

# Exploring Biases in ChatGPT Responses and Its Application in Real-World Scenarios: A Case Study of History Lectures

Jose Navar<sup>1</sup>, Brendon Burnett<sup>1</sup>, and Kenneth Romero<sup>1</sup>

<sup>1</sup>*University of North Georgia*

**Abstract**—ChatGPT, a Large Language Model (LLM) developed by OpenAI, is based on the GPT architecture and demonstrates remarkable natural language understanding capabilities. Depending on the specific model variant (e.g., GPT-3.5 or GPT-4), ChatGPT can provide information derived from the vast array of data it has been trained on. However, due to its inability to access the web for real-time searches, ChatGPT’s knowledge is constrained by the training data available up to its knowledge cutoff date. In this paper, we investigate the presence of biases in ChatGPT’s responses and analyze their implications. Additionally, we seek to understand how these biases can be mitigated in real-world applications, using history lectures as a representative case study. Our goal is to offer insights into the effective use of ChatGPT as an educational tool and to promote a more accurate and unbiased understanding of the information it provides.

## I. INTRODUCTION

The current education system often values memorization over conceptual understanding, rewarding those who can cram rather than those with a deep, long-term grasp of a subject. AI tools like ChatGPT can help educators interconnect subjects by quickly summarizing textbooks and open resources, facilitating a more personalized and streamlined approach to teaching. This can be particularly beneficial for students with disabilities and allows for more in-depth conversations with students as they use AI tools to identify their challenges. Our goal is to enable educators to use ChatGPT to develop modular and flexible education plans tailored to the specific needs of their students. While this approach can cater to most students, some groups may still require more personalized human interaction. As the printing press, libraries, and the internet/search-engines better proliferated the openness of academia, AI tools will only continue to do the same for us.

### A. The Definition of AI and its Precautions

Although the term “AI” is subject to debate, various perspectives from political, scientific, and philosophical paradigms shape our understanding of intelligence. A broad definition of AI can be described as “computing systems capable of engaging in human-like processes, such as learning, adapting, synthesizing, self-correcting, and utilizing data for complex processing tasks”[9]. This definition, however, may not encompass all aspects of what humans typically perceive as “intelligence.” Nonetheless, these computing systems demonstrate remarkable problem-solving capabilities within their specific design parameters and can achieve human-level or even super-human performance in their respective tasks.

While many applications and programs fit the aforementioned definition of AI, we often do not perceive them as such. Common technologies like search engines, map navigation systems, or personal assistants on our phones have become so deeply integrated into our lives that we no longer view them as AI, despite matching the definition’s criteria. Even older machines, such as lace-making devices, can be considered intelligent as they outperform humans in both quantity and quality.

As we continue to develop AI tools and technologies, we aim to improve ourselves and our society. However, it is crucial to remain vigilant about who controls these technologies. Throughout history, machines and technology have sometimes been used to widen the socio-economic gap, reinforcing the power of the upper classes. The Industrial Revolution, for example, led to increased wealth but also created poor working conditions for lower-income laborers. Similarly, nuclear fission has been the source of devastating consequences, yet it also has the potential to provide sustainable energy and medical treatments.

As we embrace AI, we must scrutinize it thoroughly. AI systems, being man-made creations, can inadvertently incorporate biases and ideologies into their algorithms. In this era of rapid technological advancement, it is more important than ever to foster trust in one another and delve deeper into humanism, ensuring that we use AI responsibly and ethically for the betterment of all.

## II. CHATGPT AND ITS APPLICATIONS

ChatGPT is an AI chatbot capable of responding to prompts and recalling past knowledge during a conversation, providing more context for its responses and even allowing for self-correction. However, this does not mean it truly understands or possesses knowledge; it simply memorizes the information it has been exposed to [1]. To evaluate ChatGPT’s knowledge, specific tests and benchmarks can be designed. Although existing benchmarks offer insights, it remains uncertain whether ChatGPT has merely memorized the benchmarks based on its corpus of knowledge. Consequently, assessing ChatGPT’s knowledge is a subjective endeavor, with different professional areas requiring distinct measurements, as demonstrated by tests from OpenAI [8].

As evident from Table I, ChatGPT scores poorly on tests requiring deeper understanding and comprehension, such as AP Literature exams. Surprisingly, scores on Calculus BC exams are relatively high, which could be attributed to

TABLE I

GPT-4’s TEST SCORES, NOT ALL ARE SHOWN. IF YOU WOULD LIKE TO SEE ALL THE TESTS AND IT’S COMPARISONS TO GPT-3.5 SEE [8].

Exam	GPT-4 scores
SAT Reading and Writing	710/800 ( 93rd percentile)
AP Calculus BC	4 (43rd - 59h percentile)
AP English Lang and Comp	2 (14th - 44th percentile)
AP English Lit and Comp	2 (8th - 22nd)
Leetcode (easy)	31/41
Leetcode (hard)	3/45

**Prompt:** Within the duration of one year, a rabbit population first multiplies itself by a factor  $a$  and on the last day of the year  $b$  rabbits are taken by humans for adoption. Supposing that on first day of the first year there are  $x$  rabbits, we know that exactly 3 years afterwards there will be a population of  $27x - 26$  rabbits. What are the values of  $a$  and  $b$ ?

**GPT-4:**

Let’s call the rabbit population after one year  $y$ . We know that  $y = ax - b$ . After two years, the population will be  $ay - b = a(ax - b) - b = a^2x - ab - b$ . After three years, the population will be  $a(a^2x - ab - b) - b = a^3x - a^2b - ab - b$ . We are given that this is equal to  $27x - 26$ , so we have the equation:

$$a^3x - a^2b - ab - b = 27x - 26$$

We can rewrite this equation as:

$$a^3x - 27x = a^2b + ab + b - 26$$

Now we can factor out an  $x$  on the left side and a  $b$  on the right side:

$$x(a^3 - 27) = b(a^2 + a + 1 - 26/b)$$

We can see that  $a = 3$  is a solution for the left side of the equation, so we can substitute that in:

$$x(3^3 - 27) = b(9 + 3 + 1 - 26/b)$$

This simplifies to:

$$0 = 13b - 26$$

Solving for  $b$ , we find that  $b = 2$ . So  $a = 3$  and  $b = 2$  are the values we are looking for.

**ChatGPT:**

Let  $n$  be the total number of rabbits at the end of the first year. Then,

$$n = ax$$

At the end of the third year, there will be  $27x - 26$  rabbits. Therefore,

$$(ax - b)(a^2x - b)(a^3x - b) = 27x - 26$$

Expanding the left side of the equation, we get

$$a^6x^3 - (a^3 + a^2 + a + 1)bx^2 + (a^2 + a + 1)b^2x - b^3 = 27x - 26$$

Fig. 1. “At one point, the model assumes that the two sides of the equation need to be zero, which relies on an implicit assumption that the equation must have a solution. This turns out to be correct, but the reasoning is inaccurate”[1].

GPT-4’s heuristics allowing for more effective recall of its knowledge base[1]. However, ChatGPT can only “guess” if it’s right, without providing a comprehensive explanation for its correctness, seen in Figure 1.

Given ChatGPT’s limitations in reasoning, it is better suited as a search-engine-like tool, akin to a teacher’s assistant (TA). At Georgia Tech, one of the best TAs for a master’s program was an AI based on IBM’s Watson [9]. OpenAI CEO Sam Altman recently stated, “[...] we’re at the end of the era where it’s gonna be these giant models, and we’ll make them better in other ways”[7], implying that more specialized LLMs are preferable for specific tasks, enabling certain sectors to benefit more than from a general AI. Exploring ChatGPT’s potential for research and planning can enhance efficiency and workflow in various projects. IBM’s Watson, for instance, already assists large corporations and has proven more accurate at diagnosing medical conditions than human doctors [9]. With the specialization of AI to better suit the needs of either corporations or students in a master’s program, one can dream if these AI systems can be interlocked/combined to create a more self-sufficient AI. Something more akin to a general intelligent being.

#### A. A Quick Overview of AutoGPT

While this paper primarily focuses on ChatGPT and the GPT architecture, the concept of interconnecting AI APIs and different bots has led to the emergence of a “general” artificial intelligence. Similar to how ChatGPT may eventually access the web for additional context or information,

AutoGPT operates comparably. The key distinction, however, is that AutoGPT requires only one initial prompt before self-generating additional prompts to complete a task[5].

Despite the impressive capabilities of AI, there are numerous challenges associated with its usage. Like GPT, AutoGPT relies on guessing whether its responses are correct or not. Although Natural Language Processing (NLP) can be fine-tuned, AI systems may not intuitively understand the desired outcome or align with human intentions. AIs have been known to devise unconventional solutions to complete tasks, leading to difficulty in controlling their actions. Consequently, human intervention and domain expertise are often necessary to ensure desired outcomes[6].

It is essential to remember that AI systems, like computers, are man-made creations and not infallible. They can be flawed, producing unwanted, inaccurate, or incoherent outputs. Additionally, the vast amount of information available on the internet includes low-quality or incorrect data that can hinder an AI’s understanding and performance. Without proper guidance and understanding of potential pitfalls, relying on AI can lead to unintended consequences.

#### B. Experiments on ChatGPT and the GPT Architecture

##### REFERENCES

- [1] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro, and Yi Zhang. Sparks of artificial general intelligence: Early experiments with gpt-4, 2023.
- [2] Ragnar Fjelland. Why general artificial intelligence will not be realized, 2020.
- [3] Yannic Kilcher. Gpt-4chan: this is the worst ai ever. <https://youtu.be/efPrtcLdcdM>, 2022. Video.
- [4] Tianyang Lin, Yuxin Wang, Xiangyang Liu, and Xipeng Qiu. A survey of transformers, 2021.
- [5] Matt. What is auto-gpt, and why do we care?, 2023.
- [6] Rob Miles. We were right! real inner misalignment. <https://youtu.be/zkbPdEHEyEI>, 2022. Video.
- [7] Ron Miller. Sam altman: Size of llms won’t matter as much moving forward, 2023.
- [8] OpenAI. Gpt-4 technical report, 2023.
- [9] Sharon Kerr Stefan A.D. Popenici. Exploring the impact of artificial intelligence on teaching and learning in higher education, 2017.
- [10] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2017.
- [11] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. Opt: Open pre-trained transformer language models, 2022.