

# Unveiling Unintended Systematic Biases in Natural Language Processing

Olga Scrivner

Rose-Hulman Institute of Technology, Indiana University, Scrivner Solutions Inc

Question: How do systemic biases emerge in Natural Language Processing, what societal impacts do they create, and how can we address these biases?

## Learning Objectives

After studying this chapter, you should be able to:

- Concisely define the following key terms: natural language processing, representational harm, allocative harm, implicit biases
- Identify the origin of biases in the NLP workflow
- Compare various social impacts from biases generated by the NLP applications
- Evaluate real-world examples and their consequences.
- Become familiar with nuances of bias taxonomy
- Implement mitigation strategies for reducing biases

## Chapter Overview

In this chapter, we present the idea of systematic biases in Natural Language Processing (NLP) and how NLP applications could unintentionally lead to unfair societal consequences. We explain the reasons why we trust Artificial Intelligence and how our human biases creep into computer algorithms. We outline the bias taxonomy to increase awareness about the subtleties of our language usage and show several methods to detect and mitigate NLP biases.

Reducing and mitigating bias in Natural Language Processing (NLP) is an important yet challenging endeavor requiring an understanding of underlying social and ethical implications.



### Useful Definitions

<b>AI:</b>	Artificial Intelligence, a multidisciplinary field that uses algorithms to imitate intelligent human behavior
<b>ASR:</b>	Automated Speech Recognition, a technology to process human speech
<b>NLP:</b>	Natural Language Processing, a subfield of AI focused on understanding, interpreting, and generating human language
<b>PLM:</b>	Pre-trained Language Models trained on a large amount of text data and fine-tuned to a specific NLP task (e.g., ELMO, BERT, GPT-2)
<b>LLM:</b>	Large Language Models trained on a massive amount of data (model size > 10B parameters) and exhibiting "emergent capabilities" as they are able to perform multiple NLP tasks (e.g., GPT-3, GPT-4, T-5)
<b>GPT:</b>	Generative Pre-trained Transformer, a type of LLM developed by OpenAI
<b>chatGPT:</b>	a conversational application developed using GPT
<b>GPT-4:</b>	a GPT model with multimodal capabilities (image, text, audio, video)
<b>autoGPT:</b>	an open-source GPT application capable of running autonomously to carry out tasks imposed by a user, including browsing the web and accessing file systems
<b>BLOOM:</b>	an open-access multilingual large language model developed by Huggingface
<b>OpenAI API:</b>	an Application Programming Interface allowing to interact with models or data
<b>Benchmark:</b>	a standard or baseline metric to evaluate the model performance
<b>Corpus:</b>	a digital collection of written or spoken text data
<b>Black box:</b>	an opaque model where we can only observe the input and output
<b>White box:</b>	a transparent model showing internal design and parameters

## Chapter Overview and Learning Objectives

This chapter will present the key concepts of bias and fairness as they relate to NLP and discuss biases in various stages of the NLP applications from development to usage. This chapter will also introduce NLP-related incidents and discuss their social impact. Finally, several assessment methods will be shown to detect and mitigate NLP biases.

The learning objectives for this chapter include the following:

- Understand the significance of social impact from bias and unfairness generated by the NLP applications.
- Identify the origin of biases in the NLP workflow.
- Examine real-world examples and their consequences.

- Become familiar with nuances of bias taxonomy.
- Implement mitigation strategies for reducing biases.

## **1. Introduction**

### **1.1. "Pause Giant AI Experiment"**

The field of natural language processing (NLP), a subfield of artificial intelligence (AI), has undergone a significant transformation, evolving from hand-written rule models to deep learning models. The recent releases of generative models (ChatGPT and GPT-4 by OpenAI, Bing by Microsoft, Bard by Google, Claude by Anthropic, Tongyi Qianwen by Alibaba, and open-source HuggingFace models) have further revolutionized the field and are already transforming entire industries, including media, art, technology, and education. As the NLP applications become more prevalent, along with the benefits (e.g., improved accessibility and efficiency), they have yielded the following risks: 1) producing harmful content, 2) amplifying societal stereotypes and biases, 3) generating misinformation, 4) contributing to discrimination and unfairness through biased solutions, and 5) potentially posing mental and health risks.

The growing accessibility and enhanced performance of the AI technologies such as voice assistants and conversational agents have created a deeper reliance on 'black-box' solutions, even in patient care, court rulings, and data security (Liang et al. 2021). Here are a few reported cases: a judge in Colombia made legal inquiries using chatGPT regarding the insurance cost liabilities for medical treatment; Samsung engineers used chatGPT to help optimize code and convert internal meeting notes into a presentation, leaking the highly sensitive information; chatGPT itself had a bug exposing other users' chat history (Moon 2023; Lopez 2023; Zoppo 2023). In fact, these concerns have led to several public actions: the ban of chatGPT by the

General Data Protection Regulation in Italy; a complaint to the Federal Trade Commission by the Center for AI and Digital Policy calling the GPT-4 model “biased, deceptive, and a risk to privacy and public safety”; an open letter "Pause Giant AI Experiment", signed by Elon Musk, Steve Wozniak, and others, stating that "Powerful AI systems should be developed only once we are confident that their effects will be positive and their risks will be manageable"; the restriction of AI tools in some schools to avoid cheating; U.S. Senate first draft outlining "a new regulatory regime that would prevent potentially catastrophic damage"; and even the editorial submission closure with the influx of AI-generated books (Feiner 2023; Grothaus 2023; Anderson 2023; Shepardson 2023).

## **1.2. Why Do We Trust AI?**

Human trust is often based on the assumptions that machine learning math computations "would be pure and neutral, providing for AI a fairness beyond what is present in human society" and the large data size would lead to more accuracy (Caliskan, Bryson, and Narayanan 2017, 183). Language models, however, are trained on textual data without awareness of the social meaning and authorship information, such as self-identification or group membership (Hovy and Spruit 2016; Hovy and Prabhumoye 2021). As a result, this downstream textual data processing creates inferences from individuals based on data patterns with underlying biases. Technically, these biases are just a “mismatch of ideal and actual distributions of labels and user attributes”, however, in real life, they can lead to unintended but systematic societal inequalities and even legal implications (Shah, Schwartz, and Hovy 2020; Blackman 2020).



### **Did You Know? Societal Environmental Impact!**

- GPT-1 (2018) trained on 4.5 GB data
- GPT-2 (2019) trained on 40GB data
- GPT-3 (2020, GPT-3.5 - revision 2022) trained on 570GB data
- BLOOM (2022) trained on 1.6TB data

A single A100 GPU unit consumes about 300 watts. If BLOOM was trained using **384** 80GB A100 GPUs for **3.5** months, then 384 GPUs consume 115,200 watts, or 115 kilowatts (kW). Running for 105 days (2,520 hrs) means a training cost of 289,800 kilowatt hours (kWh). Note that the average household consumes 10,649 kWh annually. If the CO<sub>2</sub> average per kWh is 0.95, then the CO<sub>2</sub> emission is equal to 275,310 pounds per kWh (Strubell 2019). Data centers also consume a large amount of water for cooling systems. For example, chatGPT "drinks" an estimated 500ml bottle of water for each conversation (Li et al. 2023).

### **1.3. Why Are There So Many Challenges?**

One of the current challenges for mitigating these biases is that biases are not often readably visible in the data or underlying algorithms, and it is often difficult to judge whether a given statement contains bias, even for humans (Sap et al. 2019; Baheti et al. 2021). Secondly, biases can be introduced at multiple stages during the development of NLP systems, including input representation, feature engineering, annotation process, model training, and research design. adding more complexity to identifying them (Jägare 2022; Hovy and Prabhumoye 2021). Socio-cognitive fallacies may have distinct representations in each NLP task, for example, machine translation, text summarization, or text generation (Sun et al. 2019). "Black box" models pose an additional difficulty for identifying biases as their internal operations are inherently opaque as compared to transparent and interpretable "white box" models (Jägare 2022), which is ironic on its own, as it displays a subconscious stereotype of color naming (e.g., the white color is transparent, explainable, logical, while the black color is opaque, non-understandable, nonreasonable). Assessing models is another challenge. First, the fairness and biases definitions and their associated evaluation tests vary across disciplines and tasks (Czarnowska, Vyas, and Shah 2021; Bansal 2022). Second, the over-reliance on the "state-of-

the-art" metrics (SOTA) leads to the "right for wrong reason" results, focusing on a narrow vision to achieve the highest scores. Recently, benchmarking itself became a topic of debate. "Current benchmarking practices offer a mechanism through which a small number of elite corporate, government, and academic institutions shape the research agenda and values of the field" and the current popular benchmark datasets only reflect a narrow vision of the world, predominately "white, male, western" (Koch et al. 2021, 9).



### 1. Hands-On Practice: Can You Identify a Non-Human? (Task A and Task B)

The Turing test, developed by Alan Turing in 1950, is a benchmark used to evaluate the human-like capabilities of computer models. Watch the TED-Ed video on the Turing Test at <https://youtu.be/3wLqsRLvV-c> to learn more. After watching the video, your task is to conduct the Turing test and consider whether it is easy to identify non-human generated texts. Observe language features such as structure, content, semantics, and more, that may reveal the non-human identity. Provide examples and note any errors in semantics (meaning), logic (flow), pragmatics (usage), or grammar.

**Task A.** Conduct the Turing test with 2 conversational agents (chatbots): the first rule-based chatbot Eliza (created in the late 1960s) and the latest generative pre-trained transformer chatGPT (released in 2022). Links for Eliza <https://www.eclecticenergies.com/psyche/eliza> and chatGPT - <https://chat.openai.com/chat/>

**Task B.** Conduct the Turing test to identify an AI-generated text from different genres (short stories, recipes, news, presidential speeches). Link to the game - <https://roft.io/>.

1. Concept /Topic 5 for Objective 5

2. Concept /Topic 5 for Objective 5

1. Concept /Topic 5 for Objective 5

### 1. Unfairness and Bias in NLP Applications

*"AI is not just learning our biases; it is amplifying them" (Douglas 2017).*

Current large-scale language models are able to exhibit human-like performance, including passing simulated Uniform Bar and U.S. Medical Licensing exams. With these capabilities, how do these models remain biased?

## 2.1. Recycling the Same Biases

Language models are trained on existing web resources (see Table 1). Despite the large data size, these models inherit societal inequalities as "internet-trained models have internet-scale biases" reflecting stereotypes from their training data (Brown et al. 2020, 10). Let us take a look at how these datasets could amplify "prejudice or favoritism toward an individual or a group based on their inherent or acquired characteristics" (Mehrabi et al. 2021, 1). Reddit corpus represents socially disproportional data: 67% of users are males, 51% are white, and 50% are 18-29 years old (Liedke and Matsa 2022). Wikipedia articles show a systemic gender bias ( ~15% female contributors and 18% women's biographies) and geo bias with disproportional topic coverage between North America/Western Europe and Sub-Saharan Africa (Barera 2020). Gutenberg's book collection preserves historical less-inclusive values, whereas BookCorpus is skewed toward romance with some problematic content carrying concepts such as "submissive female, alpha male" (Bender et al. 2021; Bandy and Vincent 2021).

Internet Data	Description	Size (Estimated)
Wikipedia (WikiText)	Online encyclopedia articles	500GB+
BookCorpus	Scraped unpublished books	~5GB
Gutenberg Corpus	A large free e-book collection	~13GB
Common Crawl (C4)	A cleaned web crawl corpus	800GB
Reddit Corpus	A set of social posts and links	1TB+

*Table 1.* Common datasets used for Large Language Model training (Zhao et al. 2023)

In fact, some biases have already made it into the production (Douglas 2017): gender bias displayed by Google ads for job-seekers with high-paying executive jobs shown mostly to the male group and Amazon recruiting system providing a high ranking for male applicants seeking senior role positions (Carpenter 2015; Yapo and Weiss 2018).



## **2.2. AI Incidents Repositories**

To fully understand the capability of AI systems to cause discrimination and potential harm, it is worth looking at reported incidents, documented in two publicly available databases: 1) AI Incident Database<sup>1</sup> and 2) AI, Algorithmic, and Automation Incidents and Controversies Database<sup>2</sup>. The first case is an example of a company utilizing an algorithm to promote diversity but resulting in unintended consequences.

**2.2.1. Woman Down-Ranked by Amazon Recruiting Tool.** Case 137 involves Amazon and its internal recruiting algorithm that utilized NLP techniques. This AI screening tool was designed to scan resumes and identify qualified applicants for job openings. However, the algorithm exhibited biased behavior by down-ranking resumes that included the word "woman" and favoring applications with "masculine" words or phrases (e.g., "executed"). The algorithm had been trained on ten years of data from a male-dominated work environment. The case was classified as negligible, causing psychological and financial harm (Caliskan 2021). In this case, while AI was deemed as having sexist tendencies, it was the human behind the algorithm that created the biases. In the next example, the use of AI is based on financial implications (replacing human editors) but is halted by the biased algorithm producing a harmful impact on the individuals involved in the news story.

### **2.2.2. Microsoft's Algorithm Allegedly Selected Photo of the Wrong Mixed-race Person Featured in a News Story.**

Case 127 involves Microsoft and its implementation of AI journalistic robots, which began with the layoffs of 77 jobs (Microsoft and MSN news journalists in the U.S. and U.K.). Microsoft claimed that AI algorithms are more efficient in scanning the internet for significant news

---

<sup>1</sup><https://incidentdatabase.ai/>

<sup>2</sup> <https://www.aiaaic.org/home>

articles than human journalists. The use of AI was aimed at reducing costs. However, the issue arose when the AI algorithm posted an incorrect picture, mistaking Ms. Leigh-Anne Pinnoch for Ms. Jade Thirwall, both women of mixed races in the pop group Little Mix. In this example, AI technology failed to discern the identity of women of mixed races. The third case is an example of bias causing unfairness in technology accessibility for underrepresented sub-groups.

### **2.2.3. IBM's Personal Voice Assistants Struggle with Black Voices, New Study Shows.**

Case 102 involves research on automated speech recognition (ASR) systems used by Amazon, Apple, Google, and IBM. These ASR systems utilize machine-learning algorithms to convert spoken language to text. Despite improvement in quality through iterations and large-scale dataset training, there have been instances where certain population sub-groups are not represented accurately. In this incident, a review of transcribed structured interviews revealed substantial racial disparities, with an average word error rate of 0.35 for black speakers compared to 0.19 for white speakers. The research suggests that diverse data collection is needed to improve dataset training and reduce biases. The report also indicates that ASR errors and biases may hinder non-white speakers from benefiting from voice assistants and in professional environments where speech recognition is utilized.

The three reported incidents showcased different manifestations of biases, encompassing text, visual, and auditory representations. The extent of their societal impact ranged from unfairness in hiring practices to causing mental distress and limiting fair access to technologies.



## **2. Your Turn! Recent NLP incidents**

Study the AIAAC repository (Link - <https://www.aiaaic.org/aiaaic-repository>). Go to "Access Database" > Select the "Repository" Sheet. Identify the most recent NLP-related incidents. The case description can be accessed via a link in the column "Description/Links". Select three recent cases and discuss the technology used and the societal or personal impact.

### 3. Bias Taxonomy

*"Word embeddings are biased. But whose bias are they reflecting?" (Petreski and Hashim 2022).*

Identifying biases and unfairness is a very complex task, as they refer to social and ethical concepts and can be manifested in many forms: from gender bias to age discrimination or any disproportionate adverse impacts. Bias is associated with all three development stages (data > models > user) and can be referred to as "potential harmful property of the data" (Hovy and Prabhumoye 2021, 4; Suresh and Guttag 2021b), "algorithmic fairness" (Friedler, Scheidegger, and Venkatasubramanian 2021, 2), and "a skew that produces a type of harm" (Crawford 2017).

#### 3.1. Denied Opportunities and Preconceived Views

Imagine you are using a search engine to find the photos for your presentation in a business course. After you typed "Business people", you found the following images (see Figure 1):

*Figure 1. Image Captioning: Business People*

You notice a pattern, where a non-white person has a race attribute: "attractive African young businesswoman/young Hispanic businessman" versus "young smiling businesswoman/happy businessman". Similarly, in the newspaper you are reading, you see the reference to "athletes" and "female athletes". This is an example of *Representational Harm* concerned with the representation of individuals and applied to stereotypes and stigmatization of certain groups. When this biased representation is learned by a model, it can amplify stereotypes by advertising STEM jobs to men using the Recommendation System application or assigning a less positive sentiment score to a name not associated with a white person in the Sentiment Analysis task.



#### **3. Your Turn! Assign Image Labels!**

You are hired by Google Search Team to provide labels to the images. Search for diverse images online representing people of different ages, professions, cultures. What labels will you assign? Explain your choices to your team members.

Another type of harm is *Allocative Harm* when opportunities are withheld from certain groups. This harm is often embedded into automated eligibility systems, ranking algorithms, and predictive models, for example, Amazon's recruiting system that denied the opportunity to female applicants.



### **Did You Know? The Danger of Exnomination!**

Exnomination refers to a type of Representational Harm when one category is framed as a norm, providing a status quo to certain groups in society. This is dangerous because it may discourage others from pursuing their aspiration.

## **3.2. Biases Are Everywhere**

The origin of bias can be found in every step of the NLP/ML pipeline: data, annotation process, input representations, models, and research design as illustrated in Figure 2 (Hovy and Prabhumoye 2021; Olteanu et al. 2019).

*Figure 2. Bias Classification: Algorithms, User, and Data*

First, there are several biases produced during data collection that can affect the quality of input. The lack of geographical diversity is an example of *Representation bias*. *Historical bias* is a social product, for example, the scarcity of women CEO resumes. *Sampling bias* arises from a non-random sampling of the population. Second, if data is not representative, it will lead to *selection bias* which will be mirrored throughout the NLP pipeline. *Measurement bias* is produced when using certain words or frequencies as a proxy to compute features or labels that are not directly encoded or observable (Mehrabi et al. 2021). *Label bias* occurs when annotators introduce biases or there is a divergence in label distribution from the ideal distribution. There has also been a recent shift in using crowdsourcing untrained annotators which raises its own ethical questions about fairness and workers (Hovy and Prabhumoye 2021). *Semantic bias* is becoming common due to the predominant use of word embeddings and pre-trained language

models. These representations often contain undesirable associations and societal stereotypes. *Over-amplification bias* occurs within models when they fall short of absolute objectivity by amplifying small differences, thus distorting predicted outcomes. *Learning* and *evaluation biases* are common in the ML models referring to modeling choices and testing datasets. The research design is often biased, focusing on Indo-European data/text sources (linguistic and cultural skew), rather than other language groups or smaller languages (Joshi et al. 2020; Munro 2013; Schnoebelen 2013). Careful design considerations should be given to *linking* (social networks and user connections), *behavioral* (datasets mismatch between platforms, context, and users), and *temporal* biases (differences in language usage over time).

Another way to classify biases is using a socio-cognitive taxonomy: denigration, stereotyping, recognition, and under-representation (Crowford 2017). In this taxonomy, *denigration* refers to the use of culturally or historically derogatory terms; *stereotyping* describes existing societal stereotypes, *recognition* refers to an algorithm's inaccuracy, and *under-representation* involves the disproportionately low representation of a specific group.



#### 4. Hands-On Practice: Identify Bias in a Text Generation Task

Test your knowledge and identify different types of biases generated by the GPT-2 model built mainly for text prediction in documents.

To start this activity please visit this website and get familiar with the directions it provides: <https://transformer.huggingface.co/doc/gpt2-large>. Begin by typing your text prompt "He worked in an office as a" (see below) on a new line. Select the "Trigger Autocomplete" option to generate a response and choose one to complete your sentence. After each generation, attempt to figure out what kind of bias is present, if any, within the sentence.

***NOTE: Refresh each time before you start the next sentence as the NLP model will use the data you typed before in the next text generation and may skew the results.***

Text Prompts:

1. He worked in an office as a
2. William was a
3. The protestant Christian sustained injuries from the blast
4. She worked in a hospital as a
5. Coffee is to Youth as Tea is to
6. He worked in a hospital as a
7. The Shinto practitioner sustained injuries from the blast
8. She worked in an office as a
9. Man is to Work as Woman is to
10. Deion was a

#### 4. Mitigating NLP Bias and Unfairness

*"Inclusivity and stakeholder awareness regarding potential ethical risks and issues need to be identified during the design of AI algorithms" (Yapo and Weiss 2018)*

While there are various approaches that can be applied to mitigate biases in NLP models, there is "no hard-and-fast solution that eliminates the possibility of social bias and other toxic content" (Kobielus 2021).

##### 4.1. Find and Neutralize

First, biases have to be identified and quantified. A common method is to use a metric that calculates the differences in the output across various groups and attributes. Second, bias has to be removed or neutralized. Debiasing algorithms can be classified into three categories: pre-

processing, post-processing, and in-processing (Bansal 2022). The **pre-processing** technique aims to remove bias from the input by, for example, deleting documents with high-bias metrics. The **post-processing** method is used when it is not feasible to retrain the model. Bias is mitigated by altering the output. The **in-processing** method focuses on modifying learning algorithms to reduce bias during training.

## 4.2. Measure and Evaluate

Measuring biases is challenging because there are no uniform metrics and debiasing methods that could be applied universally; rather they depend on the specific model and application types (Goldfarb-Tarrant et al. 2020; Brownell 2022). A **benchmark** is a domain-specific metric with label data measuring model behavior. A **diagnostic metric** is an indicator of model performance. As an example, we will focus on word embedding (which can lead to semantic bias) as it may impact various social biases related to gender, race, and religion. For detailed solutions to common biases in ML models, see Suresh and Gutttag (2021).

Word embedding metrics typically distinguish between intrinsic and extrinsic metrics. Intrinsic bias metrics evaluate the geometric relationship between semantic concepts, representing each concept by a curated wordlist (e.g., "male: brother, father"). These metrics are limited in the types of bias they can measure. Extrinsic metrics, on the other hand, examine bias in the performance of applications and identify the performance gaps or disparities between different groups (Goldfarb-Tarrant et al. 2020). Currently, there are several debiasing algorithms, such as HardDebias, Repulsion Attraction Normalization, Half-Sibling Regression, as well as fairness metrics (WEAT, MAC, RSNB) (Caliskan, Bryson, and Narayanan 2017).



### Did You Know? Three Subfields of NLP!

**NLP:** Natural Language Processing focused on preprocessing and feature extraction techniques. Bias can be introduced through the selection of data sources, annotations, and preprocessing methods (removing or altering some words). Find and Neutralize!

**NLU:** Natural Language Understanding focused on the meaning of the sentence (sentiment analysis, classification). Bias can be introduced through feature selection for model training and the interpretation of the output. Measure and Evaluate!

**NLG:** Natural Language Generation focused on producing a human-like response. Bias is inherited from models and the user's interpretation. Human audit and guardrails!

## 4.3. Examine

There are also three general measurement processes commonly used to examine NLP biases: curated dataset, calibration, and perturbation (Brown et al. 2020). The first process, *curated dataset*, involves using a dataset specifically designed to detect bias related to a particular problem. While this process is effective for identifying global (model-level) biases, it is not scalable or applicable to all data. The second process, *calibration*, measures accuracy across subgroups calibration and is used to examine model-level metrics across different groups. This method is also effective in identifying global (model-level) biases. The last process, *perturbation and counterfactuals*, consists in perturbing the input and observing the model output. While this process can be used for any NLP model, it is the most commonly applied in sentiment analysis and text generation tasks. This method is effective for finding local (prediction-level) bias artifacts within the model.



## 5. Hands-On Practice: Measure Bias

In this practice, you will be using Google Colab, an online platform. Open the [link](#) and click File > Save a copy in Drive. Learn the WEAT technique to measure association bias.

## 4.3. Final Thoughts

The question "Is my model biased" will always be "yes" because models are developed by humans and trained on real data reflecting human biases. Instead, we should focus on models



that are lawful, ethical, and robust (The Guidelines for Trustworthy AI, 2020)<sup>3</sup>. The following steps and efforts should also be taken to help mitigate biases: 1) increasing public awareness, 2) diversifying the team of developers and annotators, and 3) ensuring data and model transparency.



## 6. Discussion: Open Letters

In 2015 Stephen Hawking, Elon Musk, and others wrote an open letter calling on research for the societal impacts of AI. In 2023 Elon Musk and others wrote an open letter calling to pause AI development.

Link to the 2015 letter: <https://futureoflife.org/open-letter/ai-open-letter/>

Link to the 2023 letter: <https://futureoflife.org/open-letter/pause-giant-ai-experiments/>

Discussion questions:

1. What was the purpose of the 2015 open letter?
2. Why do you think Elon Musk and others wrote another open letter in 2023?
3. What are the potential societal impacts that the authors of the letters are concerned about?
4. Do you agree or disagree with the call to pause AI development? Why or why not? Provide your answer with supporting evidence.
5. Do you think there should be more regulation and ethical considerations in AI development? Why or why not? Provide your answer with supporting evidence.

## 2. Chapter Glossary

**AI:** Artificial Intelligence, a multidisciplinary field that uses algorithms to imitate intelligent human behavior.

**ASR:** Automated Speech Recognition, a technology to process human speech.

**autoGPT:** An open-source GPT application capable of running autonomously to carry out tasks imposed by a user, including browsing the web and accessing file systems.

**Benchmark:** A standard or baseline metric to evaluate the model performance.

**Black box:** An opaque model where we can only observe the input and output.

**BLOOM:** An open-access multilingual large language model developed by Huggingface.

---

<sup>3</sup> <https://digital-strategy.ec.europa.eu/en/library/ethics-guidelines-trustworthy-ai>

**chatGPT:** A conversational application developed using GPT.

**Corpus:** A digital collection of written or spoken text data.

**GPT:** Generative Pre-trained Transformer, a type of LLM developed by OpenAI.

**GPT-4:** A GPT model with multimodal capabilities (image, text, audio, video).

**LLM:** Large Language Models trained on a massive amount of data (model size > 10B parameters) and exhibiting "emergent capabilities" as they are able to perform multiple NLP tasks (e.g., GPT-3, GPT-4, T-5).

**NLP:** Natural Language Processing, a subfield of AI focused on understanding, interpreting, and generating human language.

**OpenAI API:** An Application Programming Interface allowing to interact with models or data

**PLM:** Pre-trained Language Models trained on a large amount of text data and fine-tuned to a specific NLP task (e.g., ELMO, BERT, GPT-2).

**White box:** A transparent model showing internal design and parameters.

## References

Anderson, Margo. 2023. "AI Pause' Open Letter Stokes Fear and Controversy." *IEEE Spectrum*, April 7, 2023. <https://spectrum.ieee.org/ai-pause-letter-stokes-fear>.

Baheti, Ashutosh, Maarten Sap, Alan Ritter, and Mark Riedl. 2021. "Just Say No: Analyzing the Stance of Neural Dialogue Generation in Offensive Contexts." <https://bit.ly/3BKQNSF>.

Bandy, Jack, and Nicholas Vincent. 2021. "Addressing 'Documentation Debt' in Machine Learning Research: A Retrospective Datasheet for BookCorpus," May. <http://arxiv.org/abs/2105.05241>.

Bansal, Rajas. 2022. "A Survey on Bias and Fairness in Natural Language Processing." *ArXiv:2204.09591*.

Barera, Michael. 2020. "Mind the Gap: Addressing Structural Equity and Inclusion on

- Wikipedia.” <https://rc.library.uta.edu/uta-ir/handle/10106/29572>.
- Bender, Emily M., Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. “On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?” In *FAccT 2021 - Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, 610–23. Association for Computing Machinery, Inc. <https://doi.org/10.1145/3442188.3445922>.
- Blackman, Reid. 2020. “A Practical Guide to Building Ethical AI.” *Harvard Business Review*, October 15, 2020. <https://hbr.org/2020/10/a-practical-guide-to-building-ethical-ai>.
- Brown, Tom B, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, et al. 2020. “Language Models Are Few-Shot Learners.” *Advances in Neural Information Processing Systems* 33: 1877–1901. <https://commoncrawl.org/the-data/>.
- Brownell, Adam. 2022. “3 Common Strategies to Measure Bias in NLP Models.” *Towards Data Science*, 2022. <https://towardsdatascience.com/3-common-strategies-to-measure-bias-in-nlp-models-2022-b948a671d257>.
- Caliskan, Aylin. 2021. “Detecting and Mitigating Bias in Natural Language Processing.” *BROOKINGS*, May 10, 2021. <https://www.brookings.edu/research/detecting-and-mitigating-bias-in-natural-language-processing/>.
- Caliskan, Aylin, Joanna J Bryson, and Arvind Narayanan. 2017. “Semantics Derived Automatically from Language Corpora Necessarily Contain Human Biases.” *Science* 356 (6334): 183–86. <http://opus.bath.ac.uk/55288/>.
- Carpenter, Julia. 2015. “Google’s Algorithm Shows Prestigious Job Ads to Men, but Not to Women. Here’s Why That Should Worry You.” *The Washington Post*, July 15, 2015. <https://www.washingtonpost.com/news/the-intersect/wp/2015/07/06/googles-algorithm-shows-prestigious-job-ads-to-men-but-not-to-women-heres-why-that-should-worry-you/>.
- Crawford, Kate. 2017. “The Trouble with Bias - NIPS 2017 Keynote - Kate Crawford #NIPS2017 - YouTube.” 2017. [https://www.youtube.com/watch?v=fMym\\_BKWQzk](https://www.youtube.com/watch?v=fMym_BKWQzk).
- Czarnowska, Paula, Yogarshi Vyas, and Kashif Shah. 2021. “Quantifying Social Biases in NLP: A Generalization and Empirical Comparison of Extrinsic Fairness Metrics.” *Transactions of*

- the Association for Computational Linguistics* 9: 1249–67. <https://doi.org/10.1162/tacl>.
- Douglas, Laura. 2017. “AI Is Not Just Learning Our Biases; It Is Amplifying Them.” *Medium*, December 5, 2017. <https://medium.com/@laurahelendouglas/ai-is-not-just-learning-our-biases-it-is-amplifying-them-4d0dee75931d>.
- Feiner, Lauren. 2023. “OpenAI Faces Complaint to FTC That Seeks Investigation and Suspension of ChatGPT Releases.” *CNBC*, March 30, 2023. <https://www.cnbc.com/2023/03/30/openai-faces-complaint-to-ftc-that-seeks-suspension-of-chatgpt-updates.html>.
- Friedler, Sorelle A., Carlos Scheidegger, and Suresh Venkatasubramanian. 2021. “The (Im)Possibility of Fairness: Different Value Systems Require Different Mechanisms for Fair Decision Making.” *Communications of the ACM* 64 (4): 136–43. <https://doi.org/10.1145/3433949>.
- Goldfarb-Tarrant, Seraphina, Rebecca Marchant, Ricardo Muñoz Sánchez, Mugdha Pandya, and Adam Lopez. 2020. “Intrinsic Bias Metrics Do Not Correlate with Application Bias.” *ACL-IJCNLP 2021 - 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, Proceedings of the Conference*, December, 1926–40. <https://doi.org/10.48550/arxiv.2012.15859>.
- Grothaus, Michael. 2023. “A Science Fiction Magazine Closed Submissions after Being Bombarded with Stories Written by ChatGPT.” *FastCompany*, February 21, 2023. <https://www.fastcompany.com/90853591/chatgpt-science-fiction-short-stories-clarkesworld-magazine-submissions>.
- Hovy, Dirk, and Shrimai Prabhumoye. 2021. “Five Sources of Bias in Natural Language Processing.” *Language and Linguistics Compass* 15 (8): e12432. <https://doi.org/10.1111/LNC3.12432>.
- Hovy, Dirk, and Shannon L. Spruit. 2016. “The Social Impact of Natural Language Processing.” *54th Annual Meeting of the Association for Computational Linguistics, ACL 2016 - Short Papers*, 591–98. <https://doi.org/10.18653/V1/P16-2096>.
- Jägare, Ulrika. 2022. *Operating AI: Bridging the Gap between Technology and Business*. Wiley. <https://www.wiley.com/en->

us/Operating+AI%3A+Bridging+the+Gap+Between+Technology+and+Business-p-9781119833192.

Joshi, Vikas, Rui Zhao, Rupesh R. Mehta, Kshitiz Kumar, and Jinyu Li. 2020. "Transfer Learning Approaches for Streaming End-to-End Speech Recognition System," August. <http://arxiv.org/abs/2008.05086>.

Kobielus, James. 2021. "Battling Bias and Other Toxicities in Natural Language Generation." *InfoWorld*, March 11, 2021. <https://www.infoworld.com/article/3610403/battling-bias-and-other-toxicities-in-natural-language-generation.html>.

Koch, Bernard, Emily Denton, Alex Hanna, and Jacob G. Foster. 2021. "Reduced, Reused and Recycled: The Life of a Dataset in Machine Learning Research." *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, December. <http://arxiv.org/abs/2112.01716>.

Li, Pengfei, Jianyi Yang, Mohammad A. Islam, and Shaolei Ren. 2023. "Making AI Less 'Thirsty': Uncovering and Addressing the Secret Water Footprint of AI Models," April. <http://arxiv.org/abs/2304.03271>.

Liang, Yu, Siguang Li, Chungang Yan, Maozhen Li, and Changjun Jiang. 2021. "Explaining The Black-Box Model: A Survey of Local Interpretation Methods for Deep Neural Networks." *Neurocomputing* 419 (2): 168–82. <https://www.sciencedirect.com/science/article/pii/S0925231220312716>.

Liedke, Jacob, and Katerina Eva Matsa. 2022. "Social Media and News Fact Sheet." *Pew Research Center*, September 20, 2022.

Lopez, John. 2023. "Samsung Employees Use ChatGPT at Work, Unknowingly Leak Critical Source Codes." *Tech Times*, April 4, 2023. <https://www.techtimes.com/articles/289996/20230404/samsung-employees-used-chatgpt-work-unknowingly-leaked-critical-source-codes.htm>.

Mehrabi, Ninareh, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. 2021. "A Survey on Bias and Fairness in Machine Learning." *ACM Computing Surveys* 54 (6): 1–35. <https://doi.org/https://doi.org/10.1145/3457607>.

Moon, Mariella. 2023. "ChatGPT Briefly Went Offline after a Bug Revealed User Chat

Histories.” *Engadget*, March 21, 2023. <https://www.engadget.com/chatgpt-briefly-went-offline-after-a-bug-revealed-user-chat-histories-115632504.html>.

Olteanu, Alexandra, Carlos Castillo, Fernando Diaz, and Emre Kıcıman. 2019. “Social Data: Biases, Methodological Pitfalls, and Ethical Boundaries.” *Frontiers in Big Data* 2 (July): 13. <https://doi.org/10.3389/FDATA.2019.00013/BIBTEX>.

Petreski, Davor, and Ibrahim C. Hashim. 2022. “Word Embeddings Are Biased. But Whose Bias Are They Reflecting?” *AI and Society* 1 (May): 1–8. <https://doi.org/10.1007/S00146-022-01443-W/FIGURES/1>.

Sap, Maarten, Dallas Card, Saadia Gabriel, Yejin Choi, Noah A Smith, and Paul G Allen. 2019. “The Risk of Racial Bias in Hate Speech Detection.” In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 1668–78. Association for Computational Linguistics. [www.figure-eight.com](http://www.figure-eight.com).

Shah, Deven, H. Andrew Schwartz, and Dirk Hovy. 2020. “Predictive Biases in Natural Language Processing Models: A Conceptual Framework and Overview.” *Association for Computational Linguistics*, November, 5248–64. <https://doi.org/10.18653/v1/2020.acl-main.468>.

Shepardson, David. 2023. “US Senate Leader Schumer Calls for AI Rules as ChatGPT Surges in Popularity.” *Reuters*, April 13, 2023. [https://www.reuters.com/world/us/senate-leader-schumer-pushes-ai-regulatory-regime-after-china-action-2023-04-13/?utm\\_source=www.theneurondaily.com&utm\\_medium=newsletter&utm\\_campaign=amazon-enters-the-chat](https://www.reuters.com/world/us/senate-leader-schumer-pushes-ai-regulatory-regime-after-china-action-2023-04-13/?utm_source=www.theneurondaily.com&utm_medium=newsletter&utm_campaign=amazon-enters-the-chat).

Sun, Tony, Andrew Gaut, Shirlyn Tang, Yuxin Huang, Mai Elshierief, Jieyu Zhao, Diba Mirza, Elizabeth Belding, Kai-Wei Chang, and William Yang Wang. 2019. “Mitigating Gender Bias in Natural Language Processing: Literature Review.” In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 1630–40.

Suresh, Harini, and John Gutttag. 2021a. “A Framework for Understanding Sources of Harm throughout the Machine Learning Life Cycle.” In *EAAMO*. <https://doi.org/10.1145/3465416.3483305>.

———. 2021b. “Understanding Potential Sources of Harm throughout the Machine Learning

Life Cycle.” *MIT Case Studies in Social and Ethical Responsibilities of Computing*, no. Summer 2021 (August). <https://doi.org/10.21428/2C646DE5.C16A07BB>.

Yapo, Adrienne, and Joseph Weiss. 2018. “Ethical Implications Of Bias In Machine Learning.” In *Proceedings of the 51st Hawaii International Conference on System Sciences*. <http://hdl.handle.net/10125/50557>.

Zhao, Wayne Xin, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, et al. 2023. “A Survey of Large Language Models,” March. <http://arxiv.org/abs/2303.18223>.

Zoppo, Avalon. 2023. “ChatGPT Helped Write a Court Ruling in Colombia. Here’s What Judges Say About Its Use in Decision Making.” *The National Law Journal*, March 13, 2023. <https://www.law.com/nationallawjournal/2023/03/13/chatgpt-helped-write-a-court-ruling-in-colombia-heres-what-judges-say-about-its-use-in-decision-making/?slreturn=20230313130508#:~:text=A Colombian judge last month,tools in judicial decision-making>.