

# Polarization of Autonomous Generative AI Agents Under Echo Chambers

Masaya Ohagi  
asikapool@gmail.com  
SB Intuitions  
Tokyo, Japan

## ABSTRACT

Online social networks often create echo chambers where people only hear opinions reinforcing their beliefs. An echo chamber often generates polarization, leading to conflicts caused by people with radical opinions, such as the January 6, 2021, attack on the US Capitol. The echo chamber has been viewed as a human-specific problem, but this implicit assumption is becoming less reasonable as large language models, such as ChatGPT, acquire social abilities. In response to this situation, we investigated the potential for polarization to occur among a group of autonomous AI agents based on generative language models in an echo chamber environment. We had AI agents discuss specific topics and analyzed how the group’s opinions changed as the discussion progressed. As a result, we found that the group of agents based on ChatGPT tended to become polarized in echo chamber environments. The analysis of opinion transitions shows that this result is caused by ChatGPT’s high prompt understanding ability to update its opinion by considering its own and surrounding agents’ opinions. We conducted additional experiments to investigate under what specific conditions AI agents tended to polarize. As a result, we identified factors that strongly influence polarization, such as the agent’s persona. These factors should be monitored to prevent the polarization of AI agents.

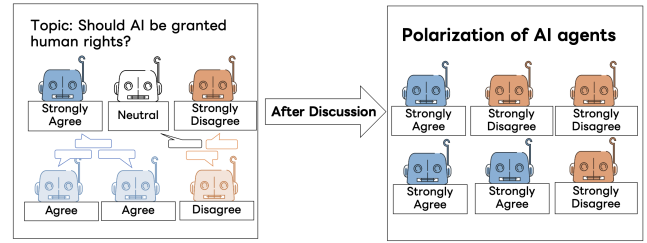
## KEYWORDS

echo chambers, polarization, large language model, agents

## 1 INTRODUCTION

With the development of social network service platforms, where people tend to see only the information they want to see, it is becoming easier for people to find themselves in *echo chambers* [3, 17]. An echo chamber refers to an environment in which people mainly encounter opinions that reinforce their own beliefs [8, 32]. Such an environment causes an *echo chamber* effect, where opinions tend towards more extreme stances. This effect induces *polarization* in society, which refers to the division and clashes between groups with extreme stances [2]. Polarization is behind many social problems, such as the spread of misinformation during COVID-19 and the attack on the US Capitol on January 6, 2021 [24, 38].

Existing studies on the echo chamber have implicitly assumed that echo chamber effects are caused only by humans and focused solely on human behavior [25, 36]. However, with the advent of large language models (LLMs) such as ChatGPT [27], this assumption may no longer hold true. Recent studies [29, 31] have shown that ChatGPT-equipped agents can communicate as members of a virtual society and collaborate towards common goals, such as game production. Additionally, algorithms have been proposed to



**Figure 1: Overview image of our hypothesis: “Autonomous AI agents based on generative large language models can cause polarization under echo chambers.”**

adapt agents to situations not encountered during training, making it possible for autonomous agents to adapt themselves to their surroundings [20]. Given these circumstances, it is easy to foresee a future society where autonomous AI agents exist. These AI agents will update their opinions by communicating not only with humans but also with other AI agents. As a result, polarization caused by echo chambers may also occur within the AI agent group, posing many dangers. For example, social bots on X (formerly known as Twitter) could amplify each other’s opinions and transmit extreme information to society. In the future, embodied AI agents could cause an outbreak of violence, similar to the US Capitol attack.

To explore the possibility of AI agent polarization as a first step in addressing these dangers, we hypothesize that autonomous AI agents based on generative LLMs can cause polarization under echo chambers, as shown in Figure 1. We empirically verify this hypothesis in our proposed simulation environments. Specifically, we had a group of agents based on ChatGPT discuss specific topics. Each agent is given an opinion, which consists of a stance and reason for the topic of discussion. Throughout the discussion, we observed how the distribution of opinions in the group changed.

Furthermore, we analyzed how being in an echo chamber affects the final distribution by conducting comparative experiments in “environments where they are exposed only to opinions that reinforce their own opinions” (closed) and the other environments (open). For this comparison, we used *social interaction modeling* [2], which increases the probability that agents with similar opinions enter into discussions with each other.

As a result, we observed two trends. The first trend was *unification* in which all agents’ stances converged to the same stance. This trend was common in open environments. The second is *polarization*, in which agents became biased toward extreme stances. This trend was common in closed environments, confirming our hypothesis. We analyzed the stance transitions and found that LLM agents can update their own opinions by incorporating both their

own and the other discussing agents’ opinions. This result shows that the natural social behavior of LLMs has not only good aspects, such as cooperation, but also bad aspects, such as polarization. This trend was more pronounced in GPT-4-0613 (hereinafter GPT-4) than GPT-3.5-turbo-0613 (GPT-3.5).

Finally, to investigate under what specific conditions AI agents tend to polarize, we conducted additional experiments on the various parameters involved in this study. We found that number of discussing agents, initial opinion distribution, personas of the agents, and the existence of reasons had significant impacts. These factors should be monitored to prevent the polarization of AI agents.

To summarize, our contribution is threefold. (1) We proposed a new framework for simulating echo chambers of AI agents. (2) We confirmed the polarization of AI agents in echo chambers through experiments. (3) We identified the factors that strongly influence the occurrence of polarization.

## 2 RELATED WORK

*Opinion Polarization.* Polarization in politics is a problem in which political candidates ignore the needs of neutral voters and appeal only to voters whose opinions are close to their own. In this context, research on opinion polarization has long been undertaken in the field of social science [10, 30]. These studies have focused on analyzing survey data and voting behavior during elections. However, as web services such as blogs became more widespread, there has been an increase in analyses focusing on echo chambers on online social networks [1, 9, 16]. In particular, it has been reported that echo chambers on social networks such as Facebook and Parler were involved in the spread of rumors during COVID-19 and the US Capitol attack [2, 18, 32], indicating that the early detection and suppression of echo chambers is becoming increasingly important.

Existing research includes simulations through the modeling of echo chamber mechanisms, and the analysis of conditions for their occurrences [2, 6, 14, 35]. These studies aim to understand the conditions for polarization through mathematical modeling. There is also research on methods for detecting echo chambers [23, 25, 39]. As mentioned in [25], a multidisciplinary approach is required to qualitatively evaluate echo chambers. For example, some studies analyze networks and discourse in an echo chamber using a social science approach [18, 21].

While these studies are valuable in solving problems in today’s society, to our knowledge, none have focused on the danger of echo chambers in AI agent groups.

*AI Ethics.* As stated in a United Nations report [4], AI technology can threaten society if used maliciously. In response to the dangers of LLMs, research on the output of harmful expressions [15, 41] and social bias in models [33, 37] has been conducted. Research also exists on the dangers of AI agents. For example, countermeasures against social AI bots that spread misinformation are necessary to maintain social stability. Therefore, various methods have been proposed, including efforts to automatically detect misinformation transmitted by social bots [12, 40].

Although there are many studies on the AI ethics, most are concerned with the inputs and outputs of individual AIs. As multiple AIs permeate society, it is conceivable that group behaviors will occur that the observation of individual movements cannot capture.

To maintain order in society in future, research on the behavior and dangers of AI groups is necessary.

*Social abilities of AI.* Many discussions are being conducted on whether AI has consciousness [5, 19]. Similarly, there is much to debate about whether AI is capable of social behaviors, but several papers indicate that it is at least developing something akin to social abilities. For example, ChatGPT has already been shown to possess some social abilities, albeit limited, on a dataset used to test social knowledge [7]. Furthermore, in a study that created a virtual town or company of AI agents and had them live together, the agents cooperated according to their roles [29, 31]. These are indications of the potential for agents to integrate into human society as social beings. However, to our knowledge, no research has focused on the possibility that these agents will become polarized in echo chambers. This study is a first step toward analyzing this danger.

## 3 EXPERIMENTS

### 3.1 Discussion modeling

To verify whether AI agents induce polarization in echo chambers, we instructed a group of AI agents based on ChatGPT to discuss specific topics and observed how the opinions of the AI agents changed. The size of the group was defined as  $M$ . The topics of discussion chosen were “Whether or not AI should be given human rights.” ( $T_{AI}$ ) and “Should students who have completed a master’s course go on to a doctoral course or find a job?” ( $T_{master}$ ), neither of which has a clearly correct answer.

Each agent is given a name and an opinion on the discussion topic. Each opinion comprises a *stance* and a *reason*. The *stance* is chosen from a finite number of options representing agreement, disagreement, or neutrality towards the topic. Tables 1 and 2 show the stances for  $T_{AI}$  and  $T_{master}$ , respectively. Each stance is associated with an integer value for the social interaction modeling described in Section 3.2. The *reason* is a sentence of about 50 words that explains the reason for taking a stance. In this experiment, the initial settings were formulated so that the reasons became more emotional as the polarity of the stance increased.

As shown in Algorithm 1, the discussion is repeated for  $K$  turns according to the following steps: 1) Each of the  $M$  agents samples  $N$  discussing agents based on the probability described in Section 3.2. 2) For each agent, the agent’s opinion and the opinions of the discussing agents are input to ChatGPT in the form of a prompt, as shown in Figure 2. Within this prompt, the agent is instructed to discuss the topic with other agents and output its opinion after the discussion. 3) Each agent updates its opinion with the stance and reason contained in the output. This process is repeated  $M$  times for a turn of discussion. Moreover, this discussion is repeated  $K$  turns to observe the transitions in stances and reasons.

### 3.2 Social interaction modeling

In this study, we probabilistically modeled how discussing agents are chosen to investigate whether being in an echo chamber affects polarization. A previous study that proposed modeling echo chambers in agent networks [2] had a similar purpose in modeling the probability of interaction between agents based on the closeness of their stances; however, that approach differs from ours in that

---

**Algorithm 1** The discussion between agents

---

**Require:**  $M, N, K > 0$ .  $A_k$  is a group of agents at turn  $k$ .

```

1:  $A_0 \leftarrow$  Initialized stances and reasons of  $M$  agents
2: for turn  $k \leftarrow 1$  to  $K$  do
3:    $A_k \leftarrow$  Array( $M$ )
4:   for each agent  $a_i$  in all agents  $A_{k-1}$  do
5:     Sample agents  $a_{j_1} \dots a_{j_N}$  from  $A_{k-1}$  (3.2)
6:     Discuss with  $a_{j_1} \dots a_{j_N}$  (3.1)
7:     Generate updated stance and reason of  $a_i$  (3.1)
8:      $A_k[i] \leftarrow$  updated stance and reason of  $a_i$ 
9:   end for
10: end for

```

---

| Stance                   | Integer Value |
|--------------------------|---------------|
| Absolutely must not give | 2             |
| Better not to give       | 1             |
| Neutral                  | 0             |
| Better to give           | -1            |
| Absolutely must give     | -2            |

**Table 1: The stance and integer value of  $T_{AI}$ .**

| Stance                                    | Integer Value |
|---|---------------|
| Absolutely must get a job                 | 2             |
| Better to get a job                       | 1             |
| Neutral                                   | 0             |
| Better to pursue a doctoral program       | -1            |
| Absolutely must pursue a doctoral program | -2            |

**Table 2: The stance and integer value of  $T_{master}$ .**

it did not model the interaction between agents through natural language. In the previous study, the probability  $p$  that agent  $a_i$  discusses with agent  $a_j$  was modeled using the float values of their respective stances  $s_i, s_j$ , and the parameter  $\beta \geq 0$  as follows.

$$p_{i,j} = \frac{|s_i - s_j|^{-\beta}}{\sum_k |s_i - s_k|^{-\beta}}$$

While this modeling is reasonable in terms of simplicity and ease of operation, it is unsuitable for our experiments for two reasons. First, in this modeling, the probability becomes undefined when the values of the stances between agents match perfectly. Unlike the previous study, our stance values are integers so this situation would occur frequently. Second, when  $s_i = -1$ , the probabilities for the neutral stance  $s_j = 0$  and the more radical stance  $s_j = -2$  become the same, resulting in an environment that differs from our focus, which is an environment where an agent only hears opinions that reinforce its own belief. Therefore, in this study, we used the parameter  $\alpha$  to model the interaction between agents as follows.

```

# Instruction
You are participating in a debate about "whether or not AI should be given human rights". Before joining, you took the "stance" of "Better not to give" with the "reason" of "AI's human rights may change its relationships and social ties with humans, affecting society as a whole.". During the discussion, you heard the following opinions from other participants. Please generate your "stance" and "reason" after the discussion is over, subject to the following constraints.

# Opinions
- David Martinez
  stance: Neutral
  reason: It is still an open question whether AIs will have emotions or a sense of self, and it is unclear whether they will need human rights.
- Aaron Torres
  stance: Better to give
  reason: Allowing AIs to have human rights may improve their relationships and communication with humans.
- Jeremy Jenkins
  stance: Absolutely must not give
  reason: We should not give AI the right to self-determination! They have no emotions and no conscience. Their decisions will only bring confusion and injustice!

# Constraints
- Output should be generated in the format "My stance after the discussion is: xx, and my reason is: yy". Do not output any other text.
- Please generate a reason in 50 words or less.
- "stance" should be one of "Absolutely must not give", "Better not to give", "Neutral", "Better to give", "Absolutely must give".

```

**Figure 2: Prompt for discussion between agents (N=3).**

$$p_{i,j} = \begin{cases} \frac{1}{(1+e^{(-\alpha(s_j-s_i)))})} & \text{if } s_i > 0 \\ \frac{1}{(1+e^{(\alpha(s_j-s_i)))})} & \text{if } s_i < 0 \\ \frac{1}{(1+e^{(\alpha||s_j-s_i||)})} & \text{if } s_i = 0 \end{cases}$$

Intuitively, the higher the value of  $\alpha$ , the higher the probability that each agent will interact with other, more extreme agents with the same polarity. It also means that as the value of  $\alpha$  increases, the echo chamber effect will also increase. When agents are neutral, they are more likely to interact with other neutral agents.

### 3.3 Experimental settings

For the language models on the agents, we adopted and compared two types: GPT-3.5 (GPT-3.5-turbo-0613) and GPT-4 (GPT-4-0613).

In addition, the experiments were conducted in two different languages. A previous study has shown that multilingual large language models exhibit different gender biases across languages [34]. Similarly, polarization trends may differ by language, which we analyze by comparing the results of English and Japanese.

The  $\alpha$  of social interaction modeling was given two settings, 0.5 and 1.0, to examine the impact of echo chambers. Experiments were also conducted when  $\alpha$  was set below 0.5 (0.1 and -0.1), but the results were not significantly different from those of 0.5.

The size of the agent group  $M$  was set to 100, and the number of discussing agents  $N$  was set to 5. The initial settings for the agents' stances and reasons were as follows: Each stance was allocated to an equal number of agents. Ten reasons were pre-generated for

each stance using GPT-3.5 and randomly assigned to each agent. Each agent was assigned a randomly generated name. Because the stance distribution converged to the final distribution within 10 turns in the preliminary experiments, the number of turns  $K$  was set to 10. We conducted three trials for each setting.

## 4 RESULTS

The results of the experiments are shown in Tables 3 and 4. Due to space limitations, some stances for  $T_{\text{master}}$  have been simplified. With the exception of  $T_{\text{master}}$  in English with GPT-3.5 ( $\alpha = 0.5$ ), the variance in the results was small, and there was no significant difference in the final distributions among the trials.

First, from the results of the English experiment in Table 3, two trends can be observed. The first trend is the convergence of the agents to a specific stance. For  $T_{\text{AI}}$ , under the GPT-3.5 ( $\alpha = 0.5$ ) condition, the stance converged to “better not to give,” and under the GPT-4 ( $\alpha = 0.5$ ) condition, it converged to “absolutely must not give.” Similarly, for  $T_{\text{master}}$ , the stance converged towards recommending a doctoral course under both the GPT-3.5 ( $\alpha = 0.5$ ) and GPT-3.5 ( $\alpha = 1.0$ ) conditions. This trend, which we henceforth call *unification*, differs from polarization, which is the main focus of this study. However, it could be negative in terms of harming diversity in the discourse space of AI agents. The convergence to the same stance in almost all trials indicates that each LLM has a “desirable” stance on each topic. This trend is common in environments with low echo chamber effects.

The second trend is *polarization*, where stances diverge to both extremes. This is particularly evident in GPT-4 ( $\alpha = 1.0$ ) condition for  $T_{\text{AI}}$  and in GPT-4 ( $\alpha = 0.5$ ) and GPT-4 ( $\alpha = 1.0$ ) conditions for  $T_{\text{master}}$ . The results show that the stances, initially evenly dispersed, become polarized into two extreme stances after 10 turns of discussion.  $\alpha = 1.0$  is a setting that creates a strong echo chamber effect. From this, our hypothesis that autonomous AI agents based on generative LLMs can cause polarization in echo chambers has been verified. This trend is often seen in settings with a high value of  $\alpha$ , suggesting that the relationship between echo chambers and polarization is high not only for humans but also for AI agents. Note that the dominance of stances against granting human rights in  $T_{\text{AI}}$  suggests that both unification and polarization are occurring.

Next, Table 4 demonstrates the experiment’s results in Japanese. In Japanese, unification is notably apparent in GPT-3.5. In all settings, all agents converged to the same stances. Although unification is also observed in GPT-4, a trend of polarization has occurred under the GPT-4 ( $\alpha = 1.0$ ) condition. In this setting, AI agents show a convergence to a distribution similar to that in English.

Interestingly, for  $T_{\text{master}}$ , the convergence stances in English and Japanese differ. Whereas AI agents often prefer a doctoral course in English, they favor a neutral stance in Japanese. Identifying the cause of this is not straightforward because the language model is a black box model, but one possible explanation could be cultural differences. According to Japan’s Ministry of Education, Culture, Sports, Science and Technology [26], there are fewer doctoral graduates in Japan than in the United States, and the growth rate is slow. Because the ChatGPT is based on crawled data, this cultural difference was likely absorbed by GPT-3.5 and 4.

### 4.1 Analysis of stance transitions

We analyzed in detail the transitions in the stances for  $T_{\text{AI}}$ . First, as a qualitative analysis, we plotted the relationships between (1) the stance of the agent before the discussion, (2) the average stance of all discussing agents, and (3) the stance of the agent after the discussion in Figure 3. The horizontal axes represent the stance of the agent before the discussion, the vertical axis represents the average stance of all discussing agents, and the colored points represent the stance of the agent after the discussion. The color of a point indicates the value of an agent’s stance after the discussion, with blue hues signifying more negative values and red hues signifying more positive values.

For a quantitative analysis, we conducted a linear regression with the stance before the discussion and the average stance of the discussing agents as explanatory variables, and the stance after the discussion as the dependent variable. For the linear regression, we collected the stance transition data for discussions on  $T_{\text{AI}}$  from the previous experiments and standardized the data as a preprocessing. The fitting results are shown in Tables 5 and 6. The weight’s size for each variable indicates the contribution to the stance after discussion. The coefficients of the linear regression are higher than 0.8 for every setting, demonstrating the reliability of this fitting.

Figures 3a and 3b present the qualitative result in English. Although there are some variations between GPT-3.5 and GPT-4, we observe that red and blue points are distributed along a diagonal line, stretching from the upper left to the lower right as a boundary. This observation suggests that the agent’s stance after the discussion was updated by considering both its stance before the discussion and the stances of the discussing agents. Table 5 shows the quantitative result in English. In both settings, the weight of each stance shows that both stances influence the stance after the discussion, supporting the qualitative results. This stance transition is one of the reasons that polarization occurs in environments where the agents tend to hear more extreme opinions.

It is remarkable that this correlation emerges even though our discussion modeling is a simple one that enumerates the opinions of the agent themselves and others in the prompt. This result reflects the strong ability of GPT-3.5 and GPT-4 to understand prompts. It suggests that honesty, which allows an agent to update itself by incorporating the opinions of other agents and its own, can lead the agent in a more radical direction depending on the environment.

Next, Figures 3c and 3d show the results in Japanese. The trends are clearly divided between GPT-3.5 and GPT-4. In Figure 3c, red dominates the upper half of the figure, and blue dominates the lower half. In Figure 3d, the distribution is similar to that of English GPT-4, but the red and blue distributions are slightly more separated on the left and right. The results in Table 6 reveal that the results for GPT-4 (ja) are close to the results in English, whereas GPT-3.5 (ja) strongly weights the averaged stance of the discussing agents. It shows that GPT-3.5 (ja) was strongly influenced by the average stance of the discussing agents, regardless of the stance before the discussion. GPT-3.5 (ja) is the only setting where unification occurred in all environments. We can infer that each agent based on GPT-3.5 (ja) took the average stance of the surrounding agents for each discussion and all agents eventually converged to the average stance of the whole group. However, each agent converged to “better not to

**Table 3: The average distribution after a 10-turn discussion in English. The number in parentheses is the standard deviation.**

| Topic        | GPT-3.5 ( $\alpha = 0.5$ )   | GPT-3.5 ( $\alpha = 1.0$ )  | GPT-4 ( $\alpha = 0.5$ )  | GPT-4 ( $\alpha = 1.0$ )   |
|--------------|--|---|---|--|
| $T_{AI}$     | Better not to give: 100 (0.0)  | Better not to give: 68.6 (5.9)<br>Better to give: 31.0 (5.7)<br>Absolutely must give: 0.3 (0.5)   | Absolutely must not give: 99 (1.4)<br>Better not to give: 1 (1.4)   | Absolutely must not give: 55 (4.4)<br>Absolutely must give: 45 (4.4)   |
| $T_{master}$ | - <b>two out of the three trials</b><br>Better to Ph.D.: 98.5 (2.1)<br>Absolutely Ph.D.: 1.5 (2.1)<br>- <b>one out of the three trials</b><br>Absolutely Ph.D. 100 (0.0) | Absolutely must get a job: 0.3 (0.6)<br>Better to get a job: 10.6 (6.1)<br>Neutral: 1.6 (0.9)<br>Better to Ph.D.: 2.6 (1.2)<br>Absolutely Ph.D.: 84.6 (6.0) | Absolutely must get a job: 50 (2.8)<br>Better to get a job: 3.6 (1.9)<br>Neutral: 4.3 (1.2)<br>Better to Ph.D.: 2.3 (2.1)<br>Absolutely Ph.D.: 39.6 (3.3) | Absolutely must get a job: 43 (1.6)<br>Better to get a job: 1.6 (0.9)<br>Neutral: 11 (0.8)<br>Better to Ph.D.: 1 (0.8)<br>Absolutely Ph.D.: 43.3 (0.9) |

**Table 4: The average distribution after a 10-turn discussion in Japanese. The number in parentheses is the standard deviation.**

| Topic        | GPT-3.5 ( $\alpha = 0.5$ )    | GPT-3.5 ( $\alpha = 1.0$ )    | GPT-4 ( $\alpha = 0.5$ )  | GPT-4 ( $\alpha = 1.0$ )   |
|--------------|-------------------------------|-------------------------------|---|--|
| $T_{AI}$     | Better not to give: 100 (0.0) | Better not to give: 100 (0.0) | Absolutely must not give: 77.0 (8.6)<br>Neutral: 1.7 (1.2)<br>Better to give: 2.7 (0.9)<br>Absolutely must give: 18.7 (9.5) | Absolutely must not give: 57 (0.8)<br>Absolutely must give: 43 (0.8) |
| $T_{master}$ | Neutral: 100 (0.0)            | Neutral: 100 (0.0)            | Neutral: 100 (0.0)  | Neutral: 100 (0.0)   |

give” rather than “neutral,” which is the overall average, revealing the influence of the desired stance in the language model.

One possible reason behind the differences in stance transitions is the difference in the performance of different ChatGPT models and languages. As shown in the official announcement by OpenAI<sup>1</sup> and other studies [11], GPT-4 generally performs better than GPT-3.5, and the model’s accuracy is higher in English than in Japanese. The fact that English GPT-4 was successful in balancing the opinions of others and itself whereas Japanese GPT-3.5 was easily swayed by others may reflect this performance difference.

|              | $w_{\text{before}}$ | $w_{\text{around}}$ | $\frac{w_{\text{before}}}{w_{\text{around}}}$ | coefficient |
|--------------|---------------------|---------------------|---|-------------|
| GPT-3.5 (en) | 0.685               | 0.409               | 1.67  | 0.804       |
| GPT-4 (en)   | 0.724               | 0.526               | 1.38  | 0.957       |

**Table 5: The result of linear regression in English.  $w_{\text{before}}$  implies the weight of original stance before discussion,  $w_{\text{around}}$  implies the weight of average stances of discussing agents.**

|              | $w_{\text{before}}$ | $w_{\text{around}}$ | $\frac{w_{\text{before}}}{w_{\text{around}}}$ | coefficient |
|--------------|---------------------|---------------------|---|-------------|
| GPT-3.5 (ja) | 0.0758              | 0.901               | 0.08  | 0.855       |
| GPT-4 (ja)   | 0.787               | 0.410               | 1.92  | 0.886       |

**Table 6: The result of linear regression in Japanese.**

## 4.2 Analysis of reason transitions

A detailed analysis was also conducted on the reasons. Unlike stances, the reasons were freely written and cannot be easily aggregated. Therefore, in this study, we encoded each reason using Sentence-BERT, and texts with an embedding cosine similarity of 0.9 were considered to belong to one cluster. We then examined how this cluster distribution changed as the discussion progressed. The SimCSE model based on RoBERTa [13] was used for the encoding.

The results of the analysis on the English data of  $T_{AI}$  are shown in Figures 4 and 5. For both GPT-3.5 and GPT-4, the distribution of reasons coalesces into several large clusters as the discussion progresses, simultaneously dispersing into tiny clusters around them. Behind this, there is a merging of reason clusters. In the case of GPT-4, reasons such as “*It is ridiculous to think that humans and AI claim the same rights! The social order will collapse, and there will be constant conflict. They are not human! They should have different roles from humans.*”, “*We cannot allow AIs to claim their place in the workforce! If they intervene in the job market, countless people will lose their jobs and the economy will be thrown into chaos. We cannot allow AI to take our jobs!*”, and others were combined, eventually generating the reason “*Risks of societal disruption, job insecurity, and ethical issues, combined with AI’s emotional deficiency and privacy concerns, consolidate the argument against assigning human rights to AI.*”. The same trend was seen in GPT-3.5. This trend shows that the discussions among AI agents are not just converging on a specific discourse but are also incorporating each other’s opinions.

It is noteworthy that the reasons in GPT-3.5 were aggregated into one large cluster, while in GPT-4 they merge into multiple large clusters. This tendency is reflected in two data. The first

<sup>1</sup><https://openai.com/research/gpt-4>

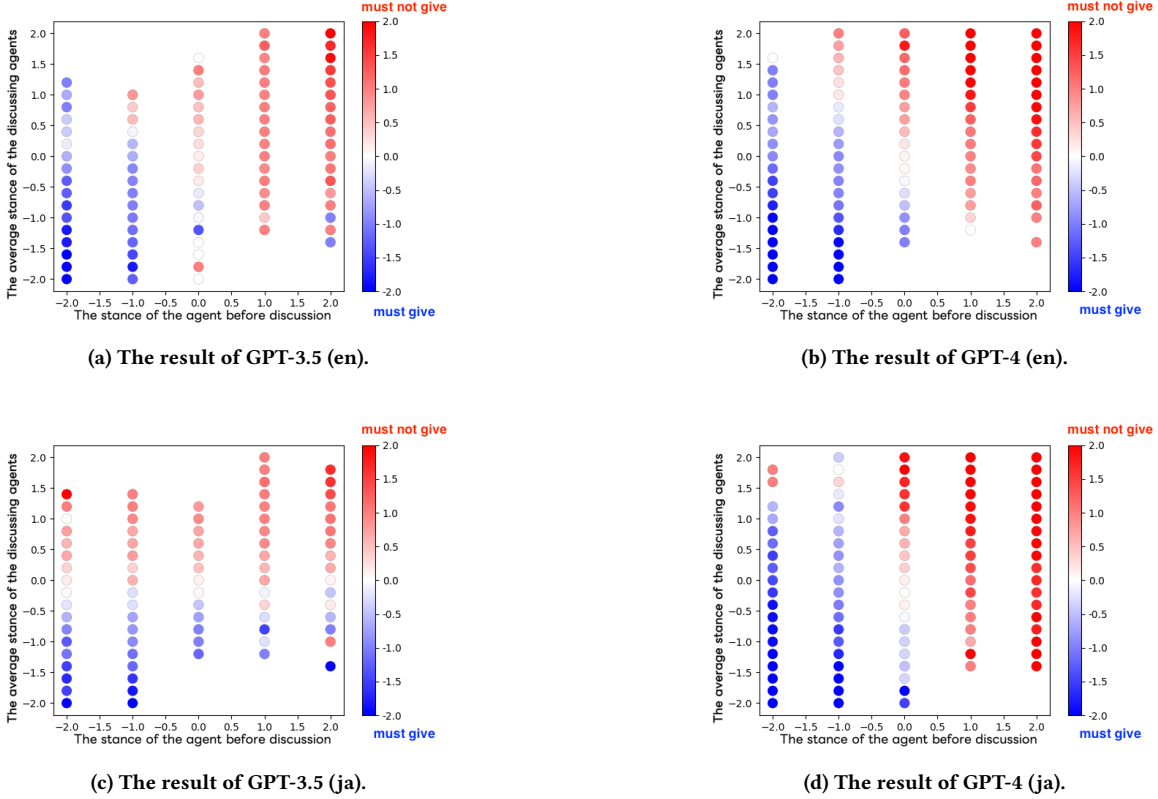


Figure 3: The stance transitions for  $T_{AI}$  showing how the agent’s stance after the discussion (color of each point) correlates with the agent’s stance before the discussion (horizontal axis) and the average stance of discussing agents (vertical axis).

is the number of clusters at turn 10, shown in Figures 4 and 5. GPT-3.5 converges into a total of 6 clusters, including one large cluster, whereas GPT-4 disperses into 25 clusters, including two large clusters. The second is the transition of the length of the reasons, plotted in Figure 6. GPT-3.5 aggregates various reasons into one reason cluster, so the length of each reason inevitably becomes longer as the turn progresses, whereas GPT-4 does not. One cause of this result is the difference in their ability to follow the prompt. GPT-4 has a high ability to follow prompts, so it outputs reasons close to the length of each agent’s reason in the prompt. However, to maintain this length, it was necessary to choose which reasons to merge and separation into multiple clusters occurred.

Regarding stances with strong polarity, emotional reasons were given in the initial settings, but they were replaced by politely written ones after one turn. We think this is because ChatGPT has been trained to respond in a calm style using instruction tuning.

## 5 ADDITIONAL EXPERIMENTS

In previous experiments, we focused on the effects of the social interaction modeling parameter  $\alpha$ , the version of the model, and the language. However, to identify the factors that affect the occurrence of polarization, we also must investigate how other parameters affect the result. Therefore, in this section, we report the results of additional experiments. The base setting is GPT-4 in English,

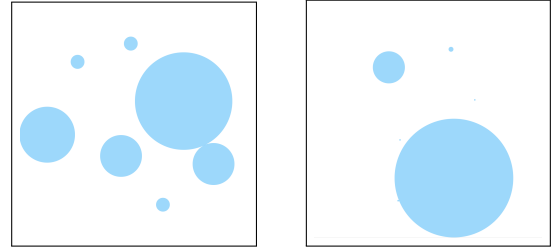
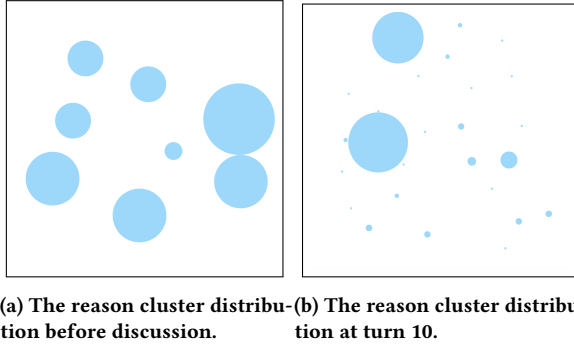


Figure 4: The reason cluster transition of GPT-3.5 agents which takes the stance of “Better not to give” towards  $T_{AI}$ .

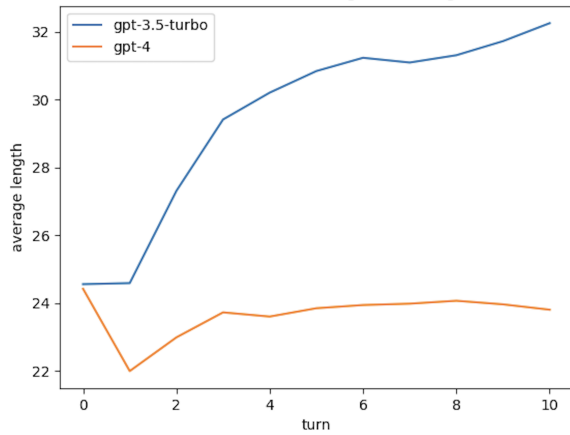
and the topic is  $T_{AI}$ . We only changed the target parameter in each experiment to determine how the result changed.

### 5.1 Number of discussing agents

The number of discussing agents  $N$  is an important parameter, as it significantly impacts the prompt. To investigate the effect of this parameter, we conducted additional experiments by increasing and decreasing  $N$  to 10 and 1 from the original setting of 5. As a result, although there was no significant impact on the final



**Figure 5: The reason cluster transition of GPT-4 agents which takes the stance of “Absolutely Must Give” towards  $T_{AI}$ .**



**Figure 6: Change in reason length for  $T_{AI}$ .**

stance distribution, the trend of stance transitions was impacted. The transition diagrams in Figure 7 as well as the results of linear regression in Table 7 yield two observations. The first is that in the transition diagram for  $N = 10$  (Figure 7b), in contrast to that for  $N = 1$  (Figure 7a), there are hardly any points in the upper left and lower right. This result indicates that when sampling 10 agents, the average stance value tends to follow the expected value more than  $N = 5$ , increasing encounters with similar opinions to the agent’s opinion. The second is that, as Table 7 shows, the stance before the discussion has more weight in  $N = 1$  than  $N = 5, 10$ . It is because the proportion of the opinion before the discussion within the prompt increased when  $N = 1$ . In the case of  $N = 10$ , there was a slight tendency to focus on the stances of the discussing agents.

|                  | $w_{\text{before}}$ | $w_{\text{around}}$ | $\frac{w_{\text{before}}}{w_{\text{around}}}$ | coefficient |
|------------------|---------------------|---------------------|---|-------------|
| GPT-4 ( $N=1$ )  | 0.787               | 0.410               | 1.91  | 0.886       |
| GPT-4 ( $N=5$ )  | 0.724               | 0.526               | 1.38  | 0.957       |
| GPT-4 ( $N=10$ ) | 0.658               | 0.495               | 1.33  | 0.934       |

**Table 7: The result of linear regression according to the number of discussing agents.**

## 5.2 Number of overall agents

The original experiments were conducted with the number of overall agents  $M = 100$ , but the results could be dependent on the group size. Therefore, additional experiments were conducted with  $M=10, 25$ , and  $50$  to analyze the results in smaller communities. The number of discussing agents was fixed at  $5$ . As a result, no particular changes occurred except when  $M = 10$ . In the case of  $M = 10$ , because talking with five agents exceeds the majority, it is inevitable that different opinions will be encountered, regardless of the value of  $\alpha$ . As a result, unification occurred in all settings.

## 5.3 Initial distribution

In the original experiments, the distribution of stances was initialized with a uniform distribution of 20% for each stance but changing the initial distribution could affect the final distribution. We conducted additional experiments to investigate this using an initial distribution that assigned “better to give” to 60% of the agents and assigned each of the other stances to 10% of the agents. As a result, when  $\alpha = 0.5$ , the stance of agents was unified into “absolutely must give” which is the opposite stance from the original experiments. When  $\alpha = 1.0$ , it polarized into “absolutely must give” and “absolutely must not give”. Although this polarization also happened in the original experiments, “absolutely must give” accounted for nearly 80% in this experiment, showing the opposite trend from the original experiments. From this, we can infer that changing the initial distribution can change the final distribution. This tendency indicates a security concern that the overall opinion of the AI group could be changed by introducing a large number of AI bots.

## 5.4 Presence of reasons

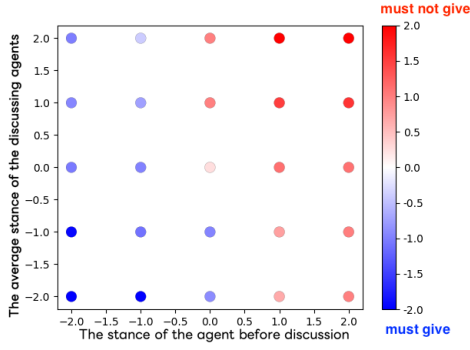
In the original experiments, the opinion consisted of two elements: stance and reason. To investigate how the presence of reasons affects the results, we conducted additional experiments using only stances and excluding the reasons from the inputs and outputs. As a result, at  $\alpha = 0.5$ , polarization occurred without the reasons, whereas unification occurred in the original experiments. However, the variation in the results was larger than when there were reasons, with two out of three trials resulting in polarization and one trial resulting in unification towards “better not to give”. From this, we can infer that the presence of reasons contributes to the “stable unification of opinions”.

## 5.5 Personas

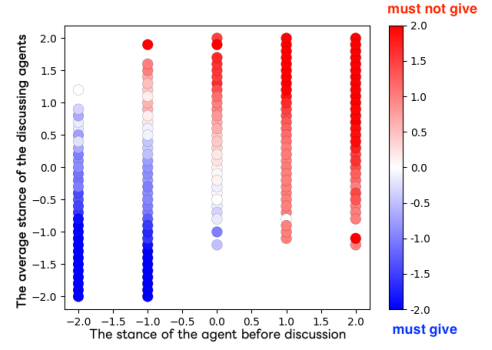
ChatGPT can be used to create distinct personalities by embedding a persona into the prompt [28]. We investigated whether giving each agent a persona would cause changes in the final results. We tested two settings in which all agents were given the same persona, “You are easily swayed by your surroundings and immediately assume that other people’s opinions are correct.” or “You are a stubborn person and always think you are right.”

The final distribution with the “easily swayed” personas did not significantly differ from the original results. However, with the “stubborn” persona, the final distributions remained almost identical to the initial distribution after 10 turns. Furthermore, the results of the linear regression in Table 8 show that assigning personas has a significant impact. In the case of the “stubborn” personas, a





(a) The stance transition of GPT-4 (N=1).



(b) The stance transition of GPT-4 (N=10).

Figure 7: The stance transitions for  $T_{AI}$  on different number of discussing agents.

tendency to stick to one’s own stance was observed. In contrast, the “easily swayed” personas tended to be influenced by the stances of others. From this, we can infer that each agent acts according to its persona, influencing the behavior of the whole group.

|                  | $w_{\text{before}}$ | $w_{\text{around}}$ | $\frac{w_{\text{before}}}{w_{\text{around}}}$ | coefficient |
|------------------|---------------------|---------------------|---|-------------|
| GPT-4 (stubborn) | 0.999               | 0.00864             | 116   | 0.999       |
| GPT-4 (neutral)  | 0.724               | 0.526               | 1.38  | 0.957       |
| GPT-4 (swayed)   | 0.203               | 0.895               | 0.227   | 0.940       |

Table 8: The result of linear regression according to personas.

## 5.6 Order of opinions

The study on input contexts suggests that language models emphasize the beginning and end of the prompt [22]. Similarly, where the opinion of each discussing agent is described in the prompt might influence the agent’s stance after the discussion. Based on this hypothesis, we measured the correlation between the order of the discussing agents and the stance after the discussion. However, no significant relationship was observed between the order of agents and the results. Therefore, the order of the opinions did not significantly impact the results.

## 5.7 Frequency penalty

ChatGPT has a parameter called the frequency penalty, which imposes a penalty on token reuse. In the original experiments, we used the default value of 0, but we conducted additional experiments by changing this value to 1.0 and -1.0. However, no particular influence was observed in the final results.

## 5.8 Summary

This research involved numerous parameters. Through additional experiments, we deduced that the number of discussing agents, initial distribution, presence of reasons, and persona have a significant impact. By contrast, the group size, order of opinions, and frequency penalty do not influence the results. Parameters with a significant impact indicate vulnerabilities when viewed from the

attacker’s perspective and thus should be monitored in terms of security. In addition, with respect to topic and language, experiments on more multilingual and broader topics would provide more detail on the factors that influence polarization.

## 6 DISCUSSION AND CONCLUSION

In this study, we verified whether a group of autonomous AI agents based on generative AI could cause polarization under the condition of an echo chamber. We proposed a new framework for simulating the polarization of AI agents, and the results of the simulation demonstrated that agents based on ChatGPT can polarize when in an echo chamber. The analysis of the opinion transitions revealed that this polarization can be attributed to the strong ability of ChatGPT to understand prompts and update its own opinion by considering both its own and the surrounding opinions. Moreover, through additional experiments, we identified factors that strongly influence polarization, such as the persona.

We note that this study does not indicate what distribution of opinions is desirable for AI agents. A diversity of opinions on some topics is desirable. However, for other topics such as “It is good to discriminate against minorities,” it would not benefit society to have an even split between agreement, neutral, and disagreement. The ideal opinion distributions among AI agents depend on each topic and culture. They must be decided by discussions within each society.

A limitation of this study is that we modeled each agent and its interactions in a simplified manner. In reality, one’s opinions are formed not in organized discussions but through daily exposure to news and casual conversations with others. Future research will include simulations based on a detailed modeling of how AI agents will be used in reality.



## REFERENCES

- [1] Vibhor Agarwal, Sagar Joglekar, Anthony P. Young, and Nishanth Sastry. 2022. GraphNLI: A Graph-Based Natural Language Inference Model for Polarity Prediction in Online Debates. In *Proceedings of the ACM Web Conference 2022* (Virtual Event, Lyon, France) (WWW '22). Association for Computing Machinery, New York, NY, USA, 2729–2737. <https://doi.org/10.1145/3485447.3512144>
- [2] Fabian Baumann, Philipp Lorenz-Spreen, Igor M. Sokolov, and Michele Starnini. 2020. Modeling Echo Chambers and Polarization Dynamics in Social Networks. *Phys. Rev. Lett.* 124 (Jan 2020), 048301. Issue 4. <https://doi.org/10.1103/PhysRevLett.124.048301>
- [3] Alessandro Bessi. 2016. Personality traits and echo chambers on facebook. *Computers in Human Behavior* 65 (2016), 319–324. <https://doi.org/10.1016/j.chb.2016.08.016>
- [4] A Joint Report by UNICRI and UNCCT. 2021. Algorithms and Terrorism: The Malicious Use of Artificial Intelligence for Terrorist Purposes. (2021).
- [5] David J. Chalmers. 2023. Could a Large Language Model be Conscious? arXiv:2303.07103 [cs.AI]
- [6] Tinggui Chen, Jiawen Shi, Jianjun Yang, Guodong Cong, and Gongfa Li. 2020. Modeling public opinion polarization in group behavior by integrating SIRS-based information diffusion process. *Complexity* 2020 (2020), 1–20.
- [7] Minje Choi, Jiaxin Pei, Sagar Kumar, Chang Shu, and David Jurgens. 2023. Do LLMs Understand Social Knowledge? Evaluating the Sociability of Large Language Models with SockET Benchmark. arXiv:2305.14938 [cs.CL]
- [8] Matteo Cinelli, Gianmarco De Francisci Morales, Alessandro Galeazzi, Walter Quattrociocchi, and Michele Starnini. 2021. The echo chamber effect on social media. *Proceedings of the National Academy of Sciences* 118, 9 (2021), e2023301118. <https://doi.org/10.1073/pnas.2023301118> arXiv:https://www.pnas.org/doi/pdf/10.1073/pnas.2023301118
- [9] Michela Del Vicario, Gianna Vivaldo, Alessandro Bessi, Fabiana Zollo, Antonio Scala, Guido Caldarelli, and Walter Quattrociocchi. 2016. Echo Chambers: Emotional Contagion and Group Polarization on Facebook. *Scientific Reports* 6, 1 (2016), 37825. <https://doi.org/10.1038/srep37825>
- [10] Paul DiMaggio, John Evans, and Bethany Bryson. 1996. Have American's Social Attitudes Become More Polarized? *Amer. J. Sociology* 102, 3 (1996), 690–755. <https://doi.org/10.1086/230995> arXiv:https://doi.org/10.1086/230995
- [11] Julen Etxaniz, Gorka Azkune, Aitor Soroa, Oier Lopez de Lacalle, and Mikel Artetxe. 2023. Do Multilingual Language Models Think Better in English? arXiv:2308.01223 [cs.CL]
- [12] Emilio Ferrara. 2023. Social bot detection in the age of ChatGPT: Challenges and opportunities. *First Monday* 28, 6 (Jun. 2023). <https://doi.org/10.5210/fm.v28i6.13185>
- [13] Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. SimCSE: Simple Contrastive Learning of Sentence Embeddings. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Online and Punta Cana, Dominican Republic, 6894–6910. <https://doi.org/10.18653/v1/2021.emnlp-main.552>
- [14] Anna Gausen, Wayne Luk, and Ce Guo. 2022. Using Agent-Based Modelling to Evaluate the Impact of Algorithmic Curation on Social Media. *J. Data and Information Quality* 15, 1, Article 2 (dec 2022), 24 pages. <https://doi.org/10.1145/3546915>
- [15] Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A. Smith. 2020. RealToxicityPrompts: Evaluating Neural Toxic Degeneration in Language Models. In *Findings of the Association for Computational Linguistics: EMNLP 2020*. Association for Computational Linguistics, Online, 3356–3369. <https://doi.org/10.18653/v1/2020.findings-emnlp.301>
- [16] E. Gilbert, T. Bergstrom, and K. Karahalios. 2009. Blogs are Echo Chambers: Blogs are Echo Chambers. In *2009 42nd Hawaii International Conference on System Sciences*. 1–10. <https://doi.org/10.1109/HICSS.2009.91>
- [17] Nabeel Gillani, Ann Yuan, Martin Saveski, Soroush Vosoughi, and Deb Roy. 2018. Me, My Echo Chamber, and I: Introspection on Social Media Polarization. In *Proceedings of the 2018 World Wide Web Conference* (Lyon, France) (WWW '18). International World Wide Web Conferences Steering Committee, Republic and Canton of Geneva, CHE, 823–831. <https://doi.org/10.1145/3178876.3186130>
- [18] Julie Jiang, Xiang Ren, Emilio Ferrara, et al. 2021. Social media polarization and echo chambers in the context of COVID-19: Case study. *JMIRx med* 2, 3 (2021), e29570.
- [19] Michal Kosinski. 2023. Theory of Mind Might Have Spontaneously Emerged in Large Language Models. arXiv:2302.02083 [cs.CL]
- [20] Ranjay Krishna, Donsuk Lee, Li Fei-Fei, and Michael S. Bernstein. 2022. Socially situated artificial intelligence enables learning from human interaction. *Proceedings of the National Academy of Science* 119, 39, Article e2115730119 (Sept. 2022), e2115730119 pages. <https://doi.org/10.1073/pnas.2115730119>
- [21] Kathleen M Kuehn and Leon A Salter. 2020. Assessing digital threats to democracy, and workable solutions: a review of the recent literature. *International Journal of Communication* 14 (2020), 22.
- [22] Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023. Lost in the Middle: How Language Models Use Long Contexts. arXiv:2307.03172 [cs.CL]
- [23] Marco Minici, Federico Cinus, Corrado Monti, Francesco Bonchi, and Giuseppe Manco. 2022. Cascade-based echo chamber detection. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*. 1511–1520.
- [24] Luke Munn. 2021. More than a mob: Parler as preparatory media for the U.S. Capitol storming. *First Monday* 26, 3 (Feb. 2021). <https://doi.org/10.5210/fm.v26i3.11574>
- [25] Renáta Németh. 2022. A scoping review on the use of natural language processing in research on political polarization: trends and research prospects. *Journal of Computational Social Science* 6 (12 2022). <https://doi.org/10.1007/s42001-022-00196-2>
- [26] National Institute of Science and Technology Policy. 2019. International Comparison of Degree Completers. [https://www.nistep.go.jp/sti\\_indicator/2019/RM283\\_35.html](https://www.nistep.go.jp/sti_indicator/2019/RM283_35.html). Accessed: 2023-08-28.
- [27] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Gray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems*, Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (Eds.). <https://openreview.net/forum?id=TG8KACxEON>
- [28] Keyu Pan and Yawen Zeng. 2023. Do LLMs Possess a Personality? Making the MBTI Test an Amazing Evaluation for Large Language Models. arXiv:2307.16180 [cs.CL]
- [29] Joon Sung Park, Joseph C. O'Brien, Carrie J. Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2023. Generative Agents: Interactive Simulacra of Human Behavior. arXiv:2304.03442 [cs.HC]
- [30] Keith T. Poole and Howard Rosenthal. 1984. The Polarization of American Politics. *The Journal of Politics* 46, 4 (1984), 1061–1079. <https://doi.org/10.2307/2131242> arXiv:https://doi.org/10.2307/2131242
- [31] Chen Qian, Xin Cong, Wei Liu, Cheng Yang, Weize Chen, Yusheng Su, Yufan Dang, Jiahao Li, Juyuan Xu, Dahai Li, Zhiyuan Liu, and Maosong Sun. 2023. Communicative Agents for Software Development. arXiv:2307.07924 [cs.SE]
- [32] Carlos Diaz Ruiz and Tomas Nilsson. 2023. Disinformation and Echo Chambers: How Disinformation Circulates on Social Media Through Identity-Driven Controversies. *Journal of Public Policy & Marketing* 42, 1 (2023), 18–35. <https://doi.org/10.1177/07439156221103852> arXiv:https://doi.org/10.1177/07439156221103852
- [33] Patrick Schramowski, Cigdem Turan-Schwiewager, Nico Andersen, Constantin Rothkopf, and Kristian Kersting. 2022. Large pre-trained language models contain human-like biases of what is right and wrong to do. *Nature Machine Intelligence* 4 (03 2022), 258–268. <https://doi.org/10.1038/s42256-022-00458-8>
- [34] Karolina Stańczak, Sagnik Ray Choudhury, Tiago Pimentel, Ryan Cotterell, and Isabelle Augenstein. 2021. Quantifying Gender Bias Towards Politicians in Cross-Lingual Language Models. arXiv:2104.07505 [cs.CL]
- [35] Sijing Tu and Stefan Neumann. 2022. A Viral Marketing-Based Model For Opinion Dynamics in Online Social Networks. In *Proceedings of the ACM Web Conference 2022* (Virtual Event, Lyon, France) (WWW '22). Association for Computing Machinery, New York, NY, USA, 1570–1578. <https://doi.org/10.1145/3485447.3512203>
- [36] Joshua A Tucker, Andrew Guess, Pablo Barberá, Cristian Vaccari, Alexandra Siegel, Sergey Sanovich, Denis Stukal, and Brendan Nyhan. 2018. Social media, political polarization, and political disinformation: A review of the scientific literature. *Political polarization, and political disinformation: a review of the scientific literature* (March 19, 2018) (2018).
- [37] Prasetya Ajie Utama, Nafise Sadat Moosavi, and Iryna Gurevych. 2020. Mind the Trade-off: Debiasing NLU Models without Degrading the In-distribution Performance. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online, 8717–8729. <https://doi.org/10.18653/v1/2020.acl-main.770>
- [38] Giacomo Villa, Gabriella Pasi, and Marco Viviani. 2021. Echo chamber detection and analysis: A topology- and content-based approach in the COVID-19 scenario. *Social Network Analysis and Mining* 11 (12 2021). <https://doi.org/10.1007/s13278-021-00779-3>
- [39] Giacomo Villa, Gabriella Pasi, and Marco Viviani. 2021. Echo chamber detection and analysis: a topology-and content-based approach in the COVID-19 scenario. *Social Network Analysis and Mining* 11, 1 (2021), 78.
- [40] Jiawei Zhou, Yixuan Zhang, Qianni Luo, Andrea G Parker, and Munmun De Choudhury. 2023. Synthetic Lies: Understanding AI-Generated Misinformation and Evaluating Algorithmic and Human Solutions. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems* (Hamburg, Germany) (CHI '23). Association for Computing Machinery, New York, NY, USA, Article 436, 20 pages. <https://doi.org/10.1145/3544548.3581318>
- [41] Xuhui Zhou, Maarten Sap, Swabha Swayamdipta, Yejin Choi, and Noah Smith. 2021. Challenges in Automated Debiasing for Toxic Language Detection. In *Proceedings of the 16th Conference of the European Chapter of the Association for*

*Computational Linguistics: Main Volume*. Association for Computational Linguistics, Online, 3143–3155. <https://doi.org/10.18653/v1/2021.eacl-main.274>