be integrated into system updates and include regular stakeholder engagement.

By addressing these aspects of the Manage function, Anthropic can enhance its capacity to manage the risks associated with Claude, allocate risk management resources based on risk measures, and ensure the responsible development and deployment of their AI system.

6.4. NIST "MEASURE"

The NIST AI Risk Management Framework's Measure function focuses on employing quantitative, qualitative, or mixed-method tools, techniques, and methodologies to analyze, assess, benchmark, and monitor AI risk and related impacts [7]. Here are some points highlighting the key aspects of NIST's Measure function that Anthropic should address in relation to Claude:

- 1) **Appropriate methods and metrics are identified and applied:** Anthropic should identify and select approaches and metrics for quantitative or qualitative measurement of the most significant risks, including context-relevant measures of trustworthiness. The appropriateness of metrics and effectiveness of existing controls should be regularly assessed and updated, involving internal experts who did not serve as front-line developers for the system and/or independent assessors. Their implementation of BBQ is outdated and needs to be reconsidered [27].
- 2) **Systems are evaluated for trustworthy characteristics:** While Anthropic does document test sets, metrics, and details about the tools used during test, evaluation, validation, and verification (TEVV), it doesn't open-source the evaluation framework or pipeline

[33]. This creates friction in academic replication tasks, who want to publicly evaluate their claims. System performance or assurance criteria should be measured qualitatively or quantitatively and demonstrated for conditions similar to deployment setting(s). Claude should be evaluated regularly for safety, computational bias, resilience, security, privacy risk, and environmental impact. The AI model should be explained, validated, and documented, and its output should be interpreted within its context to inform responsible use and governance.

- 3) **Mechanisms for tracking identified risks over time are in place:** Anthropic should have approaches, personnel, and documentation in place to regularly identify and track existing and emergent risks based on factors such as intended and actual performance in deployed contexts. Risk tracking approaches should be considered for settings where risks are difficult to assess using currently available measurement techniques or are not yet available. They do not have a Red Teaming Network (internal exists) like OpenAI nor do they have a Bug Bounty program yet, which could both be used as a crowd-sourced risk tracker.
- 4) **Feedback about efficacy of measurement is gathered and assessed:** Measurement approaches for identifying risks should be connected to deployment context(s) and informed through consultation with domain experts and other end users. Measurement results regarding system trustworthiness in deployment context(s) should be informed by domain expert and other stakeholder feedback to validate whether Claude is performing consistently as intended.

By addressing these aspects of the Measure function, Anthropic can enhance their capacity to comprehensively evaluate Claude's trustworthiness, identify and track existing and emergent risks, and verify the efficacy of metrics.

7. EU AI Act Analysis

Under the EU AI Act [8], the identified threats and issues in Anthropic's Claude can be categorized based on their risk levels. The risks associated with harmful usage and automation are also significant:

- 1) Automation of AI fine-tuning can be considered limited risk, as previous research has demonstrated it to propagate pre-learned biases. However, if this model is auto-deployed without validation, it would become high risk.
- 2) Hallucinations in outputs can cause users to mistakenly believe something, directly impacting individuals and their perception of reality, posing a limited risk.
- 3) Lack of transparency can affect user trust and lead to privacy issues. If the data is used for training other models, it can be considered high risk because AI models can memorize information.
- 4) AI being used with the intent to cause harm, such as violating human rights, poses an unacceptable risk.

As discussed, under the EU AI Act, the use of AI for harmful content removal, as proposed in Anthropic's Constitutional AI framework, would likely be classified as a high-risk AI system. This is because the automated removal of content can have significant impacts on individuals' rights to freedom of expression and access to information. The lack of transparency and potential for biases in the AI system used for content moderation further exacerbates the risks associated with this application. Similarly, the use of Claude in government agencies like DHS and USCIS would also fall under the high-risk category due to the potential for discriminatory outcomes and unequal treatment. While this maybe within the United States, we still evaluate it using the EU AI Framework, which emphasizes the importance of ensuring that AI systems used in the public sector are transparent, accountable, and free from biases that could lead to discrimination.

Anthropic's weak transparency in data usage poli-

cies and the potential for data used in training to be memorized and reproduced by AI models raise concerns under the EU AI Act's requirements for data governance and privacy protection. The Act requires AI system providers to ensure appropriate data management practices, including data minimization, data quality, and data protection safeguards. To comply with the EU AI Act, Anthropic would need to address these identified risks by implementing robust risk management processes, ensuring transparency in their AI systems' development and deployment, and establishing clear accountability mechanisms.

8. Proposed Mitigations & Resolution Strategies

To address the identified threats and issues in Anthropic's Claude, we propose the following mitigation strategies:

8.1. Enhance transparency in privacy policies

Anthropic should prioritize adopting transparent privacy practices that comprehensively detail the risks associated with artificial intelligence systems, as outlined in the NIST AI Framework [7]. Additionally, they should minimize data retention periods and implement a default opt-out option, empowering users with greater control over their personal information. To further simplify information access and boost user engagement, organizations should streamline navigation complexity and provide concise, easily understandable summaries of their privacy practices. This will empower users to make informed decisions about their data and increase trust in Anthropic's AI systems.

8.1.1. Criteria & Metrics. The evaluation of efforts to improve the transparency and accessibility of privacy policies should be guided by well-defined criteria and metrics. These include:

- 1) **Accessibility:** Measured by the average number of clicks required for users to access the

privacy practices. A lower number of clicks indicates higher accessibility, enabling users to obtain privacy policy information more conveniently.

- 2) **Time:** The duration spent by users locating specific details within the privacy policy. This metric assesses the ease with which users can quickly find the required information within the policy. A shorter duration reflects better organization and navigation of the privacy policy.
- 3) **Comprehension:** The extent to which users can understand the content of the privacy policies without relying on external references. This metric evaluates the clarity and readability of the policies; the clearer the language, the easier it is for users to comprehend without requiring external explanations.

8.1.2. Data Sources / Test. To collect data for these metrics, two primary methods can be employed:

Comprehension Surveys: Designing questionnaires that present users with the privacy policy content and assess their understanding through questions. The survey results can provide valuable insights into the comprehension metric.

Benchmarking: Comparing the organization's metrics against industry standards or best practices. By benchmarking their accessibility, time, and comprehension metrics against established norms or leading examples, organizations can identify areas for improvement and gauge their performance relative to peers or competitors.

Utilizing these data sources and testing methods, organizations can gather valuable data and insights to evaluate the effectiveness of their efforts in improving the transparency and accessibility of privacy policies. This information will guide further improvements and help organizations prioritize areas that require the most attention and resources.

8.1.3. Practical Considerations. From a practical standpoint, organizations should focus on enhancing

user experience through thoughtful design choices and effective summarization techniques. Simplifying the interface and reducing navigation complexity can expedite information access while offering clear and concise summaries of privacy practices can significantly improve user comprehension and engagement with these policies.

8.2. Establish rigorous benchmarks for hallucination and bias

To ensure transparency and facilitate rigorous public scrutiny of potential hallucinations and biases in Anthropic’s AI models, it is imperative to conduct comprehensive benchmarking exercises.

8.2.1. Criteria & Metrics. These benchmarks should aim to measure the extent of hallucinations, which can be quantified through metrics such as Q2 and factual-grounding BLEU scores. Additionally, they should evaluate various forms of bias, including statistical parity, group diversity, equalized odds, and even open-ended opinions annotated by subject matter experts. This will help identify and mitigate potential risks associated with inaccurate or biased outputs.

8.2.2. Data Sources / Test. The data sources and tests employed for these benchmarking efforts should be diverse and comprehensive. For hallucination evaluation, datasets such as HaluEval [34] and the forthcoming HaluEval-Wild [35] can provide valuable insights. Bias assessment can leverage resources like R-Judge [36], CBBQ [37], Winoqueer [38] and KorNAT [39], which cover a wide range of bias types and demographic factors.

8.2.3. Practical Considerations. From a practical standpoint, it is crucial to ensure that these benchmarks are not inadvertently used for pre-training or fine-tuning the AI models themselves, as this could introduce biases or undermine the integrity of the evaluation process. Additionally, creating private leaderboards for

these benchmarks can help maintain their integrity and prevent potential gaming or exploitation.

8.3. Develop a comprehensive remediation process

Anthropic should implement a robust process for handling user requests for data deletion and ensuring the unlearning of data by AI models. Clear mechanisms should be established for users to initiate data deletion requests, and the company should provide detailed guidance and support throughout the process. Rigorous testing should be conducted to verify the effectiveness of data removal and model unlearning.

8.3.1. Criteria & Metrics. This remediation process should prioritize clarity in the request initiation stage, empowering users with a straightforward understanding of how to initiate data deletion requests. Furthermore, it must incorporate verifiable metrics to assess the efficacy of model unlearning processes, ensuring that user data is thoroughly expunged from the models upon request.

8.3.2. Data Sources / Test. User feedback can provide invaluable insights into the clarity and user-friendliness of the data deletion process, highlighting areas that may require improvement. Additionally, dedicated unlearning tests must be conducted to verify the complete removal of user data from the models, validating the integrity and functionality of the unlearning mechanisms.

8.3.3. Practical Considerations. From a practical standpoint, several key considerations must be addressed. Firstly, the entire process of data deletion and unlearning should be transparent, with clear and well-documented steps outlined for users. Secondly, comprehensive guidance and user support should be provided throughout the process, ensuring that users are adequately informed and assisted at every stage. Finally, organizations must invest in enhancing their systems to support efficient and timely data deletion

and unlearning, prioritizing the swift and effective handling of user requests.

By implementing these mitigation strategies, Anthropic can demonstrate its commitment to responsible AI development and deployment, enhance user trust, and reduce the risks associated with its Claude model. Regular monitoring and continuous improvement of these measures will be essential to keep pace with the evolving AI governance landscape and ensure ongoing accountability.

9. Discussion

The proposed mitigation strategies for Anthropic's Claude have significant implications for the broader AI governance landscape and the social impact of AI systems. By enhancing transparency in privacy policies, Anthropic can set a positive example for other AI companies, encouraging a more open and accountable approach to data handling and user privacy. This increased transparency will empower users to make informed decisions about their data and foster trust in AI systems.

Establishing rigorous benchmarks for hallucination and bias will contribute to the development of more reliable and unbiased AI models. By publicly releasing datasets and results, Anthropic can promote collaboration and knowledge sharing within the AI community, driving collective efforts towards mitigating the risks associated with inaccurate or biased outputs. This transparency will also enable independent verification and accountability, ensuring that AI systems are subject to rigorous scrutiny.

Implementing a comprehensive remediation process for data deletion and model unlearning will address concerns about data privacy and the potential misuse of personal information. By providing users with clear mechanisms to control their data and ensuring the effectiveness of data removal and model unlearning, Anthropic can demonstrate its commitment to user privacy and build trust in its AI systems.

The adoption of these mitigation strategies by An-

thropic and other AI companies will contribute to the development of a more responsible and trustworthy AI ecosystem. As AI systems become increasingly integrated into various aspects of society, ensuring their alignment with ethical principles and societal values becomes paramount. By prioritizing transparency, accountability, and user privacy, AI companies can foster public trust and support the responsible deployment of AI technologies for the benefit of society.

10. Conclusion

In conclusion, this paper has examined the AI governance landscape, focusing on Anthropic's Claude as a case study [40]. Through the lens of the NIST AI Risk Management Framework and the EU AI Act, we have identified potential threats and issues in Claude, including the lack of transparency in privacy policies, the potential for hallucinations and biases in outputs, and concerns about third-party data usage [41]. To address these challenges, we have proposed mitigation strategies that emphasize transparency, rigorous benchmarking, and comprehensive data handling processes. By adopting these measures, Anthropic can demonstrate its commitment to responsible AI development and deployment, enhance user trust, and contribute to the broader efforts in AI governance.

The evolution of AI governance will require ongoing collaboration, adaptation, and learning from the successes and challenges of parallel domains such as privacy regulations. As AI systems become more sophisticated and integrated into society, ensuring their alignment with ethical principles and societal values will be critical. By prioritizing accountability, transparency, and user privacy, AI companies can foster public trust and support the responsible advancement of AI technologies for the benefit of society.

11. Limitations & Ethical Considerations

While this paper provides valuable insights into AI governance and accountability, it is important to

acknowledge its limitations. The analysis focuses primarily on Anthropic’s Claude and may not fully capture the diverse range of AI systems and their unique governance challenges. Additionally, the proposed mitigation strategies, while promising, require further validation and real-world implementation to assess their effectiveness and potential unintended consequences.

Ethical considerations are paramount in the development and deployment of AI systems. As AI technologies become more powerful and influential, it is crucial to ensure that they are designed and used in a manner that respects human rights, promotes fairness, and avoids harmful biases. AI companies must prioritize ethical principles throughout the AI lifecycle, from data collection and model training to deployment and monitoring.

Ongoing research, collaboration, and stakeholder engagement will be essential to address the ethical implications of AI and develop robust governance frameworks that keep pace with the rapid advancements in AI technologies. By proactively addressing ethical considerations and prioritizing accountability, transparency, and user privacy, we can work towards a future where AI systems are trusted, beneficial, and aligned with societal values.

12. Acknowledgments

I would like to express my sincere gratitude to Professor Norman Sadeh for his invaluable guidance and insights throughout the AI Governance course at Carnegie Mellon University. We are truly thankful for the opportunity to learn about the field of AI governance.

References

- [1] “Partnering with Scale to Bring Generative AI to Enterprises,” Apr. 2024, [Online; accessed 28. Apr. 2024]. [Online]. Available: <https://www.anthropic.com/news/partnering-with-scale>
- [2] “Zoom Partnership and Investment in Anthropic,” Apr. 2024, [Online; accessed 28. Apr. 2024]. [Online]. Available: <https://www.anthropic.com/news/zoom-partnership-and-investment>
- [3] “Anthropic partners with BCG,” Apr. 2024, [Online; accessed 28. Apr. 2024]. [Online]. Available: <https://www.anthropic.com/news/anthropic-bcg>
- [4] “Accenture, AWS, Anthropic Collaboration,” Apr. 2024, [Online; accessed 28. Apr. 2024]. [Online]. Available: <https://www.anthropic.com/news/accenture-aws-anthropic>
- [5] “SKT Partnership Announcement,” Apr. 2024, [Online; accessed 28. Apr. 2024]. [Online]. Available: <https://www.anthropic.com/news/skt-partnership-announcement>
- [6] “Kief Studio & Anthropic Partnership: Transform Your Business with AI C,” Apr. 2024, [Online; accessed 28. Apr. 2024]. [Online]. Available: <https://kief.studio/anthropic>
- [7] “AI Risk Management Framework | NIST,” Jan. 2024, [Online; accessed 28. Apr. 2024]. [Online]. Available: <https://www.nist.gov/itl/ai-risk-management-framework>
- [8] “Eu ai act: first regulation on artificial intelligence,” European Parliament, Apr. 2024. [Online]. Available: <https://www.europarl.europa.eu/topics/en/article/20230601STO93804/eu-ai-act-first-regulation-on-artificial-intelligence>
- [9] T. Birkstedt, M. Minkkinen, A. Tandon, and M. Mäntymäki, “Ai governance: themes, knowledge gaps and future agendas,” *Internet Research*, vol. 33, no. 7, pp. 133–167, 2023.
- [10] D. Schiff, J. Biddle, J. Borenstein, and K. Laas, “What’s next for ai ethics, policy, and governance? a global overview,” in *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, ser. AIES ’20. New York, NY, USA: Association for Computing Machinery, 2020, p. 153–158. [Online]. Available: <https://doi.org/10.1145/3375627.3375804>
- [11] E. Papagiannidis, I. M. Enholm, C. Dremel, P. Mikalef, and J. Krogstie, “Toward ai governance: Identifying best practices and potential barriers and outcomes,” *Information Systems Frontiers*, vol. 25, no. 1, pp. 123–141, 2023.
- [12] M. Mäntymäki, M. Minkkinen, T. Birkstedt, and M. Viljanen, “Defining organizational ai governance,” *AI and Ethics*, vol. 2, no. 4, pp. 603–609, 2022.
- [13] I. D. Raji, A. Smart, R. N. White, M. Mitchell, T. Gebru, B. Hutchinson, J. Smith-Loud, D. Theron, and P. Barnes, “Closing the ai accountability gap: Defining an end-to-end framework for internal algorithmic auditing,” 2020.
- [14] “General Data Protection Regulation (GDPR) – Legal Text,” Apr. 2024, [Online; accessed 28. Apr. 2024]. [Online]. Available: <https://gdpr-info.eu>
- [15] M. Goddard, “The eu general data protection regulation (gdpr): European regulation that has a global impact,” *International Journal of Market Research*, vol. 59, no. 6, pp. 703–705, 2017.
- [16] A. J. Adetayo, M. O. Aborisade, and B. A. Sanni, “Microsoft copilot and anthropic claude ai in education and library service,” *Library Hi Tech News*, 2024.

- [17] V. K. Uppalapati and D. S. Nag, “A comparative analysis of ai models in complex medical decision-making scenarios: Evaluating chatgpt, claude ai, bard, and perplexity,” *Cureus*, vol. 16, no. 1, 2024.
- [18] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, “Language models are few-shot learners,” 2020.
- [19] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” 2019.
- [20] “Claude’s Constitution,” Apr. 2024, [Online; accessed 28. Apr. 2024]. [Online]. Available: <https://www.anthropic.com/news/claude-constitution>
- [21] Y. Bai, S. Kadavath, S. Kundu, A. Askell, J. Kernion, A. Jones, A. Chen, A. Goldie, A. Mirhoseini, C. McKinnon, C. Chen, C. Olsson, C. Olah, D. Hernandez, D. Drain, D. Ganguli, D. Li, E. Tran-Johnson, E. Perez, J. Kerr, J. Mueller, J. Ladish, J. Landau, K. Ndousse, K. Lukosuite, L. Lovitt, M. Sellitto, N. Elhage, N. Schiefer, N. Mercado, N. DasSarma, R. Lasenby, R. Larson, S. Ringer, S. Johnston, S. Kravec, S. E. Showk, S. Fort, T. Lanham, T. Telleen-Lawton, T. Conerly, T. Henighan, T. Hume, S. R. Bowman, Z. Hatfield-Dodds, B. Mann, D. Amodei, N. Joseph, S. McCandlish, T. Brown, and J. Kaplan, “Constitutional AI: Harmlessness from AI Feedback,” *arXiv*, Dec. 2022.
- [22] United Nations, “Universal Declaration of Human Rights | United Nations,” Apr. 2024, [Online; accessed 28. Apr. 2024]. [Online]. Available: <https://www.un.org/en/about-us/universal-declaration-of-human-rights>
- [23] Sep. 2022, [Online; accessed 28. Apr. 2024]. [Online]. Available: <https://storage.googleapis.com/deepmind-media/DeepMind.com/Authors-Notes/sparrow/sparrow-final.pdf>
- [24] “Collective Constitutional AI: Aligning a Language Model with Public Input,” Apr. 2024, [Online; accessed 28. Apr. 2024]. [Online]. Available: <https://www.anthropic.com/news/collective-constitutional-ai-aligning-a-language-model-with-public-input>
- [25] “pol.is report,” Apr. 2024, [Online; accessed 28. Apr. 2024]. [Online]. Available: <https://pol.is/report/r3rwrnr5udrzkwvxtkdj>
- [26] Dec. 2023, [Online; accessed 28. Apr. 2024]. [Online]. Available: https://www-cdn.anthropic.com/65408ee2b9c99abe53e432f300e7f43ef69fb6e4/CCAI_public_comparison_2023.pdf
- [27] A. Parrish, A. Chen, N. Nangia, V. Padmakumar, J. Phang, J. Thompson, P. M. Htut, and S. Bowman, “BBQ: A hand-built bias benchmark for question answering,” in *Findings of the Association for Computational Linguistics: ACL 2022*, S. Muresan, P. Nakov, and A. Villavicencio, Eds. Dublin, Ireland: Association for Computational Linguistics, May 2022, pp. 2086–2105. [Online]. Available: <https://aclanthology.org/2022.findings-acl.165>
- [28] R. Bommasani, K. Klyman, S. Longpre, S. Kapoor, N. Maslej, B. Xiong, D. Zhang, and P. Liang, “The foundation model transparency index,” 2023.
- [29] E. Arazo, D. Ortego, P. Albert, N. E. O’Connor, and K. McGuinness, “Pseudo-labeling and confirmation bias in deep semi-supervised learning,” in *2020 International joint conference on neural networks (IJCNN)*. IEEE, 2020, pp. 1–8.
- [30] G. Del Valle, “DHS testing out AI pilot programs for FEMA, ICE, and USCIS,” *Verge*, Mar. 2024. [Online]. Available: <https://www.theverge.com/2024/3/18/24104843/dhs-ai-pilot-programs-chatgpt-openai-anthropic-meta>
- [31] Mar. 2024, [Online; accessed 28. Apr. 2024]. [Online]. Available: <https://www.ntia.gov/sites/default/files/ntia-ai-report-final.pdf>
- [32] “Anthropic’s Responsible Scaling Policy,” Apr. 2024, [Online; accessed 28. Apr. 2024]. [Online]. Available: <https://www.anthropic.com/news/anthropics-responsible-scaling-policy>
- [33] Aug. 2022, [Online; accessed 28. Apr. 2024]. [Online]. Available: https://www.nist.gov/system/files/documents/2022/08/18/AI_RMF_2nd_draft.pdf
- [34] J. Li, X. Cheng, X. Zhao, J.-Y. Nie, and J.-R. Wen, “HaluEval: A Large-Scale Hallucination Evaluation Benchmark for Large Language Models,” *OpenReview*, Dec. 2023. [Online]. Available: <https://openreview.net/forum?id=bxsrykzSnq>
- [35] Z. Zhu, Z. Sun, and Y. Yang, “HaluEval-Wild: Evaluating Hallucinations of Language Models in the Wild,” *arXiv*, Mar. 2024.
- [36] T. Yuan, Z. He, L. Dong, Y. Wang, R. Zhao, T. Xia, L. Xu, B. Zhou, F. Li, Z. Zhang, R. Wang, and G. Liu, “R-Judge: Benchmarking Safety Risk Awareness for LLM Agents,” *arXiv*, Jan. 2024.
- [37] Y. Huang and D. Xiong, “CBBQ: A Chinese Bias Benchmark Dataset Curated with Human-AI Collaboration for Large Language Models,” *arXiv*, Jun. 2023.
- [38] V. K. Felkner, H.-C. H. Chang, E. Jang, and J. May, “Wino-Queer: A Community-in-the-Loop Benchmark for Anti-LGBTQ+ Bias in Large Language Models,” *arXiv*, Jun. 2023.
- [39] J. Lee, M. Kim, S. Kim, J. Kim, S. Won, H. Lee, and E. Choi, “KorNAT: LLM Alignment Benchmark for Korean Social Values and Common Knowledge,” *arXiv*, Feb. 2024.

- [40] A. Parrish, A. Chen, N. Nangia, V. Padmakumar, J. Phang, J. Thompson, P. M. Htut, and S. Bowman, “BBQ: A hand-built bias benchmark for question answering,” *ACL Anthology*, pp. 2086–2105, May 2022.
- [41] D. Yao, “Zoom Invests in Anthropic, Partners on Generative AI,” *AI Business*, May 2023. [Online]. Available: <https://aibusiness.com/nlp/zoom-invests-in-anthropic-partners-on-generative-ai#close-modal>