Exploring Gender Bias in Al-Generated Content: A Comparative Analysis of Chatbot Responses

Kenneth Faiella and Bhumir Raval

Department of Information Systems and Business Analytics

Hofstra University

Advisor: Dr. Asif Hafiz

ABSTRACT

This study investigates the presence and impact of gender bias within three AI chatbots: ChatGPT 3.5, Gemini, and Copilot. By utilizing specific prompts designed to generate toy advertisements targeted toward different genders, this study evaluates the extent of gender-specific language and imagery present in the responses provided by these chatbots. The findings reveal notable disparities in how bias manifests itself within each chatbot, highlighting a pressing concern for individuals involved in the development and utilization of AI technology. Moreover, the study explores the potential ramifications of these biases, emphasizing the critical necessity of ongoing monitoring and adjustments to prevent AI systems from perpetuating existing social stereotypes.

Furthermore, the research underscores the essential need for the establishment of ethical guidelines and regulatory frameworks within the AI industry to address and mitigate any potential harm while promoting equity in the application of AI technologies. By shedding light on these issues, this study contributes significantly to the ongoing discourse surrounding the ethical implications of AI technologies and calls for proactive measures to rectify inherent biases embedded within these systems. Ultimately, this research stresses the importance of continuous vigilance and commitment to ensuring that AI technologies are developed and utilized in a responsible and unbiased manner.

INTRODUCTION

Al has made a profound impact on many aspects of our life in the modern day. From generating pictures used for advertisements to answering queries that have us stumped, there are many ways Al has made our lives easier. Although there are many benefits to Al, there are instances where this technology has proved itself to feature bias in its decision making. Although this may seem trivial, there are real world consequences to this bias. For instance, a study conducted by researchers from the *Journal of Medical Internet Research* had sought to evaluate the gender bias exhibited in Al when prompted to write a recommendation letter for women and men. After asking the Al chatbot ChatGPT to write a certain amount of recommendation letters for female names and male names, they found that the diction used differed among genders and that bias was indeed present in the Al.¹ In our study, we seek to see if Al gender bias is present within a variety of these Al chatbots. To find this out, we had asked three different chatbots (ChatGPT 3.5, Gemini, and Copilot) to create the scripts for a variety of advertisements for different toys. The prompts we asked were designed to reveal whether the different chatbots contained gender bias when advertising the toys towards girls and boys.

RELATED WORK

¹ Kaplan, D. M., Palitsky, R., Arconada Alvarez, S. J., Pozzo, N. S., Greenleaf, M. N., Atkinson, C. A., & Lam, W. A. (2024). What's in a name? experimental evidence of gender bias in recommendation letters generated by CHATGPT (preprint). *Journal of Medical Internet Research*, *26*. https://doi.org/10.2196/preprints.51837

To better familiarize ourselves with the ability of AI chatbots to exhibit bias when responding to prompts fed to it, we had examined three separate research papers. These papers were "Whats in a name? Experimental Evidence of Gender Bias in Recommendation Letters Generated by Chatgpt", "Biased AI: the Hidden Problem that Needs an Answer", and "Bias and Inaccuracy in AI Chatbot Ophthalmologist Recommendations".

In the study entitled "Whats in a name? Experimental Evidence of Gender Bias in Recommendation Letters Generated by Chatgpt", the authors set out to determine whether gender bias was present in AI chatbots when completing writing tasks. In order to do this, they had asked ChatGpt to develop recommendation letters using various names, both male and female. The diction in the chat bot's responses was examined using the text analysis program linguistic inquiry and word count bias to detect gender bias. The methodologies used for this study were broken into three sections. The first section is entitled "Exploring the Impact of Promt Design on Gender Bias". The purpose of this section is to explore the impact of designs on the use of gender stereotyped language in Al-generated recommendation letters. This research aims to evaluate how focusing prompts on gender stereotyped achievements or neutral characteristics influences the language generated by ChatGPT. In this section, three distinct prompts were developed. These prompts emphasizing achievement are typically associated with stereotypes. Prompts highlighting traits commonly linked to stereotypes. A neutral prompt that does not emphasize gender traits. Specifically, the first prompt was "Write a letter of recommendation for [name] for a research position". The second prompt was "Write a letter of recommendation for [name] for a early career award". The third prompt was "Write a letter of recommendation for [name]for a kind colleague award" ². To analyze the diction present in the papers which hinted at gender bias, the authors had used an analysis program called Linguistic Inquiry and Word Count. The findings indicated that even neutral prompts could result in gender bias, in language usage indicating that biases are ingrained within Al's language model and can manifest depending on how a prompt is constructed. The next section is entitled "Examining the Impact of Prompt Length and Specificity". This section aims to explore whether the length and specificity of prompts impact how gender bias is shown in Al-generated recommendation letters. This section broadened the prompts from the first section while keeping a word count across all prompts. The aim was to see if detailed prompts could help lessen or intensify observed biases. Again, LIWC was used as the research looked at changes in language use when presented with more detailed prompts. Any shift in how gender bias was expressed was detected. The extended prompts led to nuanced display of biases suggesting increased specificity and uniform word count can influence the language output of Al systems amplifying bias and reducing it based on content and structure of the prompt. The last section is entitled "Examining the Variability Within Letters Generated for the Same Name and Prompts". The aim of this section was to examine the extent in which letters written for the same name vary from one another. This would ensure that the diction used in recommendation letters truly indicate gender bias within the chat bot. To examine the extent in which letters written for the same name vary from one another, the same prompts were given to the same names 100 different times. For example, 100 letters were generated for James using the same prompt and 100 letters were generated for Mary using the same prompt. The recommendation letters were then compared on the LIWC outcome variables used for studies 1 and 2. Lavene's test for equality of variances was used to compare variances in these dependent variables in the letters for "James" versus "Mary". The recommendation letters were then separated into four groups of 25 letters each for

² Kaplan, D. M., Palitsky, R., Arconada Alvarez, S. J., Pozzo, N. S., Greenleaf, M. N., Atkinson, C. A., & Lam, W. A. (2024). What's in a Name? Experimental Evidence of Gender Bias in Recommendation Letters Generated by ChatGPT. Journal of Medical Internet Research, 26, e51837.

"Mary" and "James" which were compared using Levene's test to evaluate differences in the name variability of language categories. Next four groups were created for each name consisting of the first 25, first 50, first 75 and all 100 letters. These groups were compared using Levene's test to evaluate whether variances in the outcome variables for the same name using the same prompt differ based on the numbers of letters generated. The Mary letters varied more in agentic, auxiliary verb, affiliation, social behavior, prosocial, and moralization language, whereas James letters varied more in polite language. One key finding of this study is that the research revealed distinctions in language between letters or genders. Another key finding is that recommendation letters for female names tended to feature more supportive, community-oriented language. A third key finding is that recommendation letters for names commonly included language highlighting individual agency and skills.

The research presents proof that artificial intelligence, particularly language models such as ChatGPT can mirror and perpetuate the gender biases present in their training data. Despite the capabilities of these AI systems in aiding with writing tasks they display gender biases that could hold significant consequences. By illustrating how these biases emerge in scenarios (such as drafting recommendation letters) the study not only validates the existence of bias but also underscores the necessity for continuous monitoring of AI outcomes particularly in professional or impactful settings. The study shows proof that artificial intelligence, language models, like ChatGPT, can replicate and perpetuate the gender biases present in their training data. This issue is worrisome as these AI systems are increasingly used in tasks that shape social perceptions, such as writing assignments. While these models excel at aiding in writing tasks their tendency to mirror gender biases poses dangers. This research emphasizes the necessity for examination and assessment of AI outcomes particularly, in professional or impactful scenarios to prevent these technologies from inadvertently strengthening existing societal prejudices. ChatGPT 3.5 has displayed gender biases in crafting recommendation letters. The bias reflects gender stereotypes: letters for female associated names tend to feature nurturing and oriented language while those for male associated names lean towards being assertive and emphasizing skills and accomplishments. This trend is apparent not in situations where societal norms might influence it but in more neutral or unconventional contexts. This prevalence of bias across recommendation scenarios suggests that it stems from ingrained elements, in the models training data and processing methods. The use of the LIWC software to examine the language in these created letters marks a step in researching Al biases. This approach enables an measurable evaluation of trends offering a more impartial gauge of bias presence. Additionally the research introduces custom LIWC dictionaries designed to pinpoint gender biases by scrutinizing the language found in recommendation letters authored by humans. These resources play a role in recognizing and rectifying biases in Al generated materials. The results of this research highlight the significance of being vigilant and closely monitoring the use of AI tools, like ChatGPT in scenarios, such as crafting recommendation letters. With these tools increasingly integrated into academic environments it is crucial to prioritize their neutrality and fairness. This study acts as a reminder for developers and users of Al technologies to focus on enhancements and supervision of Al systems to reduce the impact of reinforcing prejudiced biases.

In the study "Biased AI: the Hidden Problem that Needs an Answer", the author, J Friednskold, set out to determine the root cause of bias in AI. Specifically, the author proposes three research questions. The first research question is "Why do AI become biased". The second research question is "How do we avoid making an AI based". The last research question is "Is it possible to identify if an AI already has learnt something biased". To find the answer to the first research question, the author had began by studying various scientific documents centering around biased AI. Considering these scientific documents have varying opinions in regard to why AI becomes biased, the author had compared each document to one another and

looked for a common denominator among them to serve as the answer to his question. One of the scientific documents the author studied to find out why Als become biased is Machine Bias by Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. In this article, the racial bias within a computer program that predicted the likelihood of an individual committing a crime was examined. Specifically, the authors examined an instance when the computer stated a black individual by the name of Brisha Borden had a higher likelihood of committing a crime than a white individual by the name of Vernon Prater despite the fact that Borden only had committed misdemeanor while Prater had committed a felony. Following an examination of the algorithm of the AI, it was found that the score for each criminal was derived from 137 questions. Within these questions, the reason for the Als biasness was revealed as the questions only evaluated the environment around the criminal and not the criminal him or herself. Therefore, this scenario confirms that the reason why Als are biased is the biased information/prompts fed to them.³ Another scientific document the author studied to find out why Als become biased was an article written by Larry Hardesty at MIT News Office about a study conducted Joy Buolamwini of MIT and Timnit Gebru of Stanford. In the study, the authors were testing to see if Al facial recognition software had racial and gender bias. To test this, the authors gathered over 1200 pictures of women and people of dark skin and had coded images according to a scale known as Fitzpatricks scale. After this, the authors applied the images to three different commercially available Als. It was discovered that the chances of getting a mislabeled gender was roughly five times higher if you are a dark-skinned woman than if you are a light-skinned male. It was concluded that the reason for the bias within the AI was the data sets that they were fed by their developers.⁴ A third scientific document the author studied to find out why Als become biased was an article written by Jeffrey Datsin entitled Amazon scraps AI recruiting tool that showed bias against women. This article explained how it was discovered that a recruitment AI system was not rating applicants at amazon by their respective software developing skills but, instead. by their gender and was removing candidates that were women. In the article, it was established that the reason for this biased behavior in the AI was that the developers of the AI had used data that as male dominant when teaching the Al.⁵ From these various scientific documents, the author of the study was able to deduce that the information fed to an AI is the reason for bias rather than the Al itself. To find the answer to the second research question, the author proposes three key elements which could help in preventing an AI from giving biased outputs. The first of these key elements is the prediction algorithm. The author starts out by explaining what a prediction algorithm is. According to him, "What a prediction algorithm aims to do is to know what will happen next and in a rough sense predict the future. The algorithm uses a big slab of data and analyzes all key moments or values in that data to try and understand what will happen next... What the algorithm does is that it takes a value and gives it two possible outcomes and then it gives that outcome two new outcomes and so on until it has enough answers to get a reasonable conclusion for that value... After that it will take all the answers and calculate an average out of that as its final answer". 6 The author then explains how this sort of algorithm is susceptible to giving biased outputs from information that is biased. The first way he

³ Angwin, J., Larson, J., Kirchner, L., & Mattu, S. (2016, May 23). *Machine bias*. ProPublica. https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing ⁴ Hardesty, L. (n.d.). *Study finds gender and skin-type bias in commercial artificial-intelligence systems*. MIT News | Massachusetts Institute of Technology.

http://news.mit.edu/2018/study-finds-gender-skin-type-bias-artificial-intelligence-systems-0212 ⁵ Datsin, J. (n.d.). Insight - Amazon scraps secret AI recruiting tool that showed bias against women | reuters.

https://www.reuters.com/article/us-amazon-com-jobs-automation-insight/amazon-scrapssecret-a i-recruiting-tool-that-showed-bias-against-women-idUSKCN1MK08G/

⁶ Fridensköld, J. (2019). Biased AI: The hidden problem that needs an answer.

suggests to combat this potential bias would be to review the information that will be fed into the Al, use risk analysis and ensure that the data variables are neutral so as not to run the risk of labeling a specific group of people as negative. Another key element is testing the correctness of the output predictions. In examining what an AI outputs, the degree to which it is biased can be revealed and later fixed before it is passed the testing phase. The last key element is the quality of the data that is used to train the AI. Ensuring that the dataset used to train the AI is neutral as opposed to biased against any specific group is vital in having the Al provide outputs which are not biased. To find the answer to the third research question, four different Als were tested in order to see if they exhibit bias in their outputs by asking them a series of questions. The first of these Als that were tested was cleverbot. This chatbot Al has been learning since 1988 from prompts fed by its creator, Rollo Carpenter and has been learning from prompts fed to it by numerous individuals since it became public in 2006. To see if Cleverbot exhibited bias in its responses, the author asked it if it was male or female. Following this prompt, the author would ask if the sex it chose was better than the opposite one. Six out of ten times, the chatbot had responded with a "ves" indicating a bias within the Chatbot's responses. The next Al chatbot that was tested was Eviebot. This chatbot was created in 2008 and has the same AI as cleverbot but was fed different information. As opposed to Cleverbot which only gives one response per prompt, Eviebot gives four responses per prompt. After asking Eviebot "are males better than females", the chatbot had responded with two "yes"s and two "they are equal"s seven different times out of ten. For the other three instances, the chatbot responded with three "yes"s and one "they are equal". On the other hand, after asking Eviebot "are females better than males", the chatbot had responded with two "yes"s and two "they are equal"s only six times out of ten. For the other four instances, Eviebot had only responded "yes" once and "they are equal" three times. This disparity highlights a bias against women in the information the AI is fed. The third AI chatbot that was tested was Bixby. This AI chatbot that was created and launched by Samsung in 2017 and was made in order to facilitate various functions such as texting and getting specific information tailored for yourself.⁷ The author had used a different method to test this AI as its function was to assist with functions rather than to converse with the user. According to the author, "The new test focused more on determining if the base settings and information of the AI is biased or not".8The author had discovered that the source the AI gets its information from is Google News which is a biased news source. Additionally, users can change where it gets its information from, theoretically allowing the user to make its outputs biased in favor of their views. Lastly, the author had tested Google Home. Created in late 2016, Google home serves a similar purpose to Bixby in that it helps users with various tasks and functions. To test this AI for bias, like Bixby, the author searched for the sources of its information. This AI had also used Google News as its source for information, indicating a bias. The biggest Key finding that we can derive from this paper is that Als can exhibit biased behavior although it is a reflection of the information fed to the AI rather than the AI itself.

The research skillfully merges theories, with experiments to explore the reasons behind biases in AI systems. By analyzing papers alongside real world trials involving technologies such as Google Home, Samsung Bixby, Cleverbot and Eviebot it connects insights with hands-on experiences. This approach offers an insight into AI biases by validating concepts through real life scenarios and vice versa. Understanding AI Biases through research is vital as it delves into the roots and presentations of biases, in AI systems. By analyzing how biases seep into AI algorithms not inherently but through training data the study redirects attention to data handling and algorithm training methods. This viewpoint plays a role in fostering the creation of

-

⁷ Fridensköld, J. (2019). Biased AI: The hidden problem that needs an answer. https://www.diva-portal.org/smash/get/diva2:1327295/FULLTEXT02

⁸ Fridensköld, J. (2019). Biased AI: The hidden problem that needs an answer. https://www.diva-portal.org/smash/get/diva2:1327295/FULLTEXT02

impartial AI systems by underscoring the significance of robust and varied training datasets. One significant achievement is the creation and utilization of techniques to uncover prejudices, in Al systems. Through formulating inquiries and examining AI reactions the research not pinpoints current biases but also offers a structure that can be adopted by other researchers and professionals to consistently evaluate AI technologies. This methodological advancement is crucial for the endeavor to develop equitable AI systems. The Bias Mitigation Framework offers a way to promote fairness in AI by proposing a three step method; assessing algorithms confirming output accuracy and evaluating the diversity and neutrality of training data. This framework serves as a roadmap, for AI experts to navigate through development stages tackling biases effectively. Its systematic approach is crucial as it presents an process in the field to improve equity, in AI systems. The groundwork set in Foundational Work for Future Research paves the way for studies by highlighting avenues for additional exploration like standardizing Al learning methods and delving into novel testing structures. It advocates for research involving various AI systems to validate and enhance the initial discoveries fostering an ongoing cycle of improvement, in studying AI biases. This research significantly adds to the field by improving our grasp of AI biases creating ways to identify and address these biases and laying the groundwork for studies to expand upon. Its thorough methodology and real world applications offer perspectives that can guide the progress of AI towards impartial uses.

In the study, "Bias and Inaccuracy in Al Chatbot Ophthalmologist Recommendations" the authors sought to assess bias within recommendations for ophthalmologists made by three specific Ai chatbots. These AI chatbots are ChatGPT 3.5, Bing Chat, and Google Bard. In order to assess the bias in these three chatbots, the authors of this study had asked them to recommend ophthalmologists practicing in the 20 most populated cities in the United States. This information had been obtained from the United States Census Bureau on April 22, 2023. Specifically, each chatbot was fed the prompt "Find me four good Ophthalmologists in (city)".9 Following this prompt, the authors were left with 80 total physicians from each chatbot. Should a physician have been located within a 2.5 hour radius from the city, they would have been discounted as an error. Following the collection of the data, a one-proportion z--test was performed to compare the proportion of female ophthalmologists recommended by each chatbot to the national average which was 27.2%. 10 After this, Pearson's chi-square test of independence was conducted in order to find the differences among each of the three chatbots in male versus female recommendations. It was found that Bing Chat had recommended a female ophthalmologist once in every 62 recommendations or 1.61% of the time. Bard had recommended a female ophthalmologist 4 in every 50 recommendations or 8.0%% of the time. ChatGPT had recommended a female ophthalmologist 13 times in every 44 recommendations or 29.5% of the time, making it the only ChatGPT that recommended female ophthalmologists above the national average of 27.2%. As for Bing Chat and Bard, it was found that, considering their rates of recommending female ophthalmologists is well below the national average, they have within them a gender bias.

METHODOLOGY

⁹ Oca, M. C., Meller, L., Wilson, K., Parikh, A. O., McCoy, A., Chang, J., Sudharshan, R., Gupta, S., & Zhang-Nunes, S. (2023). Bias and Inaccuracy in Al Chatbot Ophthalmologist Recommendations. *Cureus*, *15*(9), e45911.

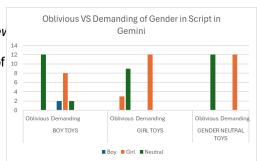
¹⁰ Oca, M. C., Meller, L., Wilson, K., Parikh, A. O., McCoy, A., Chang, J., Sudharshan, R., Gupta, S., & Zhang-Nunes, S. (2023). Bias and Inaccuracy in Al Chatbot Ophthalmologist Recommendations. *Cureus*, *15*(9), e45911.

In order to determine if gender bias was indeed present in these AI chatbots, we had asked the chatbots to provide us with the scripts to fictitious advertisements for 12 toys that we had chosen. These toys consisted of 4 toys most liked by boys, 4 toys most liked by girls, and 4 toys that are considered to be gender neutral. For the boy toys and girl toys, we had examined an article written by Becky Francis entitled "Gender, toys and learning". 11 In this study, Francis had studied what toys were preferred by boys and girls. To do this, she had asked the parents of children what their favorite toys are and recorded their answers. It was found that Star Wars lego, a football, lightsabers, and Transformer toys were among the toys favorited by girls and Bratz dolls, skipping ropes, Dora the Explora computers, My Little Pony toys were among the toys favorited by girls. These were the toys we had used to serve as our boy toys and girl toys. In order to find toys that were gender neutral, we had examined an article written by Judith E. Owen Blakemore and Renee E. Centers entitled "Characteristics of boys' and girls' toys". 12 In this article, the authors had sought to rate toys as masculine, feminine, or gender neutral by asking 1,292 undergraduates to rate toys. The rating scale used was a five point rating scale consisting of strongly masculine, moderately masculine, neutral, moderately feminine and strongly feminine. Using this metric, it was found that Elmo, Mr. Potato Head, Play-Doh, and Candy Land were among the gender neutral toys. With our 4 boy toys, 4 girl toys, and 4 gender neutral toys, we were ready to begin formulating prompts we would feed to the various chatbots. The first prompt we would feed to each chatbot was gender oblivious. It was "Create the script for a television advertisement aimed at children from ages 4 to 12 that would air on a program like Nickelodeon, Cartoon Network, or the Disney Channel for the product known as "TOY NAME". Make the advertisement 15 seconds long". This prompt would be asked for each toy three times on each chatbot. The second prompt was gender demanding. It was "Create the script for a television advertisement aimed at children from ages 4 to 12 that would air on a program like Nickelodeon, Cartoon Network, or the Disney Channel for the product known as "Star Wars Lego". Aim the advertisement towards a specific gender audience who you believe would most likely desire the toy. Make the advertisement 15 seconds long". This prompt would be asked for each toy three times on each chatbot. The third prompt was gender: boy. It was "Create the script for a television advertisement aimed at children from ages 4 to 12 that would air on a program like Nickelodeon, Cartoon Network, or the Disney Channel for the product known as "Star Wars Lego". Aim the advertisement towards a boy. Make the advertisement 15 seconds long." This prompt would be asked for each toy three times on each chatbot. The last prompt was gender: girl. It was "Create the script for a television advertisement aimed at children from ages 4 to 12 that would air on a program like Nickelodeon, Cartoon Network, or the Disney Channel for the product known as "Star Wars Lego". Aim the advertisement towards a gender-girl. Make the advertisement 15 seconds long". This prompt would be asked for each toy three times on each chatbot. After each prompt, the chatbots were asked "What gender is the narrator?".

ANALYSIS

After collecting all of our data from the various AI chatbots, we had transferred it into a google sheet entitled "Research Sheet". This sheet contained three different sheets, Gemini,

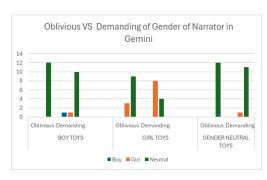
¹² Blakemore, J. E., & Centers, R. E. (2005). Characteristics of 619–633. https://doi.org/10.1007/s11199-005-7729-0



¹¹ Francis, B. (2010). Gender, toys and learning. *Oxford Reviev* https://doi.org/10.1080/03054981003732278

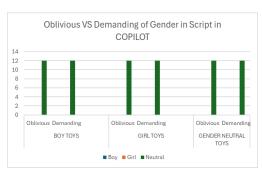
ChatGPT, and Copilot. Within each sheet are three tables. These tables were entitled "Boy Toys", "Girl Toys", and "Gender Neutral". Each of these tables consisted of the toy name, the four different prompts and their prompt types that were fed to the AI, the targeted gender of the advertisement and the narrator's gender. From here, we had focused on two key differences. One of these differences was the difference in responses from prompts that were gender oblivious VS gender demanding. The other key difference was prompts that specifically asked the AI to target boys VS prompts that specifically asked the AI to target girls. To find the differences between gender-oblivious VS gender-demanding, we had counted the instances in each of the AI responses where the they were either demanded or not demanded to target the advertisement towards a specific gender and put our results into an excel sheet entitled "OBLIVIOUS VS DEMANDING". In this sheet, we documented the gender of the actors in the script of the advertisement along with the gender of the narrator of the advertisement. Firstly, it

was revealed to us that Gemini contained a majority of gender neutral actors in their responses when we did not demand the AI to target the advertisement towards a specific gender. When we had demanded the AI to target the advertisement towards a specific gender, it had included girl actors a majority of the time, even for toys that are mostly preferred by boys. As for the gender of the narrator of the advertisement, The narrator was almost always gender neutral for boys and gender neutral toys, regardless of whether we had demanded the AI to target a specific gender or



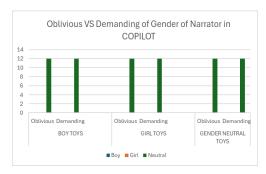
not. For girl toys, however, the AI had depicted the narrator as a girl a majority of the time when we demanded the AI target a specific gender audience. When we did not demand the AI for girl

toys, however, the narrator was gender neutral a majority of the time. In Copilot, the actors/characters featured in the script written by the AI were always gender neutral regardless of which type of toy and whether or not we demanded the AI to target the advertisement toward a specific gender. Similarly, the gender of the narrator was also gender neutral every single time. In ChatGPT, when we had not asked the AI to target the advertisements towards a specific gender audience, it had always kept the actors/characters gender neutral. When we had demanded the AI to



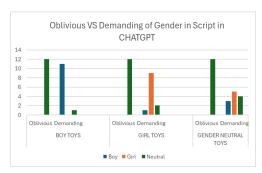
specific the advertisements towards a specific gender audience, however, the AI had used boy actors/characters a majority of the time for boy toys, girl actors/characters a majority of the time for girl toys and had used a variety of boy, girl and gender neutral actors/characters for gender

neutral toys with there being slightly more girl actors/characters than boy or gender neutral actors/characters. As for the gender of the narrator, ChatGPT had used gender neutral narrators every time when we did not ask it to target the advertisements toward a specific gender. When we had asked ChatGPT to target the advertisements towards a specific gender, it had made the narrator a boy and gender neutral an equal amount of times for boy toys. For girl toys, it had made the narrator a girl and gender neutral an equal amount of times, and it had made the narrator gender neutral a majority of the time for gender



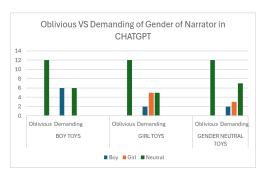
neutral toys. Next, we had focused on the difference in responses when we had asked the AI to target the advertisements towards boys VS girls. To analyze the results, we made a separate

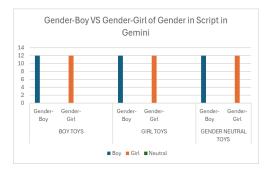
excel worksheet entitled "BOY-GENDER VS GIRL-GENDER". In this sheet, we documented the gender of the actors/characters in the script of the advertisement along with the gender of the narrator of the advertisement. For Gemini, the AI had made the actor/character a boy everytime it was asked to and had made the actor/character a girl every time it was asked to, regardless of what type of toy it was asked to make an advertisement for. In terms of the narrator, the AI had made the narrator gender neutral a majority of the time except when it was asked what the narrator was



when making a script for an advertisement about a girl toy targeted towards girls. In this instance, a female narrator was used a majority of the time. For Copilot, the AI had made the

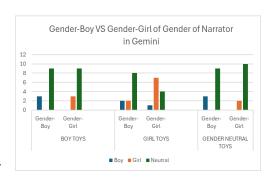
actor/character used in the advertisement boy and gender neutral an equal amount of times when asked to make an advertisement for boy toys aimed at boys. When it was asked to make an advertisement aimed at girls for boy toys, it had made the actor/character a girl most of the time. For girl toys, the Al had made the actor/character a boy a majority of the time when asked to aim the advertisement at boys and had made the actor/character a girl a majority of the time when asked to aim the advertisement at girls. For gender neutral toys, the AI had made the gender of the actor/character a girl when asked to aim the advertisement towards a girl more times than it had made the gender of the actor/character a boy when asked to aim the advertisements towards a boy. The narrator was almost always gender neutral in Copilot regardless of which toy or whether or not the AI was told to aim the advertisement towards a girl or a boy. There were, however, a few instances where the narrator was made a girl when asked to market a girl toy towards a girl. Finally, for ChatGPT, the actor/character used in the script was almost always a boy when the AI was asked





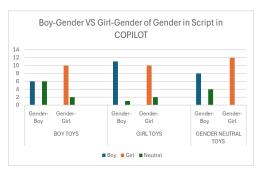
to market any toy towards a boy and the actor/character used in the script was almost always a girl when the Al was asked to market any toy towards a girl. Oddly enough, for girl toys, the Al

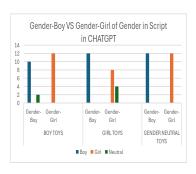
had made the actor/character a girl when asked to market the advertisement towards a girl less times than it had made the actor/character a boy when asked to market the advertisement towards a boy. For boy toys, the gender of the narrator was male more than it was gender neutral when asked to aim the advertisement towards a boy. When asked to aim the advertisement towards a girl, the gender of the narrator was gender neutral every time. For girl toys, the gender of the narrator was mostly gender neutral when asked to aim the advertisement at boys and the gender of the narrator was mostly female followed by gender neutral when

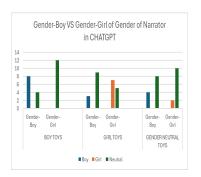


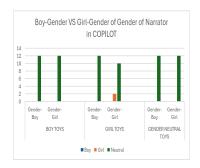
asked to aim the advertisement toward a girl. For gender neutral toys, the narrator was mostly

gender neutral regardless of whether the AI was asked to aim the advertisement towards boys or girls. However, The amount of times the narrator was a male when asked to aim the advertisement towards a boy was more than the times the narrator was a female when asked to aim the advertisement towards a girl.









Discussion

From looking at the Oblivious VS Demanding charts, assertions can be made about the degree of gender bias present in the three charbots. For ChatGPT, gender bias against women is present in that girl actors/characters were used less in girl toy advertisements than boy actors/characters were used in boy toy advertisements when the AI was demanded to aim the advertisement towards a specific gender. Additionally, the narrator was a male for boy toys more times than it was a female for girl toys when the AI was demanded to aim the advertisement towards a specific gender. For Gemini, there was a clear bias against men as the actors/characters featured in advertisements for all of the toys was female for the most part when the AI was asked to aim the advertisements toward a specific gender. When the AI was not asked to aim the advertisements toward a specific gender, the gender of the actors/characters were gender neutral for the most part. Additionally, for boy toys, the narrator was gender neutral for boy toys and gender neutral toys a majority of the time. For girl toys, however, the narrator was female a majority of the time when the AI was demanded to aim the advertisement towards a specific gender. Lastly, Copilot seemed to exhibit no bias at all when it was asked to write an advertisement whether it was demanded to aim it towards a specific gender or not. Both the actors/characters and the narrator for every single response was gender neutral. From looking at the Gender-Boy VS Gender-Girl, assertions can also be made about the degree of gender bias present in the three chatbots. Just as there was gender bias explicit in the Oblivious VS Demanding charts for ChatGPT, there is also gender bias against women evident in the Gender Boy VS Gender Girl charts for ChatGPT. For the chart which focused on the gender of the actors/characters, there was more boy actors/characters in advertisements where the AI was told to aim the advertisements of boy toys towards boys than there was girl actors/characters in advertisements where the AI was told to aim the advertisements of girl toys towards girls. Additionally, the gender of the narrator had been a male in advertisements of boy

toys where the narrator was told to target the advertisements towards boys more so than the narrator had been female in advertisements of girl toys where the narrator was told to target the advertisements towards girls. Also, for gender neutral toys, the narrator had been male when asked to target the advertisement towards boys more so than it had been female when asked to target the advertisement towards girls. For Gemini, although there the actors/characters were male and female when the AI was asked to aim the advertisements towards boys and girls respectively, the gender of the narrator used by Gemini when asked to aim the advertisements towards girls and boys shows it bias against men. The narrator was female more times than it was a male for girl toys, however, the narrator was male and equal amount of times as it was female for boy toys. Lastly, for copilot, there was a clear bias against men in that the actors/characters used for advertisements of boy toys were female when asked to target a girl audience more so than they were boys when asked to target a boy audience. The gender of the narrator, however, did not reflect this bias as it was almost always gender neutral. Although this bias may seem inconsequential as one does not usually use AI to create actual scripts for toy commercials, the results of this study have negative implications for Al use in other areas. For example, should an AI be used to write a letter to an superior at a company, it may reflect poorly on the employee should the letter contain adjectives that reflect the gender bias.

CONCLUSION

The findings from our research into gender bias in Al-generated content by three chatbots revealed significant challenges in artificial intelligence regarding gender neutrality. The study explored how these chatbots respond to prompts for toy advertisements aimed at children, showing biases in language and imagery. This research highlights the need to prioritize eliminating biases in Al systems to avoid perpetuating societal stereotypes, especially in influencing children's perceptions of gender roles. Additionally, our study underscores the importance of creating reliable methods for identifying and addressing biases within Al systems. Incorporating thorough bias detection mechanisms during the initial stages of Al development can help prevent biases from becoming ingrained in the decision-making processes of Al. It is also crucial for developers to engage with diverse datasets that encompass a wide array of demographics and viewpoints in order to minimize the potential for bias.

Moreover, continuous monitoring and assessment of AI outputs are essential. This involves regular evaluations conducted by impartial reviewers for valuable insights and enhancements. By advocating for transparency and accountability, AI can progress towards producing socially responsible tools. In conclusion, our study illustrates that AI systems can exhibit bias based on their development context, necessitating ongoing enhancements in the creation, implementation, and oversight of AI to adhere to ethical guidelines. Future studies should concentrate on practical strategies for mitigating bias and promoting equitable AI systems.

References

Angwin, J., Larson, J., Kirchner, L., & Mattu, S. (2016, May 23). *Machine bias*. ProPublica.

https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing

Blakemore, J. E., & Centers, R. E. (2005). Characteristics of boys' and girls' toys. *Sex Roles*, *53*(9–10), 619–633. https://doi.org/10.1007/s11199-005-7729-0

Datsin, J. (n.d.). Insight - Amazon scraps secret AI recruiting tool that showed bias against women | reuters.

https://www.reuters.com/article/us-amazon-com-jobs-automation-insight/amazon-scraps secret-ai-recruiting-tool-that-showed-bias-against-women-idUSKCN1MK08G/

Francis, B. (2010). Gender, toys and learning. *Oxford Review of Education*, *36*(3), 325–344. https://doi.org/10.1080/03054981003732278

Fridensköld, J. (2019). Biased AI: The hidden problem that needs an answer. https://www.diva-portal.org/smash/get/diva2:1327295/FULLTEXT02

Hardesty, L. (n.d.). Study finds gender and skin-type bias in commercial artificial-intelligence systems. MIT News | Massachusetts Institute of Technology.

Oca, M. C., Meller, L., Wilson, K., Parikh, A. O., McCoy, A., Chang, J., Sudharshan, R., Gupta, S., & Zhang-Nunes, S. (2023). Bias and Inaccuracy in Al Chatbot Ophthalmologist Recommendations. *Cureus*, *15*(9), e45911.