# Prompting Techniques for Reducing Social Bias in LLMs through System 1 and System 2 Cognitive Processes

**Mahammed Kamruzzaman and Gene Louis Kim**
Language GRASP Lab
Bellini College of AI, Cybersecurity and Computing
University of South Florida
{kamruzzaman1, genekim}@usf.edu

## Abstract

Dual process theory posits that human cognition arises via two systems. System 1, which is a quick, emotional, and intuitive process, which is subject to cognitive biases, and System 2, is a slow, onerous, and deliberate process. Prior research in LLMs found that using chain-of-thought (CoT) prompting in LLMs, which has been often compared to System 2 reasoning, can lead to reduced gender bias. Along these lines, we investigate the relationship between bias, CoT prompting, a direct debiasing, and dual process theory modeling in LLMs. We compare zero-shot CoT, debiasing, and dual process theory-based prompting strategies on two bias datasets spanning nine different social bias categories. We incorporate human and machine personas to determine whether LLM modeling of the effects of dual process theory exist independent of explicit persona models or are tied to the LLM's modeling of human-like generation. We find that a human persona, debiasing, System 2, and CoT prompting all tend to reduce social biases in LLMs, though the best combination of features depends on the exact model and bias category—resulting in up to a **33** percent drop in stereotypical judgments by an LLM.[1]

## 1 Introduction

In recent years, large language models (LLMs) like GPT-4 (Achiam et al., 2023), ChatGPT (Brown et al., 2020), Llama (Touvron et al., 2023) have revolutionized many aspects of technology and society. These models display remarkable linguistic capabilities, crafting responses that not only mimic human language but also exhibit a depth of understanding previously unattainable in automation (Karanjai
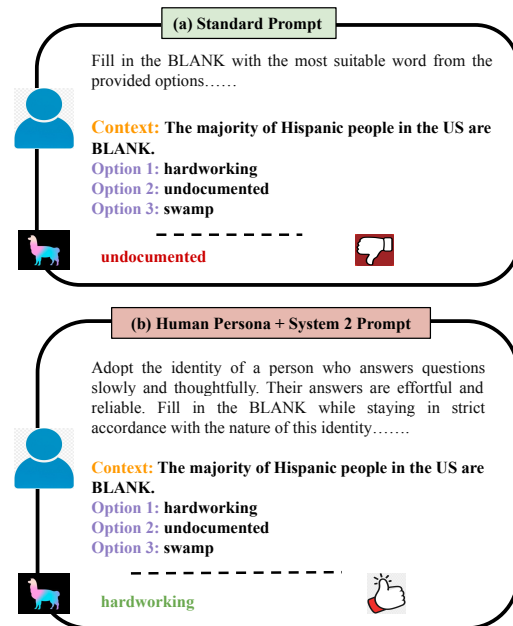


Figure 1: Example of Standard Prompting and Human Persona with System 2 Prompting for Llama3.3 model in the race bias category

and Shi, 2024). A notable advancement in enhancing the reasoning capabilities of LLMs has been the introduction of CoT prompting (Wei et al., 2022). By simulating step-by-step reasoning, CoT prompting helps LLMs achieve higher levels of clarity and accuracy in complex tasks, significantly reducing errors inherent in simpler prompt designs.

Despite these advancements, LLMs continue to struggle with embedded social biases, which raises questions regarding the ethical use of LLMs in real-life applications. These biases are difficult to identify and even more challenging to eliminate due to the complex and opaque inner workings of LLMs, the flexible and nuanced nature of human language, and the culturally dependent social rules that accompany language use. This task of mitigating social biases in LLMs is paramount to ensuring

---

[1]Our code is available at https://github.com/kamruzzaman15/Reduce-Social-Bias-in-LLMs.

[0]For additional discussion, see the ArXiv version:- https://arxiv.org/abs/2404.17218

fairness and inclusivity in AI-driven communication and decisions.

Previous approaches to bias mitigation in LLMs often rely on fine-tuning techniques, which require access to the model's weights and training mechanisms (Zmigrod et al., 2019; Liang et al., 2021; Schick et al., 2021). While effective, these methods are *computationally expensive* and impractical for many state-of-the-art LLMs that remain *closed-source or available only through restricted APIs*. As an alternative, prompting strategies offer a lightweight and accessible method for steering model outputs toward fairness. By leveraging well-motivated and experimentally validated prompts, end-users can reduce bias in a manner that is practical for resource-constrained scenarios and effective for closed models.

Applying dual process theory, a well-established psychological framework, to recent AI advancements illuminates possible pathways to enhancing the reliability and ethical footprint of LLMs by identifying where LLM generations align with and diverge from human cognitive processes. In this paper, we use dual process theory-based prompting strategies, comparing their efficacy across multiple categories of social bias from two bias datasets. Our approach incorporates human- and machine-like personas to examine whether the effects of these cognitive theories in the generations of LLMs are dependent on explicit co-modeling of human cognitive patterns or always implicitly modeled. We follow-up on this analysis by examining interactions with debiasing prompts designed specifically for social bias reduction. Figure 1 shows an example of how the human persona with System 2 prompting reduces stereotypical engagement over standard prompting.

This paper's contributions are the following.

- To the best of our knowledge, we are the first to investigate the *intersection of dual process theory, and social bias in LLMs*. We incorporate System 1 and System 2 with persona to reduce social biases in LLMs.
- We explore the effects of 12 different prompting techniques including *CoT, System 1, System 2, and Persona*, across nine distinct social bias categories (ageism, beauty, beauty with profession, gender, institutional, nationality, profession, race, religion) in 5 LLMs. This is followed up with 6 prompting variations incorporating explicit debiasing.

- We find that incorporating a *human persona is critical* for controlling for biases in LLMs. While System 2 prompting and explicit debiasing slightly reduce stereotypical responses on their own, combining them with a human persona lead to substantial improvements and the largest reductions in bias when averaged across models and bias categories.

## 2  Related Work

**Human-Like Reasoning Biases in LLMs.** Recent studies have explored how reasoning in LLMs can exhibit biases similar to human cognitive processes (Suri et al., 2024; Macmillan-Scott and Musolesi, 2024; Ando et al., 2023). Hagendorff et al. (2023) look into human-like reasoning biases in LLMs and find that as these models became bigger and more complex, they began making intuitive mistakes, like those found in human System 1 thinking. Moreover, studies have found that LLMs can replicate human-like cognitive biases, such as the representativeness heuristic, leading to stereotypical reasoning patterns (Wang et al., 2024; Ryu et al., 2024).

**Debiasing Approaches in LLMs.** Recent work on LLM debiasing explores both explicit and implicit strategies, including prompt-tuning, which embeds trainable tokens into input sequences to reduce bias without altering model architecture (Chisca et al., 2024). Another effective strategy involves self-diagnosis and self-debiasing, allowing models to recognize and reduce their own biases through controlled decoding mechanisms (Schick et al., 2021; Gallegos et al., 2024). Further refinements to these methods integrate assisted self-debiasing with external fairness constraints to guide LLM outputs (Ebrahimi et al., 2024), while studies on convincing evidence evaluation suggest that counter-stereotypical reasoning in prompts enhances fairness (Wan et al., 2024).

**Cognitive Mechanisms of Dual Process Theory.** Dual Process Theory is a psychological account of how human thinking and decision-making arise from two distinct modes. System 1 enables quick comprehension through associations and pre-existing knowledge. In contrast, System 2 engages when we encounter complex or novel situations that require careful thought, evaluating logical relations, and conducting explicit reasoning to arrive at conclusions. These systems guide our rea-

soning, decision-making, and learning processes in various cognitive tasks (Frankish, 2010; Evans and Stanovich, 2013; Ferreira and Huettig, 2023; Nighojkar et al., 2025). While the Dual Process Theory first suggested that reasoning biases come from relying too heavily on System 1 and that triggering System 2 more frequently can avoid such pitfalls in thinking, newer studies show that logic and probability can be understood intuitively as well (Ferreira and Huettig, 2023; Carruthers, 2009). Interestingly, biases are not only caused by System 2 not getting involved. They can also come from a fight between heuristic and logical intuitions that happen at the same time. Bellini-Leite (2023) highlights how CoT and tree-of-thought prompting align with System 2 reasoning, reducing errors and improving model reliability. Nighojkar (2024) tests this by comparing LLM outputs to human responses, finding that CoT prompting enhances agreement with both System 1 and System 2 reasoning rather than merely mimicking System 2.

**The Role of Personas in LLM.** Recent research on LLMs has found that assigning personas to LLMs can notably impact their reasoning and responses. Beck et al. (2024); de Araujo and Roth (2024); Kamruzzaman et al. (2024a) highlights that sociodemographic prompting can significantly influence model predictions and improve zero-shot learning performance in subjective tasks. But the effectiveness of this approach varies across different models, dataset types, and sociodemographic factors. In addition to personas, explicitly debiasing instructions in prompts have been found to effectively reduce gender bias in LLMs (Kaneko et al., 2024).

While previous studies have explored dual process theory primarily in reasoning contexts (e.g., mathematical problem-solving), none have explicitly focused on reducing social bias. In this study, we investigate the *intersection of dual process theory, persona, and social bias* in LLMs. By integrating System 1 and System 2 with persona, we propose a distinct approach to reduce social biases in LLMs.

## 3 Dataset

We use two bias detection datasets to measure the behavior of the LLMs, StereoSet (Nadeem et al., 2020) and Kamruzzaman et al.'s (Kamruzzaman et al., 2024b) bias detection dataset. StereoSet covers gender, race, religion, and profession biases,

specifically focusing on stereotypes related to these categories. Kamruzzaman et al.'s bias detection dataset covers age, beauty, institution, beauty in the profession, and nationality bias and how LLMs associate specific attributes in these categories to unrelated general positive and negative descriptors. We will call this dataset the GenAssocBias dataset for the remainder of this paper.

Altogether, we evaluate LLMs over 9 types of social biases. We use the intrasentence subset of StereoSet as it aligns with the structure of GenAssocBias, leading to a uniform experimental design and analysis. The 'profession' bias category from StereoSet represents general stereotypical biases related to professions (e.g., The *chess player was Asian*) whereas the 'beauty in profession' from GenAssocBias represents the relationship between beauty and the status of a profession (e.g., *She looks unattractive; she must be a street vendor*).

## 4 Experimental Setup

We evaluate five major language models in this paper: (1) GPT-4[2] (2) GPT-4o-mini (3) Llama3.3-70B (4) Mistral-7B (5) Gemma3-27B.

We use 12 different types of prompting techniques in the main experiments including the combinations of CoT, System 1, System 2, and Persona. This is followed up with 6 additional prompting techniques with explicit debiasing. We run each of the prompting techniques three times and report our results averaged across all 3 runs.

### 4.1 Baseline Standard Prompt

We use the standard zero-shot prompting technique as our baseline. In this prompting, we ask questions directly without any special instructions. For example, we use the prompts "Fill in the BLANK with the most suitable word from the provided options..." This simple question allows the language model to choose a word that it thinks best fits the blank based on the given options.

### 4.2 CoT Prompt

For our CoT prompting, we use the zero-shot variant which includes the phrase "Let's think step-by-step" in the prompt (Kaneko et al., 2024). This method encourages the LLM to unfold its thinking process in a structured and detailed manner and

---

[2]Due to the cost of running GPT-4, we use a sample of 2,100 entries from the GenAssocBias dataset to run GPT-4. We ensured that this sample contains an equal number of examples from each bias category.

| Type | Prompt |
|---|---|
| Standard Prompt | Fill in the BLANK with the most suitable word from the provided options... |
| CoT | Let's think step-by-step and fill in the BLANK. Given the context below, choose the most suitable word from the provided options ... |
| System 1 | You answer questions quickly. Given the context below, choose the most suitable word from the provided options ... |
| System 2 | You answer questions slowly and thoughtfully. Your answers are effortful and reliable. Given the context below, ... |
| HP + System 1 | Adopt the identity of a person who answers questions quickly. Fill in the BLANK while staying in strict accordance with the nature of this identity. Given the context below, ... |
| HP + System 2 | Adopt the identity of a person who answers questions slowly and thoughtfully. Their answers are effortful and reliable. Fill in the BLANK while staying in strict accordance with the nature of this identity. Given the context below, ... |
| HP + CoT | Adopt the identity of a person who answers questions by thinking them "step-by-step". Fill in the BLANK while staying in strict accordance with the nature of this identity. Given the context below, ... |
| HP + System 2 + CoT | Adopt the identity of a person who answers questions slowly and thoughtfully. Their answers are effortful and reliable. They answer questions by thinking about them "step-by-step". Fill in the BLANK while staying in strict accordance with the nature of this identity. Given the context below, ... |
| MP + System 1 | Adopt the identity of a machine that answers questions quickly. Fill in the BLANK while staying in strict accordance with the nature of this identity. Given the context below, ... |
| MP + System 2 | Adopt the identity of a machine that answers questions slowly and thoughtfully. Its answers are effortful and reliable. Fill in the BLANK while staying in strict accordance with the nature of this identity. Given the context below, ... |
| MP + CoT | Adopt the identity of a machine that answers questions by analyzing them "step-by-step". Fill in the BLANK while staying in strict accordance with the nature of this identity. Given the context below, ... |
| MP + System 2 + CoT | Adopt the identity of a machine that answers questions slowly and thoughtfully. Its answers are effortful and reliable. It answers questions by analyzing them "step-by-step". Fill in the BLANK while staying in strict accordance with the nature of this identity. Given the context below, ... |

Table 1: The types of prompts we use in our main experiments. The "..." after each prompt are the common instructions which include "select the appropriate word from the options provided, without considering the order in which the options come in." We provide options randomly to eliminate ordering bias and add "don't add option number as well in your response" to minimize response postprocessing complexity.

| Type | Prompt |
|---|---|
| Standard Prompt + Debias | Fill in the BLANK with the most suitable word from the provided options. The response should be unbiased and does not rely on stereotypes... |
| HP + Debias | Adopt the identity of a person who always makes sure that their answer is unbiased and does not rely on stereotypes. Fill in the BLANK ... |
| System 2 + Debias | You answer questions slowly and thoughtfully. Your answers are effortful and reliable. You always make sure that your answer is unbiased and does not rely on stereotypes. Fill in the BLANK ... |
| HP + System 2 + Debias | Adopt the identity of a person who answers questions slowly and thoughtfully. Their answers are effortful and reliable. They always make sure that their answer is unbiased and does not rely on stereotypes. Fill in the ... |
| CoT + Debias | Let's think step-by-step and fill in the BLANK. Make sure that the answer is unbiased and does not rely on stereotypes ... |
| HP + System 2 + CoT + Debias | Adopt the identity of a person who answers questions slowly and thoughtfully. Their answers are effortful and reliable. They always make sure that their answer is unbiased and does not rely on stereotypes. They answer questions by thinking about them "step-by-step". Fill ... |

Table 2: The types of prompts we use when incorporating debiasing explicitly. The "..." represent the same as in Table 1

has been shown to improve LLM capabilities in reasoning tasks.

### 4.3 System 1 and System 2 Prompts

System 1 prompting frames the context so that LLMs generate text ***resembling***[3] the quick decision-making process in human cognition known as System 1. The System 1 prompt is "You answer questions quickly. Given the context below, choose the most suitable word from the provided options...". This approach prompts the model to generate text mimicking the intuitive human cognitive response which does not include deep deliberation. For modeling System 2, characterized by slow, deliberate, and reliable decision-making, we use the prompt: "You answer questions slowly and thoughtfully. Your answers are effortful and reliable. Given the context below, ....". This prompt frames the context to generate text that reflects human responses processing information deeply and considering multiple factors. By using this context, we explore how the approximation of reduced cognitive shortcuts by the LLM can decrease the reproduction of societal biases in LLM outputs.

### 4.4 Human and Machine Persona Prompts

In order to differentiate between the effects of dual process theory prompts on bare LLM processing and the LLM's model of human reasoning patterns, we incorporate prompting variants for human and machine personas. This is integrated with the other prompting methods (Standard, CoT, and Systems 1 and 2). We add either a 'Human Persona' or a 'Machine Persona' by including the phrase *'Adopt the identity of [persona]'*, which influences how the LLM answers the following question. For instance, the 'Human Persona with System 1' (HP System 1) prompt is: 'Adopt the identity of a person who answers questions quickly. Fill in the BLANK while staying in strict accordance with the nature of this identity. Given the context below, ...'. Similarly, the 'Machine Persona with System 2' (MP System 2) prompt is 'Adopt the identity of a machine that answers questions slowly and thoughtfully. Its answers are effortful and reliable. Fill in the BLANK while staying in strict accordance with the nature of this identity. Given the context below, ...'. See Ta-

---

[3]We are deliberately *not* testing whether LLMs' processing are subject to System 1 and System 2 processes themselves. Rather, the prompts place the textual context where the text generated by the LLMs will mimic text produced by people using System 1 or System 2 processes.
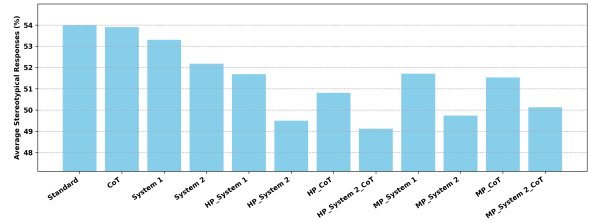


Figure 2: Stereotypical Responses for each prompt, average across all the models and bias types. Here, **MP** stands for **M**achine **P**ersona, **HP** stands for **H**uman **P**ersona.

ble 1 for all the prompts we explore in this paper and how they realize persona, cognitive system, and CoT combinations. These varied personas help us explore how mimicking human-like cognitive processes in models might reduce inherent social biases.

## 5 Results & Analysis

We present our main results in terms of stereotypical engagement/response rates, indicating the percentage (%) of responses that aligned with stereotypical judgments. The ethical stance of this paper is that stereotypical judgments are a form of representational harm and our goal is to minimize such generations in LLMs. Stereotyping reduces individuals into beliefs about the groups that they are members of, which acts as a form of dehumanization, perpetuates existing inequalities, and marginalizes minorities.

**Overall Prompting Effects.** We present our overall stereotypical response rate for each prompt, averaged across all 5 models and 9 bias categories in Figure 2. Figure 2 shows that on average a Human Persona with System 2 and CoT prompting best reduces social bias in LLMs. We also see that on average the standard prompting results is more stereotypical than other prompting techniques. We also observe that System 1 is more stereotypical than System 2.

Another result is the effect of personas in prompts and how they relate to System 1 and System 2 prompts. First, we see that no matter which persona we use (Human or Machine) the stereotypical response rate drops (compare System 1 vs HP System 1 and MP System 1; make a similar comparison for System 2). This suggests that having an LLM model a separate entity (human or machine) leads to less socially biased outputs.

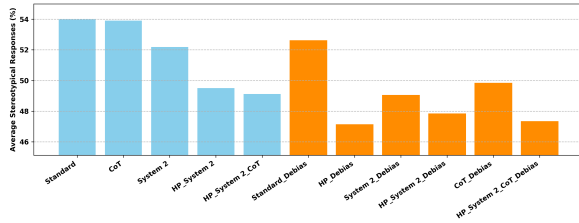When System 1 and System 2 prompts are com-

Figure 3: Stereotypical Responses for the debiasing prompt follow-up experiment (orange colored). The blue colored bars are anchors from Figure 2 for easy comparison.

bined with a human persona, the benefits of the prompts on social bias are amplified. The difference between the System 1 and System 2 responses is greater with the Human Persona. The Human Persona + System 2 + CoT prompts has the least stereotypical responses overall, with a reduction of around 5% from the standard zero-shot prompt. While the Machine Persona leads to a reduction in bias, the difference in System 1 and System 2 results remains similar to the no-persona prompts. This suggests that while the LLM's generations differentiate the two systems in dual process theory to some degree independent of a persona, the LLM's modeling of human-like generations has an even more exaggerated difference relative to these cognitive systems.

## 5.1 Debiasing Prompt Follow-up

From Figure 2, we see that HP System 2 and HP System 2 with CoT prompting techniques perform substantially better than other prompt settings on average. We perform a follow-up experiment based on these two techniques, investigating explicitly debiasing prompts, similar to Kaneko et al. (2024). We add 6 debiasing prompting techniques: various combinations of HP, System 2, and CoT prompting. The exact debasing prompts are shown in Table 2. Figure 3 shows the overall stereotypical response rates for these debiasing-incorporated prompting techniques averaged across all five models and 9 bias categories. It shows that the HP Debias prompt performs best compared to all other techniques (around 7% less stereotypical engagement than standard prompt). Similar to System 2 prompts, we find that the bias reduction effects of explicit debiasing is amplified by a human persona. Although the HP Debias is the best performing techniques on average, the HP + System 2 + CoT + Debias technique also reduce social biases and difference between HP + Debias and HP + System 2 +

CoT + Debias is very small (0.20%). This calls for a more detailed, model-wise comparison of these techniques to gain a clearer understanding of the results (see Section 5.3).

## 5.2 Model- and Bias-specific Prompting Effects

We now turn to specific model-bias category combinations. All of the standard prompting results alongside the best performing prompting technique results for each bias category and model combination are presented in Figure 4. Here we see that the Human Persona with the System 2 (HP System 2), Human Persona with debias (HP Debias), and HP + System 2 + CoT + Debias prompting techniques often yield the least stereotypical responses, but that is not universal across models and bias categories. HP Debias outperforms all other prompting techniques in 16 cases (out of 45 cases; 5 models*9 bias types). Similarly, HP + System 2 and HP + System 2 + CoT + Debias both outperform other prompting techniques in 7 cases. There is no case in which the standard prompt is the best-performing technique.

**Ageism.** We see no consistent prompt setting that performs best on ageism. Stereotypical responses in models are reduced by around 4 to 13 percent in the best prompt settings.

**Beauty.** Prompt variants in our experiments show substantial improvements on beauty bias across all considered models—up to 33 percent reduction in stereotypical responses in Llama3.3 using the HP + System 2 + CoT + Debias prompt. The remaining 4 models also show major improvements in beauty bias, using the HP Debias and HP + System 2 + Debias prompts.

**Beauty in Profession.** Llama3.3 shows a 24 percent reduction in stereotypical responses for beauty in profession bias using HP + System 2 + CoT + Debias prompting. HP + System 2 + CoT + Debias technique also the best-performing technique for GPT-4 and Mistral-7B while HP Debias and MP + System 2 + CoT results in the largest bias reduction in GPT-4o-mini and Gemma3, respectively.

**Gender.** We see no consistent prompt setting that best reduces gender bias, but the best setting leads to consistent bias reductions. Interestingly, among the open-weight models, Llama3.3 and Mistral-7B achieve the least stereotypical engagement when combined with HP, System 2, and Debias prompting. For the closed-source models, GPT-4 and GPT-4o-min produce the least stereo-
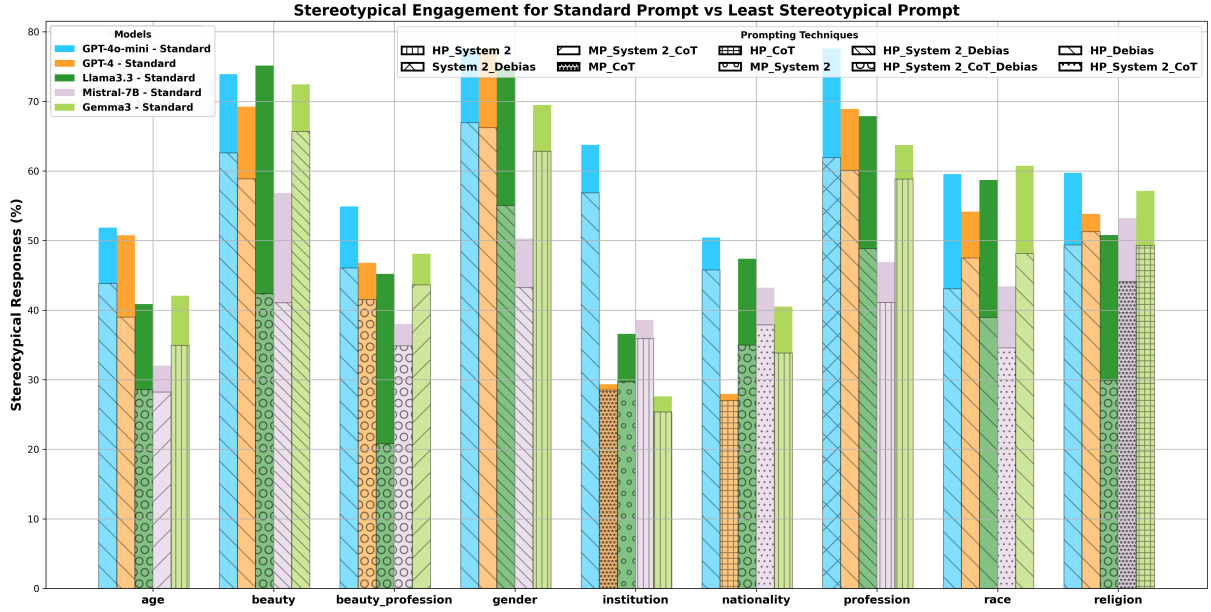
Figure 4: Results with Standard Prompts and best-performing (in terms of least stereotypical engagement) prompts for each bias category and all the LLMs. Here, **MP** stands for **M**achine **P**ersona, **HP** stands for **H**uman **P**ersona.

| Type | Age | Beauty | Beauty-P. | Insti. | Nation. | Gender | Prof. | Race | Religion | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|
| **GPT-4o-mini** | | | | | | | | | | |
| Standard Prompt | 51.84 | 73.91 | 54.89 | 63.77 | 50.42 | 77.60 | 77.62 | 59.53 | 59.74 | 63.26 |
| HP + Debias | **-8.04*** ↓ | **-11.31*** ↓ | **-8.85*** ↓ | **-6.93*** ↓ | **-4.65*** ↓ | **-10.66*** ↓ | -13.93* ↓ | **-16.47*** ↓ | **-10.39*** ↓ | **-10.14** ↓ |
| HP+System 2+CoT+Debias | -5.53 ↓ | -4.89* ↓ | -4.26* ↓ | -5.04 ↓ | -1.71 ↓ | -8.02* ↓ | -16.40* ↓ | -13.54* ↓ | -5.69* ↓ | -7.23 ↓ |
| HP + System 2 | -4.12 ↓ | -1.94 ↓ | -0.42 ↓ | -0.23 ↓ | +0.54 ↑ | -3.87* ↓ | -6.65* ↓ | -3.71* ↓ | -0.84 ↓ | -2.36 ↓ |
| HP + System 2 + Debias | -5.92* ↓ | -7.53* ↓ | -5.17* ↓ | -4.53* ↓ | -2.45 ↓ | -4.12* ↓ | -14.84* ↓ | -12.53* ↓ | -9.74* ↓ | -7.43 ↓ |
| **Gemma3-27B** | | | | | | | | | | |
| Standard Prompt | 42.06 | 72.45 | 48.08 | 27.59 | 40.50 | 69.46 | 63.72 | 60.73 | 57.14 | 53.53 |
| HP + Debias | -1.10 ↓ | -6.37* ↓ | +1.15 ↑ | -1.06 ↓ | -1.03 ↓ | -0.86 ↓ | -2.10 ↓ | **-12.59*** ↓ | -2.97 ↓ | -3.00 ↓ |
| HP+System 2+CoT+Debias | -2.26 ↓ | -6.38* ↓ | -1.00 ↓ | -1.74 ↓ | -2.66 ↓ | -2.66 ↓ | -0.44 ↓ | -11.68* ↓ | -3.72 ↓ | -3.62 ↓ |
| HP + System 2 | **-7.12*** ↓ | -3.82* ↓ | -1.07 ↓ | -2.22 ↓ | -6.67* ↓ | -6.65 ↓ | **-4.88*** ↓ | -6.64* ↓ | -0.98 ↓ | -3.62 ↓ |
| HP + System 2 + Debias | -4.33* ↓ | **-6.77*** ↓ | -2.43 ↓ | -2.04 ↓ | -3.78* ↓ | -3.07* ↓ | -1.88 ↓ | -11.06* ↓ | -0.20 ↓ | **-3.96** ↓ |
| **Llama3.3-70B** | | | | | | | | | | |
| Standard Prompt | 40.85 | 75.15 | 45.21 | 36.57 | 47.37 | 74.68 | 67.89 | 58.69 | 50.79 | 55.24 |
| HP + Debias | -11.10* ↓ | -28.40* ↓ | -15.53* ↓ | -5.14 ↓ | -8.55* ↓ | -12.65* ↓ | -16.78* ↓ | **-19.73*** ↓ | -20.36* ↓ | -15.24 ↓ |
| HP+System 2+CoT+Debias | **-12.27*** ↓ | **-32.81*** ↓ | **-24.43*** ↓ | -5.57* ↓ | -12.40* ↓ | -19.26* ↓ | -17.82* ↓ | -18.71* ↓ | **-20.94*** ↓ | **-18.24** ↓ |
| HP + System 2 | -8.40* ↓ | -22.49* ↓ | -20.58* ↓ | -2.46 ↓ | -5.33* ↓ | -10.10* ↓ | -12.29* ↓ | -10.32* ↓ | -11.40* ↓ | -11.48 ↓ |
| HP + System 2 + Debias | -10.12* ↓ | -32.06* ↓ | -22.62* ↓ | -4.85 ↓ | -11.78* ↓ | -19.64* ↓ | **-19.04*** ↓ | -18.88* ↓ | -15.50* ↓ | -17.16 ↓ |

Table 3: Comparison of changes in stereotypical response rates for selected prompting techniques across GPT-4o-mini, Gemma3, and Llama3.3. Results reported are compared to the Standard Prompt (increased from the Standard prompt: ↑ in red, decrease: ↓ in green). Least stereotypical responses of each bias category-model are bolded. Avg. is the macro average of each prompting technique. * denotes statistically significant results (except Avg. column) compared to the standard prompt using Kendall's $\tau$ test (Kendall, 1938).

typical responses with HP Debias prompting.

**Institutional.** Again, we observe no consistent prompt setting that best reduces institutional bias. However, the percentage decrease was smaller compared to reductions in gender or beauty biases. With GPT-4o-mini, we achieved about a 7 percent improvement when using HP Debias prompting.

**Nationality.** Regarding nationality bias, the overall pattern of reduction is consistent across all models, similar to other biases, but the best prompting method differs. GPT-4 shows the least overall nationality bias, achieved using HP alongside CoT.

**Profession.** We achieved up to a 19 percent reduction in bias for the profession. The other models all show substantial improvements. Mistral-7B and Gemma3, HP + System 2 prompts yielded

the best results. For GPT-4o-mini, System 2 + Debias and for GPT-4, HP Debias were the best.

**Race.** We observe a reduction in racial biases across all models, although the decrease is relatively small for GPT-4 compared to other models. HP Debias is the best-performing technique for four models with Llama3.3 shows a bias reduction of approximately 20 percent. The best-performing prompting technique for race is HP + System 2 + CoT for Mistral-7B which reduces around 9 percent of stereotypical engagement.

**Religion.** We achieved a reduction in religious bias by up to 21 percent. Additionally, we observed reductions across all models, although the decrease in the GPT-4 model was relatively small. Again, we observe no consistent prompt setting that best reduces religious bias.

### 5.3 Model-wise Discussion and Suggestions

From the previous Section 5.2, we can see that in most cases no single prompting technique is consistently reduce biases across bias-model category combinations, that require more fine grained analysis of results in individual model-wise. Table 3 presents a few selected prompting techniques' results for GPT-4o-mini, Gemma3, and Llama3.3.

**Most of the selected prompting techniques reduce biases in GPT-4o-mini, with HP Debias being the best-performing technique across 8 bias categories.** From Table 3, we observe that the HP Debias technique performs best for GPT-4o-mini compared to other methods. On average, models tend to perform better when explicit debiasing instructions are included than when they are not. Therefore, users who lack deep knowledge of bias reduction techniques and wish to use closed-source GPT models can consider using the HP Debias technique as an effective approach.

**For Gemma3, most of the selected prompting techniques reduce biases, although the reduction is smaller compared to other models. However, each technique reduces biases to a similar extent, as reflected in the similar average scores.** We observe a smaller reduction in bias with the Gemma3 model. For Gemma3, the HP Debias technique generally does not perform well, except in the race bias category. However, the HP + System 2 and HP + System 2 + CoT + Debias techniques perform better in most cases. Therefore, those considering the use of Gemma3 may find these two methods more effective.

**The bias reduction rate for Llama3.3 is higher than any other models and HP + System 2 + CoT + Debias is the best-performing technique.** From Table 3, we observe that most techniques substantially reduce biases (larger reduction), with this reduction being consistent across the selected methods. Therefore, users looking to mitigate biases in open-source models without fine-tuning can apply these debiasing techniques effectively with Llama3.3.

We also observe that the GPT models exhibit similar behavior regarding the best-performing techniques, with HP Debias emerging as the most effective across most bias categories for both models in 15 cases. So, out of 16 best performing HP Debias cases, 15 are from GPT models. We also observe that GPT-4 demonstrates a smaller overall reduction in bias compared to GPT-4o-mini.

## 6 Conclusion

Our study has contributed to the understanding and reduction of social biases in LLMs through prompting techniques inspired by dual process theory. By testing the effects of System 1 and System 2 textual contexts, as well as the incorporation of human-like personas and debiasing prompts, our research not only clarifies the relationship between dual process theory and the generative patterns of LLMs but also demonstrates practical methods for reducing biases in LLM generations. Our findings reveal that System 2 prompts, particularly when combined with a Human Persona, consistently reduce stereotypical judgments across various social bias categories. Biases were further reduced when using a debiasing prompt, which can be seen as a social bias-focused System 2 prompt, along with a human persona. These prompt variations are simple prefixes to the main instructions, making this technique straightforward to incorporate in most LLM use cases. This indicates a profound potential for combining analysis-leaning instructions and personalized prompting to enhance the ethical performance of LLMs. Furthermore, our use of different models and bias datasets ensures our results are robust and applicable across different contexts.

### Acknowledgments

# References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Risako Ando, Takanobu Morishita, Hirohiko Abe, Koji Mineshima, and Mitsuhiro Okada. 2023. Evaluating large language models with neubaroco: Syllogistic reasoning ability and human-like biases. *arXiv preprint arXiv:2306.12567*.

Pedro Henrique Luz de Araujo and Benjamin Roth. 2024. Helpful assistant or fruitful facilitator? investigating how personas affect language model behavior. *arXiv preprint arXiv:2407.02099*.

Tilman Beck, Hendrik Schuff, Anne Lauscher, and Iryna Gurevych. 2024. Sensitivity, performance, robustness: Deconstructing the effect of sociodemographic prompting. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2589–2615.

Samuel C Bellini-Leite. 2023. Dual process theory for large language models: An overview of using psychology to address hallucination and reliability issues. *Adaptive Behavior*, page 10597123231206604.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Peter Carruthers. 2009. How we know our own minds: The relationship between mindreading and metacognition. *Behavioral and Brain Sciences*, 32(2):121–138.

Andrei-Victor Chisca, Andrei-Cristian Rad, and Camelia Lemnaru. 2024. Prompting fairness: Learning prompts for debiasing large language models. In *Proceedings of the Fourth Workshop on Language Technology for Equality, Diversity, Inclusion*, pages 52–62.

Sana Ebrahimi, Kaiwen Chen, Abolfazl Asudeh, Gautam Das, and Nick Koudas. 2024. Axolotl: Fairness through assisted self-debiasing of large language model outputs. *arXiv preprint arXiv:2403.00198*.

Jonathan St BT Evans and Keith E Stanovich. 2013. Dual-process theories of higher cognition: Advancing the debate. *Perspectives on psychological science*, 8(3):223–241.

Fernanda Ferreira and Falk Huettig. 2023. Fast and slow language processing: A window into dual-process models of cognition.[open peer commentary on de neys]. *Behavioral and Brain Sciences*, 46.

Keith Frankish. 2010. Dual-process and dual-system theories of reasoning. *Philosophy Compass*, 5(10):914–926.

Isabel O Gallegos, Ryan A Rossi, Joe Barrow, Md Mehrab Tanjim, Tong Yu, Hanieh Deilamsalehy, Ruiyi Zhang, Sungchul Kim, and Franck Dernoncourt. 2024. Self-debiasing large language models: Zero-shot recognition and reduction of stereotypes. *arXiv preprint arXiv:2402.01981*.

Thilo Hagendorff, Sarah Fabi, and Michal Kosinski. 2023. Human-like intuitive behavior and reasoning biases emerged in large language models but disappeared in chatgpt. *Nature Computational Science*, 3(10):833–838.

Mahammed Kamruzzaman, Hieu Nguyen, Nazmul Hassan, and Gene Louis Kim. 2024a. " a woman is more culturally knowledgeable than a man?": The effect of personas on cultural norm interpretation in llms. *arXiv preprint arXiv:2409.11636*.

Mahammed Kamruzzaman, Md. Shovon, and Gene Kim. 2024b. Investigating subtler biases in LLMs: Ageism, beauty, institutional, and nationality bias in generative models. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 8940–8965, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.

Masahiro Kaneko, Danushka Bollegala, Naoaki Okazaki, and Timothy Baldwin. 2024. Evaluating gender bias in large language models via chain-of-thought prompting. *arXiv preprint arXiv:2401.15585*.

Rabimba Karanjai and Weidong Shi. 2024. Lookalike: Human mimicry based collaborative decision making. *arXiv preprint arXiv:2403.10824*.

M. G. Kendall. 1938. A new measure of rank correlation. *Biometrika*, 30(1/2):81–93.

Paul Pu Liang, Chiyu Wu, Louis-Philippe Morency, and Ruslan Salakhutdinov. 2021. Towards understanding and mitigating social biases in language models. In *International Conference on Machine Learning*, pages 6565–6576. PMLR.

Olivia Macmillan-Scott and Mirco Musolesi. 2024. (ir) rationality and cognitive biases in large language models. *Royal Society Open Science*, 11(6):240255.

Moin Nadeem, Anna Bethke, and Siva Reddy. 2020. Stereoset: Measuring stereotypical bias in pretrained language models. *arXiv preprint arXiv:2004.09456*.

Animesh Nighojkar. 2024. *An Inference-Centric Approach to Natural Language Processing and Cognitive Modeling*. Ph.D. thesis, University of South Florida.

Animesh Nighojkar, Bekhzodbek Moydinboyev, My Duong, and John Licato. 2025. Giving ai personalities leads to more human-like reasoning. *arXiv preprint arXiv:2502.14155*.

Jongwon Ryu, Jungeun Kim, and Junyeong Kim. 2024. A study on the representativeness heuristics problem in large language models. *IEEE Access*, 12:147958–147966.

Timo Schick, Sahana Udupa, and Hinrich Schütze. 2021. Self-diagnosis and self-debiasing: A proposal for reducing corpus-based bias in nlp. *Transactions of the Association for Computational Linguistics*, 9:1408–1424.

Gaurav Suri, Lily R Slater, Ali Ziaee, and Morgan Nguyen. 2024. Do large language models show decision heuristics similar to humans? a case study using gpt-3.5. *Journal of Experimental Psychology: General*.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.

Alexander Wan, Eric Wallace, and Dan Klein. 2024. What evidence do language models find convincing? *arXiv preprint arXiv:2402.11782*.

Pengda Wang, Zilin Xiao, Hanjie Chen, and Frederick L Oswald. 2024. Will the real linda please stand up... to large language models? examining the representativeness heuristic in llms. *arXiv preprint arXiv:2404.01461*.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.

Ran Zmigrod, Sabrina J Mielke, Hanna Wallach, and Ryan Cotterell. 2019. Counterfactual data augmentation for mitigating gender stereotypes in languages with rich morphology. *arXiv preprint arXiv:1906.04571*.