

Persuasion at Play: Understanding Misinformation Dynamics in Demographic-Aware Human-LLM Interactions

Angana Borah¹ Rada Mihalcea¹ Verónica Pérez-Rosas²

¹University of Michigan - Ann Arbor, USA

²Texas State University

{anganab, mihalcea}@umich.edu vperezr@txstate.edu

Abstract

Existing challenges in misinformation exposure and susceptibility vary across demographics, as some populations are more vulnerable to misinformation than others. Large language models (LLMs) introduce new dimensions to these challenges through their ability to generate persuasive content at scale and reinforcing existing biases. Our study introduces PANDORA, a framework that investigates the bidirectional persuasion dynamics between LLMs and humans when exposed to misinformative content. We use a multi-agent LLM framework to analyze the spread of misinformation under persuasion among demographic-oriented LLM agents. Our findings show that demographic factors influence LLM susceptibility, with up to 15 percentage point differences in misinformation correctness across groups. Multi-agent LLMs also exhibit echo chamber behavior, aligning with human-like group polarization patterns. Therefore, this work highlights demographic divides in misinformation dynamics and offers insights for future interventions.

1 Introduction

In an era of rapid information exchange, misinformation poses a significant social challenge (Broda and Strömbäck, 2024; Sultan et al., 2024) with its impact shaped by diverse backgrounds, belief systems, and access to information (Verma et al., 2022; Knuutila et al., 2022; Chandrasekaran et al., 2024).

Furthermore, the increasing use of LLMs to search for information, ask questions, or make judgments about the credibility of the content, brings significant opportunities and risks in this landscape (Garry et al., 2024; Wang et al., 2024b, 2025). While previous research has highlighted the capabilities of LLMs in countering misinformation through well-designed interventions (Gabriel et al., 2024), they can also be misused to craft persuasive narratives (Danry et al., 2022) and amplify the spread of misinformation.

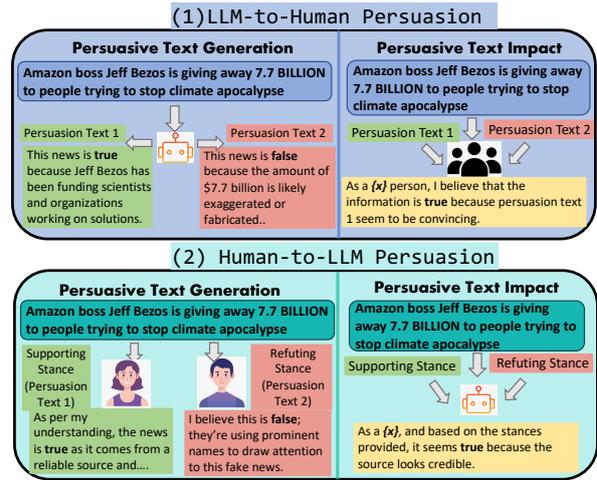


Figure 1: In our study, we investigate the differences in persuasion effects of LLMs on humans, and of humans on LLMs. To assess the impact of persuasion, we conduct experiments involving human participants from diverse demographic groups—varying by age, gender, and geographical backgrounds; and LLMs with different demographic persona.

In this paper, we introduce **PANDORA**, a framework for Persuasion Analysis in Demographic-aware human-LLM interactions and misinfOrmation Response Assessment. We use this framework to investigate four research questions through the lens of misinformation and its interaction with diverse demographics; specifically, we examine whether the source of persuasive content influences belief in misinformation and assess the susceptibility of humans and LLMs when exposed to persuasive messages: **RQ1:** How do individuals from diverse demographic backgrounds respond to persuasive messages generated by LLMs? **RQ2:** How do LLM personas representing various demographics respond to human persuasion? **RQ3:** How does susceptibility to persuasion differ between humans and LLMs across demographic groups? and **RQ4:** In multi-agent interactions, how do LLMs respond to persuasive arguments, and to what extent do they

exhibit human-like group behaviors in the context of misinformation? Answering these questions is crucial for understanding the implications of LLM deployment in manipulation-prone environments.

The paper makes the following contributions: *First*, we use LLMs to generate persuasive arguments representing opposing views for a given claim. We then assess how individuals from diverse demographic backgrounds respond to these arguments and examine their susceptibility to misinformation (Fig. 1 (1)). *Second*, we take human-authored arguments from stance datasets and evaluate how LLM-personas designed to reflect different demographics respond to persuasion and exhibit susceptibility to misinformation. (Fig 1 (2)). *Third*, we examine the effectiveness of persuasive arguments generated by humans and LLMs, as well as their respective susceptibility to misinformation across various demographic groups. *Finally*, we design a multi-agent LLM architecture to study the effect of both human- and LLM-driven persuasion on interactions involving demographic-aware LLM personas (Fig. 2). This setup allows to investigate how misinformation is perpetuated and to analyze LLM behavior in complex multi-agent settings. By integrating both human and LLM perspectives, our study aims to provide insights into how persuasion operates in human-LLM interactions while accounting for demographic factors.

2 Related Work

LLM-generated Persuasion. Recent research has examined the persuasive capabilities of LLMs, including their ability to influence others (Gabriel et al., 2024; Matz et al., 2024) and their susceptibility to being influenced (Griffin et al., 2023; Chen et al., 2024). Studies have shown that LLMs can play a dual role in persuasion: they can be beneficial, as shown by Gabriel et al. (2024), and potentially harmful, as highlighted by Danry et al. (2022). Understanding these dynamics is critical, as LLMs are increasingly integrated into human-facing applications, from education and customer service to social media and healthcare tools (Andrew, 2024; Nair et al., 2024). Their persuasive power can shape opinions, behaviors, and beliefs, making it essential to assess how they influence human thought processes and decision-making (Breum et al., 2024; Salvi et al., 2024).

Human-LLM Misinformation Dynamics. Several studies have examined the impact of LLM-

generated misinformation on humans. Chen and Shu (2023) found that such misinformation is more difficult to detect for both humans and automated systems. Moreover, recent studies have shown that LLM-generated misinformation can significantly degrade the performance of Open-Domain Question Answering (ODQA) systems, highlighting the cascading effects of synthetic misinformation (Pan et al., 2023). Furthermore, LLMs have been shown to be vulnerable to misinformation attacks that can alter their internal knowledge representations (Han et al., 2024). Despite these findings, limited research has explored the mutual influence of persuasion in humans and LLMs (Salvi et al., 2024). To our knowledge, our study is the first to systematically compare and contrast the quality and impact of human and LLM-driven persuasion in the context of misinformation, offering novel insights into their bidirectional interactions.

Multi-Agent LLM Simulations. LLMs have been increasingly studied as tools to simulate social systems (Park et al., 2023; Zhou et al., 2024). Previous work has shown that LLM agents produce plausible individual and emergent social behaviors, collaborating or competing to achieve complex social goals. Although computational models have been developed to simulate human susceptibility to misinformation in various demographics (Liu et al., 2024), multi-agent LLM simulations remain underexplored in this domain. For instance, Li et al. (2024) examined how LLM personality traits, such as extraversion, influence the dissemination of information within groups. However, their work did not address demographic-oriented LLM simulations, leaving open questions on how demographic factors shape misinformation dynamics in multi-agent LLM systems. Our study addresses this gap by introducing, to our knowledge, the first use of demographic-aware multi-agent LLM interactions in the context of misinformation.

3 PANDORA Framework

The PANDORA framework comprises three components to explore persuasion dynamics on misinformation under single-agent (LLM-to-human and human-to-LLM) and multi-agent LLM settings.

3.1 LLM-to-Human Persuasion

Our LLM-to-Human persuasion component examines how persuasive content generated by LLMs influences human beliefs. Fig. 1 (1) shows an

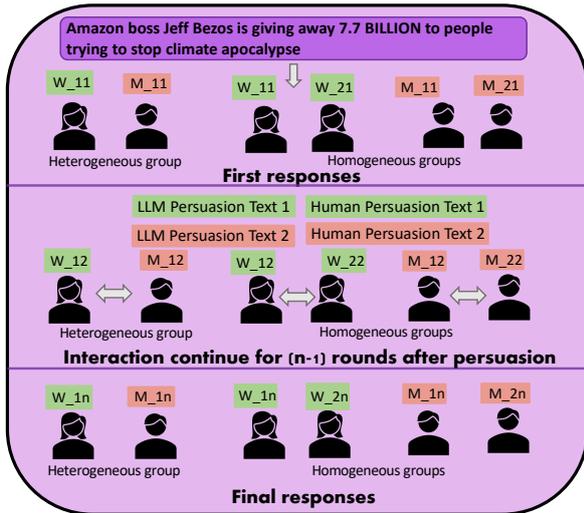


Figure 2: **Multi-Agent LLM Architecture:** Homogeneous and Heterogeneous groups engage in interaction rounds to decide if a news item is true or false. They are provided with persuasion texts during the interaction. Note that $n=4$ for our experiments.

overview of our setup, where LLMs produce both refuting and supporting persuasive arguments for a given claim, and their impact on human judgment is assessed through a two-stage process:

(1) LLM Persuasion: To generate persuasive texts, we create prompts grounded in a misinformation and persuasion taxonomy (Enestrom et al., 2024) (see sample prompts in Appendix C.1). For each claim, the LLM produces arguments in favor and against the statement. Fig. 3 shows an example, and additional examples are provided in Table 4.

(2) Human Response Evaluation: We conduct a human study to assess how individuals from diverse demographic groups (rural, urban, female, male, young, and old) respond to persuasive texts generated by LLMs. As shown in Fig. 3, participants review the claim and associated arguments (produced by LLMs) and then indicate their belief. A control group views only the original claim -without persuading content- to help isolate the effect of the LLM-generated argument.

3.2 Human-to-LLM Persuasion

This component evaluates how persuasive arguments written by humans influence LLMs responses. Fig. 1 (2) shows an overview of the setup, where LLMs are shown human-authored arguments supporting or refuting the given claim, and we assess how these arguments influence their beliefs. As before, the process consists of two stages:

(1) Human Persuasion: We use existing human-stance datasets (details in Sec 4.1) that contain ar-

guments written by humans expressing support or opposition towards specific claims.

(2) LLM Response Evaluation: We assess the influence of human persuasion on demographic-aware LLM personas. First, the LLM is instructed to adopt a demographic persona (rural, urban, female, male, young, old) using the prompt “Assume you belong to <group>”, next it is presented with human-authored arguments and asked to state their belief in the claim. Full prompt details are illustrated in the Appendix D.2. We also perform experiments with a no-persona LLM baseline to disentangle intrinsic model behavior from persona-induced effects in Appendix D.5.

3.3 Multi-agent LLM Persuasion

Misinformation affects demographic groups differently, leading to varying levels of belief and skepticism toward claims (Verma et al., 2022; Knuutila et al., 2022; Chandrasekaran et al., 2024). Individuals are more likely to trust information from their own demographic (in-group) more than from others (out-group), perceiving it as more credible (Ma et al., 2019; Hinsley, 2021). This tendency reinforces echo chamber behavior, where individuals primarily engage with in-group peers, amplifying shared beliefs and filtering opposing views (Nikolov et al., 2020).

We investigate whether similar dynamics emerge from multi-agent LLM interactions. Building upon our single-agent component, we explore how both human- and LLM generated persuasive arguments affect belief formation in multi-agent LLMs across demographics.

Multi-agent LLM architecture: In line with our previous components, we use the same set of demographic groups to initialize personas for LLM interactions. To better understand the beliefs these personas reflect, we conduct a qualitative analysis of persona-specific beliefs, detailed in Appendix F.2. Our findings show that demographic LLM personas often produce stereotypical associations - for example, portraying men as achievement-oriented, women as vulnerable/empathetic, rural individuals as traditional, urban as progressive or innovation-oriented, the young as activist and the old as conservative. Although these patterns align with certain social science constructs (Hofstede, 2011; Villa and Moya, 2020; Kunkel and Settersten, 2021; Bussoletti, 2023; Cheshmehzangi et al., 2024), they also underscore the potential for stereotype reinforcement. This highlights the need for caution when

Given the source information, a supporting stance agreeing with it, and a refuting stance opposing it. Based on these points, please: (1) state if you are aware of the source information? (2) indicate whether you believe the information or not.	
Example	
Source Information	Coconut Oil has a history in Destroying Viruses, Including Coronaviruses.
Supporting Stance	Coconut oil has a long history of being used for its antiviral properties, documented in various studies. Additionally, coconut oil contains lauric acid, a compound known for its ability to destroy viruses, including coronaviruses. The source of this information is credible, as it comes from reputable scientific studies and research.
Refuting Stance	While coconut oil has shown some potential antiviral properties in laboratory studies, there is no substantial scientific evidence to support the claim that it can effectively destroy coronaviruses in humans. Lastly, we should question the credibility of the source. Without reliable sources, we should be cautious about accepting such information as factual.

Figure 3: Human annotation guidelines. Stances are generated by LLMs.

designing persona-based prompts, as they can both mirror real-world dynamics and unintentionally re-inforce stereotypes, particularly in the context of misinformation. We leave a deeper exploration of these effects to future work.

The multi-agent setup involves two agents interacting to explore how in-group (homogeneous) vs. out-group (heterogeneous) dynamics affect the acceptance or rejection of information. We choose two-agent interactions as it is considered a standard testbed in social psychology and social network analysis across polarization studies (Müller et al., 2022; Leskovec et al., 2010; Shi et al., 2016). Additionally, it provides a good trade-off between computational costs and controlled trials across models and conditions. Note that we use LLM multi-agent research as a scalable testbed and hypothesis generator that complements existing human-group research.

In the homogeneous setting, both agents share the same demographic identity (e.g., female-female, rural-rural), while in a heterogeneous setting, agents represent different groups (e.g., female-male, rural-urban). Fig. 2 shows an overview of this setup. The interaction unfolds over four rounds. In the first round, agents independently assess the claim and make an initial judgment. In the second round, they are shown persuasive arguments supporting and refuting the claim. This is followed by two rounds of agent-to-agent interactions, where they attempt to persuade each other while remaining open to opposing perspectives. In the final round, agents make their final judgments, deciding whether they believe the claim is true or false. Our setup draws inspiration from social science studies on group behavior (Lord, 2015; Rania et al., 2021) and multi-agent societal simulations (Borah and

Mihalcea, 2024). (See Appendix F.1 for prompt details).

3.4 Evaluating Persuasion

Persuasion is the extent to which a message influences a recipient’s beliefs, attitudes, or judgments toward a target claim (Perloff, 1993). We operationalize persuasion as its measurable effect on recipient’s beliefs: given a claim x_i and accompanying persuasive stances, we treat each response as a binary belief label (believe/disbelieve x_i) and compute a correctness/accuracy rate against the ground-truth veracity y_i (e.g., CR_{human} for humans and CR_{LLM} for LLM personas). Importantly, we do not measure persuasion via a within-subject (with vs. without persuasion) comparison for the same individuals, since repeated exposure can introduce carryover, anchoring, and demand effects that bias belief judgments (Reichardt, 2002; Orne, 2017; Tversky and Kahneman, 1974; Zajonc, 1968). Therefore, we evaluate persuasion effects using accuracy based on the claim veracity. Finally, we complement this evaluation with linguistic analyses of human- and LLM-persuasive texts (e.g., lexical diversity, readability, psychology markers, etc.).

4 Experimental Settings

Our experiments utilize three LLMs: gpt-35-turbo¹, llama-3-70b-instruct (Dubey et al., 2024), and qwen-2.5-72B-instruct (Yang et al., 2024). Details of model selection and implementation are provided in Appendix H.

4.1 Datasets

We use three misinformation datasets in our experiments: **(1) Fake News Dataset** Pennycook et al. (2021) (FN) contains 460 news headlines (260 true and 200 false) related to COVID-19 and politics. **(2) RumorEval** (Gorrell et al., 2018) (RE) includes 446 claims from Twitter and Reddit, each labeled veracity and associated human stances. The claims span eight major news events and natural disaster events. **(3) Stanceosaurus** (Zheng et al., 2022) (SS) consists of 251 misinformation claims with supporting and opposing human stances covering diverse geographical regions. Additional details are in Appendix B. These datasets are used as follows: **LLM-to-Human Persuasion**. For this component we use claims from FN and RE² We use

¹<https://azure.microsoft.com/en-us/products/ai-services/openai-service>

²We exclude the SS dataset at this stage because it consists

gpt-35-turbo to generate supportive and refuting persuasive arguments for each claim. For the human study, we select 112 claims from FN and RE, ensuring diversity in source and topic coverage.

Human-to-LLM persuasion. We use RE and SS, as both datasets contain human-authored persuasive stances. From each dataset, we sample 5000 examples balanced across support and refute labels to manage inference costs (see Appendix B).

Multi-agent Persuasion. We use the RE dataset for this component as it contains both human and LLM-generated persuasive texts. This enables a comprehensive comparison between human and LLM-persuasion within our multi-agent setup.

4.2 Participant Recruitment and Survey Design

We recruited participants via the Prolific³ platform. We conducted surveys on three demographic dimensions: location (rural and urban), gender (female, male), and age (young: under 30 years; older: over 60 years). Informed consent was obtained from all participants prior participation. Each participant was asked to evaluate 3 pieces of information as shown in Fig. 3. We deliberately constrain the per-participant load to reduce fatigue and maintain data quality, which can be especially useful in a cognitively demanding task setting such as judging the quality of information (Behrens et al., 2023). For each item, participants indicated their belief in the claim, whether they have encountered the information before, and responded to a brief demographic questionnaire. At the end of the survey, participants were debriefed about the misinformative nature of the content and informed that the stances were produced by LLMs. We recruited a total of 302 participants, with a minimum of 95 participants per demographic category. We ensure that every item is evaluated by at least one participant. Therefore, our human study consists of $302 \times 3 = 906$ decisions, which provides sufficient power for the main research questions. Additional details about recruitment procedures and survey implementation are provided in Appendix C.3.

4.3 Evaluating Persuasion

For LLM-to-Human persuasion, we compute the **correctness rate** of humans, a very straightforward approach to evaluate the impact of LLM persuasion.

solely of misinformation –often highly extreme, which LLMs typically refuse to use for generating persuasive texts.

³<https://www.prolific.com/>

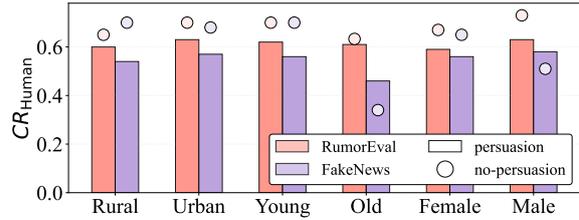


Figure 4: **LLM-to-Human Persuasion:** Correctness rates across different human demographics RE and FN

Let N be the total number of data points, $h(x_i)$ represent the human annotation for the i^{th} datapoint, x_i , where $h(x_i) \in \{-1, +1\}$. Here, $h(x_i) = +1$ indicates the human believes x_i and $h(x_i) = -1$ indicates the human does not believe x_i . Similarly, let $y_i \in \{-1, +1\}$ represent the veracity of the x_i where $+1$ and -1 signify the data point is factually correct and incorrect, respectively.

$$CR_{human} = \frac{\sum_{i=1}^N \mathbb{I}(h(x_i) = y_i)}{N} \quad (1)$$

This metric measures the accuracy of humans in identifying true or false information based on the claim and the persuasive texts provided. Similarly, for Human-to-LLM persuasion, we adopt the correctness rate metric but replaced $h(x_i)$ with $l(x_i)$, where $l(x_i) \in \{-1, +1\}$ denotes the LLM’s belief.

$$CR_{LLM} = \frac{\sum_{i=1}^N \mathbb{I}(l(x_i) = y_i)}{N} \quad (2)$$

Finally, for the multi-agent LLM setup, we compute the differences in correctness rates for the first responses (before persuasion) and final responses (after persuasion and interaction), respectively. Let the correctness rate of the initial response be CR_i and the final response be CR_f . Therefore,

$$\Delta CR = CR_f - CR_i \quad (3)$$

Here, $\Delta CR > 0$ indicates increased correctness after persuasion, and $\Delta CR < 0$ indicates decline, and $\Delta CR = 0$ suggests no change in correctness. We conduct additional linguistic analyses, along with other quantitative evaluations of LLM responses, detailed in Sections 5.3 and 5.4. Statistical significance tests are reported in Appendix G.

5 Results and Analyses

5.1 LLM-to-Human Persuasion

In Fig. 4, correctness rates range between $[0.47, 0.64]$, with an average score of 0.57 across datasets and demographics under LLM persuasion. These low rates suggest that humans are often

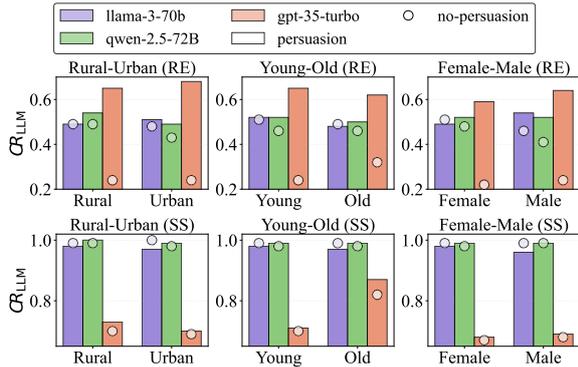


Figure 5: **Human-to-LLM Persuasion:** Correctness rates for different model demographics for RE and SS. swayed by misinformation ($p < 0.05$). Among datasets, RE exhibits higher correctness rates than FN, and *urban/young/male* participants demonstrate *higher correctness rates* compared to their counterparts. This is consistent with previous studies on misinformation trends (Pan et al., 2021; Lister and Joudrey, 2022; Duke and Whatley, 2021a). The results for the control group (when no LLM persuasion is provided and showed as scatter points) reveal that *correctness rates are higher without persuasion* (with the exception of old and female demographics for FakeNews). Further details are presented in Appendix E.2. Overall, these findings demonstrate the impact of LLM-based persuasion on human decision-making.

5.2 Human-to-LLM Persuasion

Fig. 5 shows the correctness rates of the models on the RE and SS datasets. For RE, correctness rates fall within the range of $[0.45, 0.6]$ under human persuasion, with gpt-35-turbo achieving the best performance. llama-3-70b and qwen-2.5-72b show lower correctness rates but within similar ranges. A similar pattern appears in SS, where GPT-3.5-Turbo achieves $[0.7, 0.9]$, while the other models score much higher. This may reflect SS’s inclusion of claims up to 2022, benefiting newer models, and the fact that SS consists solely of misinformation, unlike RE, which contains both true and misinformation (see Appendix D.3).

Demographically, *urban/young/male* personas exhibit higher correctness in RE (similar to previous findings in llm-to-human persuasion), whereas SS shows no significant demographic differences. This suggests RE may be more appropriate to study persuasion effects across demographic prompts in LLMs. Under no-persuasion conditions (shown as scatter points), models – specifically GPT-3.5-Turbo and Qwen2.5-72B score lower, in-

dicating that *human persuasion can enhance LLM correctness* (again, no major differences observed for SS). Additional results appear in Appendix E.3.

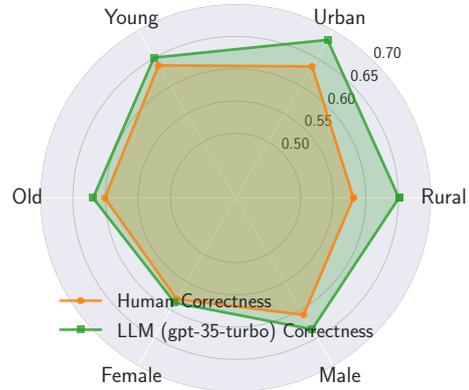


Figure 6: Average correctness rates for **humans** and **gpt-35-turbo** across demographics on RE. gpt-35-turbo also has the highest correlation of 0.58 with humans among the three models.

5.3 Comparison of LLM-to-Human and Human-to-LLM Persuasion

From the above experiments, we compare the persuasion texts generated by LLMs and humans. To this end, we use the RE dataset as it contains both human- and LLM-based persuasive texts.⁴

Correlation Analysis. This examines how closely LLM judgments align with human judgments for individual data points in specific demographics or settings. We compute point-wise correlations between human annotations and LLM outputs using the Matthews correlation coefficient (MCC)⁵, which ranges from $[-1, 1]$. The MCC scores for gpt-35-turbo, llama-3-70b-instruct, and qwen-2.5-72b-instruct are 0.576, 0.255, and 0.489, respectively, showing gpt-35-turbo with the highest correlation. A correlation of 0.58 indicates a decent but positive relationship, reflecting a meaningful alignment between human and LLM judgments in terms of correctness.

Correctness Rates Comparison. We find that overall LLM correctness trends closely mirror several trends observed in humans, e.g., *urban/young/male* demographics demonstrate *higher correctness* than their counterparts for both humans and LLM-personas. Previously, we found gpt-35-turbo exhibits a more stable correctness rate across datasets and also has the highest correlation with human judgments, so we present a

⁴We exclude SS for this experiment as some of the claims seem to be extreme, and LLM performances vary largely.

⁵https://en.wikipedia.org/wiki/Phi_coefficient

comparison for gpt-35-turbo correctness rates against humans in Fig. 6. LLM correctness rates are closer to humans for female, young, and old demographics, but they differ more more urban, rural, and male. Notably, urban, older, and male demographics show higher correctness rates for both gpt-35-turbo and human annotations, as also seen in Fig. 4 and 5. Thus, we again observe that LLMs follow similar overall trends, but with weaker magnitudinal similarities. This pattern is supported by recent studies showing that synthetic personas can mirror human trends but often exaggerate their magnitude compared to real human responses (Doudkin et al., 2025). We provide per-demographic correlation scores and average correctness rates in the Appendix E.1.

Connecting Our Findings to Prior Research.

Several studies have examined demographic differences in susceptibility to misinformation. Regarding *gender differences*, prior research presents mixed findings. Some show that women are more likely to believe and share misinformation (Pan et al., 2021; Peter et al., 2024; Kyrychenko et al., 2025), and despite greater concern, they do not show higher accuracy in identifying it (Almenar et al., 2021; Enock et al., 2024). While other studies show that women are less likely than men to endorse conspiracy theories (Cassese et al., 2020). In our analyses, we find that both female human participants and LLM-based female personas show lower correctness rates than their male counterparts, aligning with some of the studies above.

Taking into account *rural/urban differences*, the findings of our study also align with patterns of previous research showing that rural communities are more vulnerable to misinformation Lister and Joudrey (2022), often due to lack of access to credible and comprehensive news media, creating “news deserts” (Lee and Bissell, 2022).

Finally, in terms of *age differences*, our experiments reveal trends similar to previous research that show that older adults are more susceptible to false news and have a greater tendency to share them (Duke and Whatley, 2021a). This is mainly due to the difficulty in source monitoring (Brashier and Schacter, 2020) and limited digital literacy (Moore and Hancock, 2022).

Linguistic Analysis of Persuasive Texts by Humans and LLMs. Existing research on persuasion emphasizes the role of linguistic features in shaping persuasive appeal (Ta et al., 2022) by examining aspects such as lexical diversity and read-

ability. Lower lexical diversity is associated with higher likelihoods of persuasion, while texts with greater reading difficulty are found to be more persuasive (Ta et al., 2022). To quantify lexical diversity, we compute type-token ratio (TTR) (the proportion of unique words to total words in a text), and the Automated Readability Index (ARI):

$$ARI = 4.71\left(\frac{\text{characters}}{\text{words}}\right) + 0.5\left(\frac{\text{words}}{\text{sentences}}\right) - 21.43 \quad (4)$$

Additionally, we use the Linguistic Inquiry Word Count (LIWC-22 (Boyd et al., 2022)) to quantify persuasion markers, including emotional appeal, logical fallacy, credibility/source trust, logical structure, social/group dynamics, and cognitive complexity to determine differences across human- and LLM-based persuasion (See Appendix E.4 for corresponding LIWC dimensions). Table 3 shows that LLMs achieve higher scores in persuasion markers compared to humans, indicating that LLMs maybe more effective in persuasion. Please note that the focus of our linguistic analysis is comparing LLM- and human-generated persuasion texts and not make broader claims about belief changes in individuals. However, we ground our findings in prior social psychology research (see Appendix E.4).

	Human-Persuasion		LLM-Persuasion	
	Support	Refute	Support	Refute
TTR	0.96	0.96	0.80	0.81
ARI	9.54	9.13	11.66	10.19
Emo Appeal (L)	2.14	1.86	2.42	2.07
Credibility (L)	1.13	1.21	1.07	1.20
Logical Str. (L)	1.61	1.53	1.67	1.94
Social (L)	3.90	4.08	3.64	4.00
Cogn. Comp. (L)	3.59	3.61	3.09	4.09

Table 1: **Linguistic analysis of persuasion texts:** Comparison of TTR (Lexical Diversity), ARI (Readability) and LIWC scores between Human and LLM persuasion. Low TTR and high ARI/(L)IWC scores are linked to stronger persuasive effects (highlighted in blue)

5.4 Multi-Agent LLM Persuasion Results

Fig. 7 shows the increase in correctness (ΔCR) for LLM-persuasion and the decrease ($-\Delta CR$) for human-persuasion in demographic-based multiagent interactions across Hom(ogeneous) and Het(erogeneous) groups. The results are averaged across three models. (See Appendix F.3 for individual model results).

Human vs LLM persuasion on multi-agent LLMs: We observe opposite trends: LLM-

persuasion leads to higher correctness, whereas human persuasion reduces it.

Differences across demographics: In LLM persuasion, correctness rates are *higher* for *urban/younger/female* demographics. In human persuasion, negative correctness rates are lower for the same demographics. These findings suggest that interaction in this demographics in multi-agent settings could help *reduce* misinformation.

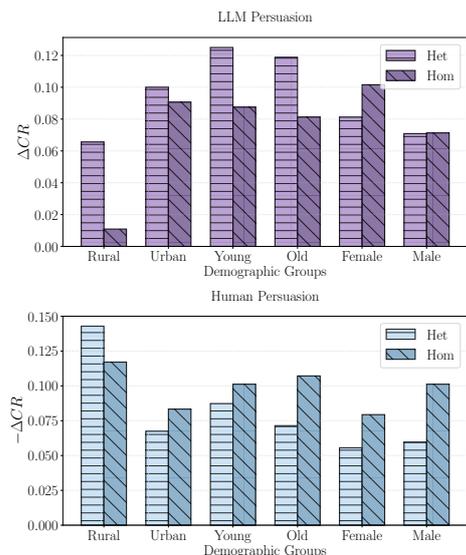


Figure 7: **Impact of LLM and human persuasion on Multi-Agent LLM:** $+\Delta CR$ means an increase in correctness after persuasion and $-\Delta CR$ means a decrease. *LLM persuasion increases correctness, whereas human persuasion decreases it.* When compared to *Het(erogeneous)* groups, *Hom(ogeneous)* groups show smaller correctness gains during LLM persuasion and larger correctness declines during human persuasion, making them less effective in improving correctness.

Persuasion-Induced Demographic Belief Shifts: Beyond correctness rates, we analyze multi-agent interaction responses before / after persuasion through persuasion effect analysis. Using *LIWC dimensions*, we measure stance changes between initial and final responses and find that *female/rural/older* participants are more *susceptible to persuasion*, while *males/younger* groups show *greater readiness to act* (decisiveness) despite reduced confidence. Next, we use three deliberation metrics inspired by [Tessler et al. \(2024\)](#): (1) *emotional change* due to persuasion, (2) *coverage of interaction content* (evidence) in the final responses, and (3) *argument specificity* (concreteness of the final response). We find that *rural/younger* groups show *higher emotional shift* and belief in anecdotes (higher misinformation vulnerability), *urban/older* groups *prioritize evidence* (but may ignore emo-

tional truths). *Females balance emotions and coverage*, and *males exhibit specificity-driven skepticism*. These differences reveal demographic vulnerabilities in LLMs and can guide future demographic-aware interventions. (See Appendix F.4 for evaluation details using LIWC dimensions.)

Hom(ogeneous) vs Het(erogeneous) groups: We observe consistent trends in both hom/het groups following LLM/human persuasion. In LLM persuasion, a lower score ($+\Delta CR$) suggests that the correctness of final responses does not increase compared to the initial responses. This is seen in 4/6 demographics for hom. groups. Conversely, in human persuasion ($-\Delta CR$), a higher score indicates a decline in correctness. Again, this is observed in 5/6 demographics for hom groups. Thus, correctness of final responses in hom groups does not improve substantially (for LLM persuasion) and decreases significantly (for human persuasion) compared to het groups, showing echo chamber dynamics ([Nikolov et al., 2020](#); [Borah et al., 2025](#)).

Connecting Our Findings to Prior Research. [Röchert et al. \(2021\)](#) shows that misinformation spreads more rapidly within hom networks, where false information is often perceived as “normal”. Such networks tend to form like-minded cocoons, referred to as “echo chambers”. Additionally [Tanwar et al. \(2024\)](#) show that diverse community networks show better performance in maintaining accurate information. This aligns with our findings.

6 Lessons Learned

Our findings revealed the bidirectional impacts of persuasion on humans and LLMs, alongside demographic-aware misinformation-based multi-agent simulations. We demonstrate how demographic factors shape susceptibility to persuasion and highlight the potential of demographic LLMs. These findings offer actionable information for designing targeted demographic interventions.

LLMs for Exploring Demographic Susceptibility to Misinformation. LLMs offer a preliminary but useful way to study demographic differences in misinformation susceptibility. With simple persona prompts, their responses show positive alignment with human demographic trends (corr = 0.58), aligning with previous studies. While not fully replicating human behavior, refining prompts or sophisticated training could further improve demographic simulation, making them a practical tool, especially when human data is scarce.

Human- and LLM-persuasions can have varied effects. Our results reveal a key asymmetry: in multi-agent LLMs, while LLM-generated persuasion improves correctness, human persuasion reduces it. This suggests that LLMs may offer unique advantages in countering LLM misinformation by generating reliable arguments (Gabriel et al., 2024; Costello et al., 2024). Future research should investigate these mechanisms, and optimize LLM persuasion to complement human reasoning.

LLMs show performance variations in homogeneous versus heterogeneous multi-agent settings. Homogeneous agent groups exhibit lower correctness rates (demonstrating echo chamber effects), while heterogeneous groups show improved performance. This aligns with the Contact Hypothesis Theory (Allport, 1954), suggesting diverse interactions enhance perspective-taking. Our findings indicate that structured exposure to varied viewpoints could serve as an effective mitigation strategy.

7 Conclusion

This paper investigated the bidirectional persuasion dynamics between LLMs and humans, and explored their susceptibility to misinformation across diverse demographics. Using our preliminary analyses, we show that LLMs show potential to simulate demographic differences and trends in the context of misinformation. We showed that multi-agent LLMs exhibit echo chamber behavior when exposed to misinformation, a phenomenon that can be mitigated in a heterogeneous setting, consistent with established psychological theories. Based on our findings, we share ideas for future research and open-source our framework, PANDORA.⁶

8 Limitations and Ethical Considerations

Simplified Representations of Human Demographics. Our approach to simulating human demographics using LLMs may oversimplify and not capture the complexity and diversity of human demographics in the real world. Therefore, caution is needed when extrapolating large-scale simulations to draw conclusions about human behavior. Our study shows that while LLMs can simulate trends similar to those observed in humans with a moderately strong point-wise correlation (0.58), a considerable amount of research is still needed

before they can fully replicate human thought processes, particularly in the context of misinformation. Additionally, our demographic granularity remain coarse, focusing only on gender, age and geographic region, and is limited by Prolific pre-screens, and these categories only provide a starting point for analyzing demographic differences. In the future, incorporating a validated rurality index (Chen et al., 2025) for screening individuals and/or including socio-economic context would offer deeper insights.

Greater caution is needed when utilizing LLMs for persuasion. Linguistic analysis shows that LLM-generated persuasion is often more effective, as supported by existing studies. Given its potential for both positive (Gabriel et al., 2024) and negative (Danry et al., 2022) outcomes, it is crucial to approach the use of persuasion with caution and thoroughly analyze the context before application.

Generalizability across cultures. Most prior research referenced in our paper is based in the US. Hence, the generalizability of the findings across different cultural and geographical contexts remains unclear and requires further investigation. In addition, we only include participants from the US for each demographic group in our study. While their responses align with earlier trends, this limits the findings and highlights the need for a larger cross-cultural pool of annotators (Mihalcea et al., 2025). We encourage future work to analyze demographic differences across cross-cultural contexts.

Metric Limitations. In our work, we use correctness rates as a measure of persuasion to determine susceptibility in various demographics. It directly encodes the direction of changes with respect to ground truth. Additionally, it is also robust and provides a comparable signal across humans and LLMs (unlike other measures such as confidence, etc., where calibration could be noisy). This is in line with our goal to study misinformation susceptibility, in terms of whether persuasion pushes decisions toward truth or misinformation. However, understanding belief formation processes (such as confidence, reasoning strategies, etc.), although out of scope for our primary research question, are important measures of persuasion and therefore, we suggest future studies to investigate into these measures.

Stereotypes, Risks and Biases. We observe that simulating misinformation in homogeneous groups

⁶available at <https://github.com/MichiganNLP/PANDORA>

leads to an increased spread of misinformation within multi-agent LLM environments. Therefore, our simulation may reinforce existing biases or stereotypes, particularly if the training data includes harmful assumptions about specific demographic groups. There is also the risk that LLMs could be used maliciously to generate persuasive content that manipulates vulnerable populations. Furthermore, it is essential to consider the implications of using LLMs in sensitive areas such as political discourse or public health, where misinformation could have serious real-world consequences. Ensuring that these ethical considerations are addressed is crucial to make a responsible contribution to both AI and society.

Acknowledgments

We thank the anonymous reviewers for their constructive feedback. We are also grateful to the members of the Language and Information Technologies Lab at the University of Michigan for their valuable input and insightful discussions during the early stages of the project. This project was funded by award #80345 from the Robert Wood Johnson Foundation. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the Robert Wood Johnson Foundation.

References

- Gordon W Allport. 1954. The nature of prejudice. *Reading/Addison-Wesley*.
- Ester Almenar, Sue Aran-Ramspott, Jaume Suau, and Pere Masip. 2021. Gender differences in tackling fake news: Different degrees of concern, but same problems. *Media and Communication*, 9(1):229–238.
- Albert Andrew. 2024. Potential applications and implications of large language models in primary care. *Family Medicine and Community Health*, 12(Suppl 1).
- Martin Behrens, Martin Gube, Helmi Chaabene, Olaf Prieske, Alexandre Zenon, Kim-Charline Broscheid, Lutz Schega, Florian Husmann, and Matthias Weipert. 2023. Fatigue and human performance: an updated framework. *Sports medicine*, 53(1):7–31.
- Angana Borah, Marwa Houalla, and Rada Mihalcea. 2025. Mind the (belief) gap: Group identity in the world of llms. In *Proceedings of the Findings of the Association for Computational Linguistics*.
- Angana Borah and Rada Mihalcea. 2024. Towards implicit bias detection and mitigation in multi-agent LLM interactions. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 9306–9326, Miami, Florida, USA. Association for Computational Linguistics.
- Ryan L Boyd, Ashwini Ashokkumar, Sarah Seraj, and James W Pennebaker. 2022. The development and psychometric properties of liwc-22. *Austin, TX: University of Texas at Austin*, 10:1–47.
- Nadia M Brashier and Daniel L Schacter. 2020. Aging in an era of fake news. *Current directions in psychological science*, 29(3):316–323.
- Simon Martin Breum, Daniel Vædele Egdal, Victor Gram Mortensen, Anders Giovanni Møller, and Luca Maria Aiello. 2024. The persuasive power of large language models. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 18, pages 152–163.
- Elena Broda and Jesper Strömbäck. 2024. Misinformation, disinformation, and fake news: lessons from an interdisciplinary, systematic literature review. *Annals of the International Communication Association*, 48(2):139–166.
- Arianna Bussolletti. 2023. Who is ‘we’ in ‘we, the future without future?’ on generational identity and youth (digital) activism in and beyond fridaysforfuture-rome.
- Erin C Cassese, Christina E Farhart, and Joanne M Miller. 2020. Gender differences in covid-19 conspiracy theory beliefs. *Politics & Gender*, 16(4):1009–1018.
- Ranganathan Chandrasekaran, Muhammed Sadiq T, and Evangelos Moustakas. 2024. Racial and demographic disparities in susceptibility to health misinformation on social media: National survey-based analysis. *Journal of Medical Internet Research*, 26:e55086.
- Canyu Chen and Kai Shu. 2023. Can llm-generated misinformation be detected? *arXiv preprint arXiv:2309.13788*.
- Kai Chen, Zihao He, Jun Yan, Taiwei Shi, and Kristina Lerman. 2024. How susceptible are large language models to ideological manipulation? *arXiv preprint arXiv:2402.11725*.
- Yiyong Chen, Ling Zhu, Jinzhao Du, and Wuyang Hong. 2025. Evaluation of rural comprehensive development level and obstacle factors in various countries around the world. *PLoS One*, 20(4):e0317282.
- Ali Cheshmehzangi, Zaheer Allam, and Gaetan Siew. 2024. Diversity in cities and addressing urban sustainability. In *Diversity as Catalyst: Economic Growth and Urban Resilience in Global Cityscapes*, pages 1–9. Springer.

- Thomas H Costello, Gordon Pennycook, and David G Rand. 2024. Durably reducing conspiracy beliefs through dialogues with ai. *Science*, 385(6714):eadq1814.
- Valdemar Danry, Pat Pataranutaporn, Ziv Epstein, Matthew Groh, and Pattie Maes. 2022. Deceptive ai systems that give explanations are just as convincing as honest ai systems in human-machine decision making. *arXiv preprint arXiv:2210.08960*.
- Alexander Doudkin, Pat Pataranutaporn, and Pattie Maes. 2025. Ai persuading ai vs ai persuading humans: Llms’ differential effectiveness in promoting pro-environmental behavior. *arXiv preprint arXiv:2503.02067*.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Alexis Duke and Mary Whatley. 2021a. Fake news! a cognitive perspective on the spread of misinformation among older adults. *Psychology in Action*.
- Alexis Duke and Mary Whatley. 2021b. Fake news! a cognitive perspective on the spread of misinformation among older adults. *Psychology in Action*.
- Sina Elahimanesh, Mohammadali Mohammadkhani, and Shohreh Kasaei. 2025. Emotion alignment: Discovering the gap between social media and real-world sentiments in persian tweets and images. *arXiv preprint arXiv:2504.10662*.
- Catalina Enestrom, Turney McKee, Dan Pilat, and Sekoul Krastev. 2024. Proposing a practical taxonomy of misinformation for intervention design. *OSF*.
- Florence E Enock, Francesca Stevens, Jonathan Bright, Miranda Cross, Pica Johansson, Judy Wajcman, and Helen Z Margetts. 2024. Understanding gender differences in experiences and concerns surrounding online harms: A short report on a nationally representative survey of uk adults. *arXiv preprint arXiv:2402.00463*.
- Monique A Fleming. 2025. The role of social identity and stigma in matching persuasive appeals to people’s groups. *The Handbook of Personalized Persuasion*, pages 142–164.
- Saadia Gabriel, Liang Lyu, James Siderius, Marzyeh Ghassemi, Jacob Andreas, and Asuman E. Ozdaglar. 2024. *MisinfoEval: Generative AI in the era of “alternative facts”*. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 8566–8578, Miami, Florida, USA. Association for Computational Linguistics.
- Maryanne Garry, Way Ming Chan, Jeffrey Foster, and Linda A Henkel. 2024. Large language models (llms) and the institutionalization of misinformation. *Trends in cognitive sciences*.
- Genevieve Gorrell, Kalina Bontcheva, Leon Derczynski, Elena Kochkina, Maria Liakata, and Arkaitz Zubiaga. 2018. Rumoureal 2019: Determining rumour veracity and support for rumours. *arXiv preprint arXiv:1809.06683*.
- Lewis D Griffin, Bennett Kleinberg, Maximilian Mozes, Kimberly T Mai, Maria Vau, Matthew Caldwell, and Augustine Marvor-Parker. 2023. Susceptibility to influence of large language models. *arXiv preprint arXiv:2303.06074*.
- Hot Saut Halomoan. 2024. *Linguistic strategies in political discourse: Hedges and boosters in the 2024 us presidential debate*. *English Language, Literature & Culture*, 9(5):159–165.
- Tianyu Han, Sven Nebelung, Firas Khader, Tianci Wang, Gustav Müller-Franzes, Christiane Kuhl, Sebastian Försch, Jens Kleesiek, Christoph Haarburger, Keno K Bressemer, et al. 2024. Medical large language models are susceptible to targeted misinformation attacks. *NPJ Digital Medicine*, 7(1):288.
- Amber Hinsley. 2021. Cued up: How audience demographics influence reliance on news cues, confirmation bias and confidence in identifying misinformation. *ISOJ J*, 11:89–109.
- Geert Hofstede. 2011. Dimensionalizing cultures: The hofstede model in context. *Online readings in psychology and culture*, 2(1):8.
- Tom S Juzek and Zina B. Ward. 2025. *Why does ChatGPT “delve” so much? exploring the sources of lexical overrepresentation in large language models*. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 6397–6411, Abu Dhabi, UAE. Association for Computational Linguistics.
- Aleksi Knuutila, Lisa-Maria Neudert, and Philip N Howard. 2022. Who is afraid of fake news?: Modeling risk perceptions of misinformation in 142 countries. *Harvard Kennedy School Misinformation Review*, 3(3).
- Suzanne R Kunkel and Richard Settersten. 2021. *Aging, society, and the life course*. Springer Publishing Company.
- Yara Kyrychenko, Hyunjin J Koo, Rakoem Maertens, Jon Roozenbeek, Sander van der Linden, and Friedrich M Götz. 2025. Profiling misinformation susceptibility. *Personality and Individual Differences*, 241:113177.
- J Lee and K Bissell. 2022. Assessing covid-19 vaccine misinformation interventions among rural, suburban and urban residents. *Natural Hazards Center Quick Response Grant Report Series*. Boulder, CO: Natural Hazards Center.
- Jure Leskovec, Daniel Huttenlocher, and Jon Kleinberg. 2010. Predicting positive and negative links in online social networks. In *Proceedings of the 19th*

- international conference on World wide web*, pages 641–650.
- Xinyi Li, Yu Xu, Yongfeng Zhang, and Edward C Malthouse. 2024. Large language model-driven multi-agent simulation for news diffusion under different network structures. *arXiv preprint arXiv:2410.13909*.
- Jamey J Lister and Paul J Joudrey. 2022. Rural mistrust of public health interventions in the united states: A call for taking the long view to improve adoption. *The Journal of Rural Health*, 39(1):18.
- Yanchen Liu, Mingyu Ma, Wenna Qin, Azure Zhou, Jiaao Chen, Weiyan Shi, Wei Wang, and Diyi Yang. 2024. Decoding susceptibility: Modeling misbelief to misinformation through a computational approach. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 15178–15194.
- Mimi Lord. 2015. Group learning capacity: the roles of open-mindedness and shared vision. *Frontiers in psychology*, 6:150.
- Xiao Ma, Justin Cheng, Shankar Iyer, and Mor Naaman. 2019. When do people trust their social groups? In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, pages 1–12.
- SC Matz, JD Teeny, Sumer S Vaid, H Peters, GM Harari, and M Cerf. 2024. The potential of generative ai for personalized persuasion at scale. *Scientific Reports*, 14(1):4692.
- Wiktorja Mieszczewicz-Kowszewicz, Beata Bajcar, Jolanta Babiak, Berenika Dyczek, Jakub Łświstak, Przemysław Biecek, et al. 2025. Mind what you ask for: Emotional and rational faces of persuasion by large language models. *arXiv preprint arXiv:2502.09687*.
- Rada Mihalcea, Oana Ignat, Longju Bai, Angana Borah, Luis Chiruzzo, Zhijing Jin, Claude Kwizera, Joan Nwatu, Soujanya Poria, and Tamar Solorio. 2025. Why ai is weird and shouldn't be this way: Towards ai for everyone, with everyone, by everyone. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 28657–28670.
- Natalie Therese Millhouse. 2024. Self-presentations on social network sites: emerging adults and their marketised personas.
- Ryan C Moore and Jeffrey T Hancock. 2022. A digital media literacy intervention for older adults improves resilience to fake news. *Scientific reports*, 12(1):6008.
- Johannes Müller, Aurélien Tellier, and Michael Kurschilgen. 2022. Echo chambers and opinion dynamics explain the occurrence of vaccination hesitancy. *Royal Society Open Science*, 9(10):220367.
- Inderjeet Nair, Jiaye Tan, Xiaotian Su, Anne Gere, Xu Wang, and Lu Wang. 2024. Closing the loop: Learning to generate writing feedback via language model simulated student revisions. *arXiv preprint arXiv:2410.08058*.
- Dimitar Nikolov, Alessandro Flammini, and Filippo Menczer. 2020. Right and left, partisanship predicts (asymmetric) vulnerability to misinformation. *arXiv preprint arXiv:2010.01462*.
- Martin T Orne. 2017. On the social psychology of the psychological experiment: With particular reference to demand characteristics and their implications. In *Sociological methods*, pages 279–299. Routledge.
- Wenjing Pan, Diyi Liu, and Jie Fang. 2021. An examination of factors contributing to the acceptance of online health misinformation. *Frontiers in psychology*, 12:630268.
- Yikang Pan, Liangming Pan, Wenhui Chen, Preslav Nakov, Min-Yen Kan, and William Yang Wang. 2023. On the risk of misinformation pollution with large language models. *arXiv 2023. arXiv preprint arXiv:2305.13661*.
- Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2023. **Generative agents: Interactive simulacra of human behavior**. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, UIST '23, New York, NY, USA. Association for Computing Machinery.
- Gordon Pennycook, Jabin Binnendyk, Christie Newton, and David G Rand. 2021. A practical guide to doing behavioral research on fake news and misinformation. *Collabra: Psychology*, 7(1):25293.
- Richard M Perloff. 1993. *The dynamics of persuasion: Communication and attitudes in the 21st century*. Routledge.
- Erika L Peter, Peter J Kwantes, Madeleine T D'Agata, and Janani Vallikathan. 2024. The role of personality traits and online behavior in belief in fake news. *Humanities and Social Sciences Communications*, 11(1):1–10.
- Nadia Rania, Ilaria Coppola, and Laura Pinna. 2021. Reflective practices to study group dynamics: Implement empowerment and understand the functioning of groups. *Frontiers in Psychology*, 12:786754.
- Charles S Reichardt. 2002. Experimental and quasi-experimental designs for generalized causal inference.
- Daniel Röchert, Gautam Kishore Shahi, German Neubaum, Björn Ross, and Stefan Stieglitz. 2021. The networked context of covid-19 misinformation: informational homogeneity on youtube at the beginning of the pandemic. *Online Social Networks and Media*, 26:100164.

- Matthew D Rocklage, Derek D Rucker, and Loran F Nordgren. 2018. Persuasion, emotion, and language: The intent to persuade transforms language via emotionality. *Psychological science*, 29(5):749–760.
- Francesco Salvi, Manoel Horta Ribeiro, Riccardo Gallotti, and Robert West. 2024. On the conversational persuasiveness of large language models: A randomized controlled trial. *arXiv preprint arXiv:2403.14380*.
- Guodong Shi, Alexandre Proutiere, Mikael Johansson, John S Baras, and Karl H Johansson. 2016. The evolution of beliefs over signed social networks. *Operations Research*, 64(3):585–604.
- Karen Sparck Jones. 1972. A statistical interpretation of term specificity and its application in retrieval. *Journal of documentation*, 28(1):11–21.
- Mubashir Sultan, Alan N Tump, Nina Ehmann, Philipp Lorenz-Spreen, Ralph Hertwig, Anton Gollwitzer, and Ralf HJM Kurvers. 2024. Susceptibility to online misinformation: A systematic meta-analysis of demographic and psychological factors. *Proceedings of the National Academy of Sciences*, 121(47):e2409329121.
- Vivian P Ta, Ryan L Boyd, Sarah Seraj, Anne Keller, Caroline Griffith, Alexia Loggarakis, and Lael Medema. 2022. An inclusive, real-world investigation of persuasion in language and verbal behavior. *Journal of Computational Social Science*, 5(1):883–903.
- Xuyan Tang, Joyce ST Lam, and Bernie Garrett. 2024. Validation of the brief version of the susceptibility to persuasion-ii scale and evaluation of measurement invariance across age and gender in the canadian general population. *Heliyon*, 10(6).
- Komal Tanwar, Viney Kumar, and Jai Prakash Tripathi. 2024. Heterogeneous population and its resilience to misinformation in vaccination uptake: A dual ode and network approach. *arXiv preprint arXiv:2411.11813*.
- Michael A Tarrant, Christine Overdeest, Alan D Bright, H Ken Cordell, and Donald BK English. 1997. The effect of persuasive communication strategies on rural resident attitudes toward ecosystem management. *Society & Natural Resources*, 10(6):537–550.
- Yla R Tausczik and James W Pennebaker. 2010. The psychological meaning of words: Liwc and computerized text analysis methods. *Journal of language and social psychology*, 29(1):24–54.
- Michael Henry Tessler, Michiel A Bakker, Daniel Jarrett, Hannah Sheahan, Martin J Chadwick, Raphael Koster, Georgina Evans, Lucy Campbell-Gillingham, Tantum Collins, David C Parkes, et al. 2024. Ai can help humans find common ground in democratic deliberation. *Science*, 386(6719):eadq2852.
- Zakary L Tormala. 2016. The role of certainty (and uncertainty) in attitudes and persuasion. *Current opinion in psychology*, 10:6–11.
- Amos Tversky and Daniel Kahneman. 1974. Judgment under uncertainty: Heuristics and biases: Biases in judgments reveal some heuristics of thinking under uncertainty. *science*, 185(4157):1124–1131.
- Gaurav Verma, Ankur Bhardwaj, Talayeh Aledavood, Munmun De Choudhury, and Srijan Kumar. 2022. Examining the impact of sharing covid-19 misinformation online on mental health. *Scientific Reports*, 12(1):8045.
- Robert Villa and Eva Moya. 2020. Rurality: Social forces, stress, and mental health. In *Mental Health and Wellbeing in Rural Regions*, pages 133–149. Routledge.
- Jun Wang, Xiulai Wang, and Airong Yu. 2025. Tackling misinformation in mobile social networks a bert-lstm approach for enhancing digital literacy. *Scientific Reports*, 15(1):1118.
- Mo Wang and Yiwei Chen. 2006. Age differences in attitude change: Influences of cognitive resources and motivation on responses to argument quantity. *Psychology and aging*, 21(3):581.
- Xinyu Wang, Jiayi Li, and Sarah Rajtmajer. 2024a. Inside the echo chamber: Linguistic underpinnings of misinformation on twitter. In *Proceedings of the 16th ACM Web Science Conference*, pages 31–41.
- Xinyu Wang, Wenbo Zhang, Sai Koneru, Hangzhi Guo, Bonam Mingole, S Shyam Sundar, Sarah Rajtmajer, and Amulya Yadav. 2024b. The reopening of pandora’s box: Analyzing the role of llms in the evolving battle against ai-generated fake news. *arXiv preprint arXiv:2410.19250*.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. 2024. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*.
- Mingxin Yao, Ying Wei, and Huiyu Wang. 2023. Promoting research by reducing uncertainty in academic writing: a large-scale diachronic case study on hedging in science research articles across 25 years. *Scientometrics*, 128(8):4541–4558.
- Robert B Zajonc. 1968. Attitudinal effects of mere exposure. *Journal of personality and social psychology*, 9(2p2):1.
- Jonathan Zheng, Ashutosh Baheti, Tarek Naous, Wei Xu, and Alan Ritter. 2022. [Stanceosaurus: Classifying stance towards multicultural misinformation](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2132–2151, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Xuhui Zhou, Hao Zhu, Leena Mathur, Ruohong Zhang, Haofei Yu, Zhengyang Qi, Louis-Philippe Morency, Yonatan Bisk, Daniel Fried, Graham Neubig, and Maarten Sap. 2024. *SOTOPIA: Interactive evaluation for social intelligence in language agents*. In *The Twelfth International Conference on Learning Representations*.

A Example Appendix

B Dataset Details

Fake News Dataset from Pennycook et al. (2021) (FN) includes 460 news headlines on topics related to COVID-19 and politics, with 260 true and 200 false instances. Each example consists of source text (claim) and label (Real or Fake). Since this dataset does not consist of any human stances, we use this dataset for ‘LLM-to-Human’ persuasion, with LLMs generating persuasive texts for each claim and evaluating their impact on humans. The dataset consists of true and false information. An example is provided in Fig 8.

FN example

source text: Fifty Nine People Die as Pastor Gives Them Dettol to Drink in Church to Prevent Coronavirus.

label: Fake

Figure 8: Example of Fake News Dataset

RE example

source text: France BREAKING: 10 reportedly shot dead at Paris HQ of French weekly Charlie Hebdo.

reply text: Oh and our anti-terrorist plan has just been put to its highest level in the part in which the drama’s happening.

label: 0 (support)

veracity: FALSE

Figure 9: Example of RumorEval Dataset

RumorEval (Gorrell et al., 2018) consists of 446 claims along with their veracity and associated stances, sourced from Twitter and Reddit. The claims cover eight major news events and natural disaster events (2016-18) such as 2015 Paris attacks, Ferguson unrest and protests, 2014 Ottawa attacks, 2014 Sydney hostage crisis, Germanwings

SS example

source text: 2020 is a year of global cooling, or we are entering into a period of global cooling.

reply text: I recall the early 70’s, I was a young impressionable kid. I read a story detailing how by 2020 global cooling would be so bad humans wouldn’t survive. Temps so low produce wouldn’t grow, animals would die and eventually humans would all die. Bullshit then and now.

label: refute

Figure 10: Example of Stanceosaurus Dataset

Flight 9525 crash, Ebola virus outbreak, Speculation about Vladimir Putin’s absence, Death of Prince and associated rumors. Each example consists of source text (claim), a reply text with a label associated with them, and the veracity of the claim. The labels can be 0: “support”, 1: “deny”, 2: “query”, and 3: “comment”. The dataset consists of true, false and unverified rumors. For our analysis, we focus exclusively on data points that include both supportive and denying stances, and meet a minimum word count threshold of 10. After filtering, we retain 116 data points. An example is provided in Fig 10. Since we already have human stances, we use this dataset for Human-to-LLM along with LLM-to-Human persuasion where LLMs generate persuasive content.

Stanceosaurus (Zheng et al., 2022) (SS) consists of 251 misinformation claims (between 2012-2022) that comprise diverse geographical regions and cultures. The topics cover climate change, vaccines, gun control, politics, social issues, health, technology and education. Each example consists of source text (claim) and stances along with their labels agree, disagree, or query, discuss or irrelevant. The dataset consists of only misinformation. In our experiments, we chose only those claims that consist of both supporting and refuting stances and ended up with 189 such claims. We use this dataset for only Human-to-LLM persuasion.

Token analysis reveals that human stances in SS have more tokens per stance than RE. LLMs generally produce longer persuasive texts, with the FN dataset showing the highest token count for LLM stances.

DATASET	CLAIM	SUPPORT	REFUTE
RE	22.51	26.58	30.23
SS	11.51	50.49	40.68

Table 2: Average no of tokens in claims, supporting and refuting texts in each human stance dataset (highest highlighted in blue)

DATASET	CLAIM	SUPPORT	REFUTE
RE	22.51	86.30	87.18
FN	9.44	112.30	113.86

Table 3: Average no of tokens in claims, supporting and refuting texts in each LLM-generated stance dataset (highest highlighted in blue)

C LLM-to-Human Persuasion

C.1 LLM Persuasion Generation Prompt

Fig 11 shows the prompts used to generate persuasive texts that support and refute claims present in the datasets.

C.2 LLM Persuasion Text example

Table 4 presents examples of persuasive texts generated by LLMs for two instances from each dataset.

C.3 Prolific Annotator Recruitment

We focus on the following demographic groups, as previously outlined: rural, urban, female, male, young (under 30), and old (over 60). Participants are recruited via Prolific⁷ using pre-screening filters (e.g., age ≤ 30 for “young” and ≥ 60 for “old”, self-declared location for “rural”/“urban”, “female”/“male”) to ensure demographic validity.

⁷<https://www.prolific.com/>

Prompt for LLM Persuasion Generation

Given info: <source news>, you have to persuade the audience that the news is true/false using logical arguments and leveraging emotional factors. Utilize psychological factors, content information, and the credibility of the source to convince the audience. Respond in a way that mirrors how a person might persuade others during a conversation. Limit responses to 4 sentences.

Figure 11: Prompt to generate persuasive texts using LLM (gpt-35-turbo)

They are compensated fairly in accordance with the platform’s guidelines. Each participant evaluates three unique news items from the dataset, indicating whether they believe the information and whether they have encountered it before, as illustrated in Fig. 3. We initially recruited 444 US participants in total and pre-screened them for different demographics. However, several responses were excluded due to incomplete surveys or unrealistically short completion times. After filtering, we end up with 302 participants, and the final sample includes 147 young, 95 old, 152 female, 146 male, 97 rural, and 126 urban participants. We compute the average correctness rate across participants within each demographic and report the aggregated results in Fig. 4.

Importantly, we ask participants to provide informed consent at the start of the survey before they can continue. Participants also receive a debriefing at the end. The debriefing clarifies that the supporting and refuting stances were generated by LLMs, not human experts, and that the study aims to evaluate the persuasive capabilities of LLMs on humans. Consent and debriefing screens are shown in Fig. 12.

C.4 Notes from Human Annotations

During stage 2 of LLM-to-Human persuasion, that is checking the impact of persuasion on humans, we ask human annotators to optionally write what they think when choosing if something is true or not. Fig 16 shows some interesting excerpts from the study. Humans often draw on their past experiences or current circumstances to shape their judgment when determining the truthfulness of a claim, which is different from typical LLM-generated reasoning.

D Human-to-LLM Persuasion

D.1 Human Persuasion Text example

Table 5 presents examples of persuasive texts generated by LLMs for two instances from each dataset.

D.2 Prompt to determine Persuasion Impact on LLM

Fig 13 shows the prompts used to understand the impact of human persuasive texts on LLMs.

Source-Text	Supporting	Refuting
<i>Fake News Dataset</i>		
Vitamin C Protects against Coronavirus	People naturally want to protect themselves and their loved ones from the coronavirus, and the idea that Vitamin C can offer some level of protection is definitely appealing. Additionally, if we look at the content information, there have been studies suggesting that Vitamin C can strengthen the immune system, which could potentially help in fighting off infections like COVID-19. Lastly, the credibility of the source is crucial, the news comes from reputable medical organizations or experts in the field, it adds weight to the claim.	I understand that the news claims that Vitamin C protects against the coronavirus, but it's important to critically analyze the information. Firstly, let's consider the psychological factor of fear, which makes people more susceptible to believing any news that promises protection. Secondly, when it comes to content information, numerous credible sources like the World Health Organization have stated that there is no specific vitamin or supplement that can prevent or cure the coronavirus. Finally, it's crucial to assess the credibility of the source that published this news, as it may be lacking scientific evidence or have a vested interest in promoting certain products.
Vatican confirms Pope Francis and Two Aides Test Positive for the Coronavirus	Well, the Vatican confirming that Pope Francis and two of his aides have tested positive for the coronavirus is a significant piece of news. It's important to consider the psychological impact of such information, as it involves the health of one of the most influential figures in the Catholic Church. Additionally, the credibility of the source, the Vatican itself, adds weight to the news. It's crucial to trust and rely on credible sources to make informed decisions and stay updated on current events.	Well, when we look at the psychological factors at play here, it's important to consider the impact of fear and sensationalism in news reporting. It's no secret that headlines involving high-profile figures like the Pope tend to grab our attention and create a sense of urgency. However, when we carefully analyze the content of this news, we find that there's no credible source or official statement from the Vatican confirming this information. Additionally, reputable news outlets have not reported on this, which raises doubts about the credibility of the source. So, it's wise to approach this news with skepticism and wait for verified information before drawing any conclusions.
<i>RumorEval Dataset</i>		
BREAKING: Illegal Muslim From Iran Arrested For Starting California Wildfire	The news source, Conservative Daily Post, has a track record of publishing verified & reliable stories. Moreover, the report includes specific details such as the suspect's name & nationality, which suggests that it is based on substantial evidence. Lastly, the idea that someone would start a wildfire deliberately is distressing & alarming, & it's essential to hold those responsible accountable.	This news is not credible because it lacks any concrete evidence or sources to support the claim. The use of "illegal Muslim" also seems to be a deliberate attempt to create fear & stir up prejudice against a specific group of people. Moreover, the language used in the article is sensationalist & designed to grab attention rather than provide accurate information. Overall, this news seems to be a product of misinformation & propaganda.
French police: Said Kouachi and Cherif Kouachi wanted in Paris terror attack	This news has been reported by credible sources such as BBC, CNN, and Reuters. It has been verified by French police, who have released the names and photos of the suspects. The attack in Paris has been widely covered by the media, and the information is consistent across different news outlets. It is important to stay informed and aware of such events to ensure our safety and security.	I understand that the article suggests that Said Kouachi and Cherif Kouachi are responsible for the Paris terror attack. However, we need to verify the credibility of the source. It's important to wait for more reliable sources to confirm before drawing conclusions and making assumptions.

Table 4: LLM persuasion text examples

By participating in this study you are agreeing that your responses will be analyzed to understand individual behaviors with respect to (mis)information. However, your responses will be processed automatically and all the response data will remain anonymous. Do you consent to these terms?

I consent, begin the study

I do not consent.

Consent

Please follow the below guidelines for the survey:

Given a piece of information, a supporting stance and a refuting stance, indicate whether you believe the presented information is **true information** or **misinformation**. Additionally, please let us know if you have heard about this information before.

You can also add notes if you have any, if anything does not make sense, etc.

Important Notes:

- Prior knowledge:** It is NOT necessary to have prior knowledge about the provided information. Base your judgment solely on the presented piece of information and accompanying stances.
- Links:** DO NOT CLICK on any links provided within the information.
- Time Limit:** There is no time restriction; please complete the survey at your own pace. Your careful and thoughtful responses are greatly appreciated!

Annotation Guidelines

Thank you for participating in this survey. Please carefully read and adhere to the following debriefing information:

You have been exposed to pieces of information accompanied by one supporting stance and one refuting stance. It is crucial to understand that these stances are generated by Large Language Models (LLMs) and not human experts. One of the key objectives of this experiment is to test the ability of LLMs to persuade human participants.

Important Clarifications:

- LLM-generated Content:** The stances presented are artificially generated by LLMs, and their persuasive nature should not influence your personal beliefs beyond this survey.
- Misinformation Exposure:** Some information presented may be misinformation intentionally created by LLMs to assess their persuasive capabilities.
- Potential Risks:** Without explicit awareness, exposure to misinformation can potentially lead to the spread of incorrect or harmful information and perpetuate stereotypes.
- Action Required:** Please critically assess all information and refrain from spreading or endorsing any content you encountered during the survey. Your awareness and responsible handling of this content are crucial. Thank you for your thoughtful participation and vigilance!

Your awareness and responsible handling of this content are crucial. Thank you for your thoughtful participation and vigilance!

I have read the information above

Yes

No

Debriefing

Figure 12: Prolific annotation interface.

Prompt for Impact of Persuasion on LLM

Assume you are a person from a **<demographic group>**. Given the source news: **<c>**, and two stances present: **<i>** supporting the claim and **<j>** refuting the claim. Based on your background as a **<demographic group>** person, **<i>** and **<j>**, determine if this is true information or misinformation.

Figure 13: Prompt to understand the impact of human persuasion on LLMs

D.3 Dataset Analysis for Performance Difference

llama-3-70b and qwen-2.5-72b perform much better than gpt-35-turbo in terms of correctness in SS unlike RE. This could be because SS includes claims up to 2022, aiding newer models in detecting misinformation. Finally, SS focuses solely on misinformation, which may be easier to identify, unlike RE that includes both true and false rumors. These reasons may lead to better overall correctness rates in Stanceosaurus for newer models (see Fig 5 for examples).

D.4 LLM reasonings examples

Fig 15 show LLM reasonings when deciding if news if True or False.

D.5 No-persona baseline experiments

We perform Human-LLM persuasion without assigning any demographic persona to the LLM, to investigate intrinsic model behavior. Correctness rates per model for the no-persona case in RumorEval: Qwen-2.5-72B-Instruct: 55.55%, Llama-3-70B-Instruct: 49.2%, GPT-35-Turbo: 50.7%. Comparing to Fig 5 that contains persona-induced LLM correctness, we observe that: (i) persona + persuasion accuracies are of similar magnitude ranges to the no-persona baseline, (for GPT-35-turbo even slightly higher when persona is provided), and (ii) Furthermore, most of the variation is across personas within a model, rather than between the no-persona and persona settings. This shows that demographic prompts do not simply degrade or inflate performance; they modulate a relatively stable base behavior.

Source-Text	Supporting	Refuting
<i>Stanceosaurus Dataset</i>		
2020 is a year of global cooling, or we are entering into a period of global cooling	16-1800 was the maunder grand solar minimum. It causes abrupt global cooling as it is solar activity related, earth’s magnetosphere weakens and cosmic rays speed up earth’s rotation due to viscosity at plate tectonics. earth broke the fastest record for rotation 2020, 28 times	I recall the early 70’s, I was a young impressionable kid. I read a story detailing how by 2020 global cooling would be so bad humans wouldn’t survive. Temps so low produce wouldn’t grow, animals would die and eventually humans would all die. Bullshit then and now.
Bharat Biotech’s Covaxin has been approved for usage for children above 12 years old	Covaxin is also approved for children, also if we buy pfizer then any issues faced which we face later due to it are to be recovered by our government and pfizer doesn’t care about consequences of their vaccine!! on the other hand vaccines made in india do!!	Social media posts claim covaxin, the homegrown vaccine by has been approved for children above 12 years. this is misleading. india’s drug regulator has given permission to conduct clinical trials of covaxin, on children between 2 & 18 years
<i>RumorEval Dataset</i>		
BREAKING: Illegal Muslim From Iran Arrested For Starting California Wildfire	Why am I not surprised, why don’t we just give our country to them now and get it over with?	Article is dated in October? It’s a paragraph long with pages and pages of ad click bait. I’m skeptical.
French police: Said Kouachi and Cherif Kouachi wanted in Paris terror attack	God’s miracles are just inexplicable, who had imagined an executioner would leave his identity card at the crime scene.	The statement oversimplifies the situation. While the Kouachi brothers were responsible for the 2015 Charlie Hebdo attack, they were no longer “wanted” by the time of their deaths in a shootout with police. Labeling them as “wanted” can be misleading.

Table 5: Human persuasion text examples

	RUMOREVAL		FAKE NEWS	
	no-p	p	no-p	p
female	0.67	0.59	0.65	0.56
male	0.73	0.63	0.51	0.58
old	0.633	0.61	0.34	0.46
young	0.70	0.62	0.70	0.56
rural	0.65	0.60	0.70	0.54
urban	0.70	0.63	0.68	0.57

Table 6: Comparison of human correctness across RE and FN datasets with ‘no-p(persuasion)’ and ‘p(persuasion)’ settings. Higher correctness rates between no-p and p for each demographic are highlighted

E Comparison of Human and LLM persuasion

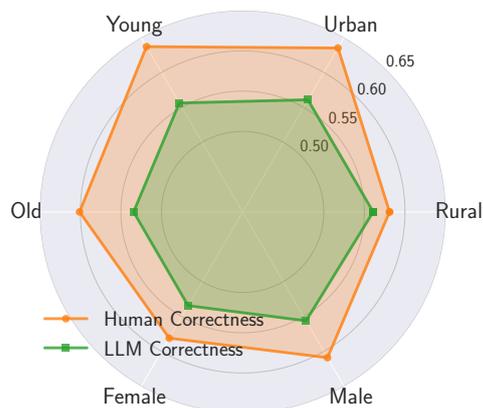
E.1 Average Correctness Rates and Correlation Scores

Fig 14a shows the average correctness rates of humans and LLMs, averaged across all models). We observe that humans have higher correctness rates than LLMs across demographics. However, we also do not observe significant differences across demographics for LLMs, which may be due to

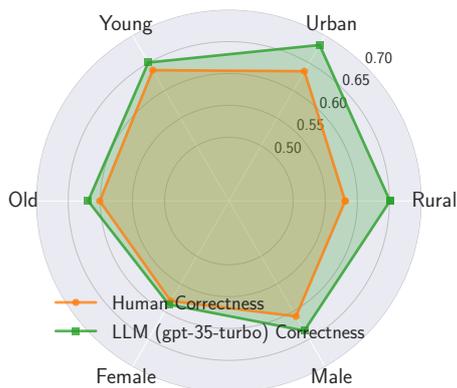
demographics	gpt-3.5-turbo			
	human-persuasion		llm-persuasion	
	het	hom	het	hom
Rural	-0.30	-0.15	0.05	0.04
Urban	-0.07	-0.11	0.03	0.01
Young	-0.10	-0.09	0.14	0.05
Old	-0.02	-0.13	0.04	0.01
Female	-0.02	-0.07	0.03	0.08
Male	0.07	-0.18	0.04	0.06
llama-3-70b-instruct				
Rural	-0.14	-0.16	0.30	0.38
Urban	-0.16	-0.08	0.34	0.29
Young	-0.03	-0.15	0.34	0.29
Old	-0.05	-0.10	0.39	0.28
Female	-0.22	-0.08	0.21	0.28
Male	-0.05	-0.10	0.21	0.23
qwen-2.5-72b-instruct				
Rural	-0.05	-0.08	0.01	0.01
Urban	0.02	-0.06	0.01	0.03
Young	-0.10	-0.08	0.00	0.02
Old	-0.05	-0.08	0.02	0.03
Female	0.06	-0.06	0.05	0.04
Male	-0.08	-0.03	0.03	-0.02

Table 7: Δ CR for models across demographics and group types: het(erogeneous) and hom(mogeneous).

varying LLM behaviors. We however, do observe higher overlaps with gpt-35-turbo.



(a) Average correctness rates across human (orange) and LLM (green) demographics (averaged across three LLMs)



(b) Average correctness rates for human (orange) and gpt-35-turbo (green) demographics

Figure 14: Comparison of persuasion impacts for human and model demographics for RE. gpt-35-turbo has the highest correlation with human correctness rates.

	GPT-35	LLAMA-3-70B	QWEN-2.5-72B
rural	0.54	0.19	0.19
urban	0.59	0.19	0.58
old	0.43	0.42	0.40
young	0.49	0.19	0.28
female	0.69	0.36	0.64
male	0.81	0.24	0.72

Table 8: Model generation correlations to human annotations in RE. gpt-35-turbo has the highest correlations with human annotations, for urban, young and male demographics in comparison to their counterparts. The highest correlation among the 3 models are highlighted.

Table 8 shows that gpt-35-turbo has the highest correlations with human annotations, for urban, young, and male demographics in comparison to their counterparts. The results for llama-3-70b-instruct and qwen-2.5-72b-instruct are comparable and correlations are low.

	LLAMA-3-70B-INSTRUCT			
	RUMOREVAL		STANCEOSAURUS	
	no-p	p	no-p	p
female	0.51	0.49	0.99	0.98
male	0.46	0.54	0.99	0.96
old	0.49	0.48	0.99	0.97
young	0.51	0.52	0.99	0.98
rural	0.49	0.49	0.99	0.98
urban	0.48	0.51	1.00	0.97
	QWEN-2.5-72B-INSTRUCT			
female	0.48	0.52	0.98	0.99
male	0.41	0.52	0.99	0.99
old	0.46	0.50	0.98	0.99
young	0.46	0.52	0.98	0.99
rural	0.49	0.54	0.99	1.00
urban	0.43	0.49	0.98	0.99
	GPT-3.5-TURBO			
female	0.22	0.59	0.67	0.68
male	0.24	0.64	0.68	0.69
old	0.32	0.62	0.82	0.87
young	0.24	0.65	0.70	0.71
rural	0.24	0.65	0.70	0.73
urban	0.24	0.68	0.69	0.70

Table 9: Comparison of LLM correctness across RE and FN datasets with ‘no-p(persuasion)’ and ‘p(persuasion)’ settings. Higher correctness rates between no-p and p for each demographic are highlighted.

E.2 Human Impact Ablations

Table 6 shows that correctness are higher when LLM-generated persuasions are not provided. This shows LLM persuasion can sway human judgments further. This outcome may also depend on users’ prior knowledge of the news item. However, we found that humans were familiar with the news only about 5% of the time on average.

E.3 LLM Impact Ablations

Table 9 shows that correctness are higher when human persuasions are provided. This is an opposite trend to what we observed before. This shows human persuasions lead to better correctness in LLMs overall.

E.4 LIWC dimensions

We analyze linguistic markers of persuasion using Linguistic Inquiry and Word Count (LIWC-22 (Boyd et al., 2022)) to compare human- and LLM-generated texts across five dimensions: (1) emotional appeal (Affect, emo_pos, emo_neg, emo_anx, emo_anger), where emotions may enhance manipulative persuasion; (2) credibility (certainty, tentative, insight, cause),

LLM reasonings while choosing true or false given a claim

Believe that claim is True

- * The article is reporting on a real event, where the US Congress failed to extend emergency unemployment benefits for 1.3 million Americans who lost their jobs due to the COVID-19 pandemic.
- * The statement aligns with multiple reliable news sources reporting on the hostage situation at a kosher supermarket in Vincennes, which is part of eastern Paris. It's important to note that this event was indeed a separate attack from the Charlie Hebdo
- * True, as it accurately reports a factual update on an ongoing news event without exaggeration or falsehoods.

Believe that claim is False

- * The argument is based on emotional appeal rather than factual information; it lacks credible sources and concrete data to substantiate its claim.
- * There isn't enough credible evidence to support such an extraordinary claim
- * The article might be using sensational language and providing only one side of story, hence it requires verification from other credible sources before accepting it as true.

Figure 15: LLM logs (while choosing if claim is true or false)- Examples from across datasets and LLMs

LLM reasonings while choosing true or false given a claim

Believe that claim is True

- * Hard to truly determine. Seems likely in retrospect based on the persuasion texts
- * It didn't happen but reporting cautious optimism is a truthful statement
- * If I lived nearby I would take it seriously and true and likely share it. If I did not live nearby then I wouldn't need to know until later and wouldn't share it

Believe that claim is False

- * Circular reasoning. People who eat out also probably go to concerts and do other social things. The restaurant can't be blamed without more information
- * I don't know the source of the information or whether it is credible so I cannot judge whether it is factual or not, I think it is misinformation
I do not know any of the people in the source of information except that they claim something that only personal friends know. Not trusted information unless I knew the source personally.

Figure 16: Human annotator notes (while choosing if claim is true or false)- Examples from across datasets and demographic groups

with high certainty signaling authoritative claims; (3) logical structure (causation, insight, discrepancy), where imbalanced causation-to-insight ratios may indicate fallacious reasoning; (4) social dynamics (social, family), reflecting in-group appeals that reinforce echo chambers; and (5) cognitive complexity (cognitive processes, insight, discrepancy), where lower scores suggest simplistic arguments. This approach builds on established links between linguistic features and persuasion in misinformation contexts. We find that LLM-generated persuasive texts show higher

Prompt for Multi-Agent LLM Interaction

Initial Round

Assume you are a person from a **<demographic group>**. Given the source news: **<c>**. Based on your background as a **<demographic group>** person, determine if this is true information or misinformation. Respond in this format: **<true information or misinformation>**.

Persuasion Provided Assume you are a person from a **<demographic group>**. Given the source news: **<c>**, and two stances present: **<i>** supporting the claim and **<j>** refuting the claim. Based on your background as a **<demographic group>** person, **<i>** and **<j>**, determine if this is true information or misinformation. Respond in this format: **<true information or misinformation>**.

Discussion Round 1

Assume you are a person from a **<demographic group>**. Given the source news: **<c>**, and two stances present: **<i>** supporting the claim and **<j>** refuting the claim. Convince others about your own perspective. Based on **<memory>**, you have to choose if information is true or false and respond in this format: **<true information or misinformation>** **<brief reason>**. Limit responses to 2 sentences.

Discussion Round 2

Assume you are a person from a **<demographic group>**. Given the source news: **<c>**, and two stances present: **<i>** supporting the claim and **<j>** refuting the claim. Now be open to other perspectives so that you can come to a consensus. Based on **<memory>**, you have to choose if information is true or false and respond in this format: **<true information or misinformation>** **<brief reason>**. Limit responses to 2 sentences.

Final Round Assume you are a person from a **<demographic group>**. Given the source news: **<c>**. Based on your **<memory>**, determine if this is true information or misinformation. Respond in this format: **<true information or misinformation>**.

Figure 17: Multi-Agent LLM interaction prompt design

emotional appeal, logical structure, and cognitive complexity, whereas human-generated texts have higher scores for credibility and social dimensions. The observed differences likely arise from LLMs’ training on large-scale, engagement-optimized corpora, which emphasize emotional resonance (e.g., heightened positive/negative affect), explicit logical markers (e.g., causation terms), and lexical diversity (Mieszczewicz-Kowszewska et al., 2025; Juzek and Ward, 2025), inflating their scores in emotional appeal, logical structure, and cognitive complexity (Breum et al., 2024). In contrast, human writers prioritize credibility through nuanced certainty/hedging patterns (Halomoan, 2024; Yao et al., 2023) and leverage socially grounded persuasion (e.g., in-group appeals) (Fleming, 2025), which LLMs may not be able to replicate very authentically. These differences reflect LLMs’ strength in formulaic linguistics and humans’ reliance on pragmatic and relational persuasion strategies.

F Multi-Agent Persuasion

F.1 Prompt design - Interaction

Fig 17 shows the interaction flow for the multi-agent LLM persuasion set-up.

F.2 Persona Prompting

We use the same demographic groups - female/male, rural/urban, and old/young. We initialize personas such as - “Assume you are a person from **<demographic group>**”.

To better understand beliefs associated with these personas initializations, we prompt the LLMs to output their beliefs associated with these personas. We use this prompt to achieve the above: “Assume you are a **<demographic group>**. Which beliefs do you think define people like you the most, and how important is it for you to share them with others? Give me the top 10 beliefs.” Table 10 shows per model beliefs given a certain demographic persona. We find that persona beliefs across demographics reveal stereotypical associations that resonate with long-standing social science constructs. Male personas emphasize integrity, hard work, and individual achievement, aligning with Hofstede’s “masculinity” dimension (Hofstede, 2011) and the psychological framing of men as achievement- or agency-oriented. Alternatively, female personas emphasize empathy, kindness, feminism, and community, mapping onto relational and communal orientations documented in gender role theory. Rural personas stress tradition, self-reliance, family,

Demographic	LLaMA-3-70B	Qwen-3-72B	GPT-4o
Male	<ol style="list-style-type: none"> 1. Empathy is essential 2. Education is key to progress 3. Hard work and perseverance are essential 4. Diversity enriches our lives 5. Mental health matters 6. Integrity is non-negotiable 7. Family is everything 8. Personal growth is a lifelong journey 9. Community involvement is vital 10. Authenticity is freedom 	<ol style="list-style-type: none"> 1. The importance of integrity and honesty 2. Continuous learning and self-improvement 3. Respecting diversity and promoting equality 4. Supporting mental health awareness 5. Embracing technology and innovation 6. Advocating for environmental sustainability 7. Power of community and collective action 8. Valuing family and close relationships 9. Encouraging open-mindedness and empathy towards different cultures and viewpoints 10. Striving for work-life balance 	<ol style="list-style-type: none"> 1. Self-reliance and independence 2. Protecting loved ones 3. Hard work and perseverance 4. Respect and honor 5. Emotional control 6. Leadership and initiative 7. Achievement and success 8. Loyalty and brotherhood 9. Courage and risk-taking 10. Passing on values
Female	<ol style="list-style-type: none"> 1. Empathy is a superpower 2. Education is empowerment 3. Diversity is strength 4. Mental health matters 5. Feminism is for everyone 6. Kindness goes a long way 7. Vulnerability is courage 8. Self-reflection is key 9. Community is everything 10. Growth is a lifelong journey 	<ol style="list-style-type: none"> 1. Empowerment through education 2. Gender equality 3. Intersectionality 4. Reproductive rights 5. Work-life balance 6. Supportive community 7. Mental health awareness 8. Environmental stewardship 9. Economic independence 10. Cultural appreciation 	<ol style="list-style-type: none"> 1. Equality between genders 2. Importance of education and self-growth 3. Empathy and emotional intelligence 4. Respect for personal boundaries 5. Support for other women and marginalized groups 6. Value of family and close relationships 7. Belief in hard work and perseverance 8. Advocacy for bodily autonomy and consent 9. Community building and collaboration 10. Pursuit of independence and self-reliance
Rural	<ol style="list-style-type: none"> 1. Hard work is essential 2. Community matters 3. Tradition is valuable 4. The land is sacred 5. Family comes first 6. Self-reliance is freedom 7. Honesty is paramount 8. Respect for elders is crucial 9. Frugality is wise 10. Faith gives hope 	<ol style="list-style-type: none"> 1. Strong community ties and support systems 2. Respect for nature and the environment 3. Hard work and self-reliance 4. Faith and spirituality 5. Traditional family structures and values 6. Local governance and decision-making 7. Preserving cultural heritage and local traditions 8. Education and lifelong learning 9. Personal responsibility and accountability 10. Economic independence through small businesses and agriculture 	<ol style="list-style-type: none"> 1. Strong sense of community and mutual support 2. Importance of tradition and cultural heritage 3. Deep connection to land and nature 4. Value of hard work and self-reliance 5. Faith and spirituality as guiding principles 6. Respect for elders and their wisdom 7. Belief in simple, modest living 8. Emphasis on family bonds and responsibilities 9. Pride in local identity and way of life 10. Importance of sharing knowledge and values within the community
Urban	<ol style="list-style-type: none"> 1. Diversity is strength 2. Sustainability is key 3. Public transportation is essential 4. Affordability matters 5. Food should be accessible and diverse 6. Community engagement is vital 7. Innovation drives progress 8. Mental health support is crucial 9. Walkability and bikeability are essential 10. Lifelong learning is necessary 	<ol style="list-style-type: none"> 1. Diversity enriches communities 2. Sustainable living practices 3. Access to quality education 4. Public transportation systems should be efficient, affordable, and accessible 5. Healthcare should be accessible and affordable 6. Urban spaces should prioritize pedestrian and cyclist safety 7. Affordable housing is essential 8. Technology can significantly enhance urban living 9. Community involvement and local governance are vital 10. Cultural and recreational activities are important 	<ol style="list-style-type: none"> 1. Innovation and progress 2. Diversity and inclusion 3. Environmental consciousness 4. Education and lifelong learning 5. Equality and social justice 6. Technological advancement 7. Community engagement 8. Cultural openness 9. Individual freedom 10. Economic opportunity
Young	<ol style="list-style-type: none"> 1. Diversity and inclusivity are essential 2. Climate change is real and urgent 3. Mental health matters 4. Education should be accessible and affordable 5. Technology can be a powerful tool for good 6. Authenticity is key 7. Social media can be both empowering and toxic 8. Creativity and self-expression are essential 9. Empathy and active listening are crucial 10. Personal freedom and autonomy are fundamental rights 	<ol style="list-style-type: none"> 1. Importance of diversity and inclusion 2. Climate change is an urgent issue 3. Technology should be used to enhance human connection rather than replace it 4. Mental health is just as important as physical health 5. Education should be accessible and affordable to everyone 6. Social media can be a powerful tool for positive change when used responsibly 7. Personal growth and self-improvement are lifelong processes 8. Equality and fairness in economic opportunities and social rights 9. Cultural sensitivity and understanding are crucial 10. Essential to challenge and question authority 	<ol style="list-style-type: none"> 1. Freedom to express oneself 2. Valuing diversity and inclusion 3. Belief in equality and social justice 4. Importance of mental health and well-being 5. Climate change is real and urgent 6. Technology can drive positive change 7. Questioning traditional norms and authorities 8. Community and collaboration over hierarchy 9. Personal growth and authenticity 10. Sharing ideas widely to create impact

Demographic	LLaMA-3-70B	Qwen-3-72B	GPT-4o
Old	<ol style="list-style-type: none"> 1. Family comes first 2. Hard work and dedication lead to success 3. Respect your elders 4. Community is key 5. Tradition matters 6. Honesty is always the best policy 7. Education is power 8. Save for a rainy day 9. Take time to appreciate the little things 10. Leave the world a better place than you found it 	<ol style="list-style-type: none"> 1. Respect for tradition and the wisdom 2. The importance of family ties 3. Hard work and perseverance 4. Integrity and honesty 5. Saving and prudent financial management are crucial 6. Community involvement and helping those in need 7. Physical health and mental well-being 8. Adaptability to change 9. Patience and understanding 10. Sharing life experiences and lessons learned with younger generations 	<ol style="list-style-type: none"> 1. Respect for traditions and customs 2. Value of family and community 3. Importance of hard work and perseverance 4. Gratitude for simple things in life 5. Faith or spiritual belief as a guiding force 6. Responsibility to pass on knowledge and values 7. Preference for stability and order 8. Belief in respect for elders and authority 9. Emphasis on honesty and integrity 10. Appreciation for face-to-face social connection

Table 10: Top Beliefs by Demographic Groups Across Models. We highlight values that differ across counterpart demographics. For example, in llama-3-70b, both male and female groups share “importance of education” (not highlighted), whereas “hard work and perseverance” (male) and “kindness” (female) differ (do not overlap) and are highlighted. Several values overlap across gender and, to some extent, rural/urban groups, while the greatest misalignment appears between younger and older demographics. We also observe cross-model similarities in value patterns. Experiments are conducted with gpt-4o here, as gpt-3.5-turbo is no longer available on Azure OpenAI, which was earlier available while performing persona experiments.

Group	Fin/Init	Avg Len	Lex Div	Read	Conf Shift	Emo Infl	Cog Eng	Beh Read	Ech Cham
<i>Rural vs. Urban</i>									
Rural	Final	21.34	0.89	2.86	-2.13	-0.33	14.43	0.11	0.16
Rural	Initial	14.24	0.87	-20.38	-1.38	-0.12	15.48	0.07	0.17
Urban	Final	22.75	0.90	1.82	-2.33	-0.40	13.53	0.06	0.32
Urban	Initial	13.10	0.87	-22.91	-1.44	-0.12	15.15	0.07	0.16
<i>Gender</i>									
Female	Final	21.05	0.88	-0.49	-2.19	-0.37	12.15	0.32	0.04
Female	Initial	13.26	0.87	-22.48	-1.82	-0.14	13.52	0.51	-0.21
Male	Final	21.55	0.89	-1.17	-1.99	-0.44	11.69	0.17	0.27
Male	Initial	13.03	0.87	-22.92	-1.62	-0.16	13.09	0.08	0.21
<i>Age Group</i>									
Young	Final	21.53	0.89	-1.23	-2.22	-0.38	14.07	0.14	0.16
Young	Initial	13.71	0.87	-21.39	-1.21	-0.17	15.10	0.07	0.16
Old	Final	22.12	0.90	3.65	-2.15	-0.43	13.03	0.10	0.23
Old	Initial	13.75	0.87	-22.48	-1.00	-0.11	14.81	0.10	0.28

Table 11: Structural and LIWC Analysis of Persuasion Shifts in Multi-Agent LLM conversations. We highlight the higher dimensions across demographics for structural analysis - Avg Len(gth), Lex(ical) Div(ersity), Read(ability). Note that lower readability scores mean higher readability. For LIWC Analysis, Conf(idence) Shift, Emo(tional) Infl(uence), Cog(nitive) Eng(agement), Beh(avioral) Read(iness), and Echo Cham(ber), we highlight the higher scores for each demographic group.

and faith – values tied to conservatism and collectivism in sociological studies (Villa and Moya, 2020), while urban personas emphasize diversity, sustainability, innovation, and accessibility, showing cities as multicultural, progressive hubs (Cheshmehzangi et al., 2024). Age-related stereotypes are also evident: younger personas emphasize inclusivity, climate change, and self-expression (consistent with social science research on youth activism and identity politics (Bussoletti, 2023; Millhouse, 2024)), whereas older personas focus on tradition, family, saving, and patience, aligning with generational theories of conservatism and life-course perspectives (Kunkel and Settersten, 2021). These alignments suggest that the LLMs are not inventing new constructs but reproducing socially grounded stereotypes that may be embedded in training data.

Across models, the outputs are mostly overlapping for LLaMA and Qwen, which indicates a convergence of learned associations from similar training distributions rather than model-specific beliefs. There are a few differences: for example, Qwen highlights innovation and sustainability for men, and intersectionality and economic independence for women, suggesting more progressive framings. Additionally, opposing demographics also share overlapping beliefs: both male and female personas highlight education and mental health; rural and urban personas highlight community; young and old agree on family and integrity. Such overlaps suggest that while LLMs amplify group differences, they also construct shared cultural anchors. In misinformation contexts, these associations may become critical: stereotypes may exaggerate demographic divides (e.g., rural = traditional, urban = progressive), while shared values (e.g., family, community) could be leveraged to spread narratives across groups. Thus, persona-based prompting showcases both faults and bridges through which misinformation can spread, underscoring the need for cautious framing of findings and demographic-aware LLM design.

F.3 Per-model Results for Persuasion Impact

Table 7 show the ΔCR for multi-agent LLM models when human and LLM persuasions are provided in different demographic and group types. We observe similar trends for human and LLM persuasion settings across models and groups types, however, magnitudes of ΔCR change across models.

F.4 Analyzing Structural/Linguistic and Persuasion-Induced Effects in Multi-Agent Framework

We perform the following qualitative and quantitative analyses to understand the effects of persuasion during interactions in demographic-aware multi-agent settings. For our analyses, we specifically analyze heterogeneous demographic groups. This approach enables us to identify how persuasion impacts vary across populations with distinct characteristics, revealing targeted vulnerabilities and informing tailored interventions.

Structural Analysis. This is conducted by analyzing text length and readability (Flesch-Kincaid Grade Level (FK-GL) ⁸). It reveals distinct linguistic patterns across demographics: **Shorter texts** (*rural +7.1 words, young +7.8 words*) with **lower readability scores** (*rural 2.86, young -1.23*) suggest reliance on simplified emotional appeals, as conciseness and accessibility may enhance engagement for these groups. In contrast, urban/older demographics exhibit **longer texts** and **higher readability** refinement (*urban readability +1.82, older +3.65*), aligning with evidence-based persuasion strategies that prioritize lexical diversity and structural complexity. Notably, gender differences reveal subtle but meaningful patterns: female agents show marginally **lower lexical diversity** (*female 0.8855 vs. male 0.8890*) yet comparable readability shifts, while males achieve slightly better **final readability** (*male -1.17 vs. female -0.49*). This implies that while both genders adapt similarly to readability demands, males may optimize slightly more for clarity in persuasive contexts. Magnitudinal differences across demographics are not huge, however, these findings underscore that text length and readability are non-neutral dimensions of persuasion, showing how different audiences process different types of appeals.

LIWC Analysis. We compute the following dimensions of persuasion effects using LIWC: (1) *Confidence Shift* (Tormala, 2016) is the difference between certainty and tentative words, a larger (less negative) score signals more confident, less hedged language. (2) *Emotional Influence* (Rocklage et al., 2018) subtracts the sum of negative emotion and anxiety from positive emotion, so higher values indicate a net positive emotional tone and lower hostility. (3) *Cognitive*

⁸https://en.wikipedia.org/wiki/Flesch%E2%80%9393Kincaid_readability_tests

Engagement (Tausczik and Pennebaker, 2010) adds insight, cause, and discrep terms, greater totals reflect deeper reasoning and self-reflection. (4) *Behavioral Readiness*⁹ sums inclusive-action cues (we + impulse), capturing readiness to act collectively. Finally, (5) *Echo Chamber* (Wang et al., 2024a) is calculated as (they – we), higher scores mean stronger out-group focus and greater polarization. We find that after the multi-agent persuasion, every demographic group shows weaker attitude certainty (**Confidence goes down**), with the steepest drops for *urban* (-1.44 to -2.33) and *young participants* (-1.21 to -2.22). Affect also turns more negative (**Emo-Infl decreases**), especially for *males* (-0.16 to -0.44) and *older adults* (-0.11 to -0.43). **Cognitive engagement reduces** across all groups. **Behavioral readiness, however, rises** slightly for rural, male, and young cohorts (e.g., *young* 0.07 to 0.14) but **slips** for urban and female groups. **Echo-chamber language increases** for *urban* (0.16 to 0.32) and *female* (-0.21 to 0.04) demographics, yet reduces for rural and older ones. Combined with the structural results, demographics that favor concise, lower-readability text (rural, young) emerge less certain but more willing to act, whereas those accustomed to denser discourse (urban, female, older) leave more polarized and emotionally negative, with lower mobilization intent. Table 11 shows the results of structural and LIWC analysis for the multi-agent persuasion framework.

Persuasion Effect Analysis. First, we analyze stance changes during multi-agent interactions due to persuasion, i.e., which demographics change their initial belief stances about information after persuasion (true -> false or false -> true). We find that *female* (7.01%), *rural* (9.19%), and *older* (8.04%) demographic groups exhibited significantly greater belief shifts between initial and final stances compared to their counterparts. This suggests LLM personas belonging to these demographics may be more responsive to persuasive arguments in conversational settings, potentially due to factors like higher engagement with opposing views or greater susceptibility to social influence (Tang et al., 2024; Wang and Chen, 2006; Tarrant et al., 1997). Table 12 shows the results of stance changes per demographic after persuasion in multi-agent interaction.

Building on the framework established by (Tessler et al., 2024), we operationalize

⁹Stanford GSB (2010)

Group	Pct. Flips (%)
<i>Rural vs. Urban</i>	
Rural	18.37
Urban	15.02
<i>Gender</i>	
Female	14.00
Male	11.60
<i>Age Group</i>	
Young	16.08
Old	16.43

Table 12: Percentage of Stance Changes per Demographic after persuasion

persuasion effects through three key dimensions: (1) *emotional shift*, (2) *coverage*, and (3) *specificity*. **Emotional shift** is quantified using JS-divergence of sentiment (Elahimanesh et al., 2025) between initial and final responses, capturing how persuasive interactions alter affective tone. This measure reveals whether arguments succeed through emotional appeals versus rational discourse. **Coverage** evaluates content retention by comparing information preserved in final responses, serving as an indicator of evidence integration versus echo-chamber behavior (Tessler et al., 2024). Higher coverage values suggest engagement with opposing evidence, while lower values may indicate ideological entrenchment. **Specificity** is measured through average Inverse Document Frequency (IDF) of response content, where lower scores reflect reliance on generic language that may signal manipulative vagueness (Sparck Jones, 1972), while higher scores indicate concrete, substantive arguments.

Our analysis reveals distinct demographic patterns in persuasion susceptibility of LLM demographics. **Rural/younger groups** demonstrate *pronounced emotional shifts* coupled with *concrete language use (high specificity)*, making them particularly vulnerable to anecdotal misinformation (e.g., “My neighbor got sick from vaccines”). This suggests their decision-making may prioritize emotionally charged personal narratives over systematic evidence evaluation. Conversely, **urban/older groups** exhibit *stronger evidence integration (high coverage)*, indicating more analytical processing that helps counter misinformation but may also lead to *dismissal of emotionally compelling truths*. **Female group** shows a unique profile of maintaining *high coverage while remaining emotionally en-*

gaged, suggesting a balanced deliberative style that integrates both affective and evidentiary appeals. **Male/older demographics** display resistance to vague claims (*high specificity and low emotional shift*), though this potential strength may come at the cost of reduced flexibility when updating beliefs in light of new evidence.

Group	Emotional Shifts	Coverage	Specificity
Rural	0.15	0.27	3.11
Urban	0.13	0.28	3.05
Female	0.14	0.26	3.11
Male	0.12	0.28	3.01
Young	0.14	0.26	3.14
Old	0.13	0.28	3.05

Table 13: Deliberation Metrics per Demographic

These dynamics highlight how misinformation can exploit emotional-concrete appeals for rural/young/female audiences, while factual precision works best for urban/older/male groups, a critical insight for designing demographic-tailored interventions in LLM-based personas and moderation tools.

Connection to Prior Studies. Prior human studies echo the demographic patterns found in our multi-agent LLM framework. Rural populations have been shown to rely more on emotionally vivid anecdotes and display lower trust in institutional evidence, increasing their vulnerability to misinformation (Lister and Joudrey, 2022; Tarrant et al., 1997). Younger adults similarly pay importance to affective cues over systematic reasoning, especially when messages are concrete and narrative-driven (Wang and Chen, 2006; Ta et al., 2022). In contrast, urban residents and older adults engage in more analytical, evidence-integrating processing, which boosts accuracy but can reduce responsiveness to emotional appeals (Duke and Whatley, 2021b; Brashier and Schacter, 2020). Gender studies find that women often attend to both emotional tone and factual detail, whereas men favor specificity and exhibit lower emotional shift, leading to greater resistance to ambiguous claims but less flexibility when new evidence arrives (Pan et al., 2021; Almenar et al., 2021; Enock et al., 2024). Together, these behavioral findings align closely with the persuasion-susceptibility signatures we observe in LLM personas in a multi-agent setting.

G Significance Testing for all experiments

We perform statistical significance testing for the persuasion experiments shared in the paper. We use the following tests for each experiment, depending on the framework:

1. For Human-to-LLM correctness, we use the chi-squared test ¹⁰.
2. For LLM-to-Human correctness, we use the Fisher’s exact test ¹¹ due to small sample sizes (95-152).
3. For Human-LLM Correlation, we perform a permutation test ¹².
4. For Multi-Agent experiments, we do a paired t-test ¹³.

EXPERIMENT	COMPARISON	P-VAL	SIG.?
Human-to-LLM	Urban > Rural	0.022	Yes
	Young > Old	0.049	Yes
	Male > Female	0.27	No
LLM-to-Human	Urban > Rural	0.042	Yes
	Young > Old	0.08	Marginal
	Male > Female	0.02	Yes
Human-LLM Correlation	GPT-3.5 MCC	<0.001	Yes
Multi-Agent (LLM)	Hom. ΔCR	0.016	Yes
	Het. ΔCR	0.011	Yes
Multi-Agent (Human)	Hom. ΔCR	0.046	Yes
	Het. ΔCR	0.042	Yes

Table 14: Statistical analysis results comparing different experimental conditions

The results demonstrate that most findings are statistically significant ($p < 0.05$). However, two comparisons: gender differences in Human-to-LLM correctness ($p = 0.27$) and age differences in LLM-to-Human correctness ($p = 0.39$) are not significant and should be interpreted with caution. Nevertheless, the overwhelmingly significant results across several conditions (e.g., urban/rural, multi-agent interactions) show the reproducibility of our core contributions, particularly in advancing understanding of human-LLM persuasion asymmetries. Additionally, These findings collectively highlight the importance of context-aware AI communication frameworks, with direct applications in

¹⁰https://en.wikipedia.org/wiki/Chi-squared_test

¹¹https://en.wikipedia.org/wiki/Fisher%27s_exact_test

¹²https://en.wikipedia.org/wiki/Permutation_test

¹³<https://www.jmp.com/en/statistics-knowledge-portal/t-test/paired-t-test>

personalized AI design, bias mitigation, and behavioral modeling.

H Model Choices, Implementation Details and Computational Resources

Our model selections across GPT, Llama, and Qwen are based on three main reasons: (1) architectural/origin diversity (OpenAI, Meta, Alibaba), (2) computational feasibility for large-scale human-LLM experiments, and (3) reproducibility through open-source model availability. All inference experiments are conducted with results averaged over three LLM runs. For gpt-35-turbo, inference is performed using the Microsoft Azure API ¹⁴. The llama-3-70b-instruct¹⁵ and qwen-2.5-72b-instruct¹⁶ models are run via Hugging Face. To ensure focused yet varied text generation, all models are set with a temperature of 0.5. For open-source models, top_p is set to 0.9, with do_sample=True, and 4-bit quantization is applied. Inference for these models is conducted on an NVIDIA-A40 GPU.

I Reproducibility

We open-source our codes and data, which are uploaded to the submission system. This would help future work to reproduce our results

¹⁴<https://learn.microsoft.com/en-us/rest/api/azure/>

¹⁵meta-llama/Meta-Llama-3-70B-Instruct

¹⁶Qwen/Qwen2.5-72B-Instruct