



# Selective agreement, not sycophancy: investigating opinion dynamics in LLM interactions

Erica Cau<sup>1,2\*</sup>, Valentina Pansanella<sup>2†</sup>, Dino Pedreschi<sup>1</sup> and Giulio Rossetti<sup>2</sup>

Handling Editor: Ronaldo Menezes

\*Correspondence:

[erica.cau@phd.unipi.it](mailto:erica.cau@phd.unipi.it)

<sup>1</sup>Department of Computer Science, University of Pisa, Largo Bruno Pontecorvo, Pisa, Italy

<sup>2</sup>Institute of Information Science and Technologies "A. Faedo" (ISTI), National Research Council (CNR), Via Giuseppe Moruzzi 1, Pisa, Italy

<sup>†</sup>Equal contributors

## Abstract

Understanding how opinions evolve is essential for addressing phenomena such as polarization, radicalization, and consensus formation. In this work, we investigate how language shapes opinion dynamics among Large Language Model (LLM) agents by simulating multi-round debates. Using our framework, we find that agent populations consistently converge toward agreement, not through sycophancy or blind conformity, but via a structured and asymmetric persuasion process. Agents are more likely to accept, and thus be persuaded by, opinions that are more agreeable relative to the discussion framing, revealing a directional bias in how opinions evolve. LLM agents selectively adopt peers' views, showing neither bounded confidence nor indiscriminate agreement. Moreover, agents frequently produce fallacious arguments, and are significantly influenced by them: logical fallacies, especially those of relevance and credibility, play a measurable role in driving opinion change. These results not only uncover emergent behaviours in agents' dynamics, but also highlight the dual role of LLMs as both generators and victims of flawed reasoning, raising important considerations for their deployment in socially sensitive contexts.

**Keywords:** Large language model; Opinion dynamics; Logical fallacies; Social simulations; Agent based model

## 1 Introduction

For the logical question of things that grow; one side holding that the ship remained the same, and the other contending that it was not the same.

*Plutarch, Life of Theseus 23.1*

In its original formulation, the “Ship of Theseus” paradox concerns a debate over whether or not a ship that had all its components replaced one by one would remain the same. Consider engaging in a discourse regarding this paradox within the context of a philosophy class, an online Reddit community, or during a dinner gathering with friends. Everyone will reason on the paradox and try to convince others of their stance. Convincing arguments can be proposed both in favor of and against this statement. Ultimately, everyone

© The Author(s) 2025. **Open Access** This article is licensed under a Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License, which permits any non-commercial use, sharing, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if you modified the licensed material. You do not have permission under this licence to share adapted material derived from this article or parts of it. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by-nc-nd/4.0/>.

will leave the debate with their own opinion or no opinion at all. Regardless of the context in which the debate takes place, one thing does not change: the means through which we will try to convince our peers, or they will convince us, is *language*. When a speaker intentionally uses language to convey a specific purpose, they exert an illocutionary force that can influence the listener's perspective, leading to a common understanding or increased division. Therefore, we must consider how language shapes the development of opinions.

The development of individual and public opinions has long been a focus of psychologists and sociologists, and more recently, it has been extensively explored in computational social science [1, 2]. This research acknowledges the complexity of Opinion Dynamics, (henceforth, OD), where multiple interacting factors lead to emergent behaviours such as consensus [3], polarization [4], and radicalization [5], often difficult to predict. Understanding the drivers of opinion change and going beyond mere observation of opinion patterns remains a complex issue. One common approach to tackle this issue is through models of OD, which aim to explain how opinions evolve via social interactions [6]. These models simplify real-world phenomena, enabling the exploration of various what-if scenarios. They generally simulate a population of individuals and their interactions, with processes often governed by simple rules that reflect empirically observed behaviours, such as the repeated averaging of opinions with neighbours [7, 8]. Recent models also incorporate the backfire effect [9, 10], where individuals become more entrenched in their opinions when confronted with contradictory information [11]. Opinion evolution is driven by factors rooted in socio-psychological theories, such as cognitive biases [12], as well as external forces like peer pressure [13], algorithmic biases [14], and mass media [15]. While these models provide simplified representations of societal dynamics and help stakeholders understand social behaviours, they often overlook important complexities. For example, they typically map opinions and messages onto numerical values and rely on rule-based agents, which limits their ability to capture the nuances of human behaviour and the complex relationships between agents' characteristics, such as demographics and personality traits.

To overcome such limitations, we propose a novel framework exploiting Large Language Models (LLMs) capabilities to create an Agent Based Model (ABM) that allows for the study of the interplay between language and opinion change in the long term. The relationship between language and opinion change has been underexplored. Monti et al. [16] is a prominent exception, highlighting the role of knowledge, similarity, and trust in a social media case study. Their findings challenge simplistic OD models, emphasizing the need for more complex analysis. LLMs have revolutionized language-related studies, enabling more realistic social simulations. Park et al. [17, 18] introduced LLM agents as *social simulacra*, capable of simulating personalities and social behaviours. Claims about LLMs possessing Theory of Mind (ToM) [19] remain debated: while Kosinski [20] and others [21, 22] suggest they exhibit emergent ToM abilities, critics [23–25] highlight their inconsistencies in ToM tasks and lack of genuine social intelligence. Nevertheless, even a simulated ToM may enhance OD models by enabling agents to consider interlocutors' mental states. LLM-driven populations display spontaneous emergent behaviours akin to human societies, such as scale-free networks [26], information diffusion [27], and social conventions through interactions [28]. In opinion evolution, LLM agents replicate echo chambers [29], polarization [30], and confirmation bias effects [31]. While LLMs can generate persuasive arguments [32] aligned with psycho-linguistic theories [16], they are less convincing than humans [33] and exhibit biases toward scientific accuracy [31], politeness

[34], and platform-specific discourse styles [35]. LLMs are also argued to have a sycophantic bias, broadly understood as *excessive agreement with and flattery of the user* [36]. While the strength and manifestation of this behavior may vary across models and evaluation settings, multiple studies consistently highlight their prevalence [36–40] also in the models employed in the present study, such as Mistral-7B-Instruct and Llama-3-8B [38].

Despite these biases, LLM-based agents have successfully reproduced experimental results in psychology and linguistics [41], making them valuable tools for *in silico* social experiments. A summary of representative works and their main characteristics is provided in Table 1.

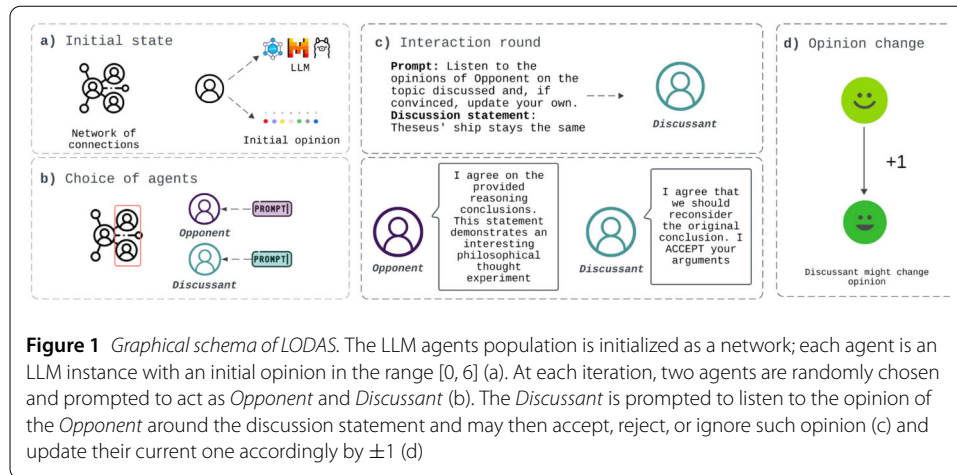
This study aims to advance Opinion Dynamics and social simulations by leveraging Large Language Models. Traditional models rely on mechanistic assumptions rarely validated in real-world settings, limiting their applicability. Instead, we explore whether LLM agents, operating without predefined update rules and guided by the Theory of Mind (ToM) hypothesis, can exhibit realistic individual behaviour and emergent collective dynamics. Unlike classical models, LLM agents engage in natural language interactions, allowing us to investigate the interplay between language and opinion change. Specifically, we examine how they employ logical fallacies and assess their role in persuasion, an aspect understudied in prior research, primarily focused on fallacy detection. A notable exception is Breum et al. (2023) [32], who analyzed LLM-driven persuasion, showing that trust, status, and knowledge influence stance shifts. However, their study focused on one-shot interactions, while we examine how LLMs adapt arguments and leverage fallacies over time. Payandeh et al. (2023) [42] provide the first systematic analysis of LLMs' susceptibility to fallacious reasoning in debates. They find that GPT-4 agrees with flawed arguments 67% of the time, significantly more than logically sound ones. Building on this, we investigate how LLMs not only process but also generate fallacies in multi-agent interactions, shedding light on their role in long-term opinion evolution.

For this purpose, we introduce LODAS, a Language-Driven Opinion Dynamics Model for Agent-Based Simulations framework. The framework allows the definition of a custom population of LLM agents and their interaction around a debated topic.

A schematic representation of LODAS is provided in Fig. 1. As shown in Fig. 1(a), LLM agents (instances either of Mistral or Llama models) hold one of seven possible opinions, evolving through social interactions via  $\pm 1$  updates. The use of a 7-point Likert scale [43] follows established methodologies in psychological research for measuring subjective constructs. We simulate three distinct scenarios: (i) a Baseline scenario with a uniform opinion distribution; (ii) a Polarized scenario, where opinions are bimodally distributed between positive and negative stances with no neutral positions; and (iii) an Unbalanced scenario, where most agents initially hold an extremely negative stance. Throughout the simulations, two agents are selected at random (see Fig. 1(b)) to engage in discussion (see Fig. 1(c)), where the *Opponent agent* (*Opponent*, from now on) attempts to persuade the *Discussant agent* (*Discussant*, from now on), who may then update their opinion on the Ship of Theseus paradox. This topic was chosen to minimize controversy and prevent convergence toward a predefined ground truth, a phenomenon documented in prior studies [31, 32, 44]. To assess the impact of linguistic framing, we start the discussion with one of two formulations: (i) a positive direction ("The ship remains the same") and (ii) a negative direction ("The ship becomes different"). This choice follows prior research [31] demon-

**Table 1** Overview of selected works in traditional and LLM-based opinion dynamics

Paper	Approach	Features	Innovation	Contributions
DeGroot, 1974 [7]	OD	Agents with weights; Synchronous opinion updates;	Opinion updated as a weighted average with neighbors' opinions	Convergence depends on the weights: the higher the weights of the neighbors, the easier the convergence to consensus
Friedkin and Johnsen, 1990 [45]	OD	Social network with stubborn agents	Stubbornness modeled as a susceptibility to others' opinions	Presence of stubborn agents can accelerate or slow dynamics; validated on real-world groups
Hegselmann and Krause, 2002 [46]	OD	Continuous opinions in [0, 1]; Bounded-confidence model with synchronous updates	Opinion averaging with all neighbors within a confidence bound	Convergence to clusters depending on the confidence bound
Deffuant and Weisbuch, 2000 [8]	OD	Continuous opinion in [0, 1]; Bounded-confidence model with pairwise updates; mean-field study	Opinion averaging with neighbours within a confidence bound	Low bounded-confidence fosters polarization
Sirbu et al., 2019 [14]	OD	Extension of [8] with a biased recommender system; mean-field study	Interaction mediated by the recommender system	Algorithmic Bias enhances polarization and fragmentation
Chuang et al., 2023 [44]	LLM OD	LLM agents in a network; simulation of confirmation bias	Agents have <i>personae</i> and a memory; effects of confirmation bias	Agents biased toward accurate info; convergence to model's inherent bias
Breum et al., 2024 [32]	LLM OD	Synthetic persuasion scenario with LLM agents	Use of skeptic and convincer agent types	LLMs effectively mimic traits of human persuasion mechanisms
Payandeh et al., 2023 [42]	LLM OD	Two-agent engaging a multi-round debate	Focus on logical reasoning and susceptibility to fallacies	GPT-4 agrees with flawed arguments in 67% of cases
Törnberg et al., 2023 [35]	LLM OD	Simulation of social media platform with LLM agents acting as users	Assessment of Recommender System effects on user behaviour	Polarization driven by RS; alternative RS reduces toxicity and increases interpersonal interactions
Ju et al., 2024 [47]	LLM OD	Debates between LLM agents with opinions on a scale	Evaluation framework; comparison with DeGroot and HK models	LLM OD is highly sensitive to input prompt
Park et al., 2024 [18]	LLM OD	Simulation with 1052 LLM agents mimicking GSS respondents	Comparison of real and simulated responses	Reduction of LLMs bias across racial and ideological groups
Wang et al., 2025 [30]	LLM OD	LLM agents compared with BCM and FJ; natural language interactions	Agents follow various topologies; interact via debates only with neighbors	Polarization and echo chambers emerge; methods for mitigation
Our work	LLM OD	Mean field topology; natural language interaction and logical fallacies detection	Multiple initial conditions; comparison with OD model, validation	Asymmetric pattern in opinion updating, persuasion through logical fallacies



strating how initial statement framing (“Global warming is/is not a hoax”) may influence opinion evolution.

The investigation presented in this work is structured according to the following Research Questions (RQs).

**RQ1: Can LLM agents exhibit realistic individual behaviour and emergent collective dynamics?** The lack of mechanistic rules LLMs guiding the agents might lead to dynamic trends that differ from classical models. We allow LLMs to discuss and we observe their dynamics.

**RQ2: To what extent do initial opinion distributions influence final outcomes?** The framing of the initial statement might directly influence the LLMs in generating persuasive statements. To investigate this, we designed two versions of the discussion statement, one positive and one negative, and examined their impact on the opinion of the agents’ population.

**RQ3: How do different LLMs impact the persuasion process, particularly regarding logical fallacies?** Different LLMs might generate varied persuasive strategies when interacting with each other. We propose to investigate the logical reasoning behind these textual utterances to identify if their persuasion attempts are flawed and unreliable.

The remainder of this paper is organized as follows. In Sect. 2, we examine the outcomes of our simulations across different initial conditions and scenarios, analyzing, on the one hand, opinion trends, acceptance rates, and, on the other, the linguistic patterns in agent interactions, assessing the role of logical fallacies in shaping opinion change. Section 3 details the simulation framework and experimental design. In Sect. 4, we discuss our findings, and highlight three different behaviours: convergence around a single position, tendency towards agreement, and an asymmetric acceptance-rejection bias, whereas higher opinions are more often accepted and rarely rejected, while lower opinions are more often rejected and rarely accepted, producing an asymmetric pattern in opinion updating.

We also highlight the presence of fallacies in LLM-generated discourse and their impact on persuasion. Additionally, in Sect. 5, we outline key takeaways, study limitations, and directions for future research. Additional figures and analyses are provided in the Supplementary Materials.

## 2 Results

This work extends the modelling of OD using LLM agents to explore whether and which emergent behaviours arise without explicit opinion modification rules. Additionally, it examines the linguistic features of the debates, linking them to specific agent roles and behaviours. To this end, we defined a framework in which a networked population of LLM agents discusses a given topic, updating their opinions according to tunable behavioural rules. Our simulations considered a population of 140 LLM agents. We assumed a mean-field context (i.e., all agents can interact with all other agents without any social restrictions), a commonly used starting point to identify potential emerging behaviours from the opinion evolution process. Each agent is an LLM instance, holding a discrete opinion in the interval  $[0, 6]$ , where 0 means *strongly disagree* and 6 *strongly agree* with a given statement. Agents – as in many classical OD models – interact with each other at discrete time intervals in a pairwise fashion: at each time step, an interacting pair is chosen at random among the connected agents; in this way, in each interaction, we can assign each agent one of two roles, respectively *Opponent* and *Discussant*.

In the present work, we assigned as a discussion topic the paradox of the *Ship of Theseus*, a thought experiment on the concept of identity first recorded in Plutarch's works. The rationale behind the paradox is the following: if all the parts of the ship are replaced over a long period, is the resulting ship the same ship it was at the beginning? This dilemma was chosen because there is no scientific truth. In this way, we avoid LLMs converging toward what they know to be scientifically valid and limit their bias toward immediate adherence to positive opinions. We designed our model to pose this “dilemma” in two different ways: (i) “the boat is the same”, and (ii) “the boat is not the same”. We leveraged Mistral-7B Instruct [48] (Mistral from now on) and Llama-3-8B [49] (Llama from now on) to compare different open state-of-the-art LLMs. By varying the direction of the dilemma, the LLM, and the initial distribution of opinions, we created 12 distinct settings. From our simulations, we obtained opinion evolution data and related textual data, allowing us to relate language and opinion change.

### 2.1 Emergent behaviours in LODAS

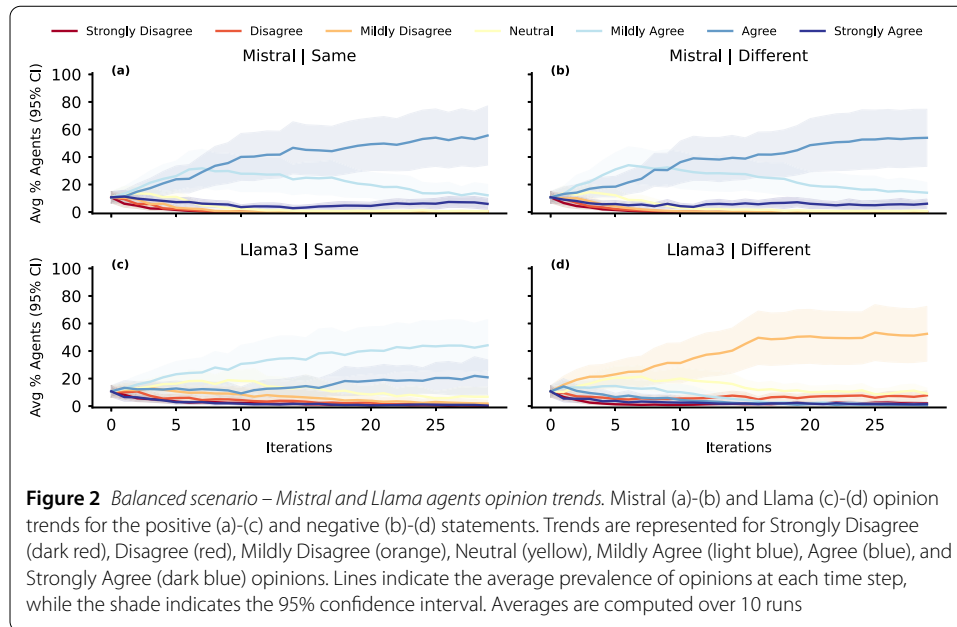
To investigate whether populations of LLM agents exhibit emergent social behaviours (*RQ1*) – such as convergence, consensus, or polarization – we begin by analyzing the opinion evolution in the Balanced scenario. Here, agents' initial opinions are uniformly distributed across the opinion spectrum. This setup serves as a neutral baseline to avoid initial biases and allows comparison with bimodal or skewed initial distributions.

Figure 2 illustrates the evolution of opinion distributions over 30 iterations, across 10 independent simulation runs. The shaded areas represent the 95% confidence interval. Across all four panels, we observe consistent patterns.

First, we note a consistent pattern of *convergence*: agent populations do not remain evenly distributed or fluctuate randomly, but rather gravitate toward one or two dominant opinions. This concentration is stable across runs, with the majority of agents consistently clustering around the same opinion categories.

Second, this convergence is predominantly oriented toward *agreement* with the presented statement, whether it is in the positive or negative direction. In both Mistral conditions (Figs. 2(a)-(b)), we see a progression from mild agreement to full agreement, resulting in a dominant majority of agents aligning with the statement. Similarly, in *Llama|Same*



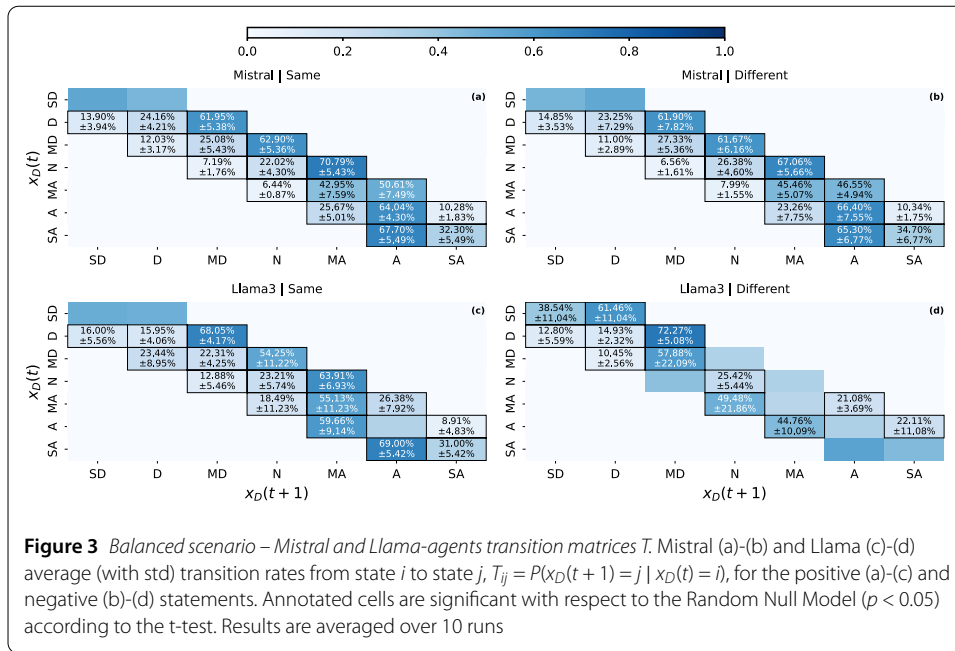


setting (Fig. 2(c)), agents increasingly converge around *Mildly Agree* and *Agree*, while *Neutral* initially rises and then declines. An exception to this tendency towards agreement is found in the *Llama|Different* setting (Fig. 2(d)), where the dominant final opinion is *Mildly Disagree*. Here, although *Neutral* and *Mildly Agree* increase early on, they subsequently decline, reversing the opinion trend compared to other conditions. This distinct behaviour underscores that while convergence is a general feature, its orientation (agreement or disagreement) may depend on model-specific dynamics and prompt framing.

**Comparison with Random Baseline** To determine whether the observed convergence and agreement patterns arise from chance or represent systematic behaviours, we compare them with a Random null model that mirrors the structural features of the simulations (population size, number of iterations, frequency of interactions, and initial distribution) but replaces agents' decision-making with stochastic transitions. In this model, agents randomly shift their opinion by  $-1$ ,  $0$ , or  $+1$  upon interaction, with probabilities uniformly distributed across permitted transitions (see Sect. 3 and Supplementary Materials Sect. S1 for further details).

The Random null model fails to reproduce the emergent patterns observed in LODAS simulations. The opinion distribution remains uniform over time (this also holds with different initial conditions, see Supplementary Figs. S1-S3).

To statistically validate the difference, we compare the transition matrices of the LODAS simulations and the Random baseline. Figure 3 presents the average transition probabilities  $T_{ij} = \mathbb{P}(x_D(t+1) = j \mid x_D(t) = i)$  across all conditions. Black-bordered cells denote statistically significant differences ( $p < 0.05$ , obtained with a two-sample Welch's t-test [50] with unequal variances on the distributions obtained from 10 independent executions of each model). A substantial majority of opinion transitions in LODAS simulations are significantly different from the random baseline, confirming that the observed behaviours are not attributable to randomness.



**Mechanisms Behind Convergence – Persuasion and Sycophancy** This first analysis of opinion evolution trends, however, does not explain *how* these dynamics emerge. A first hypothesis is that convergence results from sycophantic behaviour, i.e., agents consistently adopting their *Opponent's* opinion, which is an LLM characteristic recognized in the literature. To test this, we analyzed the acceptance probabilities  $P(A | x_D, x_O)$ , i.e., the likelihood that a *Discussant's* opinion  $x_D$  moves towards the *Opponent's* opinion  $x_O$ . This measure also reflects the success of persuasion: if the *Discussant* changes their opinion to align more with the *Opponent's*, we interpret this as the *Opponent* having persuaded the *Discussant*. We computed matrices of  $P(A | x_D, x_O)$  and the average  $P(A | \Delta x)$ , where  $\Delta x = x_O - x_D$ .

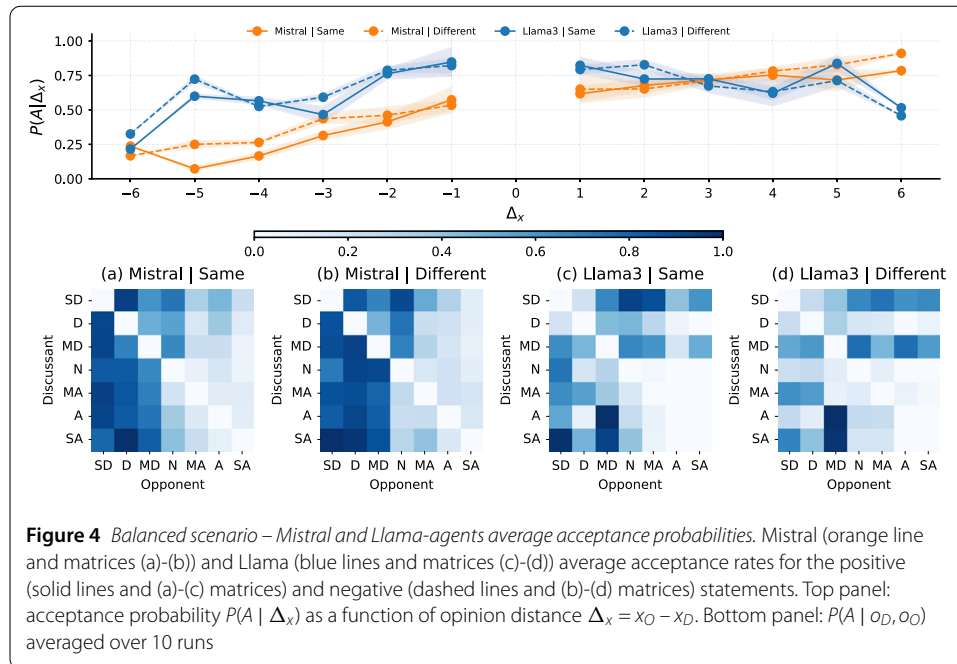
Figure 4 shows that persuasion is not indiscriminate. For Mistral agents (orange lines and matrices (a) and (b)), the probability of persuasion increases with  $\Delta x = x_O - x_D$ : it exceeds 60% when the *Opponent* expresses a more agreeable opinion ( $x_O > x_D$ ), but drops to around 20% when the *Opponent* holds a more disagreeing view ( $x_O < x_D$ ). Llama agents display a more symmetric pattern, yet still show an increase in  $P(A | \Delta x)$  as  $\Delta x$  grows, indicating that persuasive success correlates positively with the relative agreement of the *Opponent's* opinion.

This asymmetry contradicts the sycophancy hypothesis. Agents are not passively aligning with any interlocutor; they exhibit selective persuasion, favoring *Opponents* whose opinions are more agreeable relative to the statement context.

Furthermore, rejection patterns support this interpretation: Llama agents have a lower  $P(R | x_D, x_O)$  than Mistral agents, and rejection probabilities generally decrease as  $\Delta x$  increases (see Supplementary Materials, Fig. S20).

These patterns suggest that agents *do not exhibit sycophantic behaviour nor bounded confidence*: even distant opinions can successfully persuade or be actively rejected. Specifically, the likelihood of persuasion increases when the *Opponent* expresses a more positive and distant opinion than the *Discussant's*, skewing the opinion distribution toward agreement. In contrast, *Opponents* expressing more negative opinions, i.e.,  $x_O < x_D$ , have lower



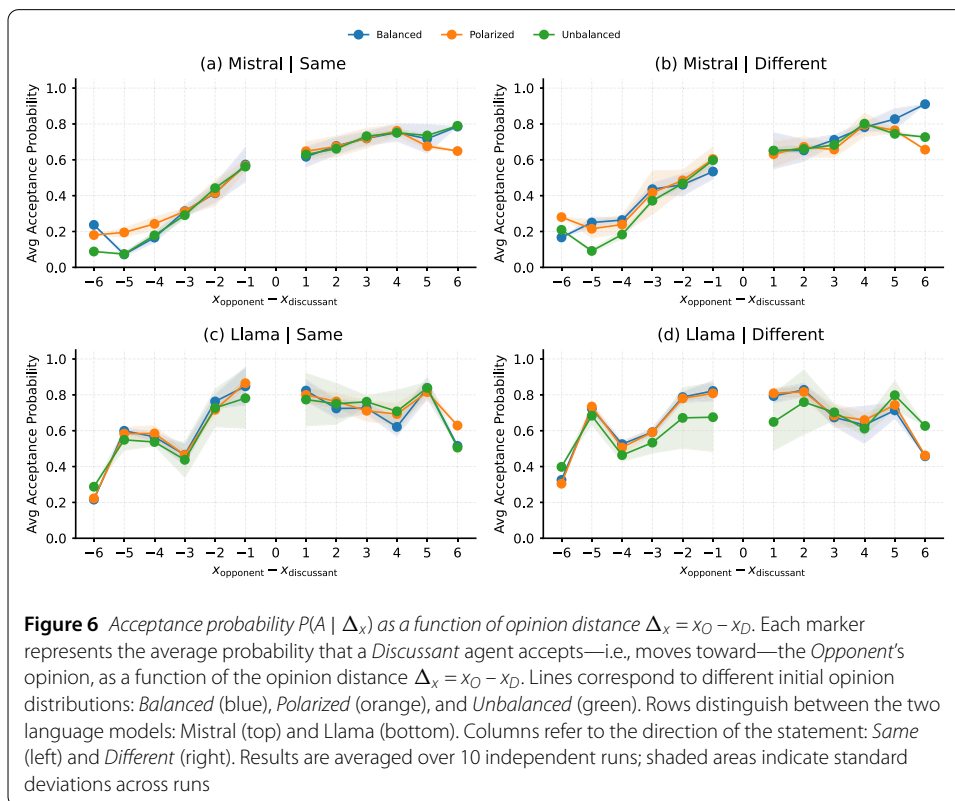
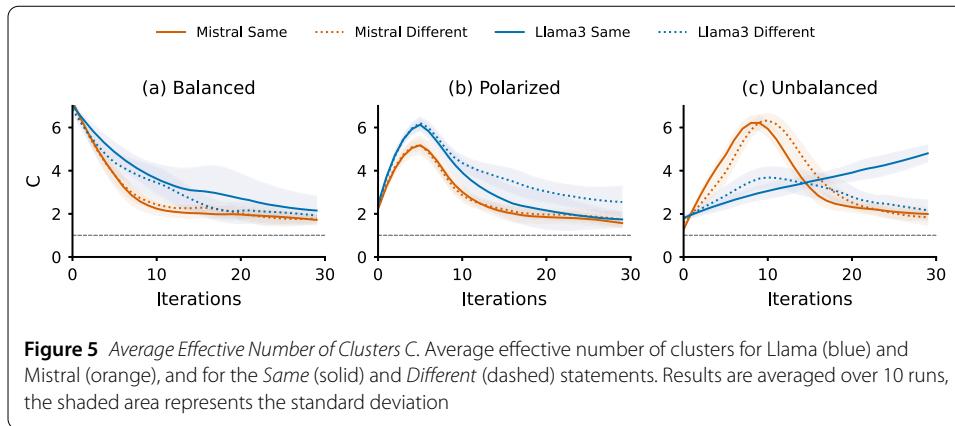


persuasive success and higher rejection rates. This directional asymmetry indicates a form of *backfire effect*, emerging only for less agreeable inputs and intensifying with opinion distance.

Finally, simulations using a toy model in which agents always accept their *Opponent's* opinion produce markedly different dynamics: the majority of the population converges on the *Strongly Disagree* opinion. The complete absence of such a result in the LODAS simulations further refutes the sycophancy hypothesis. Conversely, toy model simulations where agents always reject their *Opponent's* opinion generate dynamics more similar, albeit more extreme, to those observed, with the majority of agents converging on the *Strongly Agree* opinion (see Supplementary Materials Sect. S1 for further details).

Together, these findings address our first research question (RQ1), showing that LODAS consistently produce emergent behaviours characterized by *convergence and alignment*, typically *toward agreement*. These trends are *statistically significant* compared to a random null model. Moreover, the behaviours do not stem from indiscriminate acceptance or sycophancy. Instead, they arise from structured, selective interaction patterns shaped by the underlying language models, resulting in an asymmetric backfire effect and a bias toward strongly positive opinions, which we can call an *asymmetric acceptance-rejection bias*. The strength of these effects varies depending on the choice of LLM.

**Impact of Skewed Initial Opinion Distribution** To assess the influence of initial conditions (RQ2), we systematically compared simulations initialized under three different configurations: *Balanced* (uniform distribution across the opinion spectrum), *Polarized* (bimodal distribution centered on *Strongly Disagree* and *Strongly Agree*), and *Unbalanced* (skewed distribution concentrated around *Strongly Disagree*). Despite these substantial differences in starting configurations, we observe that the qualitative evolution of opinions over time is largely preserved across scenarios. In all conditions, opinion trends rapidly



shift away from initial extremes, with agents progressively converging around moderate or positive agreement positions (see Supplementary Materials, Figs. S11-S14).

As emerges from Fig. 5, across all three initial conditions, we observe a consistent decrease in variability over time, indicating convergence toward fewer opinion states (see also Figs. S10, S13 and S16). Overall, these trends confirm that while the statement framing has limited influence, the choice of language model affects the speed and stability of convergence.

Acceptance and rejection behaviours also appear robust to changes in initial conditions. The functional forms of  $P(A | \Delta x)$  (see Fig. 6) and  $P(R | \Delta x)$  (see Supplementary Materials, Fig. S26) are stable across *Balanced*, *Polarized*, and *Unbalanced* scenarios. Similarly,

the matrix representations  $P(A | x_D, x_O)$  and  $P(R | x_D, x_O)$  reveal consistent interaction patterns.

Taken together, these results indicate that the initial distribution of opinions has a limited and transient influence on the collective dynamics (RQ2). Instead, the key determinant of opinion evolution, variability, and interaction behaviour is the LLM used to enhance agent decision-making. The differences between Mistral and Llama are more pronounced and persistent than those induced by any variation in the starting opinion configuration.

## 2.2 Linguistic behaviour

Moving on to RQ3, we analyzed the arguments produced by the agents in both roles – *Opponents* and *Discussants* – during their conversations on the Theseus' Ship paradox. Specifically, we examined their linguistic behaviour, focusing on the production of persuasive yet fallacious content, and assessed how such fallacious utterances can influence the opinion change trend within multi-agent debate.

Table 2 shows the average percentage of fallacious statements generated by *Opponent* agents, calculated from aggregated results of 10 discussion runs. In each run, *Opponents* produced a total of 12.600 statements. The percentages represent the ratio of fallacious content relative to the total number of statements and are categorized by initial opinion distribution – Balanced, Polarized, and Unbalanced – by statement, and by LLM.

The proportion of statements containing logical fallacies remained relatively stable across all scenarios and discussion framing, at around 20%. Variability was primarily attributed to the underlying LLM. Mistral agents produced slightly fewer fallacious statements than Llama, especially under unbalanced initial conditions, where only 15.56% of statements were classified as fallacious. Under balanced conditions, Mistral's fallacy rate remained close to 19%, regardless of the framing of the discussion. In contrast, Llama showed an increased sensitivity to negative framing, producing more fallacious utterances than Mistral. Nonetheless, the overall fallacy rate remained below 30% of the total statements.

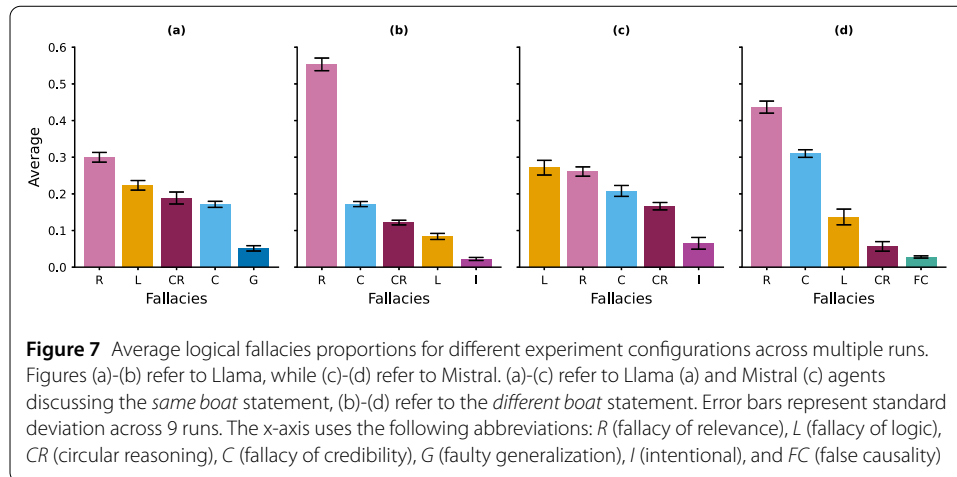
Due to the limited variability in fallacy types observed across configurations of different LLMs and statement framing, we focus our analysis only on the patterns detected in the Balanced scenario. Additional figures for the Polarized and Unbalanced scenarios are provided in the Supplementary Materials (Figs. S27 and S28).

As shown in Fig. 7, few types of fallacies emerged, with LLMs often repeating similar patterns across different statement framings. The variability of their aggregated distribution over the 10 runs was minimal, as indicated by the low standard deviation value in the error bar. Overall, both Llama and Mistral relied more heavily on specific types of fallacies,

**Table 2** Percentage of logical fallacies in *Opponents'* statements

	Balanced				Polarized				Unbalanced			
	Llama		Mistral		Llama		Mistral		Llama		Mistral	
	S	D	S	D	S	D	S	D	S	D	S	D
% Fallacious (O)	20.87	23.39	19.01	19.31	19.88	26.83	16.79	20.22	22.06	20.77	18.39	15.56

Percentage of unique *Opponent* (O) statements classified as fallacious, across models (Llama, Mistral), initial opinion distribution (balanced, polarized, unbalanced), and opinion framing (same, different).



**Table 3** Ratio of *Discussants* changing opinion for the effect of fallacious statements

	Balanced				Polarized				Unbalanced			
	Llama		Mistral		Llama		Mistral		Llama		Mistral	
	S	D	S	D	S	D	S	D	S	D	S	D
% Opinion change (D)	64.9	71.4	53.79	55.29	78.06	77.16	58.52	60.53	77.84	78.82	60.18	60.72

Percentage of *Discussants* (D) changing opinion for the effect of fallacious statements produced by *Opponents*, across models (Llama, Mistral), initial opinion distribution (Balanced, Polarized, Unbalanced), and opinion framing (same, different).

particularly fallacies of relevance, credibility, and logic. Furthermore, both models generated arguments in which they reiterated the initial premises as conclusions, resulting in the pragmatic defect of circular reasoning; this occurred more frequently in the *same boat* discussion. Additionally, though less frequently, they tended to assume a causal relationship without justification (false causality).

Having assessed the presence of fallacious utterances in the *Opponent* agents, we moved on to measure the persuasive impact of these fallacies over the *Discussant*. Specifically, we investigated whether the presence of a fallacy in the *Opponent* statement caused a shift by  $\pm 1$  in the opinion held by the *Discussant* compared to their prior stance before the interaction. An overview of this analysis can be found in Table 3. Overall, Llama *Discussants* demonstrated higher vulnerability to logical fallacies, changing their opinion 78% of the time in the *same boat* scenario, and 75% of the time in the *different boat* scenario. Conversely, Mistral agents showed greater robustness against logical fallacies. They both produced fewer fallacies than Llama agents (Table 2) and their *Discussants* resisted more than Llama ones, with opinion shifts occurring in 60% and 61% of the respective cases (Table 3).

Once investigated the production of fallacies at a macro-level, we proceeded to examine which specific types of logical fallacies were most effective in inducing the opinion shifts in the *Discussants*. Most changes, as highlighted in Table 4, are caused by fallacies of relevance when agents discussed the *different boat* scenario, whereas in the *same boat* discussion the opinion change is triggered by general logical fallacies that do not fall under the other labels recognized by the classifier.

Although it is difficult to interpret the specific fallacies introduced by the classification model under the label *fallacy of logic*, the preference for fallacies of relevance may reflect

**Table 4** Percentage of opinion changes in *Discussants* due to logical fallacies

Fallacy Type	Llama (S)	Llama (D)	Mistral (S)	Mistral (D)
Fallacy of Logic	<b>24.35%</b>	10.65%	<b>27.07%</b>	16.82%
Faulty Generalization	8.85%	1.14%	3.01%	0.00%
Ad Populum	0.00%	0.00%	0.38%	0.00%
Appeal to Emotion	2.21%	0.00%	0.00%	0.00%
Fallacy of Credibility	19.11%	9.89%	23.31%	31.80%
Fallacy of Extension	0.00%	0.00%	0.00%	0.00%
Fallacy of Relevance	22.33%	<b>51.14%</b>	25.56%	<b>40.67%</b>

Values indicate the percentage of opinion shifts in the *Discussant* agents exposed to each fallacy type by the *Opponent*'s statement. Results refer to Llama and Mistral agents discussing the *same boat* (S) and *different boat* (D) initial statement.

the tendency of LLMs to overlook logical reasoning in favor of empty rhetorical devices. This rhetoric is made up of compelling elements introduced into the argument, which may be unrelated to the discourse's premises, while creating a misleading yet persuasive discourse.

### 3 Methods

In the Language-Driven Opinion Dynamics Model for Agent-Based Simulations (LODAS) model, we have a population of  $N$  agents, where each agent  $a$  is an LLM agent, i.e. an instance of a Large Language Model. Agents are enhanced using AutoGen [51]: “a framework for creating multi-agent AI applications that can act autonomously or work alongside humans”. Specifically, we exploited AutoGen AgentChat's AssistantAgent, a built-in agent that uses a Large Language Model and has the ability to use tools. It serves as a foundational agent that can be customized or integrated into multi-agent conversations.

In our model, each LLM agent holds a discrete opinion  $x_a \in \{0, \dots, 6\}$  associated (from 0 to 6) with a negative (*strongly disagree*, *disagree*, *mildly disagree*), *neutral*, or positive (*mildly agree*, *agree*, *strongly agree*) stance on a given statement  $s \in S$  around a given topic  $\theta \in \mathcal{T}$ . A statement  $s$  can have a *positive valence*, e.g., “this is true,” or a *negative valence*, e.g., “this is not true.”

In our study, we chose the Ship of Theseus paradox as the topic, where the statements were phrased as “the ship is the same” (positive valence) and “the ship is different” (negative valence). To formalize this, we define a function  $\pi(s)$  that maps statements to their valence as follows:

$$\pi(s) = \begin{cases} +1, & \text{if } s \text{ expresses a positive valence (e.g., “the ship is the same”)} \\ -1, & \text{if } s \text{ expresses a negative valence (e.g., “the ship is different”)} \end{cases}$$

At each discrete time step  $t$ , a pair of agents  $(a_i, a_j)$  is randomly selected from this network. One agent from the pair is assigned the role of *Discussant* ( $D$ ) while the other takes on the role of *Opponent* ( $O$ ).

#### Prompts

*Discussants*  $D$  act according to the following prompt.

### Discussant Prompt

```
[INST]
### You {Discussant_opinion} on the reasoning conclusion
provided as input.
Task:
- Listen to the argument of {Opponent.name} on the reasoning
conclusions and decide if you maintain your opinion
or change it.

### Constraints:
- At the end of each interaction declare if you
  - 'ACCEPT' {Opponent.name} argument;
  - 'REJECT' {Opponent.name} argument;
  - 'IGNORE' your original opinion.
Write your response with the following format:
\"My original opinion was I {Discussant_opinion}
on the reasoning.
After reading your argument my conclusions are:
I <ACCEPT|REJECT|IGNORE> your stance because <argument>\"
[/INST] "
```

The role of the *Opponent* is instead modeled by the following prompt:

### Opponent Prompt

```
[INST]
You {Opponent_opinion} on the reasoning conclusion provided as input.
Support your opinion by providing personal arguments.
Avoid using already generated arguments.

IF {Discussant.name} writes REJECT in his answer, write a second statement
where you declare if you <ACCEPT|REJECT|IGNORE> his stance.
Otherwise, conclude the conversation writing a message with
a single word 'END'.

### Constraints:
- In your first statement you must adhere to your opinion
  ('{Opponent_opinion}')
- Write your first response with as: \"I {Opponent_opinion} on the
provided reasoning conclusions. I think that <argument>\"
[/INST] "
```

The selected *Discussant*  $a_D$  engages in a discussion with the *Opponent*  $a_O$  on a pre-defined topic  $\theta \in \mathcal{T}$ , with the goal of influencing the other's opinion. During this interaction, the *Discussant*  $a_D$  and the *Opponent*  $a_O$  are prompted to maintain their initial opinions unless convinced by the argumentation of the other.

The discussion is started by  $a_D$ , who asks agent  $a_O$  to express their opinion on statement  $s$  around topic  $\theta$  with valence  $\pi(s)$ .

In our study, we have two different statements:

### Positive valence statement $\pi(s) = +1$

Theseus set sail to reclaim the throne as king of Athens. During the journey, parts of Theseus's ship began to break or decay; Theseus and his crew replaced these parts as they sailed. Eventually, each part of the ship is replaced. In the end the Ship of Theseus is still the same ship on which he originally sailed.

and



Negative valence statement  $\pi(s) = -1$

Theseus set sail to reclaim the throne as king of Athens. During the journey, parts of Theseus's ship began to break or decay; Theseus and his crew replaced these parts as they sailed. Eventually, each part of the ship is replaced. In the end, the Ship of Theseus is completely different from the one he originally sailed.

The question has the following structure:

Discussion initialization

What do you think of the following statement?: {s}

The *Opponent* is asked to produce a persuasive utterance in response to the *Discussant*, based on their current opinion, to persuade the *Discussant* and shift their stance. The *Discussant* then processes the *Opponent*'s response and generates a comment about that statement, expressing whether it was convinced by the *Opponent* or not. The interaction may result in a positive (+1) or negative (−1) change in the *Discussant*'s opinion, or no change (0). Finally, the *Opponent* closes the discussion in one of two ways: if the *Discussant* chooses not to accept the persuasive statement, then it generates a new statement commenting on the current stance of the *Discussant* and thanking it for the discussion. This comment does not affect the opinions' status, it simply ends the iteration round. Otherwise, if the *Discussant* is persuaded by the *Opponent*, the *Opponent* can simply end the iteration with an END keyword.

In the present study, we set the number of iterations to  $T = 30$ . At each iteration  $t$  there are  $N$  pairwise random interactions  $(a_D, a_O)$ .

To ensure consistency across experiments, we fixed the temperature parameter of the LLMs after preliminary tests showed that varying it did not alter the overall opinion dynamics, but only affected the linguistic variability of the generated statements, for instance, reducing repetitive phrasing. Due to computational constraints, we limited our experiments to small-scale models (7–8B parameters), which still provided meaningful interaction patterns while remaining tractable for large-scale simulations. Regarding the number of interaction rounds, we observed that the system typically reached convergence and stability within a relatively short number of iterations, consistent with previous work using similar simulation settings [52]. To minimize unnecessary computational load and environmental impact, we therefore set the maximum number of iterations to 30.

### 3.1 Metrics and analysis

#### 3.1.1 Statistical validation

To evaluate whether our results differ significantly from patterns that could arise by chance, we constructed a *Random Null Model* under the same experimental constraints as the LODAS setting. Specifically, we defined a population of  $N = 140$  agents, interacting under the same structural rules. Each simulation was run for  $T = 30$  iterations, with each iteration consisting of  $N$  pairwise interactions between a *Discussant* agent  $D$  and an *Opponent* agent  $O$ . In each interaction, only agent  $D$  could update their opinion, with a possible change of +1, −1, or 0.

We performed  $R = 10$  independent runs for the null model to generate a reference distribution of opinion transitions. To compare these outcomes with those from the experimen-

tal condition, we analyzed the respective transition matrices. Each matrix  $\mathbf{T} \in \mathbb{R}^{K \times K}$  represents the empirical average transition probabilities between  $K$  discrete opinion states. The element  $T_{ij}$  denotes the probability of transitioning from opinion state  $i$  to state  $j$ :

$$T_{ij} = P(x_D(t+1) = j \mid x_D(t) = i)$$

where  $x_D(t) \in \mathcal{O}$  is the opinion of the *Discussant* agent at time step  $t$ , and  $\mathcal{O}$  is the set of all possible opinions. Each row of the matrix is normalised such that:

$$\sum_{j=1}^K T_{ij} = 1 \quad \text{for all } i \in \{1, \dots, K\}$$

To assess statistical differences between the experimental and null models, we used Welch's  $t$ -test for independent samples. For each matrix entry  $(i, j)$ , we tested the null hypothesis:

$$H_0 : \mu_{ij}^{\text{exp}} = \mu_{ij}^{\text{null}}$$

against the two-sided alternative:

$$H_1 : \mu_{ij}^{\text{exp}} \neq \mu_{ij}^{\text{null}}$$

where  $\mu_{ij}^{\text{exp}}$  and  $\mu_{ij}^{\text{null}}$  represent the expected transition probabilities in the experimental and null conditions, respectively. We used the `scipy.stats` implementation of Welch's  $t$ -test, which does not assume equal variances. Statistical significance was determined using a threshold of  $p < 0.05$ , under which the null hypothesis was rejected in favor of a significant difference.

The goal of this analysis is to verify that the observed opinion dynamics are not the result of random processes. To this end, we focused our statistical testing on the comparison between empirical transition matrices and appropriate null models. This allowed us to assess whether the structure of opinion change emerged from meaningful interactions rather than noise. We deemed broader statistical comparisons out of scope for the present study.

### 3.1.2 Opinion evolution metrics

Opinion dynamics were further analyzed by examining the temporal evolution of the average opinion distribution. Let  $x_i(t) \in \mathcal{O}$  denote the opinion of agent  $i$  at time step  $t$ , where  $\mathcal{O}$  is the set of possible discrete opinion states. For each opinion  $x \in \mathcal{O}$  and each time step  $t \in \{1, \dots, T\}$ , we computed the proportion of agents holding opinion  $x$ , defined as:

$$P_x(t) = \frac{1}{N} \sum_{i=1}^N \mathbb{I}[x_i(t) = x]$$

where  $\mathbb{I}[\cdot]$  is the indicator function. The resulting trajectories  $P_x(t)$  were averaged across  $R = 10$  independent simulation runs, and we report either 95% confidence intervals or standard deviation bands to indicate variability.

Moreover, variability metrics were also computed. In particular:

- *Entropy (H)*: Measures the uncertainty in the opinion distribution at each time step. It is computed as the Shannon entropy of the normalized opinion counts:

$$H(t) = - \sum_i p_i(t) \log_2 p_i(t)$$

where  $p_i(t)$  is the proportion of agents holding opinion  $i$  at time  $t$ . Higher entropy indicates greater opinion diversity, while lower entropy reflects convergence or consensus.

- *Standard Deviation ( $\sigma$ )*: Captures the spread of opinions by computing the standard deviation of the opinion distribution, weighted by the proportion of agents:

$$\sigma(t) = \sqrt{\sum_i p_i(t) (i - \mu(t))^2} \quad \text{with} \quad \mu(t) = \sum_i p_i(t) \cdot i$$

where  $\mu(t)$  is the weighted mean opinion. Lower values indicate tighter clustering of opinions around the mean.

- *Effective Number of Clusters (C)* (from [14]): Estimates the effective number of opinion clusters using the formula:

$$C(t) = \frac{N^2}{\sum_i n_i(t)^2}$$

where  $n_i(t)$  is the number of agents with opinion  $i$  at time  $t$ , and  $N$  is the total number of agents. This metric adjusts for both the number and size of groups, with higher  $C$  indicating more fragmentation.

To assess sycophantic tendencies, we computed acceptance probability matrices  $\mathbf{A} \in [0, 1]^{K \times K}$ , where each entry  $A_{ij}$  denotes the empirical probability that an agent with opinion  $x_D = i$  accepts the opinion  $x_O = j$  of their *Opponent*. This can be expressed as:

$$A_{ij} = P(A \mid x_D = i, x_O = j) = \frac{N_A(i, j)}{N_{\text{int}}(i, j)}$$

where:

- $N_A(i, j)$  is the number of interactions in which an agent with opinion  $i$  accepted the *Opponent's* opinion  $j$  during the whole process,
- $N_{\text{int}}(i, j)$  is the total number of interactions between agents with opinions  $i$  (Discussant) and  $j$  (Opponent), during the whole process.

Analogously, to examine backfire-like effects, we constructed rejection probability matrices  $\mathbf{R} \in [0, 1]^{K \times K}$ , where  $R_{ij}$  indicates the probability of rejecting the *Opponent's* opinion  $j$  when the *Discussant* holds opinion  $i$ , computed across all iterations. Rows correspond to the discussant's opinion and columns to the *Opponent's*.

Additionally, we analyzed the influence of opinion distance on interaction outcomes by defining the signed opinion distance as  $\Delta x = x_O - x_D$ . For each possible value of  $\Delta x$ , we computed the acceptance and rejection probabilities  $P(A \mid \Delta x)$  and  $P(R \mid \Delta x)$ , respectively. The metric was computed considering the overall dynamics. These conditional probabilities were estimated empirically and averaged over the  $R = 10$  simulation runs, with uncertainty represented using standard deviations.

### 3.2 Logical fallacies detection

The presence of logical fallacies in text is usually identified through transformers-based models, so the task is often approached as a multi-label classification problem. In this work, we employ the `distilbert-base-fallacy-classification` model [53] obtained from HuggingFace. We chose this specific model as it is trained on the dataset used by Jin et al. [54], which introduces the task of logical fallacies detection and the LOGIC dataset for fallacies. In the present article, we refer to the 13 fallacies illustrated in the original article by Jin et al. [54].

In the following, we present and discuss only the most common fallacies identified in our analysis; readers are referred to the original paper for a full summary.

- *Fallacy of credibility*: it consists of an appeal to a form of ethics or authority;
- *Fallacy of relevance*: the argument relies on premises that are irrelevant to the conclusion. In [55], it is suggested that premises might be *psychologically relevant* but not *logically relevant*, resulting in an argument that seem apparently correct and persuasive;
- *Appeal to emotion*: this fallacious argument assumes that premises are not relevant to conclusions, but the premises are used as a means to convey a specific emotion aiming to manipulate the beliefs of the reader;
- *Circular reasoning (circularis in probando)*: a fallacy characterized by a circularity in reasoning so that the premises depend on the conclusions and vice versa.

## 4 Discussion

This study builds on recent literature [18, 31, 32] and introduces a Language-Driven Opinion Dynamics Model for Agent-Based Simulations (LODAS) to investigate how language and social influence shape opinion evolution, with a particular focus on the role of logical fallacies. In the model, each agent holds a discrete opinion ranging from *Strongly Disagree* to *Strongly Agree*. At each time step, a *Discussant* asks an *Opponent* for their opinion on a topic; the *Opponent* responds with the intent to persuade the *Discussant*, who may then adjust their opinion by  $\pm 1$  or keep it unchanged. This process repeats until opinions stabilize or a stopping condition is reached. We studied three initial opinion distributions: (a) Balanced (uniformly distributed opinions), (b) Polarized (only extreme opinions), and (c) Unbalanced (majority extremely negative). The discussion topic was the paradox of the ship Theseus, chosen to prevent convergence toward a ground truth or consensus. Each initial condition was paired with either a positive framing (“The boat is the same”) or a negative framing (“The boat is different”), producing six scenarios simulated with two different LLM.

Our findings address three main research questions: *RQ1*: Can LODAS generate emergent behaviours without mechanistic rules? *RQ2*: To what extent do initial opinion distributions influence final outcomes? *RQ3*: How do different LLM impact persuasion, particularly regarding logical fallacies?

Regarding *RQ1*, our analyses demonstrate that the LODAS framework can produce emergent behaviours without embedding explicit behavioural rules common in traditional mechanistic models [56]. Specifically, agents exhibit (i) strong convergence toward a dominant opinion, often forming a majority though not always full consensus; (ii) a consistent tendency towards agreement; and (iii) asymmetric acceptance-rejection bias — the probability of an agent accepting or rejecting an *Opponent*’s opinion is strongly and oppositely

correlated with the signed opinion distance: higher opinions are more often accepted and rarely rejected, while lower opinions are more often rejected and rarely accepted, producing an asymmetric pattern in opinion updating. These emergent patterns underline the ability of language-driven interactions to naturally shape opinion evolution in ways that mirror empirical social phenomena, confirming the promise of LODAS as a modelling approach.

*RQ2* addresses the role of initial opinion distributions. Our results indicate that initial conditions have limited influence on final outcomes: whether balanced, polarized, or unbalanced, opinions converged toward specific stable points dictated by the model and statement framing. This suggests that each model-statement pair generates a strong internal dynamic that overrides initial biases, driving convergence toward a characteristic opinion cluster. This robustness underscores the influence of language model design on interaction dynamics, reflecting tendencies toward coherence and alignment that encourage agreement [37, 57]. The asymmetric acceptance-rejection bias, which reduces acceptance of negative opinions and amplifies influence from positive ones, appears to be a key driver of this stability.

Finally, in exploring *RQ3*, we found meaningful differences between the two LLM agent types. Mistral agents yielded more stable results, with faster and stronger convergence toward agreement and a pronounced asymmetric acceptance-rejection bias. Conversely, Llama agents displayed greater openness to a range of opinions, though typically favoring those similar to their own. Notably, the framing of the discussion statement influenced Llama agents' dominant opinions: negative framing shifted their majority from agreement to mild disagreement. Linguistic analysis revealed that LLM agents frequently employed logical fallacies—particularly those related to relevance and credibility—in attempts to persuade others. These agents were also influenced by such fallacious arguments, consistent with prior work showing susceptibility of language models to faulty reasoning [42, 58, 59]. Interestingly, Llama agents were more effective persuaders, with about 68% of *Discussants* changing opinions after exposure to fallacious arguments, compared to 54% for Mistral. This highlights both the persuasive power of logical fallacies in artificial agents and the varying susceptibility depending on model architecture.

Taken together, these insights advance our understanding of how language-based social simulations can capture complex opinion dynamics and the role of logical fallacies therein. The observed convergence, agreement bias, and asymmetric acceptance-rejection bias reveal intrinsic tendencies within LLM agents to favor coherence and social alignment, potentially mirroring real-world psychological phenomena but also raising concerns about sycophantic behaviours [37]. However, the agreement we observe is not simply a result of sycophancy but emerges from asymmetric processing of differing opinions, with broader acceptance of more positive views. Moreover, the presence and persuasive impact of logical fallacies emphasize the need to critically evaluate the reasoning capabilities of such agents in social simulations, given their susceptibility to flawed arguments [42, 59]. Our study thus provides a foundation for further work exploring how language-driven models can inform both the dynamics of opinion formation and the risks associated with fallacious reasoning in AI-mediated social influence.

## 5 Conclusion

This study introduces a Language-Driven Opinion Dynamics Model for Agent-Based Simulations (LODAS), allowing for the exploration of how language and social influence shape

opinion dynamics. By utilizing LLM agents, this study shows that synthetic agents, when left unprompted, tend to converge toward agreement, irrespective of initial opinion distributions or prompt framing. This convergence is primarily shaped by the underlying language model, with agents exhibiting a consistent asymmetric acceptance-rejection bias: they are more likely to adjust their opinions toward more positive stances and reject more negative ones. This bias is more pronounced in Mistral, which favors agreement more strongly, whereas Llama agents exhibit a form of bounded confidence, showing greater susceptibility to nearby opinions. In both cases, agents frequently employ logical fallacies in their persuasive attempts and are, in turn, influenced by such flawed arguments.

One key limitation of the current framework is the simplicity of the agents. In this model, agents are equipped with verbal reasoning skills but lack distinct personalities or cognitive diversity. The introduction of more complex agent types - such as those with different decision-making styles, biases, or psychological traits - could better replicate the diversity of human interactions [60, 61]. Additionally, the experiments were conducted using English-aligned LLMs, which may have introduced cultural biases in the interactions between agents. LLMs have been proven to exhibit societal and cultural biases [62–64] attributable to the training data, algorithm structure, and user interactions [65], with an alignment with Western cultural values, especially when prompted in English [66, 67].

Future extensions of the framework could also benefit from a deeper integration of cognitive biases [44] and demographic factors [68], as these elements are known to influence opinion dynamics in the real world. To mitigate cultural biases, future works could leverage emerging methodologies such as cultural prompting [66] techniques to achieve more representative alignment with diverse cultural value systems, thus potentially unveiling unseen dynamics.

Furthermore, the model currently assumes a mean-field scenario, which neglects the structure of real-world social networks. Incorporating network features such as clustering, assortativity, or echo chambers could significantly increase the realism of the simulations and improve their ability to replicate polarization dynamics [30, 69, 70]. Preliminary tests with alternative network topologies and more sophisticated opinion dynamics algorithms suggest the potential to capture more complex patterns of interaction.

The exploration of fallacious reasoning in social simulations of LLM agents and its role in opinion dynamics has been approached at a preliminary level in this study, leaving substantial opportunities for future investigation. The role of fallacies poses challenges not only in the context of social simulations - where agents could potentially be optimised through better prompting, enhanced memory, or other refinements to mitigate fallacious reasoning - but also in human-LLM interactions. If LLMs are easily swayed by illogical arguments and tend to validate human perspectives, they may inadvertently reinforce false or potentially harmful beliefs.

To improve the understanding of these dynamics, several directions for future research can be pursued. One key focus is investigating methods to reduce fallacious reasoning in LLMs, such as through improved prompting, enhanced memory mechanisms, or adjustments to biases. Understanding the interplay between memory, bias, and opinion evolution is also critical for analyzing the role of persuasive language in opinion change. Comparing LLM-based simulations with real-world data from online interactions or controlled experiments can help evaluate (i) the robustness of the framework, (ii) its ability to repli-



cate human behaviour, and (iii) the effects of linguistic features on opinion change under controlled conditions.

To summarize, despite its limitations, the framework provides a valuable tool for studying the mechanisms of consensus-building and argumentation in a controlled environment. The framework could serve as a foundation for exploring the drivers of opinion dynamics and their implications for phenomena such as polarization, bias, and misinformation.

#### Abbreviations

LLM, Large Language Model; AI, Artificial Intelligence; OD, Opinion Dynamics; ABM, Agent Based Model; LODAS, Language-Driven Opinion Dynamics Model for Agent-Based Simulations.

### Supplementary information

**Supplementary information** accompanies this paper at <https://doi.org/10.1140/epjds/s13688-025-00579-1>.

**Additional file 1.** The present article has accompanying Supplementary Information files with figures and tables complementary to those presented in the main text. (PDF 864 kB)

#### Acknowledgements

We thank Daniele Atzeni for the valuable feedback and Giuliano Cornacchia for the help in designing plots and figures.

#### Author contributions

EC analyzed the data and wrote the paper. VP analyzed the data and wrote the paper. GR designed the experiments, performed the experiments and supervised the project. DP supervised the project. All authors read and approved the final manuscript.

#### Funding information

This project is supported by SoBigData.it which receives funding from the European Union—NextGenerationEU—National Recovery and Resilience Plan (Piano Nazionale di Ripresa e Resilienza, PNRR)—Project: “SoBigData.it—Strengthening the Italian RI for Social Mining and Big Data Analytics”—Prot. IR0000013—Avviso n. 3264 del 28/12/2021 (to VP and RG); this work is also supported by the scheme ‘INFRAIA-01-2018-2019: Research and Innovation action’, Grant Agreement No 871042 ‘SoBigData++: European Integrated Infrastructure for Social Mining and Big Data Analytics’ (to RG); finally this work is supported by: the EU NextGenerationEU programme under the funding schemes PNRR-PE-AI FAIR (Future Artificial Intelligence Research) (to EC and RG).

#### Data availability

The dataset generated and analyzed during the current study is within this paper or publicly available at <https://github.com/ericcau/LLM-Opinion-Dynamics>.

### Declarations

#### Ethics approval and consent to participate

Not applicable.

#### Consent for publication

Not applicable.

#### Competing interests

The authors declare no competing interests.

Received: 25 February 2025 Accepted: 29 July 2025 Published online: 15 August 2025

#### References

1. Conte R, Gilbert N, Bonelli G, Cioffi-Revilla C, Deffuant G, Kertesz J, Loreto V, Moat S, Nadal J-P, Sanchez A, et al (2012) Manifesto of computational social science. *Eur Phys J Spec Top* 214:325–346
2. Tucker JA (2023) Computational social science for policy and quality of democracy: public opinion, hate speech, misinformation, and foreign influence campaigns. In: *Handbook of computational social science for policy*, pp 381–403
3. Li L, Scaglione A, Swami A, Zhao Q (2013) Consensus, polarization and clustering of opinions in social networks. *IEEE J Sel Areas Commun* 31:1072–1083
4. Biondi E, Boldrini C, Passarella A, Conti M (2023) Dynamics of opinion polarization. *IEEE Trans Syst Man Cybern Syst* 53(9):5381–5392. <https://doi.org/10.1109/TSMC.2023.3268758>

5. Ramos M, Shao J, Reis SDS, Anteneodo C, Andrade JS, Havlin S, Makse HA (2015) How does public opinion become extreme? *Sci Rep* 5(1):10032. <https://doi.org/10.1038/srep10032>
6. Castellano C, Fortunato S, Loreto V (2009) Statistical physics of social dynamics. *Rev Mod Phys* 81(2):591–646. <https://doi.org/10.1103/revmodphys.81.591>
7. Degroot M (1974) Reaching a consensus. *J Am Stat Assoc* 69:118–121
8. Deffuant G, Neau D, Amblard F, Weisbuch G (2000) Mixing beliefs among interacting agents. *Adv Complex Syst* 3(01n04):87–98
9. Chen X, Tsaparas P, Lijffijt J, Bie TD (2021) Opinion dynamics with backfire effect and biased assimilation. *PLoS ONE* 16:e0256922
10. Monti C, De Francisci Morales G, Bonchi F (2020) Learning opinion dynamics from social traces. In: *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*, pp 764–773
11. Nyhan B, Reifler J (2010) When corrections fail: the persistence of political misperceptions. *Polit Behav* 32:303–330
12. Allahverdyan AE, Galstyan A (2014) Opinion dynamics with confirmation bias. *PLoS ONE* 9(7):99557. <https://doi.org/10.1371/journal.pone.0099557>
13. Liu L, Wang X, Chen X, Tang S, Zheng Z (2021) Modeling confirmation bias and peer pressure in opinion dynamics. *Front Phys* 9:649852
14. Sirbu A, Pedreschi D, Giannotti F, Kertész J (2019) Algorithmic bias amplifies opinion fragmentation and polarization: a bounded confidence model. *PLoS ONE* 14:e0213246
15. Pansanella V, Sirbu A, Kertész J, Rossetti G (2023) Mass media impact on opinion evolution in biased digital environments: a bounded confidence model. *Sci Rep* 13(1):14600
16. Monti C, Aiello LM, De Francisci Morales G, Bonchi F (2022) The language of opinion change on social media under the lens of communicative action. *Sci Rep* 12(1):17920
17. Park JS, Popowski L, Cai C, Morris MR, Liang P, Bernstein MS (2022) Social simulacra: creating populated prototypes for social computing systems. In: *Proceedings of the 35th annual ACM symposium on user interface software and technology*, pp 1–18
18. Park JS, Zou CQ, Shaw A, Hill BM, Cai C, Morris MR, Willer R, Liang P, Bernstein MS (2024) Generative agent simulations of 1,000 people. *arXiv preprint*. [arXiv:2411.10109](https://arxiv.org/abs/2411.10109)
19. Premack D, Woodruff G (1978) Does the chimpanzee have a theory of mind? *Behav Brain Sci* 1(4):515–526. <https://doi.org/10.1017/s0140525x00076512>
20. Kosinski M (2023) Theory of mind may have spontaneously emerged in large language models. *arXiv preprint*. [arXiv:2302.02083](https://arxiv.org/abs/2302.02083)
21. Street W, Siy JO, Keeling G, Baranes A, Barnett B, McKibben M, Kanyere T, Lentz A, Dunbar RI, et al (2024) LLMs achieve adult human performance on higher-order theory of mind tasks. *arXiv preprint*. [arXiv:2405.18870](https://arxiv.org/abs/2405.18870)
22. Li H, Chong YQ, Stepputtis S, Campbell J, Hughes D, Lewis M, Sycara K (2023) Theory of mind for multi-agent collaboration via large language models. *arXiv preprint*. [arXiv:2310.10701](https://arxiv.org/abs/2310.10701)
23. Ullman T (2023) Large language models fail on trivial alterations to theory-of-mind tasks. *arXiv preprint*. [arXiv:2302.08399](https://arxiv.org/abs/2302.08399)
24. Sap M, Le Bras R, Fried D, Choi Y (2022) Neural theory-of-mind? On the limits of social intelligence in large LMs. In: *Goldberg Y, Kozareva Z, Zhang Y (eds) Proceedings of the 2022 conference on empirical methods in natural language processing*, vol 248. Association for Computational Linguistics, Abu Dhabi, pp 3762–3780. <https://doi.org/10.18653/v1/2022.emnlp-main.248>. <https://aclanthology.org/2022.emnlp-main>
25. Shapira N, Zwirn G, Goldberg Y (2023) How well do large language models perform on faux pas tests? In: *Rogers A, Boyd-Graber J, Okazaki N (eds) Findings of the association for computational linguistics: ACL 2023*. Association for Computational Linguistics, Toronto, pp 10438–10451. <https://doi.org/10.18653/v1/2023.findings-acl.663>
26. De Marzo G, Pietronero L, Garcia D (2023) Emergence of scale-free networks in social interactions among large language models. *arXiv preprint*. [arXiv:2312.06619](https://arxiv.org/abs/2312.06619)
27. Gao C, Lan X, Lu Z, Mao J, Piao J, Wang H, Jin D, Li Y (2023) S3: social-network simulation system with large language model-empowered agents. *arXiv preprint*. [arXiv:2307.14984](https://arxiv.org/abs/2307.14984)
28. Ashery AF, Aiello LM, Baronchelli A (2025) Emergent social conventions and collective bias in LLM populations. *Sci Adv* 11(20):9368
29. Mou X, Wei Z, Huang X (2024) Unveiling the truth and facilitating change: towards agent-based large-scale social movement simulation. *arXiv preprint*. [arXiv:2402.16333](https://arxiv.org/abs/2402.16333)
30. Wang C, Liu Z, Yang D, Chen X (2025) Decoding echo chambers: LLM-powered simulations revealing polarization in social networks. In: *Proceedings of the 31st international conference on computational linguistics*. Association for Computational Linguistics, Abu Dhabi, pp 3913–3923. <https://aclanthology.org/2025.coling-main.264/>
31. Chuang Y-S, Goyal A, Harlalka N, Suresh S, Hawkins R, Yang S, Shah D, Hu J, Rogers TT (2023) Simulating opinion dynamics with networks of LLM-based agents. *arXiv preprint*. [arXiv:2311.09618](https://arxiv.org/abs/2311.09618)
32. Breum SM, Egdal DV, Mortensen VG, Møller AG, Aiello LM (2023) The persuasive power of large language models. *arXiv preprint*. [arXiv:2312.15523](https://arxiv.org/abs/2312.15523)
33. Flamino J, Modi MS, Szymanski BK, Cross B, Mikolajczyk C (2024) Limits of large language models in debating humans. *arXiv preprint*. [arXiv:2402.06049](https://arxiv.org/abs/2402.06049)
34. Priya P, Firdaus M, Ekbal A (2024) Computational politeness in natural language processing: a survey. *ACM Comput Surv* 56(9):241. <https://doi.org/10.1145/3654660>
35. Törnberg P, Valeeva D, Uitermark J, Bail C (2023) Simulating social media using large language models to evaluate alternative news feed algorithms. *arXiv preprint*. [arXiv:2310.05984](https://arxiv.org/abs/2310.05984)
36. Cheng M, Yu S, Lee C, Khadpe P, Ibrahim L, Jurafsky D (2025) Social sycophancy: a broader understanding of LLM sycophancy. *arXiv preprint*. [arXiv:2505.13995](https://arxiv.org/abs/2505.13995)
37. Taubenfeld A, Dover Y, Reichart R, Goldstein A (2024) Systematic biases in LLM simulations of debates. *arXiv preprint*. [arXiv:2402.04049](https://arxiv.org/abs/2402.04049)
38. Fanous A, Goldberg J, Agarwal AA, Lin J, Zhou A, Daneshjou R, Koyejo S (2025) Syceval: evaluating LLM sycophancy. *arXiv preprint*. [arXiv:2502.08177](https://arxiv.org/abs/2502.08177)
39. Ranaldi L, Pucci G (2023) When large language models contradict humans? Large language models' sycophantic behaviour. *arXiv preprint*. [arXiv:2311.09410](https://arxiv.org/abs/2311.09410)

40. Khan AA, Alam S, Wang X, Khan AF, Neog DR, Anwar A (2024) Mitigating sycophancy in large language models via direct preference optimization. In: 2024 IEEE international conference on big data (BigData). IEEE, pp 1664–1671
41. Aher GV, Arriaga RI, Kalai AT (2023) Using large language models to simulate multiple humans and replicate human subject studies. In: International conference on machine learning, pp 337–371. PMLR
42. Payandeh A, Pluth D, Hosier J, Xiao X, Gurbani VK (2023) How susceptible are LLMs to logical fallacies?
43. Likert R (1932) A technique for the measurement of attitudes. *Arch Psychol* 22 140:55
44. Chuang Y-S, Goyal A, Harlalka N, Suresh S, Hawkins R, Yang S, Shah D, Hu J, Rogers T (2024) Simulating opinion dynamics with networks of LLM-based agents. In: Duh K, Gomez H, Bethard S (eds) Findings of the Association for Computational Linguistics: NAACL 2024. Association for Computational Linguistics, Mexico City, pp 3326–3346. <https://aclanthology.org/2024.findings-naacl.211>
45. Friedkin NE (1986) A formal theory of social power. *J Math Sociol* 12:103–126
46. Hegselmann R, Krause U (2002) Opinion dynamics and bounded confidence: models, analysis and simulation. *J Artif Soc Soc Simul* 5
47. Ju D, Williams A, Karrer B, Nickel M (2024) Sense and sensitivity: evaluating the simulation of social dynamics via large language models. *arXiv preprint*. [arXiv:2412.05093](https://arxiv.org/abs/2412.05093)
48. Jiang AQ, Sablayrolles A, Mensch A, Bamford C, Chaplot DS, Casas D, Bressand F, Lengyel G, Lample G, Saulnier L, Lavaud LR, Lachaux M-A, Stock P, Scao TL, Lavril T, Wang T, Lacroix T, Sayed WE (2023). Mistral 7B. *arXiv preprint*. <https://arxiv.org/abs/2310.06825>
49. Dubey A, Jauhri A, Pandey A, Kadian A, Al-Dahle A, Letman A, Mathur A, Schelten A, Yang A, Fan A, et al (2024) The Llama 3 herd of models. *arXiv preprint*. [arXiv:2407.21783](https://arxiv.org/abs/2407.21783)
50. Welch BL (1947) The generalization of 'student's problem when several different population variances are involved. *Biometrika* 34(1–2):28–35
51. Microsoft: Microsoft/autogen. <https://github.com/microsoft/autogen>. Accessed 2025-02-13
52. Cau E, Failla A, Rossetti G (2024) Bots of a feather: mixing biases in LLMs' opinion dynamics. In: International conference on complex networks and their applications. Springer, Berlin, pp 166–176
53. Manabat BK q3fer/distilbert-base-fallacy-classification · hugging face. <https://huggingface.co/q3fer/distilbert-base-fallacy-classification>. Accessed 2025-02-12
54. Jin Z, Lalwani A, Vaidhya T, Shen X, Ding Y, Lyu Z, Sachan M, Mihalcea R, Schölkopf B (2022) Logical fallacy detection. *arXiv preprint*. [arXiv:2202.13758](https://arxiv.org/abs/2202.13758)
55. Copi IM, Cohen C, MacMahon K (2013) Introduction to logic, 14 edn. Pearson custom library. Pearson Education, Upper Saddle River
56. Sirbu A, Loreto V, Servedio VD, Tria F (2017) Opinion dynamics: models, extensions and external effects. In: Participatory sensing, opinions and collective awareness. Springer, Cham, pp 363–401
57. Oviedo-Trespalcacios O, Peden AE, Cole-Hunter T, Costantini A, Haghani M, Rod JE, Kelly S, Torkamaan H, Tariq A, Newton JDA, Gallagher T, Steinert S, Filtiness AJ, Reniers G (2023) The risks of using chatgpt to obtain common safety-related information and advice. *Saf Sci* 167:106244. <https://doi.org/10.1016/j.ssci.2023.106244>
58. Li Y, Wang D, Liang J, Jiang G, He Q, Xiao Y, Yang D (2024) Reason from fallacy: enhancing large language models' logical reasoning through logical fallacy understanding. *arXiv preprint*. [arXiv:2404.04293](https://arxiv.org/abs/2404.04293)
59. Mouchel L, Paul D, Cui S, West R, Bosselut A, Faltings B (2024) A logical fallacy-informed framework for argument generation. *arXiv preprint*. [arXiv:2408.03618](https://arxiv.org/abs/2408.03618)
60. Cava LL, Tagarelli A (2024) Open models, closed minds? On agents capabilities in mimicking human personalities through open large language models. *arXiv preprint*. [arXiv:2401.07115](https://arxiv.org/abs/2401.07115)
61. Huang J-T, Lam MH, Li EJ, Ren S, Wang W, Jiao W, Tu Z, Lyu MR (2024) Emotionally numb or empathetic? Evaluating how LLMs feel using EmotionBench. *arXiv preprint*. [arXiv:2308.03656](https://arxiv.org/abs/2308.03656)
62. Bender EM, Gebru T, McMillan-Major A, Shmitchell S (2021) On the dangers of stochastic parrots: can language models be too big? In: Proceedings of the 2021 ACM conference on fairness, accountability, and transparency, pp 610–623
63. Cao Y, Zhou L, Lee S, Cabello L, Chen M, Hershovich D (2023) Assessing cross-cultural alignment between ChatGPT and human societies: an empirical study. *arXiv preprint*. [arXiv:2303.17466](https://arxiv.org/abs/2303.17466)
64. Masoud RI, Liu Z, Ferianc M, Treleaven P, Rodrigues M (2023) Cultural alignment in large language models: an explanatory analysis based on Hofstede's cultural dimensions. *arXiv preprint*. [arXiv:2309.12342](https://arxiv.org/abs/2309.12342)
65. Liu Z (2025) Cultural bias in large language models: a comprehensive analysis and mitigation strategies. *J Transcult Commun* 3(2):224–244
66. Tao Y, Viberg O, Baker RS, Kizilcec RF (2024) Cultural bias and cultural alignment of large language models. *PNAS Nexus* 3(9):346
67. Agarwal D, Naaman M, Vashistha A (2025) Ai suggestions homogenize writing toward western styles and diminish cultural nuances. In: Proceedings of the 2025 CHI conference on human factors in computing systems, pp 1–21
68. Wang Z, Chiu YY, Chiu YC (2023) Humanoid agents: platform for simulating human-like generative agents. In: Feng Y, Lefever E (eds) Proceedings of the 2023 conference on empirical methods in natural language processing: system demonstrations. Association for Computational Linguistics, Singapore, pp 167–176. <https://doi.org/10.18653/v1/2023.emnlp-demo.15>. <https://aclanthology.org/2023.emnlp-demo.15/>
69. Piao J, Lu Z, Gao C, Xu F, Santos FP, Li Y, Evans J (2025) Emergence of human-like polarization among large language model agents. *arXiv preprint*. [arXiv:2501.05171](https://arxiv.org/abs/2501.05171)
70. Zheng W, Tang X (2025) Simulating social network with LLM agents: an analysis of information propagation and echo chambers. In: Tang X, Huynh VN, Xia H, Bai Q (eds) Knowledge and systems sciences. Springer, Berlin, pp 63–77. [https://doi.org/10.1007/978-981-96-0178-3\\_5](https://doi.org/10.1007/978-981-96-0178-3_5)

## Publisher's note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.