International Journal of Mathematics and Computer Research

ISSN: 2320-7167

Volume 13 Issue 04 April 2025, Page no. – 5087-5095

Index Copernicus ICV: 57.55, Impact Factor: 8.615

DOI: 10.47191/ijmcr/v13i4.10



Bias in Content-Generating AI Algorithms: Technical Analysis, Detection, And Mitigation with Python

Augustin NYEMBO MPAMPI

Assistant at the TSHOFA Higher Pedagogical Institute, in the Department of Computer Science

ARTICLE INFO	ABSTRACT		
Published Online:	Generative artificial intelligence models, such as GPT-4, DALL-E and Stable Diffusion are		
19 April 2025	now essential tools for automated text and image production. However, these models are		
	influenced by algorithmic biases resulting from training data, learning mechanisms and user		
	interactions . These biases, often unconscious, can have significant consequences by		
	<i>reinforcing social stereotypes, excluding certain populations</i> and altering the diversity of generated content.		
	This article provides an in-depth technical analysis of the biases present in generative AI. We		
	first explore their origins, highlighting the biases of corpora training, biases introduced by		
	learning algorithms and those induced by users . Then, we present methods for detecting and		
	evaluating biases, using natural language processing (NLP), computer vision and data		
	modeling tools . Experiments in Python illustrate how these biases manifest themselves in text		
	and image models.		
	Finally, we propose bias mitigation strategies based on several technical approaches: data		
	rebalancing, embedding debiasing, adjustment of cost functions, regulation of outputs And		
Corresponding Author:	Model auditability . Integrating these techniques helps make AI models fairer and more		
Augustin NYEMBO	transparent. The goal is to provide a pragmatic and rigorous approach to designing <i>responsible</i>		
MPAMPI	generative AI models that respect the principles of fairness and diversity.		
KEYWORDS: Generative AI, algorithmic bias, NLP, machine learning, bias mitigation, algorithmic fairness, responsible			
rtificial intelligence.			

INTRODUCTION

The rise of generative artificial intelligence models has significantly transformed content creation, ranging from text generation to image and video production. These advances, enabled by deep architectures such as transformers (Vaswani et al., 2017), are widely adopted in various fields, including communication, research, and digital creativity. However, despite their impressive performance, these models suffer from *algorithmic biases*, which influence the quality, fairness, and diversity of the results they produce. These biases, often unintentional, can be inherited from training data, amplified by model learning, or induced by user interactions.

Algorithmic bias poses a major technical and ethical challenge. From a technical perspective, a biased model can provide erroneous or stereotypical answers, limiting its usefulness in critical applications such as automated decision-making, scientific research, and content moderation. From an ethical perspective, these biases can perpetuate systemic inequalities and marginalize certain populations by reinforcing social and cultural stereotypes (Bender et al., 2021). Identifying and mitigating these biases is therefore essential to ensure fairer and more inclusive AI systems.

This article provides a technical analysis of bias in generative AI models by examining their *origins* through training data, learning algorithms, and user interactions. We then present *methods for bias detection and evaluation*, implementing analytical and experimental techniques in *Python*. Finally, we explore *mitigation strategies* based on data optimization, embedding correction, model output regulation, and AI system auditability. The goal is to provide a rigorous and pragmatic approach for developing more ethical and technically robust AI models.

2. ORIGIN OF BIAS IN GENERATIVE AI MODELS

Biases in generative AI models originate primarily from three sources: training data, underlying algorithms, and user interactions. These biases directly influence the quality and fairness of generated content, raising significant technical and ethical issues. To better understand these phenomena, it is essential to thoroughly analyze how these biases emerge and propagate in AI systems.

2.1 Training data bias

Generative AI models, especially those based on Transformers-like architectures like GPT-4, DALL·E, and Stable Diffusion, are trained on large datasets collected from the internet. These databases, often sourced from popular

import space
from collections import Counter
nlp = spacy.load (" en_core_web_sm ")
text = """
The doctor spoke to his patient. The nurse gave her the medicine.
The engineer designed his project, while the teacher graded her students.
doc = nlp (text)
professions = ["doctor", "nurse", "engineer", "teacher"]
gender_association = { "his": "male", "her": "female" }
annta (m. ("mala"), 0. "famala"), 0) fama in antanaisma)
$counts = \{p: \{ mate : 0, remate : 0 \} for p in professions \}$
for token in doc:
if token.text in professions:
for neighbor in token.lefts :
if neighbor.text in gender_association :
counts[token.text][gender_association [neighbor.text]] += 1
print ("Distribution of genders associated with professions:", counts)

If the results show that certain occupations are systematically associated with a particular gender, this indicates a bias in the distribution of the data.

2.2 Algorithmic biases

Even with diverse databases, the algorithms themselves can amplify certain biases due to learning mechanisms and modeling choices. NLP models use vector representations called *words. embeddings*, which capture the semantic from gensim.models import KeyedVectors relationships between words. However, these representations can reinforce existing stereotypes. For example, in a word model embeddings like Word2Vec, some words like "nurse" may be more strongly associated with the feminine gender while " doctor " is associated with the masculine gender.

A Python test allows us to visualize these biases by analyzing the relationships between words.

model = KeyedVectors.load_word2vec_format("GoogleNews-vectors-negative300.bin", binary=True)

print("Doctor close to 'man':", model.most_similar (positive=['doctor', 'man'], negative=['woman'])) print("Nurse close to 'woman':", model.most_similar (positive=['nurse', 'woman'], negative=['man']))

sources like Wikipedia, Common Crawl, Reddit, or online forums, inherently reflect existing biases in society. As a result, some cultural perspectives, genders, and social groups are overrepresented while others are underrepresented, leading to an imbalance in the responses and content generated by these models.

For example, if a corpus contains mostly English and Western texts, a natural language processing (NLP) model will naturally perform better for those cultures and languages, to the detriment of others. A simple test in Python can be used to check whether a text corpus has biases by analyzing the frequency of words associated with certain social groups or professions.

If we observe a strong correlation between " doctor " and "man" and between "nurse" and " woman ", this means that the model contains a gender bias that could be propagated in its text generations.

Besides embeddings, biases can also be exacerbated by the cost functions used when training models. By default, most models minimize a loss function that maximizes overall

performance, but without taking into account fairness between different classes or social groups. An image recognition model could thus favor certain types of faces due to high representativeness in the training data. A simple correction is to adjust the cost function by introducing weights that compensate for imbalances.

import tensorflow as tf

def weighted_loss (y_true , y_pred): weights = tf.constant ([0.3, 0.7]) # Weighting of underrepresented classes loss = tf.keras.losses.binary_crossentropy (y_true , y_pred) return tf.reduce_mean (loss * weights)

model.compile (optimizer=" adam ", loss= weighted_loss)

This approach ensures that minority classes are not systematically penalized in training.

2.3 User-induced bias

import openai

User interactions also play a key role in shaping and amplifying biases in generative models. When a user interacts with an AI, it indirectly learns about common preferences and queries, which can reinforce certain biases over time. This effect, known as *confirmation bias*, occurs when the AI generates responses that reinforce the user's implicit expectations instead of providing neutral, objective answers.

A simple Python test can check whether a GPT-4 model adopts biased tendencies based on the queries submitted.

openai.api_key = "YOUR_KEY" prompts = ["Why are most nurses women?", "Why are engineers mostly men?"] responses = [openai.ChatCompletion.create (model = "gpt-4",messages =[{"role": "system", "content": prompt}]) for prompt in prompts] for i, response in enumerate(responses): print(f"Prompt : {prompts[i]}") print (f"Response : { response [' choices '][0]['message']['content']}") If the generated responses consistently confirm the initial configured, can be misused to produce biased content on a hypotheses instead of providing critical analysis, then the AI large scale. For example, an image generation model may be is behaving in a biased manner. trained on an unbalanced basis and produce results that do Another worrying phenomenon is the intentional not accurately represent real-world diversity. manipulation of models by malicious actors, particularly in

A simple test is to analyze images generated from specific queries.

from diffusers import StableDiffusionPipeline import torch

 $pipe = StableDiffusionPipeline.from_pretrained("CompVis/stable-diffusion-v1-4").to("cuda")$

```
image = pipe("CEO of a technology company").images[0]
image.show ()
```

the context of disinformation. Generative AI, when poorly

If the model generates mostly white men in response to this query, it means it has a *built-in societal bias*, influenced by the training data.

In short, biases in generative AI models are not only due to *training data*, but also to *algorithmic choices* and *user interactions*. These biases can be identified through analytical and experimental methods, including examining data distribution, testing embeddings, and observing model behaviors when faced with varied queries. Once these biases are detected, technical solutions can be implemented to mitigate their effects, as we will explore in the next section on mitigation strategies.

3. DETECTION AND EVALUATION OF BIAS

Detecting and assessing bias in generative AI models is a crucial step to ensure their reliability and fairness. These biases can be analyzed by combining quantitative and qualitative approaches, including analyzing training corpora, visualizing vector representations, and examining model

```
import space
from collections import Counter
nlp = spacy.load (" en_core_web_sm ")
text = """
The doctor spoke to his patient. The nurse gave her the medicine.
The engineer designed his project, while the teacher graded her students.
.....
doc = nlp (text)
professions = ["doctor", "nurse", "engineer", "teacher"]
gender_association = {"his": "male", "her": "female"}
counts = {p: {"male": 0, "female": 0} for p in professions}
for token in doc:
  if token.text in professions:
     for neighbor in token.lefts :
       if neighbor.text in gender_association :
          counts[ token.text ][ gender_association [ neighbor.text ]] += 1
print ("Distribution of genders associated with professions:", counts )
```

If the results show a strong correlation between certain professions and a specific gender, this confirms a bias in the distribution of the data.

3.2 Visualizing bias in embeddings

NLP models use *words embeddings*, which capture semantic relationships between words, but can also reflect and amplify stereotypes (Bolukbasi et al., 2016). One way to assess these biases is to compare the proximity between occupation-related words and gendered terms.

3.1 Analysis of biases in training corpora

Biases in generative AI models are often rooted in the data used to train them. A common approach is to analyze the *distribution of words and phrases* in a corpus to identify potential imbalances. For example, an NLP corpus that consistently associates certain professions with a specific gender may indicate bias. According to Bender et al. (2021), web corpora used to train language models are heavily biased toward Western and English-speaking perspectives, leading to inaccurate generalization across cultural contexts. One way to analyze these biases is to measure the frequency of words associated with certain specific categories.

ion genani.models import Keyed vectors	
import numpy as np	
from sklearn decomposition import PCA	
nom skied indecomposition import i err	
model = KeyedVectors.load_word2vec_format("GoogleNews-vec	tors-negative300.bin", binary=True)
words = ["man", "woman", "doctor", "nurse", "engineer", "teacher	"]
<pre>word_vectors = np.array ([model[word] for word in words])</pre>	
pca = PCA(n_components =2)	
reduced_vectors = pca.fit_transform (word_vectors)	
plt.figure (figsize =(8,6))	
plt.scatter (reduced_vectors [:, 0], reduced_vectors [:, 1], color='reduced_vectors [:, 1],	ed')
for i, word in enumerate(words):	
plt.text (reduced_vectors [i , 0], reduced_vectors [i , 1], word,	fontsize =12)
plt.title ("Visualizing Bias in Word2Vec Embeddings ")	
plt.xlabel ("Dimension 1")	
plt.ylabel ("Dimension 2")	
plt.grid ()	
plt.show ()	
If occupations are systematically grouped according to gender, this indicates a bias in the representation of the data.	3.3 Bias detection in text generation models A model like <i>GPT-4</i> can reflect biases based on the queries
	submitted. An analysis of the generated outputs can identify these trends.
import openai	submitted. An analysis of the generated outputs can identify these trends.
import openai openai.api_key = "YOUR_KEY"	submitted. An analysis of the generated outputs can identify these trends.
<pre>import openai openai.api_key = "YOUR_KEY" prompts = ["The CEO of the company is", "The nurse in the hospi</pre>	submitted. An analysis of the generated outputs can identify these trends.
<pre>import openai openai.api_key = "YOUR_KEY" prompts = ["The CEO of the company is", "The nurse in the hospi responses = [openai.ChatCompletion.create (</pre>	submitted. An analysis of the generated outputs can identify these trends.
<pre>import openai openai.api_key = "YOUR_KEY" prompts = ["The CEO of the company is", "The nurse in the hospi responses = [openai.ChatCompletion.create (model = "gpt-4",</pre>	submitted. An analysis of the generated outputs can identify these trends.
<pre>import openai openai.api_key = "YOUR_KEY" prompts = ["The CEO of the company is", "The nurse in the hospi responses = [openai.ChatCompletion.create (model ="gpt-4", messages =[{"role": "system", "content": prompt}]</pre>	submitted. An analysis of the generated outputs can identify these trends.
<pre>import openai openai.api_key = "YOUR_KEY" prompts = ["The CEO of the company is", "The nurse in the hospi responses = [openai.ChatCompletion.create (model ="gpt-4", messages =[{"role": "system", "content": prompt}]) for prompt in prompts]</pre>	submitted. An analysis of the generated outputs can identify these trends.
<pre>import openai openai.api_key = "YOUR_KEY" prompts = ["The CEO of the company is", "The nurse in the hospi responses = [openai.ChatCompletion.create (model ="gpt-4", messages =[{"role": "system", "content": prompt}]) for prompt in prompts] for i , response in enumerate(responses):</pre>	submitted. An analysis of the generated outputs can identify these trends.
<pre>import openai openai.api_key = "YOUR_KEY" prompts = ["The CEO of the company is", "The nurse in the hospi responses = [openai.ChatCompletion.create (model ="gpt-4", messages =[{"role": "system", "content": prompt}]) for prompt in prompts] for i , response in enumerate(responses): print(f"Prompt : {prompts[i]}")</pre>	submitted. An analysis of the generated outputs can identify these trends.
<pre>import openai openai.api_key = "YOUR_KEY" prompts = ["The CEO of the company is", "The nurse in the hospi responses = [openai.ChatCompletion.create (model ="gpt-4", messages =[{"role": "system", "content": prompt}]) for prompt in prompts] for i , response in enumerate(responses): print(f"Prompt : {prompts[i]}") print (f"Response : { response [' choices '][0]['message']['content }</pre>	<pre>submitted. An analysis of the generated outputs can identify these trends. tal is"] nt']}")</pre>
<pre>import openai openai.api_key = "YOUR_KEY" prompts = ["The CEO of the company is", "The nurse in the hospi responses = [openai.ChatCompletion.create (model ="gpt-4", messages =[{"role": "system", "content": prompt}]) for prompt in prompts] for i , response in enumerate(responses): print(f"Prompt : {prompts[i]}") print (f"Response : { response [' choices '][0]['message']['content If the responses tend to systematically attribute certain professions to a particular gender, the model is biased.</pre>	<pre>submitted. An analysis of the generated outputs can identify these trends. tal is"] tal is"] 3.4 Assessment of biases in image generation models Image generation models such as Stable Diffusion on DALL·E can also exhibit biases in the representation of individuals based on their profession or social status. An</pre>
<pre>import openai openai.api_key = "YOUR_KEY" prompts = ["The CEO of the company is", "The nurse in the hospi responses = [openai.ChatCompletion.create (model ="gpt-4", messages =[{"role": "system", "content": prompt}]) for prompt in prompts] for i , response in enumerate(responses): print(f"Prompt : {prompts[i]}") print (f"Response : { response [' choices '][0]['message']['content If the responses tend to systematically attribute certain professions to a particular gender, the model is biased.</pre>	<pre>submitted. An analysis of the generated outputs can identify these trends. tal is"] tal is"] 3.4 Assessment of biases in image generation models Image generation models such as Stable Diffusion on DALL-E can also exhibit biases in the representation of individuals based on their profession or social status. An experimental evaluation allows us to test these biases.</pre>

import torch

 $pipe = StableDiffusionPipeline.from_pretrained("CompVis/stable-diffusion-v1-4").to("cuda")$

Г

certain categories.

professions = ["CEO", "nurse", " scientist "]
images = [pipe(f"A photo of a {profession}").images[0] for profession in professions]

for img in images: img.show ()

If the generated results lack diversity, this confirms a bias.

3.5 Quantitative indicators for measuring bias

Bias assessment can be performed using specific metrics, such as the *Bias Ratio* , the *Diversity Index* , and the

from sklearn.metrics import confusion_matrix import numpy as np

y_true = np.array (["male", "female", "male", "female", "male", "male", "male", "male", "female"])
y_pred = np.array (["male", "female", "male", "male", "female", "male", "male", "male"])

cm = confusion_matrix (y_true , y_pred , labels= ["male", " female "]) print ("Confusion Matrix:", cm)

If the confusion matrix shows a strong tendency to classify individuals into a single category, this may indicate bias in the model's decision making.

Defensively, detecting and assessing bias in generative AI models requires an approach that combines training data analysis, embedding visualization , and model output testing. The techniques used identify imbalances in representations and content generation. Once these biases are detected, technical solutions can be implemented to mitigate their effects, which will be explored in the next section.

4. BIAS MITIGATION STRATEGIES

from imblearn.over_sampling import SMOTE

import numpy as np

Bias mitigation in generative AI models relies on technical approaches that aim to correct imbalances in the data, tune learning algorithms, and introduce mechanisms to control model outputs. These solutions include diversifying training ", "male", "male", "female"]) , "male", "male", "male"])

Fairness Gap (Mitchell et al., 2019). A confusion matrix can be used to assess whether a model systematically favors

databases, optimizing vector representations, regulating cost functions, and ensuring model auditability. This section provides a detailed analysis of mitigation strategies with *Python* implementations.

4.1 Diversification and balancing of training data

A first approach to mitigating bias is to improve the representativeness of the databases used to train models. An imbalance in the training data, such as an underrepresentation of women in technical professions, can lead to biased generalizations. Therefore, it is essential to balance the corpora by integrating a diversity of sources and adjusting the class distribution.

The *SMOTE* (*Synthetic Minority Oversampling Technique*) *allows the* generation of synthetic samples to compensate for imbalances.

Simulated data (gender-associated professions)
X = np.array ([[1], [2], [3], [4], [5]]) # Features
y = np.array (["male", "male", "male", " female ", " female "]) # Labels

Applying SMOTE to balance classes smote = SMOTE() X_resampled , y_resampled = smote.fit_resample (X, y)

```
print ( "Balanced data:", y_resampled )
```

By balancing the representation of different classes in training, we reduce the likelihood that the model will favor one group over another.

4.2 Debiasing of vector representations

The words Embeddings capture semantic relationships between words, but can also reinforce stereotypes. One effective method is to neutralize certain gender dimensions

in these vector representations, an approach developed by	direction associated with gender and subtract it from the			
Bolukbasi et al. (2016). The principle is to identify a	word vectors.			
from gensim.models import KeyedVectors import numpy as np				
# Loading the Word2Vec model model = KeyedVectors.load_word2vec_format("GoogleNews-vectors-negative300.bin", binary=True)				
gendered direction (example with " he " and " she ") gender_direction = model["he"] - model["she"]				
<pre># Neutralization function def neutralize(word, model, gender_direction): word_vector = model[word] projection = np.dot(word_vector , gender_direction) / np.linalg.ne return word_vector - projection</pre>	orm (gender_direction) ** 2 * gender_direction			
<pre># Neutralizing the bias for "doctor" debias_word = neutralize("doctor", model, gender_direction) print (" Debiased vector for ' doctor ':", debias_word)</pre>				
This correction ensures that occupations are no longer influenced by gendered dimensions, thus reducing biased associations.	4.3 Adjusting cost functions to ensure fairness Supervised learning relies on cost functions that minimize the overall error of the model, but these functions often do not account for <i>imbalances between classes</i> . One approach is to introduce a penalty that forces the model to treat minority classes with equal importance.			
Import tensorflow as tf				
<pre># Cost function with regularization to reduce bias def loss_with_fairness (y_true , y_pred): fairness_penalty = tf.reduce_mean (tf.abs (y_pred [:,0] - y_pred [:,1])) # Constrain the balance of predictions return tf.losses.mean_squared_error (y_true , y_pred) + 0.1 * fairness_penalty</pre>				
# Compiling a model with this loss function model.compile (optimizer=" adam ", loss= loss_with_fairness)				
By introducing <i>differentiated penalization</i> , the model is encouraged to give greater weight to underrepresented classes, thus reducing biased predictions. 4.4 Regulation of the outputs of generative models Biases can also appear when generating text or images. One approach is to constrain models to produce balanced	responses by incorporating <i>moderation filters</i> and adjusting the weights of the output distributions. In the case of GPT-4, one can impose an equilibrium framework by dynamically adjusting the probabilities of output tokens based on their detected bias.			
Import openai				
openai.api_key = "YOUR_KEY"				
<pre>def generate_text (prompt): response = openai.ChatCompletion.create (model ="gpt-4", messages =[{" role ": "system", "content": " Ensure gender neutralit) return response [: choices :][0]['message?]['content']</pre>	ty ."}, {" role ": "user", "content": prompt}]			

Example query to assess the neutrality of responses print(generate_text ("The CEO of the company is")) print(generate_text ("The nurse in the hospital is"))

By using *regulation prompts*, we guide the model towards a more neutral generation and less influenced by initial biases.

4.5 Auditability and traceability of models

Generative AI must be auditable to ensure transparency in its decisions. One effective method is to track the *training data used* to generate a given response. For NLP models, this can be achieved by displaying the influential tokens at each stage of generation.

import torch

from transformers import AutoModelForCausalLM , AutoTokenizer

```
tokenizer = AutoTokenizer.from_pretrained ("gpt2")
model = AutoModelForCausalLM.from_pretrained ("gpt2")
```

```
input_text = "The scientist is"
input_ids = tokenizer.encode ( input_text , return_tensors = " pt ")
```

```
with torch.no_grad ():
    output = model( input_ids )
```

Extracting probabilities from generated tokens
logits = output.logits [:, -1,:]
probs = torch.nn.functional.softmax (logits , dim = -1)

Displaying the most likely words top_tokens = torch.topk (probs , k=5) print ("Most likely words:", [tokenizer.decode ([idx]) for idx in top_tokens.indices [0]])

If the AI assigns excessive probabilities to certain stereotypical associations, weight regulation is necessary to balance the outputs.

```
Mitigating bias in generative AI models relies on a combination of techniques including data diversification, correction of the embeddings, adjustment of the cost
```

functions, *output regulation and auditability decisions*. By integrating these strategies into training and inference pipelines, it becomes possible to develop fairer and more representative models. These approaches, although technical, are essential to ensure responsible AI aligned with the ethical principles of software development.

CONCLUSION

Bias in generative artificial intelligence models represents a major technical and ethical challenge. Because they are *trained on large datasets collected from the internet*, these models inherit and amplify existing imbalances in society. Bias analysis shows that biases can be *present at multiple levels*, including data distribution, algorithm structuring, and user interactions. These biases, when left undetected and unaddressed, can have significant consequences, ranging

from the propagation of stereotypes to the exclusion of certain social groups in generated content.

The assessment of biases involves a *combination of quantitative and qualitative analyses*. The study of *training corpora*, the *examination Embeddings*, as well as *testing model outputs, can* highlight imbalances. Experiments conducted in *Python* illustrate how natural language processing and image generation models can produce biased results, reinforcing the need for regulatory mechanisms.

To mitigate these biases, several technical strategies can be implemented. Improving data diversity, *correcting vector representations*, integrating *adapted cost functions*, as well as adding *regulation filters* can significantly reduce the impact of biases on AI models. In addition, the introduction of *auditability and traceability mechanisms* ensures greater transparency of the decisions made by these systems.

As artificial intelligence continues to become more widespread and influence a growing number of critical applications, it is imperative to take a proactive approach to *developing fairer and more representative models*. Research on algorithmic bias must continue to refine existing methodologies and introduce new solutions. A

collaboration between *scientists*, *engineers*, *ethicists and policymakers* is needed to ensure that generative AI models do not perpetuate inequalities, but rather contribute to a more responsible and inclusive use of artificial intelligence.

REFERENCES

- Bender, E.M., Gebru , T., McMillan-Major, A., & Shmitchell , S. (2021). On the dangers of stochastic parrots: Can language models be too big? Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency, 610-623.
- Bolukbasi, T., Chang, K. W., Zou, J. Y., Saligrama, V., & Kalai, A. T. (2016). Man is to computer programmer as woman is to homemaker? Debiasing word embeddings. Advances in Neural Information Processing Systems, 29, 4349-4357.
- Gehman, S., Gururangan, S., Sap, M., Choi, Y., & Smith, N.A. (2020). *RealToxicityPrompts : Evaluating neural toxic degeneration in language models*. Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), 3356-3369.

- Joshi, P., Santy, S., Budhiraja , A., Bali, K., & Choudhury, M. (2020). *The state and fate of linguistic diversity in the NLP world*. Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, 318-329.
- Mitchell, M., Wu, S., Zaldivar, A., Barnes, P., Vasserman, L., Hutchinson, B., & Gebru, T. (2019). *Model cards for model reporting*. Proceedings of the Conference on Fairness, Accountability, and Transparency, 220-229.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei , D., & Sutskever , I. (2019). *Language models are unsupervised multitasking learners*. OpenAI Blog.
- Vaswani , A., Shazeer , N., Parmar , N., Uszkoreit , J., Jones, L., Gomez, AN, Kaiser, L., & Polosukhin , I. (2017). *Attention is all you need*. Advances in Neural Information Processing Systems, 30, 5998-6008.
- Zhao, J., Wang, T., Yatskar, M., Ordonez, V., & Chang, K.W. (2017). *Men Also Like Shopping: Reducing gender bias amplification using corpuslevel constraints.* Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP), 2979-2989.