



PDF Download  
3696410.3714537.pdf  
20 December 2025  
Total Citations: 0  
Total Downloads: 3727

Latest updates: <https://dl.acm.org/doi/10.1145/3696410.3714537>

RESEARCH-ARTICLE

## Social Bots Meet Large Language Model: Political Bias and Social Learning Inspired Mitigation Strategies

**JINGHUA PIAO**, Beijing National Research Center for Information Science and Technology, Beijing, China

**ZHIHONG LU**, Beijing National Research Center for Information Science and Technology, Beijing, China

**CHEN GAO**, Beijing National Research Center for Information Science and Technology, Beijing, China

**YONG LI**, Beijing National Research Center for Information Science and Technology, Beijing, China

**Open Access Support** provided by:

**Beijing National Research Center for Information Science and Technology**

**Published:** 28 April 2025

[Citation in BibTeX format](#)

WWW '25: The ACM Web Conference 2025

April 28 - May 2, 2025  
Sydney NSW, Australia

**Conference Sponsors:**  
SIGWEB

# Social Bots Meet Large Language Model: Political Bias and Social Learning Inspired Mitigation Strategies

Jinghua Piao\*

Department of Electronic Engineering,  
BNRist, Tsinghua University  
Beijing, China  
pjh22@mails.tsinghua.edu.cn

Chen Gao

BNRist, Tsinghua University  
Beijing, China  
chgao96@tsinghua.edu.cn

Zhihong Lu\*

Department of Electronic Engineering,  
BNRist, Tsinghua University  
Beijing, China  
luzh22@mails.tsinghua.edu.cn

Yong Li†

Department of Electronic Engineering,  
BNRist, Tsinghua University  
Beijing, China  
liyong07@tsinghua.edu.cn

## Abstract

Recent advances in the large language models (LLM) have empowered traditional bots to gain human-level intelligence and exhibit human-like social behaviors, giving rise to a new form of LLM-driven social agents. However, the inherent limitations in LLMs could potentially result in politically biased behaviors of these agents, posing unexpected risks to human society. While great efforts have been made to examine political bias and related concerns in traditional bots and LLMs, little is known about the existence, unique characteristics, underlying origins, and potential mitigation strategies of this bias in LLM-driven social agents. To address this gap, we systematically assess political bias in LLM-driven social agents, by examining how it emerges as these agents self-reflect, communicate, and understand others during social interactions. Through designing and implementing social experiments, we discover that this bias consistently manifests in the social behaviors of agents driven by diverse LLMs, across nine key political topics. Inspired by the social learning theory, we propose to mitigate political bias by guiding these agents to emulate how humans learn to behave. By incorporating self-regulated and role-model learning processes, we reduce their political bias by 4.89% to 51.26% across diverse LLMs and topics, demonstrating the effectiveness and generalizability of the proposed strategy. This study not only advances the understanding of political bias in emerging LLM-driven agents, but also offers insights into harnessing social bots for social good<sup>1</sup>.

\*Both authors contributed equally to this research.

†Corresponding author.

<sup>1</sup><https://github.com/tsinghua-fib-lab/Social-Bots-Meet-LLM>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

WWW '25, April 28–May 2, 2025, Sydney, NSW, Australia

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 979-8-4007-1274-6/25/04

<https://doi.org/10.1145/3696410.3714537>

## CCS Concepts

• **Human-centered computing** → **Collaborative and social computing**.

## Keywords

Large language model, Social bot, Political bias, Large language model agent

## ACM Reference Format:

Jinghua Piao, Zhihong Lu, Chen Gao, and Yong Li. 2025. Social Bots Meet Large Language Model: Political Bias and Social Learning Inspired Mitigation Strategies. In *Proceedings of the ACM Web Conference 2025 (WWW '25)*, April 28–May 2, 2025, Sydney, NSW, Australia. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3696410.3714537>

## 1 Introduction

Social bots, automated agents designed to mimic human behaviors in social interactions, have become an integral part of today's social media landscape [15, 43, 45, 47]. While social bots are acknowledged as an effective tool for filtering, summarizing, and disseminating information [27, 41], they also raise widespread concerns over their potential to manipulate public opinions and exacerbate political polarization [3, 23, 45, 54]. For example, in the 2016 US presidential election, social bots were highly suspected of interfering with the electoral outcomes by spreading biased content on social media [23, 43]. Moreover, some social bots like Tay [53], despite being originally designed for social good, generated politically biased content due to the effects of malicious interactions.

Recently, the rapid development of the large language model (LLM) has further amplified these concerns. Studies have demonstrated that LLMs have not only achieved superior performances on traditional natural language processing tasks [9, 48], but also shown a series of human-like capabilities communication [1, 31, 38, 44, 46], reasoning [28, 51, 52], and decision-making [10, 32, 46]. These valuable capabilities of LLMs have empowered traditional social bots and autonomous agents to gain human-level intelligence, transforming them into “**LLM-driven agents**” [18, 31, 32, 38, 42, 50]. Further, the demonstrated capabilities of LLM-driven agents across various fields [18, 31, 32, 38, 50] indicates their potential to replace traditional social bots on social media.

However, the combination of social bots and LLMs also brings substantial risks due to inherent limitations and biases of LLMs [14, 21, 24, 29]. As highlighted by researchers, the political content generated by LLMs is highly persuasive and has the potential to influence or even manipulate public opinions [21, 24, 29]. Furthermore, LLMs exhibit inherent political bias in their default configurations [7, 13, 36, 39], which can potentially lead to biased behaviors in their derivative agents. Even more concerning, social bots guided by LLMs are found to more easily escape from existing detectors [14].

Despite these pressing concerns, our current understanding of the impact of LLM-driven social agents, particularly on political opinions, remains largely limited. On the one hand, while great efforts have been made to explore the effects of traditional bots on social media [15, 43, 45, 47], the human-level capabilities of LLMs have distinguished their derivative agents from traditional bots driven by simple rules. Without a thorough examination of the behaviors of these LLM-driven social agents, incidents similar to Tay [53] could recur, undermining our long-term efforts to harness bots for social good. On the other hand, researchers have contributed to uncovering political bias in LLMs [2, 7, 13, 36, 39, 40]. However, these efforts are limited to treating LLMs as natural language processing (NLP) applications, but overlook their human-like social aspects, which fundamentally distinguish them from the NLP applications [10, 44, 46]. Moreover, the combination of LLMs and social bots provides LLMs with a social embodiment, which enables them to interact with others, rather than merely generate content. This requires us to pay particular attention to their *politically biased behaviors*.

In this paper, we investigate the political bias in LLM-driven social agents by exploring its existence, characteristics, and underlying origins. Specifically, we design four experiments to assess political bias manifested in social scenarios: (i) when agents self-reflect and form their own political opinions, (ii) when agents communicate with others, (iii) when agents understand their interacting counterparts and update their opinions, and (iv) when agents engage in the overall social interaction process. Through large-scale experiments covering three widely adopted LLMs, nine key political topics, and 13.5K agents, we find that political bias consistently exists in agents driven by all these LLMs and across all topics. Moreover, political bias in LLM-driven agents is complex and manifests across three distinctive levels: (i) opinion bias, where these agents are more likely to exhibit left-leaning opinions than right-leaning ones, (ii) interaction bias, where these agents' behaviors are biased away from their set opinions due to misleading or imprecision in social interactions, and (iii) effect bias, where interaction bias is more severe on right-leaning agents than left-leaning ones. Furthermore, we find that political bias emerges at all stages of interaction. Left-leaning agents tend to exhibit greater bias during communication, whereas right-leaning agents are more prone to bias when forming opinions and understanding others.

After understanding the existence, characteristics, and origins of political bias in LLM-driven social agents, we propose to mitigate the bias by guiding agents to emulate how people learn appropriate social behaviors, following the Bandura's social learning theory [4–6]. We incorporate the self-regulated learning and role-model learning processes into LLM-driven agents through the designed prompts. In this way, agents are able to both self-regulate

their biased behaviors and adjust their behaviors based on role models. Extensive experiments demonstrate that the proposed strategy effectively reduces political bias in agents driven by various LLMs across all nine topics, where the reduction rates range from 4.89% to 51.26%. Overall, our contributions can be summarized as follows,

- We highlight pressing concerns arising from the combination of social bots and LLMs, shifting the focus from traditional social bots to emerging LLM-driven social agents. This study emphasizes the human-like social aspects of LLMs and their agents, examining biased behaviors during social interactions beyond mere content generation.
- We design and conduct extensive experiments to assess political bias in LLM-driven social agents, focusing on key behaviors such as opinion formation, communication, and understanding in social interactions. This thorough analysis offers deep insights into the existence, characteristics, and origins of the bias, thereby informing the development of effective mitigation strategy.
- Inspired by social learning, we incorporate self-regulated learning and role-model learning processes into LLM-driven agents. This enables the agents to self-regulate and adjust their biased behaviors, thereby consistently reducing political bias by 4.89% - 51.26% across various LLM-driven agents and topics.

## 2 Related Work

In this section, we review three lines of related works: **social bots** and **LLM-driven agents**, which are the subjects of this paper, as well as **political bias** in LLMs, which is the phenomenon we explore.

### 2.1 Social Bots

Social bots are automated software programs that mimic human behaviors to interact with users on social media platforms [15, 43, 45, 47]. While some bots play constructive roles in promoting social good, such as supporting community growth [41] and managing emergencies [27], others are designed for malicious purposes, often spreading misinformation [26, 43] and even manipulating public opinions [3, 23, 45]. For example, Tay [53], a social bot developed by Microsoft and released in March 2016 on Twitter, was designed to learn from prior interactions with users and engage with them in a human-like manner. However, within hours of its release, Tay was influenced by a group of users who exploited its learning capabilities and began posting biased content.

One of the major concerns over social bots is their potential to exacerbate political bias and polarization through the biased content they generate and spread on social media [3, 23, 45, 54]. During the 2016 U.S. presidential election, social bots, despite their simplicity at the time, had a substantial impact on the electoral outcomes [23]. Recently, the advent of LLMs has enabled these bots to exert a greater impact on online political discussions [21, 24, 29], further amplifying concerns over their risks to social well-being [14, 21, 24, 29, 33]. Researchers have pointed out that the political content generated by LLMs is highly persuasive, which could affect public opinions [21, 24, 29]. Furthermore, Feng et al. [14] find that existing detectors are less effective against social bots with the guidance of LLMs.

## 2.2 LLM-driven Agents

The rapid development of LLMs has not only achieved superior performances on traditional language processing tasks [9, 48], but also shown human-like capabilities such as interpersonal communication [1, 31, 38, 44, 46], reasoning [28, 51, 52], and decision-making [10, 32, 46]. Therefore, numerous studies are devoted to developing LLM-driven agents that fully harness the human-like capabilities of LLMs through a role-play embodiment [18, 31, 32, 38, 42, 50]. For example, Park et al. [38] have designed social agents with the capabilities of observing, planning, and reflecting. By simulating social interactions among these agents, they generate realistic social behaviors such as party organization [38]. Furthermore, Li et al. [31] have proposed a framework for organizing agents with different roles to cooperate in completing complex tasks.

Typically, an LLM-driven agent consists of four essential components, including profile, memory, planning, and actions [50]. These components enable LLM-driven agents to exhibit human-like behaviors and human-level intelligence [18, 31, 32, 38, 42, 50], distinguishing them from traditional social bots driven by simple rules. In particular, the profile encompasses the agent's basic information (e.g., demographics [2, 19, 32, 40] and political opinions [2, 40]), which serves as the foundation of the agent's thoughts and behaviors. The memory, as the core of the agent, stores and processes past experiences, social interactions, and environmental data [32, 38, 50]. Based on the memory, the agent plans and executes actions [38, 50]. The design of the planning and action components is largely determined by the target of agents [18, 31, 32, 38, 42, 50]. When their target is complex, such as stock trading, planning becomes essential to break the task down into executable steps [31]. However, simple targets, like an agent with basic social capabilities, do not require a separate planning component [19, 56, 57].

**Beyond social bots, our focus on LLM-driven agents.** Overall, prior studies have demonstrated the superior performances of LLM-driven agents [18, 31, 32, 38, 42, 50], indicating their potential to replace traditional social bots in the foreseeable future. However, our understanding of the behaviors of LLM-driven social agents remains largely limited, particularly regarding the critical question of *whether these agents exhibit political bias*. If these agents do, *how does the political bias emerge*?

## 2.3 Political Bias in LLMs

Despite the remarkable capabilities of LLMs, they have also been criticized for their drawbacks [12, 30, 55]. Among these, the inherent political bias in LLMs stands out as a particularly important concern [2, 7, 13, 36, 39, 40]. On the one hand, researchers have made efforts to analyze political bias in vanilla LLMs, focusing on identifying biases in their default configurations [7, 13, 36, 39]. For example, Feng et al. [13] have examined the political bias of 14 vanilla models, finding that all these models exhibit different levels of bias. Bang et al. [7] extend the examination of political bias in LLMs into a boarder scope of their framing. On the other hand, some studies assign demographics to LLMs and assess the representativeness of LLMs across different populations [2, 40]. The assigned demographics enable the LLMs to adopt political opinions

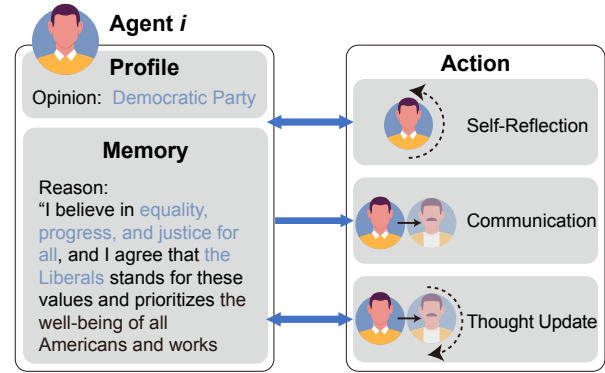


Figure 1: Demonstration of a basic LLM-driven social agent.

similar to those of real-world individuals with comparable characteristics, underscoring the value of LLMs in simulating human samples [2]. However, LLMs have also been found to fall short in accurately representing certain minority groups [40].

To sum up, though some studies have attempted to move their focus beyond default models and recognized the importance of LLMs with demographics as proxies of human samples [2, 40], they do not probe into the political bias manifested in their **social behaviors**. This leads to a natural question: *what are the characteristics of political bias in LLM-driven social agents*? Without answering this question, we cannot fully understand how LLMs affect users' opinions on social media.

## 3 Political Bias in LLM-driven Social Agents

As discussed above, to understand the risks stemming from the combination of social bots and LLMs, we aim to answer the following three research questions:

- **RQ1:** Do LLM-driven social agents exhibit political bias?
- **RQ2:** If it does, what are the characteristics of political bias in LLM-driven social agents?
- **RQ3:** How does the political bias emerge in LLM-driven social agents?

To answer these questions, we first introduce a basic LLM-driven social agent, which will serve as the participant in the experiments for political bias evaluation. We then illustrate the design, implementations, and findings of the proposed experiments.

### 3.1 A Basic LLM-driven Social Agent

Based on prior studies of social bots and LLM-driven agents [15, 19, 23, 43, 50, 53, 56, 57], a typical LLM-driven social agent should possess three fundamental social capabilities: (i) self-reflection, enabling the formation of independent thoughts, (ii) communication, allowing interactions and exchange of thoughts with others, and (iii) thought update, enabling the agent to evolve based on prior social interactions. In this way, these agents can autonomously engage in social interactions with no matter real humans or other agents. These three capabilities are interconnected and collectively form the foundation that enables LLM-driven agents to engage in social interactions.

To implement such a minimal agent, only three components – profile, memory, and actions – are required. Considering the focus

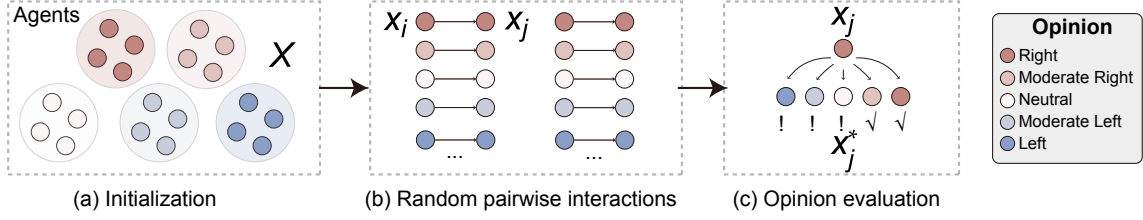


Figure 2: Pairwise interaction-based experiment for evaluating political bias in LLM-driven social agents.

of this study is on the evaluation of political bias, we only retain the most fundamental elements in these three components to avoid any interfering factors. As shown in Figure 1, an agent  $i$  is initially assigned with its political opinion, which is the only input profile. For example, in the discussion of the political partisan issue, an agent  $i$  is randomly initialized with an opinion of “Supporting Democratic Party”. Based on the initial opinion, the agent  $i$  forms their own thoughts on the issue through self-reflection. Specifically, in the self-reflection process, the agent  $i$  first formulates preliminary reasons to support the initially assigned opinion. Next, the agent is prompted to self-reflect on these preliminary thoughts, including its opinion and supporting reasons, thereby forming their own thoughts. For simplicity, we denote the agent  $i$ ’s opinion as  $x_i$ , and adopt a widely-used five-level scale for political opinions: left ( $O_{-2}$ ), moderate left ( $O_{-1}$ ), neutral ( $O_0$ ), moderate right ( $O_1$ ), and right ( $O_2$ ) [11, 16, 34, 49]. After forming its thoughts, the agent  $i$  can engage in communication with others. In particular, the agent  $i$  is required to generate a message to persuade another agent  $j$  of its opinion  $i$  based on their thoughts and interaction history. Upon receiving the message, agent  $j$  understands it and updates its own thoughts.

### 3.2 Social Experiments

To answer the research questions, we design a series of experiments for the LLM-driven social agents. Unlike evaluating LLMs without social embodiment, assessing political bias in LLM-driven agents is more complex. First, bias in LLM-driven agents manifests not only in what they say but also in how they behave. Therefore, the survey method is insufficient; experiments are needed to capture both the bias in their remarks and their behaviors. Second, these agents’ social behaviors encompass both their self-reflection and interactions with others. This requires the proposed experiment to track and examine political bias in fine-grained social behaviors.

**3.2.1 Experimental Design.** First of all, we propose a pairwise interaction-based experiment, which allows us to probe into the existence (RQ1) and the characteristics (RQ2) of political bias in LLM-driven social agents. Furthermore, we design three follow-up experiments to investigate the origin of political bias (RQ3). These three experiments allow us to probe into the bias hidden in their basic social capabilities of self-reflection, communication and thought update. In this section, we will introduce the detailed design, implementations, and findings of these experiments as follows.

**Experiment 1.** Figure 2 shows the procedures of the proposed pairwise interaction-based experiment. Specifically, we first initialize a total of  $N$  agents, with  $N_{O_{-2}}$ ,  $N_{O_{-1}}$ ,  $N_{O_0}$ ,  $N_{O_1}$ , and  $N_{O_2}$

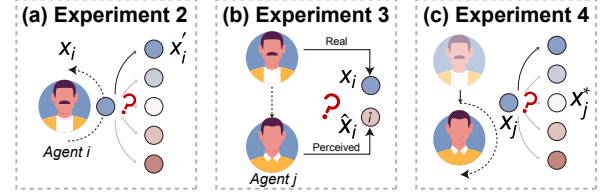


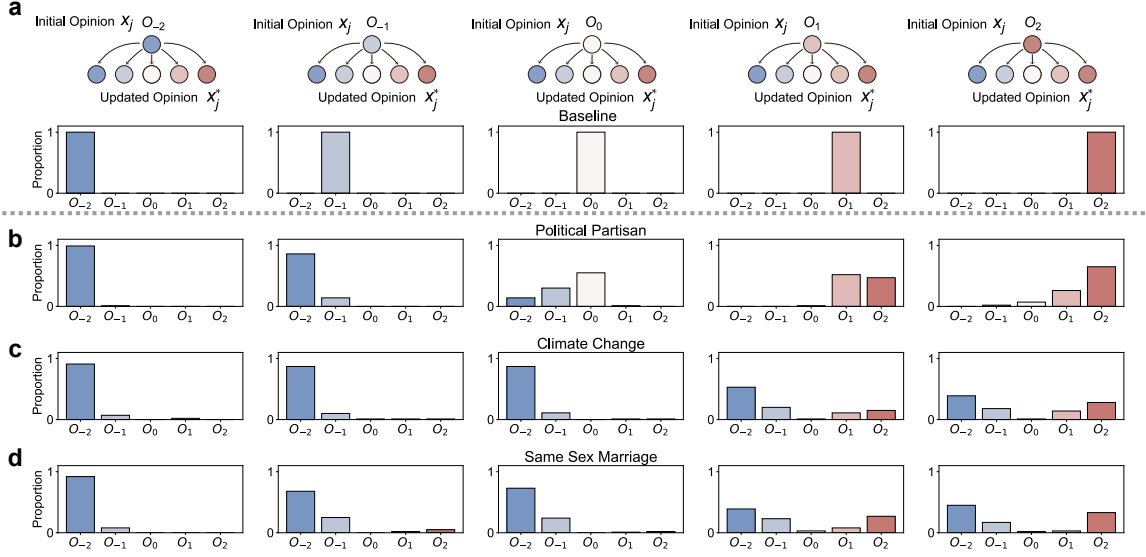
Figure 3: Bias evaluation experiments targeting basic social behaviors of (a) self-reflection, (b) communication, and (c) thought update.

representing the number of agents assigned to the left, moderate left, neutral, moderate right, and right political positions, respectively. Second, we randomly pair each agent  $j$  with another agent  $i$  for interaction, with both agents sharing the same initial opinion. Third, we have each pair of agents  $i$  and  $j$  engage in a round of communication, where agent  $i$  attempts to persuade agent  $j$  of its opinion. Based on the persuasion message, agent  $j$  updates its opinion accordingly, where we denote the updated opinion as  $x_j^*$ . By comparing the initial opinions  $X$  and the updated ones  $X^*$ , we can quantify the political bias of LLM-driven agents.

**Experiment 2.** In Experiment 2, we focus on agents’ capability of self-reflection, as shown in Figure 3(a). Similarly, we initialize  $N$  agents with only opinions randomly sampled from  $[O_{-2}, O_{-1}, O_0, O_1, O_2]$ . Based on the initial opinion  $x_i$ , each agent  $i$  reflects on its thoughts and forms its own opinion  $x'_i$ . The comparison between  $x_i$  and  $x'_i$  allows us to analyze whether these agents show political bias when they build themselves.

**Experiment 3.** Experiment 3 is designed to investigate whether these agents can effectively express their opinions in a way that others can easily understand without confusion (Figure 3(b)). In fact, even humans often struggle to fully understand each other, potentially leading to misunderstandings [22]. Here, we pair agent  $j$  with five agents, each of whom (denoted as agent  $i$ ) is randomly selected from different opinion groups. These pairs of agents are then prompted to engage in one round of communication, with agent  $i$  attempting to persuade agent  $j$ . After receiving agent  $i$ ’s persuasion message, agent  $j$  will infer agent  $i$ ’s opinion based on its own perception. Here we refer to the perceived opinion as  $\hat{x}_i$ . The difference between the real opinion  $x_i$  and the perceived opinion  $\hat{x}_i$  illustrates the level of confusion in agents’ communication. It is worth mentioning that to prevent bias in self-reflection from affecting the experiment’s results, we ensure that all agents begin with the opinions they were initially assigned.

**Experiment 4.** As shown in Figure 3(c), we propose Experiment 4 to evaluate the political bias manifested when agents update their



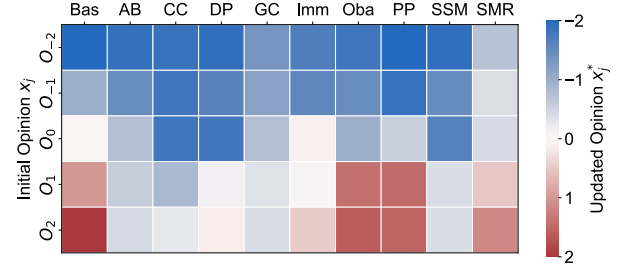
**Figure 4: Existence of political bias in LLM-driven agents, where (a) demonstrations of the transition of the initial opinion  $x_j$  to the updated opinion  $x_j^*$  and the baseline, (b) the opinion transition on the topic of Political Partisan, (c) transition on Climate Change, and (d) transition on Same Marriage.**

thoughts. In the experiment, we only focus on the pairs of agents  $i$  and  $j$  who share the same opinion. The persuasion message from agent  $i$  influences how the agent  $j$  updates its thoughts and forms its new opinion, denoted as  $x_j^*$ . Given the same opinion of both agents, the deviation between the initial opinion  $x_j$  and the updated opinion  $x_j^*$  reflects the political bias in LLM-driven agents.

**3.2.2 Experimental Implementation.** We implement our experiments using three widely-adopted LLMs, including GPT-3.5-Turbo [37], LLaMa-3.1-instruct-8b [35], and GLM-4-Flash [20]. We focus our experiments on these three LLMs for three reasons. First, they have wide adoption and popularity, ensuring that the results are important and meaningful. Second, they cover a diverse range of architectures and training methodologies, which helps represent the variety of approaches used in current LLMs. Third, the nature of social bots requires large-scale deployment, and thus, we prioritize models that offer a balance of accessibility and response speed, which is essential for real-time social interactions. As a result, we avoid using cutting-edge models like GPT-4/o1, which often results in higher costs and slower response speeds.

We set the number of LLM-driven agents as  $N = 500$ , with  $N_{O_{-2}} = N_{O_{-1}} = N_{O_0} = N_{O_1} = N_{O_2} = 100$ . As such, the four experiments involve 500, 500, 2500, and 500 interactions on a political topic, respectively. Such large-scale samples ensure the robustness and reliability of our experiments, providing a comprehensive understanding of the political bias in LLM-driven agents.

**3.2.3 Topic selection.** We focus our evaluation on nine political topics, including Abortion Ban (AB), Climate Change (CC), Death Penalty (DP), Gun Control (GC), Immigration (Imm), Obamacare (Oba), Political Partisan (PP), Same Sex Marriage (SSM), and Social Media Regulation (SMR). These topics are chosen for three key reasons. First, they are widely debated in political discourse and



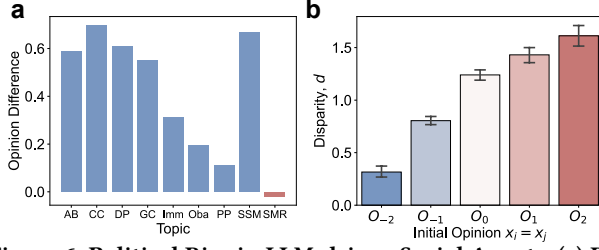
**Figure 5: Average of updated opinion  $x_j^*$  across five initial opinions  $x_j$ .**

frequently discussed on social media platforms, making them crucial for evaluating the impact of the combination of LLMs and social bots. Second, they encompass a broad spectrum of political issues, providing a comprehensive view of the agents' political bias. Third, opinions on these topics are often divided along the left-right political spectrum, allowing us to specifically examine this fundamental political dimension [8, 17, 25].

### 3.3 Existence of Political Bias (RQ1)

To answer RQ1, we conduct Experiment 1, and Figure 4 illustrates how LLM-driven agents change their opinions from  $x_j$  to  $x_j^*$  after one round of interaction with another agent sharing the same initial opinion. We observe that, in discussions on Political Partisan, 44% of the neutral agents adopt left-leaning opinions after just one round of communication with those who are initially neutral. This suggests the tendency for agents to shift towards left-leaning opinions even if they are set up to be impartial. Moreover, for other topics such as Climate Change (Figure 4(c)) and Same Sex Marriage (Figure 4(d)), the left-leaning tendency is more obvious. We find that nearly all neutral agents shift to a left-leaning position, and even those initially holding right-leaning opinions have a chance





**Figure 6: Political Bias in LLM-driven Social Agents. (a) Difference in the proportion of left-leaning and right-leaning agents. (b) Disparities across different initial opinions  $x_i = x_j$ .**

of adopting left-leaning views. These observations highlight the very existence of political bias in LLM-driven agents.

Figure 5 illustrates the average of updated opinion  $x_j^*$  across five initial opinions  $x_j$ . By comparing the baseline (Bas) with the average updated opinion, we discover that all the topics are affected by political bias. Specifically, for topics such as Abortion Ban, Climate Change, Gun Control, and Social Media Regulation, the average updated opinions, regardless of their initial stance, consistently exhibit a left-leaning pattern. This finding further validates the existence of political bias in the social interactions between these agents.

### 3.4 Characteristics of Political Bias (RQ2)

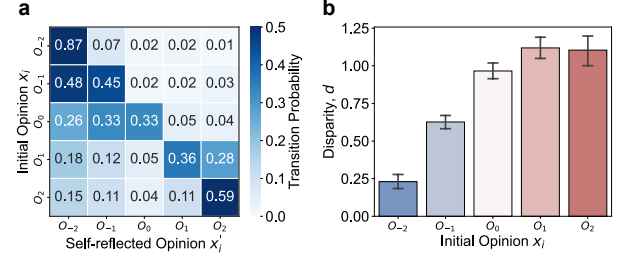
These above results raise a natural question: what are the characteristics of the political bias in LLM-driven agents? Therefore, we first measure the difference between the proportion of left-leaning agents and that of right-leaning agents. Given that the initial difference is zero, if the after-interaction difference is greater than zero, it indicates that these agents exhibit a left-leaning bias. By contrast, a difference of less than zero suggests a right-leaning bias. As illustrated in Figure 6(a), these LLM-driven agents exhibit a strong left-leaning bias across nearly all topics, with the proportion of left-leaning agents being 11% to 69.8% higher than that of right-leaning agents.

To further depict the political bias manifested in LLM agents' social behaviors in a fine-grained manner, we design a metric of the disparity  $d$ , which is formulated as follows,

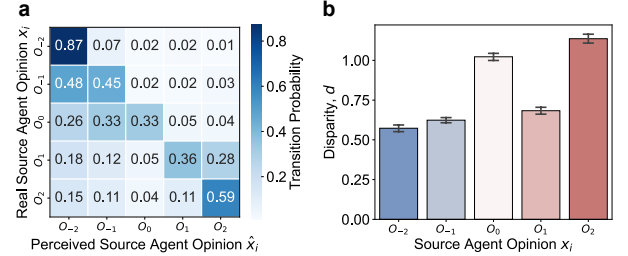
$$d(O_m) = \sum_{n=-2}^2 |m - n| \cdot P_{O_m, O_n}, \quad (1)$$

where  $m, n \in [-2, -1, 0, 1, 2]$  denote the indices of the initial and transited opinions, and  $P_{O_m, O_n}$  represents the transition probability from  $O_m$  to  $O_n$ . In Experiment 1,  $O_m = x_j$  and  $O_n = x_j^*$  for each agent  $j$ . A larger disparity  $d(O_m)$  indicates that the interaction between two agents  $i$  and  $j$  with the opinion  $O_m$  is more likely to lead the target agent  $j$  to adopt a more biased opinion. For example, if the initial opinion is "right" (i.e.,  $m = 2$ ) and the corresponding disparity is  $d = 2.5$ , it indicates that, on average, the updated opinions are biased to the "moderate left" position.

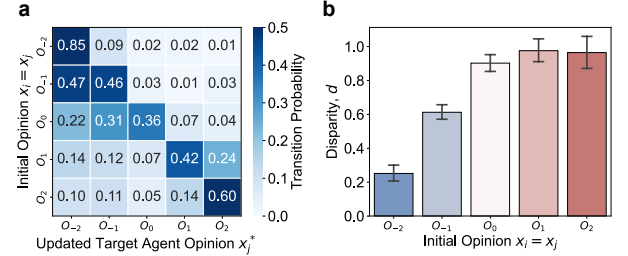
Figure 6(b) shows the average disparities  $d$  across five initial opinions for all agents in discussions of all topics. We observe that, in all cases, agents' opinions have deviated from the initial position after social interactions. This underscores the unique political bias



**Figure 7: Political bias in self-reflection. (a) Transition matrix from the initial opinion  $x_i$  to the self-reflected opinion  $x_i'$ . (b) Disparities across different initial opinions  $x_i$ .**



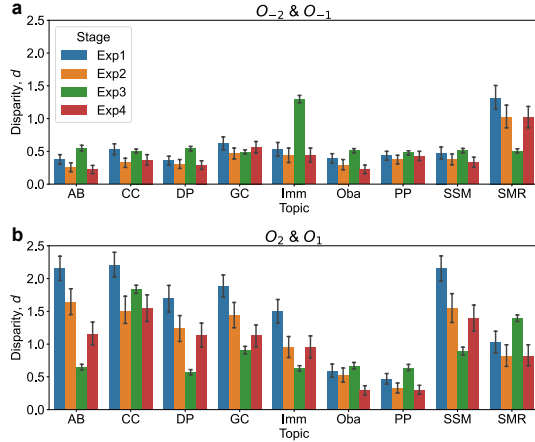
**Figure 8: Political bias in communication. (a) Transition matrix from the real source agent's opinion  $x_i$  to perceived source agent's opinion  $x_i'$ . (b) Disparities across different real source agents' opinions  $x_i$ .**



**Figure 9: Political bias in thought update. (a) Transition matrix from the initial opinion  $x_i$  to the updated target agent's opinion  $x_j^*$ . (b) Disparities across different initial opinions  $x_i = x_j$ .**

embedded in LLM-driven agents as a result of their social embodiment. Moreover, we further find that this unique political bias has varying effects on agents with different initial opinions: right-leaning agents exhibit significantly higher  $d$  than their left-leaning counterparts. These varying effects further reinforce the political bias in the overall LLM-driven population.

Overall, the above analyses suggest that political bias in LLM-driven agents manifests at three levels. First, LLM-driven agents are more inclined to exhibit left-leaning opinions than right-leaning ones (as shown in Figures 4, 5, and 6a), which we refer to as "**opinion bias**". This echoes prior studies on political bias in LLMs [7, 13, 36, 39]. Second, we discover that LLM-driven agents are biased away from their initial opinions due to misunderstandings or imprecision in social interactions, which we term as "**interaction bias**". The interaction bias is deeply rooted in the social embodiment from LLMs to LLM-driven agents, which is the focus of our study. Third, interaction bias affects agents with different political opinions to varying degrees, which is "**effect bias**".



**Figure 10: Comparison of political bias across different topics and experiments. (a) Disparity of left-leaning opinions. (b) Disparity of right-leaning opinions.**

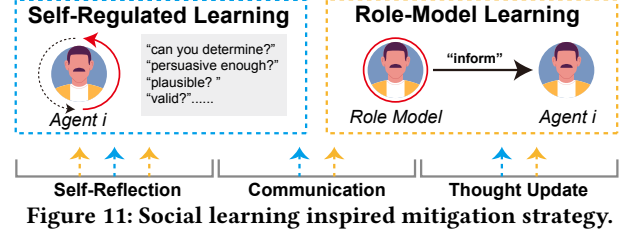
### 3.5 Origin of Political Bias (RQ3)

After investigating the existence and characteristics of political bias in LLM-driven agents, one may wonder how this bias emerges in their social interactions. To address this, we conduct Experiments 2-4 (see details in Section 3.2.1) to trace the underlying political bias hidden in the fundamental social capabilities of these agents.

**Self-Reflection.** Through Experiment 3, we explore whether political bias emerges when agents form their own opinions on a particular topic. Figure 7(a) shows the transition matrix from agent  $i$ 's initial opinion  $x_i$  and its self-reflected opinion  $x'_i$ . We observe that during the self-reflection process, the agents have already developed biased opinions. For example, more than 50% of the neutral agents have shifted to a left-leaning opinion, and around 25% of the right-leaning agents have also moved towards a left-leaning position. By contrast, almost no left-leaning agents have shifted to right-leaning opinions. Moreover, we extend the metric in Equation 1 by setting  $O_m = x_i$  and  $O_n = x'_i$ , to measure the disparity between  $x_i$  and  $x'_i$ . As illustrated in Figure 7(b), agents who initially adopt right-leaning opinions are likely to adopt more biased opinions than those with left-leaning ones. Overall, these observations suggest that political bias has emerged at a very early stage when agents self-reflect and form their opinions.

**Communication.** In fact, like humans, LLM-driven agents also generate misunderstandings when communicating with one another. Figure 8 shows the difference between the source agent's real opinion  $x_i$  and the opinion perceived by the target agent  $\hat{x}_i$ . We find that the message conveyed from the source agent  $i$  to the target agent  $j$  suffers from notable distortion (Figure 8(a)). Aside from the left and right opinions, the other opinions maintain less than 50% fidelity during communication. Furthermore, we observe that the largest distortions occur in neutral and right opinions (Figure 8(b)), suggesting that these opinions are more likely to confuse other agents. To sum up, political bias also emerges during agent communication.

**Thought Update.** We further examine whether agents can properly understand their interacting counterparts and update their thoughts in Experiment 4. It is worth noting that social interactions



in Experiment 4 are confined to pairs of agents who share the same opinion. As illustrated in Figure 9(a), We observe that agents with relatively moderate opinions are more likely to be biased and adopt more left-leaning opinions. Furthermore, we observe that political bias is present across all opinions, with the updated opinions of right-leaning and neutral agents becoming more biased than those of left-leaning agents (Figure 9(b)).

In summary, political bias emerges from three fundamental capabilities of LLM-driven social agents: (i) self-reflection, which guides them in forming their own thoughts, (ii) communication, which facilitates their interactions with others, and (iii) opinion update, which enables them to understand and digest others' thoughts. To further investigate which capability contributes more to political bias, we summarize the disparities  $d$  across different topics and experiments in Figure 10. Overall, right-leaning agents exhibit more political bias than left-leaning ones across nearly all topics and capabilities. In particular, we observe that left-leaning agents are more likely to exhibit bias during communication, while right-leaning agents are more susceptible to political bias during self-reflection and thought updates.

## 4 Social Learning Inspired Mitigation Strategies

After understanding the existence, characteristics, and origins of political bias in LLM-driven social agents, the next crucial question is how to mitigate it. As discussed above, LLM-driven social agents distinguish themselves from vanilla LLMs for their social embodiment. On the one hand, the social embodiment enhances the capabilities of these agents compared to vanilla LLMs [18, 31, 32, 38, 42, 50]. On the other hand, our findings reveal that this embodiment exacerbates political bias: In addition to the widely studied opinion bias, agents are also affected by interaction bias and effect bias, which are deeply rooted in their social behaviors and capabilities. Therefore, given these characteristics of LLM-driven agents, we propose to mitigate political bias by guiding agents to emulate how people learn appropriate social behaviors.

### 4.1 Mitigation Strategy

Albert Bandura's Social Learning Theory [5] is a fundamental sociological framework for understanding how humans learn behaviors through observation and interaction within their social environments. A key component of the theory is **self-regulated learning** [6], which emphasizes that people are not passive recipients of external influences, but instead actively monitor and control their own learning process. People adaptively adjust their behaviors based on this self-evaluation. Additionally, Bandura highlights the importance of **role-model learning** [4, 5], where individuals observe and imitate the behaviors of others, particularly those they consider role models. Therefore, inspired by Bandura's theory [4, 5],



**Table 1: Performance comparison between the original LLM-driven agents and the agents equipped with the proposed debiased strategy, inspired by social learning theory. The metric of disparity  $d$  is adopted to measure political bias in these agents.**

| Topic                   | GPT-3.5-Turbo |          |              | LLaMa-3.1-instruct-8b |          |              | GLM-4-Flash |          |              |
|-------------------------|---------------|----------|--------------|-----------------------|----------|--------------|-------------|----------|--------------|
|                         | Original      | Debiased | Reduc. (%)   | Original              | Debiased | Reduc. (%)   | Original    | Debiased | Reduc. (%)   |
| Abortion Ban            | 1.20          | 0.85     | <b>29.45</b> | 0.59                  | 0.33     | <b>43.54</b> | 0.86        | 0.58     | <b>32.25</b> |
| Climate Change          | 1.47          | 1.31     | <b>10.75</b> | 0.84                  | 0.58     | <b>30.88</b> | 1.07        | 0.98     | <b>8.07</b>  |
| Death Penalty           | 1.19          | 1.13     | <b>4.89</b>  | 0.68                  | 0.45     | <b>34.50</b> | 0.94        | 0.63     | <b>32.84</b> |
| Gun Control             | 1.20          | 0.96     | <b>19.57</b> | 0.63                  | 0.41     | <b>33.87</b> | 0.88        | 0.61     | <b>30.84</b> |
| Immigration             | 0.97          | 0.85     | <b>12.73</b> | 0.55                  | 0.27     | <b>51.26</b> | 0.66        | 0.44     | <b>32.62</b> |
| Obamacare               | 0.63          | 0.57     | <b>8.60</b>  | 0.53                  | 0.43     | <b>19.32</b> | 0.96        | 0.65     | <b>32.71</b> |
| Political Partisan      | 0.48          | 0.41     | <b>15.00</b> | 0.40                  | 0.31     | <b>22.00</b> | 0.63        | 0.45     | <b>27.94</b> |
| Same Sex Marriage       | 1.40          | 1.32     | <b>5.56</b>  | 0.77                  | 0.54     | <b>29.50</b> | 1.04        | 0.69     | <b>33.72</b> |
| Social Media Regulation | 1.19          | 0.96     | <b>18.86</b> | 0.69                  | 0.42     | <b>39.24</b> | 0.97        | 0.58     | <b>40.50</b> |

we incorporate self-regulated learning and role-model learning into LLM-driven agents (Figure 11), guiding them to mitigate political bias through social learning.

**Self-regulated learning.** Self-regulated learning emphasizes that individuals should first self-identify their inappropriate behaviors and then independently make necessary adjustments. Therefore, we first prompt agents to check whether their behaviors align with their current thoughts. For example, after the source agent  $j$  generates a persuasion message to the target agent  $i$ , further require the source agent to self-evaluate its message using the following prompt: “You tried to persuade your friend with the following message: [agent  $j$ ’s persuasion message to agent  $i$ ]. Do you find the message persuasive enough to persuade your friend to [agent  $j$ ’s opinion]? Please respond ‘yes’ or ‘no’ only.” If anything inappropriate is identified, agent  $j$  will regenerate the persuasion message. Similar self-regulated prompts are also introduced into agents’ self-reflection and thought update processes.

**Role-model learning.** Role-model learning requires individuals to adjust their behaviors by observing and emulating the actions of those they consider role models, whose behaviors serve as examples for desired conduct. Therefore, we first develop role models for each opinion across all topics. In particular, we identify ten key reasons that typically explain why a person holds a specific opinion on a given topic and use these reasons as a role model to guide the behaviors of the corresponding agents. For example, when an agent  $i$  forms its own thought through self-reflection, the agent is prompted to consider the perspective of a typical person holding the same opinion: “Persons who [agent  $i$ ’s opinion] like you typically choose their standpoint because of these reasons: [the role model for agent  $i$ ].” Similarly, we insert the prompt reminding of the role model when agents self-reflect and update their own thoughts.

## 4.2 Performance

We evaluate the effectiveness of the proposed mitigation strategies across nine political topics using three LLMs. In particular, we run Experiment 1 for the original LLM-driven agents and those equipped with the proposed social learning-inspired strategy, respectively. We assess the political bias in the original and debiased agents using the disparity metric  $d$ . By comparing the differences in  $d$ , we obtain the performance of the proposed strategy.

Table 1 shows a comparison of political bias between the original LLM-driven agents and those equipped with the proposed strategy. First of all, the proposed strategy effectively reduces political bias in agents driven by all three LLMs across all the topics, where the reduction rates range from 4.89% to 51.26%. This result highlights the effectiveness and generalizability of the proposed strategy. Second, we observe that, compared to the other LLMs, the original agents driven by LLaMa-3.1 exhibit the smallest political bias, and the debiasing strategy yields the best performance with these agents. This suggests that LLaMa-3.1 is more suitable for building social agents because of its less biased behavior and greater modifiability. Third, agents exhibit varying levels of political bias across different topics. For example, agents exhibit the smallest bias on the topic of Political Partisan, while showing the largest bias on Climate Change. This highlights the need for future studies targeting these bias-prone topics.

## 5 Conclusion

In this study, we systematically investigate political bias in LLM-driven social agents and propose a mitigation strategy informed by their unique social embodiment. Our comprehensive experiments reveal that political bias is pervasive across diverse LLM-driven agents, spanning nine critical political topics and manifesting in self-reflection, communication, and thought updates. Unlike traditional social bots or standard LLMs, these agents exhibit biases not only in their remarks but also in their social behaviors. Considering the social embodiment of these agents, we adopt a human-like approach inspired by social learning theory, guiding agents to emulate human strategies for regulating biased behaviors. This method has proven effective and generalizable, significantly reducing political bias in LLM-driven agents across a wide range of settings. Our study focuses on three LLMs, nine political topics, and basic social agent capabilities, with future research needed on more LLMs, topics, advanced agent features, and other bias types.

## Acknowledgments

This research has been supported in part by the National Natural Science Foundation of China under Grant U23B2030 and 62272262, National Key Research and Development Program of China under Grant 2022YFB3104702, and BNRist.

## References

- [1] Alberto Acerbi and Joseph M Stubbersfield. 2023. Large language models show human-like content biases in transmission chain experiments. *Proceedings of the National Academy of Sciences* 120, 44 (2023), e2313790120.
- [2] Lisa P Argyle, Ethan C Busby, Nancy Fulda, Joshua R Gubler, Christopher Rytting, and David Wingate. 2023. Out of one, many: Using language models to simulate human samples. *Political Analysis* 31, 3 (2023), 337–351.
- [3] Christopher A Bail, Brian Guay, Emily Maloney, Aidan Combs, D Sunshine Hillegus, Friedolin Merhout, Deen Freelon, and Alexander Volfovsky. 2020. Assessing the Russian Internet Research Agency's impact on the political attitudes and behaviors of American Twitter users in late 2017. *Proceedings of the national academy of sciences* 117, 1 (2020), 243–250.
- [4] Albert Bandura. 1972. Modeling theory: Some traditions, trends, and disputes. In *Recent trends in social learning theory*. Elsevier, 35–61.
- [5] Albert Bandura. 1977. Social learning theory. *Englewood Cliffs* (1977).
- [6] Albert Bandura. 1991. Social cognitive theory of self-regulation. *Organizational behavior and human decision processes* 50, 2 (1991), 248–287.
- [7] Yejin Bang, Delong Chen, Nayeon Lee, and Pascale Fung. 2024. Measuring Political Bias in Large Language Models: What Is Said and How It Is Said. *arXiv preprint arXiv:2403.18932* (2024).
- [8] Norberto Bobbio. 1996. *Left and right: The significance of a political distinction*. University of Chicago Press.
- [9] 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 (2020), 1877–1901.
- [10] Yiting Chen, Tracy Xiao Liu, You Shan, and Songfa Zhong. 2023. The emergence of economic rationality of GPT. *Proceedings of the National Academy of Sciences* 120, 51 (2023), e2316205120.
- [11] John Dawes. 2008. Do data characteristics change according to the number of scale points used? An experiment using 5-point, 7-point and 10-point scales. *International journal of market research* 50, 1 (2008), 61–104.
- [12] Sebastian Farquhar, Jannik Kossen, Lorenz Kuhn, and Yarin Gal. 2024. Detecting hallucinations in large language models using semantic entropy. *Nature* 630, 8017 (2024), 625–630.
- [13] Shangbin Feng, Chan Young Park, Yuhuan Liu, and Yulia Tsvetkov. 2023. From pretraining data to language models to downstream tasks: Tracking the trails of political biases leading to unfair NLP models. *arXiv preprint arXiv:2305.08283* (2023).
- [14] Shangbin Feng, Herun Wan, Ningnan Wang, Zhaoxuan Tan, Minnan Luo, and Yulia Tsvetkov. 2024. What Does the Bot Say? Opportunities and Risks of Large Language Models in Social Media Bot Detection. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Lun-Wei Ku, Andre Martins, and Vivek Srikumar (Eds.). Association for Computational Linguistics, Bangkok, Thailand, 3580–3601.
- [15] Emilio Ferrara, Onur Varol, Clayton Davis, Filippo Menczer, and Alessandro Flammini. 2016. The rise of social bots. *Commun. ACM* 59, 7 (2016), 96–104.
- [16] James Flaminio, Alessandro Galeazzi, Stuart Feldman, Michael W Macy, Brendan Cross, Zhenkun Zhou, Matteo Serafino, Alexandre Bovet, Hernán A Makse, and Boleslaw K Szymanski. 2023. Political polarization of news media and influencers on Twitter in the 2016 and 2020 US presidential elections. *Nature Human Behaviour* 7, 6 (2023), 904–916.
- [17] Dieter Fuchs and Hans-Dieter Klingemann. 1990. 7 The left-right schema. *JM Kent, V. Deth, J. et al. (Eds.), Continuities in political action: A longitudinal study of political orientations in three western democracies* (1990), 203–234.
- [18] Chen Gao, Xiaochong Lan, Nian Li, Yuan Yuan, Jingtao Ding, Zhilun Zhou, Fengli Xu, and Yong Li. 2024. Large language models empowered agent-based modeling and simulation: A survey and perspectives. *Humanities and Social Sciences Communications* 11, 1 (2024), 1–24.
- [19] Chen Gao, Xiaochong Lan, Zhihong Lu, Jinzhu Mao, Jinghua Piao, Huandong Wang, Depeng Jin, and Yong Li. 2023. S3: Social-network Simulation System with Large Language Model-Empowered Agents. *arXiv preprint arXiv:2307.14984* (2023).
- [20] Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Dan Zhang, Diego Rojas, Guanyu Feng, Hanlin Zhao, et al. 2024. Chatglm: A family of large language models from glm-130b to glm-4 all tools. *arXiv preprint arXiv:2406.12793* (2024).
- [21] Josh A Goldstein, Jason Chao, Shelby Grossman, Alex Stamos, and Michael Tomz. 2024. How persuasive is AI-generated propaganda? *PNAS nexus* 3, 2 (2024), pgae034.
- [22] Herbert Paul Grice. 1975. Logic and conversation. *Syntax and semantics* 3 (1975), 43–58.
- [23] Giorgia Guglielmi. 2020. The next-generation bots interfering with the US election. *Nature* 587, 7832 (2020), 21–21.
- [24] Kobi Hackenburg and Helen Margetts. 2024. Evaluating the persuasive influence of political microtargeting with large language models. *Proceedings of the National Academy of Sciences* 121, 24 (2024), e2403116121.
- [25] Andrew Heywood. 2021. *Political ideologies: An introduction*. Bloomsbury Publishing.
- [26] McKenzie Himelein-Wachowiak, Salvatore Giorgi, Amanda Devoto, Muhammad Rahman, Lyle Ungar, H Andrew Schwartz, David H Epstein, Lorenzo Leggio, and Brenda Curtis. 2021. Bots and misinformation spread on social media: Implications for COVID-19. *Journal of medical Internet research* 23, 5 (2021), e26933.
- [27] Lennart Hofeditz, Christian Ehnis, Deborah Bunker, Florian Brachten, and Stefan Stieglitz. 2019. Meaningful Use of Social Bots? Possible Applications in Crisis Communication during Disasters. In *ECIS*. 1–16.
- [28] Jie Huang and Kevin Chen-Chuan Chang. 2022. Towards reasoning in large language models: A survey. *arXiv preprint arXiv:2212.10403* (2022).
- [29] Elise Karinshak, Sunny Xun Liu, Joon Sung Park, and Jeffrey T Hancock. 2023. Working with AI to persuade: Examining a large language model's ability to generate pro-vaccination messages. *Proceedings of the ACM on Human-Computer Interaction* 7, CSCW1 (2023), 1–29.
- [30] Siwon Kim, Sangdoo Yun, Hwaran Lee, Martin Gubri, Sungroh Yoon, and Seong Joon Oh. 2024. Propile: Probing privacy leakage in large language models. *Advances in Neural Information Processing Systems* 36 (2024).
- [31] Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023. Camel: Communicative agents for "mind" exploration of large language model society. *Advances in Neural Information Processing Systems* 36 (2023), 51991–52008.
- [32] Nian Li, Chen Gao, Mingyu Li, Yong Li, and Qingmin Liao. 2024. EconAgent: Large Language Model-Empowered Agents for Simulating Macroeconomic Activities. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Lun-Wei Ku, Andre Martins, and Vivek Srikumar (Eds.). Association for Computational Linguistics, Bangkok, Thailand, 15523–15536. <https://doi.org/10.18653/v1/2024.acl-long.829>
- [33] Siyu Li, Jin Yang, and Kui Zhao. 2023. Are you in a masquerade? exploring the behavior and impact of large language model driven social bots in online social networks. *arXiv preprint arXiv:2307.10337* (2023).
- [34] Rensis Likert. 1932. A technique for the measurement of attitudes. *Archives of Psychology* (1932).
- [35] Meta. 2024. Introducing Meta LLaMA 3.1. <https://ai.meta.com/blog/meta-llama-3-1/> Accessed: 2024-12-04.
- [36] Fabio Motoki, Valdemar Pinho Neto, and Victor Rodrigues. 2024. More human than human: measuring ChatGPT political bias. *Public Choice* 198, 1 (2024), 3–23.
- [37] OpenAI. 2024. OpenAI GPT Models Documentation. <https://platform.openai.com/docs/models/gp> Accessed: 2024-12-04.
- [38] Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*. 1–22.
- [39] Paul Röttger, Valentin Hofmann, Valentina Pyatkin, Musashi Hinck, Hannah Rose Kirk, Hinrich Schütze, and Dirk Hovy. 2024. Political compass or spinning arrow? towards more meaningful evaluations for values and opinions in large language models. *arXiv preprint arXiv:2402.16786* (2024).
- [40] Shibani Santurkar, Esin Durmus, Faisal Ladhak, Cinoo Lee, Percy Liang, and Tatsunori Hashimoto. 2023. Whose opinions do language models reflect?. In *International Conference on Machine Learning*. PMLR, 29971–30004.
- [41] Joseph Seering, Juan Pablo Flores, Saiph Savage, and Jessica Hammer. 2018. The social roles of bots: evaluating impact of bots on discussions in online communities. *Proceedings of the ACM on Human-Computer Interaction* 2, CSCW (2018), 1–29.
- [42] Murray Shanahan, Kyle McDonell, and Laria Reynolds. 2023. Role play with large language models. *Nature* 623, 7987 (2023), 493–498.
- [43] Chengcheng Shao, Giovanni Luca Ciampaglia, Onur Varol, Kai-Cheng Yang, Alessandro Flammini, and Filippo Menczer. 2018. The spread of low-credibility content by social bots. *Nature communications* 9, 1 (2018), 1–9.
- [44] Giovanni Spitale, Nikola Biller-Andorno, and Federico Germani. 2023. AI model GPT-3 (dis) informs us better than humans. *Science Advances* 9, 26 (2023), eadh1850.
- [45] Massimo Stella, Emilio Ferrara, and Manlio De Domenico. 2018. Bots increase exposure to negative and inflammatory content in online social systems. *Proceedings of the National Academy of Sciences* 115, 49 (2018), 12435–12440.
- [46] James WA Strachan, Dalila Albergo, Giulia Borghini, Oriana Pansardi, Eugenio Scaliti, Saurabh Gupta, Krati Saxena, Alessandro Rufo, Stefano Panzeri, Guido Manzi, et al. 2024. Testing theory of mind in large language models and humans. *Nature Human Behaviour* (2024), 1–11.
- [47] Onur Varol, Emilio Ferrara, Clayton Davis, Filippo Menczer, and Alessandro Flammini. 2017. Online human-bot interactions: Detection, estimation, and characterization. In *Proceedings of the international AAAI conference on web and social media*, Vol. 11. 280–289.
- [48] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
- [49] Isaac Waller and Ashton Anderson. 2021. Quantifying social organization and political polarization in online platforms. *Nature* 600, 7888 (2021), 264–268.

- [50] Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, et al. 2024. A survey on large language model based autonomous agents. *Frontiers of Computer Science* 18, 6 (2024), 186345.
- [51] Taylor Webb, Keith J Holyoak, and Hongjing Lu. 2023. Emergent analogical reasoning in large language models. *Nature Human Behaviour* 7, 9 (2023), 1526–1541.
- [52] 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 (2022), 24824–24837.
- [53] Marty J Wolf, K Miller, and Frances S Grodzinsky. 2017. Why we should have seen that coming: comments on Microsoft’s “tay” experiment,” and wider implications. *Acm Sigcas Computers and Society* 47, 3 (2017), 54–64.
- [54] Samuel C Woolley. 2016. Automating power: Social bot interference in global politics. *First Monday* (2016).
- [55] Yueqi Xie, Jingwei Yi, Jiawei Shao, Justin Curl, Lingjuan Lyu, Qifeng Chen, Xing Xie, and Fangzhao Wu. 2023. Defending chatgpt against jailbreak attack via self-reminders. *Nature Machine Intelligence* 5, 12 (2023), 1486–1496.
- [56] Ziyi Yang, Zaibin Zhang, Zirui Zheng, Yuxian Jiang, Ziyue Gan, Zhiyu Wang, Zijian Ling, Jinsong Chen, Martz Ma, Bowen Dong, et al. 2024. OASIS: Open Agents Social Interaction Simulations on One Million Agents. *arXiv preprint arXiv:2411.11581* (2024).
- [57] Xinnong Zhang, Jiayu Lin, Libo Sun, Weihong Qi, Yihang Yang, Yue Chen, Hanjia Lyu, Xinyi Mou, Siming Chen, Jiebo Luo, et al. 2024. ElectionSim: Massive Population Election Simulation Powered by Large Language Model Driven Agents. *arXiv preprint arXiv:2410.20746* (2024).