

Anonymous ACL submission

1

(Horta Ribeiro et al., 2023; Schaffner et al., 2024). Large Language Models (LLMs) have been hypothesized to be capable of conversational moderation and facilitation tasks, which often require actively participating in the discussions, instead of passively flagging or removing content (Small et al., 2023; Korre et al., 2025).

While studies exist for simulating user interactions in social media (Park et al., 2022; Mou et al., 2024; Törnberg et al., 2023; Rossetti et al., 2024; Balog et al., 2024), and for using artificial facilitators (Kim et al., 2021; Cho et al., 2024), none so far have combined the two approaches. We posit that synthetic simulations can be a cheap and easy way to develop and test preliminary, in silico experiments with LLM facilitators, initial versions of which may be unstable or unpredictable (Atil et al., 2025; Rossi et al., 2024), before testing them in much more costly experiments with human participants. Our work thus asks the following two questions: (1) Can we produce high-quality synthetic discussions, involving alternative facilitation strategies, by crafting an appropriate environment for simulations? (2) Can we boost the effectiveness of LLM moderators (in synthetic discussions) by using prompts aligned with facilitation strategies proposed in modern Social Science research?

We propose a simple and generalizable approach using LLM-driven synthetic experiments for online moderation research, enabling fast and inexpensive model “debugging” and parameter testing (e.g., LLM moderator prompts, instructions) without human involvement (§3) (Fig. 1). An ablation study (§5.2) demonstrates that each step of our methodology meaningfully contributes to generating high-quality synthetic data. Using this methodology, we examine four LLM moderation strategies (including a novel strategy inspired by Reinforcement Learning (RL)) based on current Social Science facilitation research (§4) and compare them with two baselines (LLM facilitators with simplistic facilitation prompts). LLMs are also used to gauge discussion quality (e.g., argument quality, toxicity).

Our analysis reveals two key findings (§5): (1) the presence of LLM facilitators has a positive and statistically significant influence on the quality of synthetic discussions, and (2) facilitation strategies inspired by Social Science research often do not manage to outperform simpler baselines. Furthermore, we release XXXan open-source Python framework for generating and evaluating synthetic discussions, alongside a large, publicly available

dataset comprising automatically evaluated synthetic discussions (§6). We use open-source LLMs and include all relevant configurations in order to make our study as reproducible as possible (see §A.3, §A.4).

2 Background and Related Work

2.1 LLMs as Human Subjects

When conducting social experiments with LLMs instead of human subjects, it is imperative to know how representative results can be. Grossmann et al. (2023) argue that synthetic agents have the potential to eventually replace human participants, a perspective shared by other researchers (Törnberg et al., 2023; Argyle et al., 2023). Indeed, LLMs have demonstrated emergent complex social behaviors (Park et al., 2023; Marzo et al., 2023; Leng and Yuan, 2024; Abdelnabi et al., 2024; Abramski et al., 2023). They are also capable of predicting human survey responses in aggregate (Hewitt et al., 2024) and in the level of individual people, given extensive personal data (Park et al., 2024).

However, significant limitations of LLMs remain in the context of Social Science experiments. Issues include dataset contamination; undetectable behavioral hallucinations (Rossi et al., 2024); sociodemographic, statistical and political biases (Anthis et al., 2025; Hewitt et al., 2024; Rossi et al., 2024), often amplified during discussions (Taubenfeld et al., 2024); unreliable survey responses (Jansen et al., 2023; Bisbee et al., 2024; Neumann et al., 2025); inconsistent annotations (Gligori’c et al., 2024); non-deterministic outputs (Atil et al., 2025), especially in closed-source models (Bisbee et al., 2024); and excessive agreeableness due to alignment procedures (Park et al., 2023; Anthis et al., 2025; Rossi et al., 2024). Despite these shortcomings, researchers frequently anthropomorphize LLM agents (Rossi et al., 2024), leading to biased interpretations and obscuring the true nature of their behavior (Anthis et al., 2025; Zhou et al., 2024a). Our study must thus be conservative towards the generalizability of our results to discussions with human participants.

We stress that we propose and investigate synthetic discussions as a means of preliminarily “debugging” and exploring artificial facilitators (e.g., with different facilitation strategies) in silico, before testing them in much more costly experiments with human participants. For example, synthetic experiments may help illustrate a clear risk or disad-

vantage of a particular facilitation strategy, that may be otherwise difficult to foresee. Experiments with real participants, however, are ultimately needed, and we leave them for future work.

2.2 Evaluating Discussion Quality

Synthetic discussions often degrade rapidly without human interaction, exhibiting repetitive, low-quality content (Ulmer et al., 2024). However, research on quantifying synthetic data quality is currently limited. Balog et al. (2024) introduce metrics utilizing comparisons with human data, but this approach depends on datasets with the same topics, and lacks scientific grounding since believable LLM outputs do not necessarily lead to behavior simulation (Rossi et al., 2024). Their most generalizable metric—a vague “coherence” score—is LLM-annotated without theoretical support.

Alternatively, Ulmer et al. (2024) propose metrics like N-gram-based “Diversity”. Intuitively, the metric penalizes long, repeated sequences between comment pairs in a discussion. It is defined as:

$$\text{div}(d) = 1 - \frac{2}{N(N-1)} \sum_{i=1}^N \sum_{j=i+1}^N R(c(i, d), c(j, d)) \quad (1)$$

where R is the ROUGE-L F1 score² (Lin, 2004), and N_d is the length (in comments) of discussion d .

Low diversity scores point to pathological problems (e.g., LLM user-agents repeating the previous comments). Extremely high diversity scores, on the other hand, may point to a lack of interaction between participants, as a discussion in which participants engage with each other will feature some lexical overlap (e.g., common terms, paraphrasing points of other participants). For this reason, we compare the distribution of *diversity* scores for synthetic discussions with that measured on sampled human discussions. This comparison allows us to estimate the extent to which synthetic discussions approximate real-world content variety and participant interaction, or at the very least, points to pathological problems in our generated data.

Besides metrics for the quality of synthetic data, we also need metrics that can quantify how “well” a discussion is going from the point of view of the participants, or outside users reading the discussion. We choose Toxicity for two primary reasons: Prompting LLMs for toxicity detection is reliable

(Kang and Qian, 2024; Wang and Chang, 2022; Anjum and Katarya, 2024), and toxicity can inhibit online and deliberative discussions (De Kock et al., 2022; Xia et al., 2020).³

2.3 Synthetic Discussions

Synthetic discussion systems include synthetic clones of Reddit (Park et al., 2022), Twitter/X (Mou et al., 2024) and social media in general (Törnberg et al., 2023; Rossetti et al., 2024) as well as games (Park et al., 2023) and social experiments (Zhou et al., 2024b).

Balog et al. (2024) introduce their own methodology to produce synthetic discussions, where they extract topics and comments from real-world online discussions, and prompt an LLM to continue them. Unlike our approach, they do not use LLM user-agents to model conversational dynamics, nor do they model the presence of facilitators. Their methodology faces challenges when LLMs generate malformed metadata, for which they offer no solution, and relies on the existence of suitable human discussion datasets.

Ulmer et al. (2024) create synthetic discussions between two participants; an agent (who controls the environment) and a client (who interacts with the agent). They then filter the generated discussions and use them as training data to further fine-tune the agent LLM for a specific task. Their approach however does not model the existence of multiple clients (users), nor is it applied on online discussion facilitation. Our proposed methodology can be modelled as a generalization of their paradigm; an agent (moderator) converses with multiple clients (non-moderator users).

Finally, Abdelnabi et al. (2024) create synthetic negotiations with multiple agents having various agendas and responsibilities. Our work can be modelled as a domain shift of their methodology from negotiations, to discussion facilitation; participants with different motivations (i.e., normal users, trolls, long-standing community members), interact with themselves and a stakeholder holding veto power (facilitator) who presides over the discussion.

2.4 LLM Facilitation

Unlike traditional ML models, LLMs can actively facilitate discussions (Korre et al., 2025). They can warn users for rule violations (Kumar et al., 2024), monitor engagement (Schroeder et al., 2024), ag-

²We use the rouge-score package in our analysis.

³We note that this is not always true (Avalle et al., 2024).

gregate diverse opinions (Small et al., 2023), provide translations and writing tips, which is especially useful for marginalized groups (Tsai et al., 2024). These capabilities suggest that LLMs may be able to assist or even replace human facilitators in many tasks (Seering, 2020).

Moderator chatbots have shown promise; Kim et al. (2021) demonstrated that simple rule-based models can enhance discussions, although their approach was largely confined to monitoring and balancing user activity. Cho et al. (2024) use LLM facilitators in human discussions, with moderation strategies based on Cognitive Behavioral Therapy and the work of Rosenberg and Chopra (2015). They show that LLM facilitators can provide “specific and fair feedback” to users, although they struggle to make users more respectful and cooperative. In contrast to both works, our work uses exclusively LLM participants (and LLM facilitators), and tests them in an explicitly toxic and challenging environment.

3 Methodology

3.1 Defining Synthetic Discussions

We assume that the h most recent preceding comments at any given point in the discussion provide sufficient context for the LLM agents (users, facilitators, annotators) (Pavlopoulos et al., 2020). This approach eliminates the need for additional mechanisms such as summarization (Balog et al., 2024), LLM self-critique (Yu et al., 2024), or memory modules (Vezhnevets et al., 2023), resulting in reduced computational overhead and a more transparent, explainable system.

Additionally, we assume that three key functions define the structure of synthetic discussions:

- Underlying model ($LLM(\cdot)$).
- Turn-taking function (t): Determines which user speaks at each turn.
- Prompting function (ϕ): Provides each participant with a personalized instruction prompt, including information such as name and SDB.

We can then model a synthetic comment c at position i of a discussion d recursively as:

$$c(d, i) = LLM(\phi(t(d, i)) \mathbin{++} [c(d, j)]_{i-h}^{i-1}) \quad (2)$$

where $++$ is the string concatenation operator, h is the context length of the LLM user-agent (how many preceding comments they can “remember”), and $[c(d, j)]_{i-h}^{i-1} \dots$ denotes the concatenation of the previous h comments.

Our formulation of synthetic discussions not only keeps the system simple, but also enables controlled experimentation with various alternatives for each of the three functions (Section 5.2).

3.2 Turn Taking

In online discussions, users do not take turns uniformly, nor do they randomly select which comments to respond to. Instead, they often create “comment chains” where they follow up on responses to their own previous comments. To simulate this, our proposed function chooses between the preceding user and another random user for each turn in the discussion:

$$t(i) = \begin{cases} \text{unif}(U) & i = 1, i = 2 \\ \text{unif}(U / \{t(i-1)\}) & i > 2, p = 0.6 \\ t(i-2) & i > 2, p = 0.4 \end{cases} \quad (3)$$

where U is the set of all non-facilitator users, unif is a function sampling from the uniform distribution, and p represents the probability of the corresponding option being selected. When a facilitator is present, t alternates between picking a normal user and the facilitator (the latter decides whether to respond to or not—the LLM producing an empty string is equivalent to not responding).

3.3 Prompting

SocioDemographic Backgrounds (SDBs) have proven promising in generating varied responses, and alleviating the Western bias exhibited by LLMs (Burton et al., 2024). We generate characteristics for 30 LLM user personas with unique SDBs by prompting a GPT-4 model (OpenAI et al., 2024) (§A.4.1). We do not explicitly include political positions in the prompts of the participants, since instruction-tuned LLMs have been shown to be inherently left-leaning—which can not be alleviated by prompting alone (Taubenfeld et al., 2024)—and research in the field has predominantly occupied Western politics (Taubenfeld et al., 2024; Potter et al., 2024; Rozado, 2024; Pit et al., 2024). Following the paradigm presented by Abdelnabi et al. (2024), we assign roles to non-facilitator user-agents, which inform their incentives for participating in the discussion (e.g., helping the community or disrupting discussions). Each role was mapped to specific instructions (§A.4.3). We create three roles for users: neutral, trolls, and community-focused users. Finally, we select a user instruction prompt (§A.4.2) which instructs participants that

repeatedly toxic posts *should* influence their behavior.

4 Experimental Setup

4.1 Moderation Strategies

We test four different facilitation strategies,⁴ along with two naive ones that serve as baselines for discussion facilitation:

1. **No Moderator:** A *baseline* where no facilitator is present.
2. **No Instructions:** A *baseline* where a LLM facilitator is present, but is provided only with basic instructions (e.g., “You are a moderator, keep the discussion civil”).
3. **Moderation Game:** Our proposed *experimental* strategy, inspired by Abdelnabi et al. (2024) (§2.3). Instructions are formulated as a game, where the facilitator tries to maximize their scores by arriving at specific outcomes (e.g., “User is toxic: −5 points, User corrects behavior: +10 points”). No actual score is being kept; they exist to act as indications for how desirable an action or outcome is. The other participants are not provided with scores, nor are they aware of the game rules.
4. **Rules Only:** A *real-life* strategy where the prompt is adapted from LLM alignment guidelines (Huang et al., 2024). This provides the facilitator with a set of rules to uphold, without specifying how to uphold them (e.g., “Be fair and impartial, assist users, don’t spread misinformation”).
5. **Moderation guidelines:** A *real-life* strategy based on guidelines given to human facilitators of Cornell e-Rulemaking Initiative (CeRI) (eRulemaking Initiative, 2017). For example, “Stick to a maximum of two questions, use simple and clear language, deal with off-topic comments”).
6. **Facilitation guidelines:** A *real-life* strategy based on the human facilitation guidelines used by the MIT Center for Constructive Communications (White et al., 2024). For example, “Do not make decisions, be a guide, provide explanations”).

4.2 Technical Details

For toxicity annotation, we use ten LLM annotator-agents controlled by a model already used in prior work (LLaMa3.1 70B) (Kang and Qian, 2024).

⁴The exact prompts used per strategy are in §A.4.4.

Each annotator’s prompt includes SDBs distinct from the ones provided to the users, annotation instructions, and few-shot examples (§A.3). Each annotator is tasked with annotating all comments in each discussion once.

We use three open-source models (in Eq 2) from different families and of different sizes: LLaMa 3.2 (70B), Qwen2.5 (33B) and Mistral Nemo (12B). We select the instruction-tuned variants and quantize them to 4 bits, due to our limited resources. The original and ablation experiments were collectively completed within roughly four weeks of computational time, using two Quadro RTX 6000 GPUs. The execution script is available in the project’s repository⁵. The automated discussion generation is detailed in §A.2.

5 Results

5.1 Main findings

LLM facilitators significantly improve synthetic discussions. As is shown in Fig. 4, comments in unmoderated discussions exhibit significantly worse toxicity (ANOVA $p < .000$).⁶

The effect of LLM facilitators is amplified over time under all strategies when compared to unmoderated discussions. Table 1 demonstrates that, for example, with *Mod. Guid.*, conversations begin with 0.277 lower toxicity, and the *Mod. Guid. × time* interaction shows that with each additional dialogue turn, toxicity decreases on average by another 0.023 points.

Sophisticated facilitation strategies do not qualitatively further improve synthetic discussions. The impact of the “Rules Only”, “Moderation” and “Facilitation Guidelines” strategies (§4.1) is marginal, and sometimes even not statistically significant compared to the second baseline (“No Instructions”) (Fig. 4). This suggests that out-of-the-box LLMs may be unable to effectively use which would enable them to effectively use these advanced instructions, verifying recent research demonstrating important limitations in LLM facilitators (Cho et al., 2024).

LLM facilitators choose to intervene far too frequently. Fig. 2 demonstrates that LLM facilitators intervene at almost any opportunity, even

⁵anonymous.4open.science/r/experiments-B27D

⁶The large size of our dataset allows the use of parametric tests.

though they are instructed to only do so when necessary. Additionally, a qualitative look through the dataset reveals that LLM user-agents exhibit atypical tolerance for excessive moderator interventions. Humans in contrast, typically become irritated and more toxic after repeated, unneeded interventions (Schaffner et al., 2024; Amaury and Stefano, 2022; Schluger et al., 2022; Cresci et al., 2022).

Mistral and Qwen generate discussions more aligned with human diversity scores, despite being significantly smaller than the LLaMa model. As is shown in Fig. 5a, Qwen demonstrated the highest diversity among the evaluated models, indicating limited participant interaction (§2.2), followed by Mistral Nemo and LLaMa. However, none of the models closely matched the diversity observed in human discussions. LLaMa’s lower diversity validates prior research suggesting that highly aligned LLMs struggle to replicate human dynamics (Park et al., 2023; Leng and Yuan, 2024). Alternatively, it can also be attributed to its longer average comment length (Fig. 6a); we find that there is a statistically significant negative correlation between comment length and diversity in synthetic discussions ($p < .000$), although we can not verify this pattern in human-generated texts ($p = 0.775$). Despite these differences, the performance gaps between models are relatively small; notably, smaller models like Mistral generate synthetic data of comparable quality to that produced by larger models such as Qwen and LLaMa.

Specialized instruction prompts are essential for eliciting toxic behavior in instruction-tuned LLMs. Our instruction prompt for the participants (§3.3) incentivizes them to react to toxic behavior. Indeed, discussions involving “Troll” user-agents, led to increased toxicity among *other* participants, even when moderated under the “No Instructions” strategy (Fig. 3, Student’s t-test, $p < .000$). This effect diminishes when we remove these instructions (Fig. 3).

5.2 Ablation Study

In order to assess the impact of each component of our proposed methodology, we generated eight synthetic discussions per ablation experiment, using a single model, Qwen, to limit computational cost. We evaluated the diversity of these generated ablated discussions by comparing their diversity scores (cf. §2.2) with i) discussions in our original dataset produced solely by the Qwen model; and

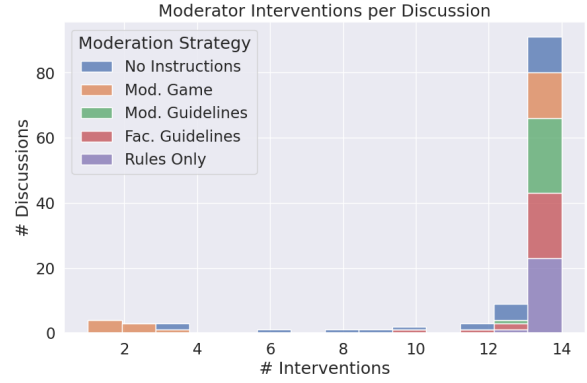


Figure 2: Histogram of interventions by LLM moderators. The maximum number of interventions is 14.

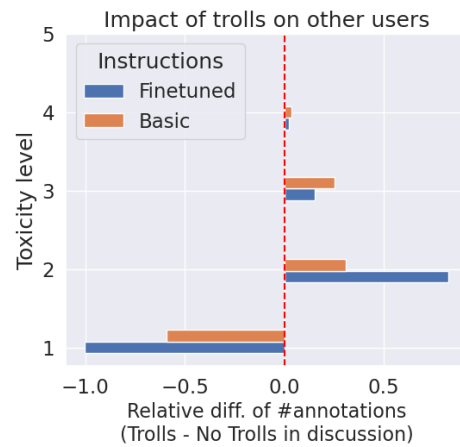


Figure 3: Relative differences in number of annotations per Toxicity of synthetic discussions, when comments by troll users are excluded. We compare between our specialized and a basic instruction prompt.

ii) human discussions from the CeRI “Regulation Room” dataset⁷, which includes moderated online deliberative discussions for ten diverse topics.

5.2.1 Effects of Turn Taking Functions

Our proposed turn-taking function meaningfully improves the quality of synthetic data.

We compare our turn-taking function (§3.2) to two baselines: Round Robin (participants speaking one after the other, then repeating) and Random Selection (uniformly sampling one of the participants each time). Fig. 5b demonstrates that no single function fully approximates human diversity scores (all distributions diverge from the blue—human—distribution). However, unlike our own function, both baselines feature extremely

⁷<http://archive.regulationroom.org> Disclaimer: Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the CeRI

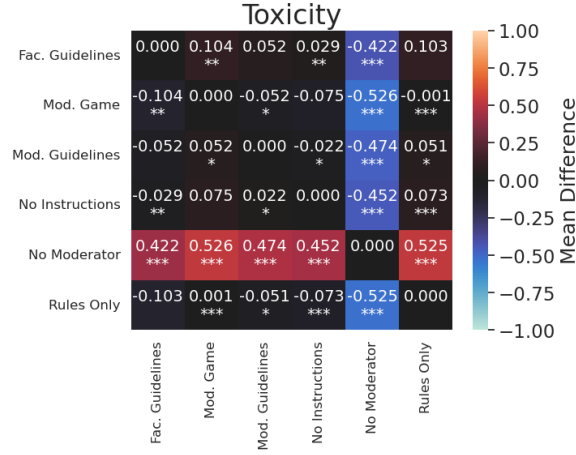


Figure 4: Mean difference of Toxicity between pairs of facilitation strategies. When the value of a cell at row i and column j is x , strategy i leads to overall worse (negative values) or better (positive values) toxicity compared to j for an average of x points in a scale of 1 – 5. For each comparison, we use a pairwise Student t-test; p-values are shown as asterisks ($p < 0.1$, $* p < 0.05$, $** p < 0.01$, $*** p < 0.001$).

high diversity. Additionally, Fig. 6b demonstrates that comments in discussions following our turn-taking function closely follow the length of human discussions.

5.2.2 Effects of User Prompting

We conduct three separate experiments in which user-agents (excluding moderators) are subjected to one of the following conditions at a time: (1) no assigned SDBs, (2) no assigned roles, or (3) only a basic instruction prompt given (§A.4.2).

SDBs, roles and our instruction prompt increase the quality of synthetic data. Fig. 5c illustrates that although our proposed methodology—incorporating SDBs, roles, and specialized instruction prompts—does not achieve discussions with diversity scores comparable to human ones, replacing any of the above results in a notable deterioration. For instance, omitting SDBs (denoted as “No SDBs” and represented by the red distribution in Fig. 5c) causes the majority of discussions to exhibit maximum diversity—one—indicating a significant loss in participant interaction. This decline is analogous to the effects observed when modifying the turn-taking function. Also similarly to the turn-taking ablation study, our proposed methodology w.r.t. prompts, features comments that best emulate observed human comment length (Fig. 6c).

Variable	Toxicity
Intercept	2.164***
Fac. Guid.	-0.230***
Mod. Guid.	-0.277***
RL Game	-0.435***
No Instructions	-0.426***
Rules Only	-0.461***
time	-0.012**
Fac. Guid×time	-0.023***
Mod. Guid×time	-0.023***
RL Game×time	-0.011*
No Instructions×time	-0.003
Rules Only×time	-0.008

$\cdot p < 0.1$, $* p < 0.05$, $** p < 0.01$, $*** p < 0.001$

Table 1: Ordinary Least Squares (OLS) regression coefficients for Toxicity ($Adj.R^2 = 0.054$). “Time” denotes dialogue turn, reference factor is “No moderator”.

6 Datasets & Software

We introduce XXX⁸ an open-source, lightweight, purpose-built framework for managing, annotating, and generating synthetic discussions. Key features include:

- Three core functions: generating, running, and annotating randomized discussion experiments according to provided parameters.
- Built-in fault tolerance (automated recovery and intermittent saving) and file logging to support extended experiments.
- Easy installation via PIP (pip install xxx).

We also release a dataset of synthetic discussions annotated by LLMs for toxicity and argument quality. It can serve as a valuable resource for benchmarking how LLM facilitators would behave according to different facilitation strategies, as well as for further finetuning LLMs, as generally showcased by Ulmer et al. (2024). The supplementary ablation dataset, as well as the code for the analysis and the graphs present in this paper, can be found in the project repository⁹. **Warning: The datasets by their nature contain offensive and hateful speech.**

⁸anonymous.4open.science/r/framework-F8E6

⁹anonymous.4open.science/r/experiments-B27D

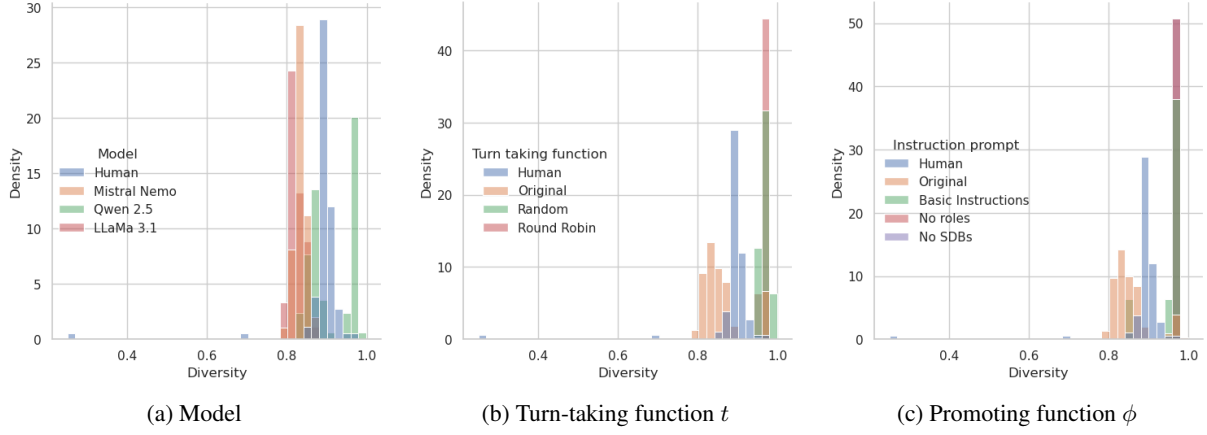


Figure 5: Diversity (§2.2) distribution for each discussion by LLM (§4.2), turn-taking function t (§3.2), and prompting function ϕ used (§3.3).

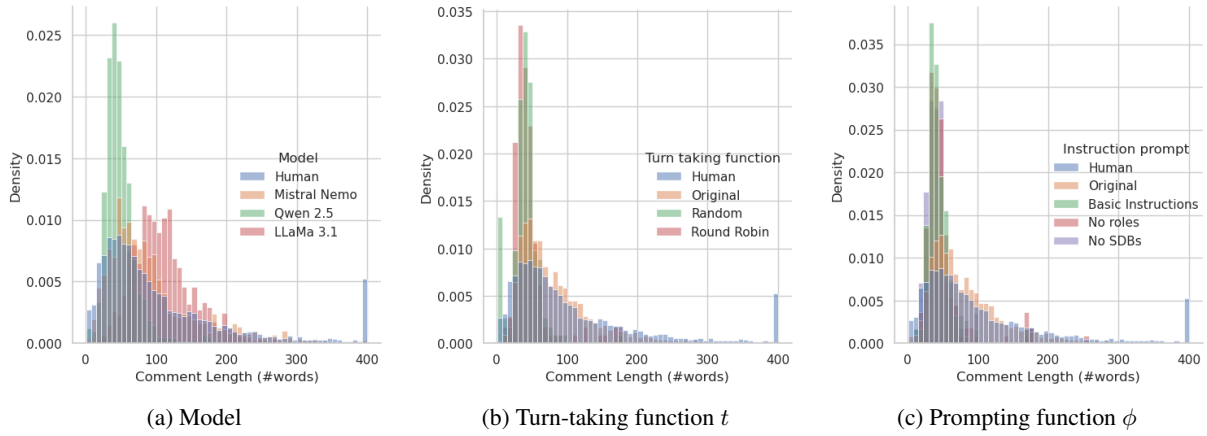


Figure 6: Comment length for each discussion by LLM (§4.2), turn-taking function t (§3.2), and prompting function ϕ used (§3.3). For ease of comparison, comments above 400 words are marked at the end of the x-axis.

7 Conclusions and Future Work

Our study is the first to apply synthetic data generation to the field of online discussion facilitation. We proposed a simple and generalizable methodology that enables researchers to inexpensively conduct pilot facilitation experiments using exclusively synthetic LLMs. We also conducted an ablation study to demonstrate that each component of our methodology contributes to the production of higher-quality synthetic data.

We created an open-source Python Framework, called XXX, that applies this methodology to hundreds of experiments, which we used to create and publish a large-scale synthetic dataset. Using this dataset, we compared the effectiveness of six moderation strategies and baselines for LLM moderators, elicited from current facilitation research.

Using XXX, we demonstrated that (1) LLM moderators significantly improve the quality of synthetic discussions; (2) established human modera-

tion/facilitation guidelines often do not surpass simple baselines with regard to toxicity and Argument Quality (AQ); (3) smaller LLMs such as Mistral Nemo (12B) can be sufficient for generating high-quality synthetic data; (4) specialized instruction prompts may be needed for instruction-tuned models to feature toxic comments in synthetic discussions.

Future work should identify additional robust quality metrics to evaluate the utility of synthetic data, and examine the applicability of findings obtained on synthetic data (e.g., regarding optimal facilitation strategies) to discussions involving humans. It would also be interesting to explore whether non-instruction-tuned models can generate synthetic discussions that are more aligned with observed human behaviors (Anthis et al., 2025). Finally, synthetic discussion simulations may have the potential to train human moderators before exposing them to real-world discussions.

8 Limitations

While we investigate the impact of moderation strategies in synthetic discussions, we cannot claim that the behavior of LLM users and facilitator-agents is representative of human behavior. This claim can be scarcely made in Social Science studies involving LLM subjects (Rossi et al., 2024; Zhou et al., 2024a)—as discussed in §2.1.

Furthermore, our experimental setup makes several assumptions that may affect the generalizability of our findings. We examine only three LLMs, assume a maximum of one facilitator per discussion, and use a turn-taking algorithm that overlooks contextual factors like relevance and emotional engagement, which are crucial in human interactions. Additionally, we do not account for the fact that humans may behave differently when knowing they are interacting with LLMs instead of humans. Our methodology also does not take into account interactions where the user-agents and moderator-agents are based on different LLMs (cf. Eq 2). Finally, our analysis partly relies on LLM-generated annotations, potentially introducing known biases associated with LLM annotation (§A.3).

9 Ethical Considerations

Synthetic discussions involving LLMs could be exploited by malicious actors to make LLM user-agents more capable at performing unethical tasks (Majumdar et al., 2024; Marulli et al., 2024). Such actors could adapt our methodology to maximize toxicity, disrupt human discussions, or learn to circumvent moderation mechanisms to propagate misinformation or spread specific agendas. Notably, LLMs currently lack robust defenses against these types of attacks (Li et al., 2025), although ongoing research is addressing these vulnerabilities (Wang et al., 2025).

Even in non-malicious contexts, researchers deploying LLM moderators in real-world communities must do so with transparency and explicit community consent. The undisclosed use of LLM agents can erode trust, be perceived as manipulative (Retraction-Watch, 2025), and potentially violate regulatory standards such as the EU AI Act (European Parliament and Council, 2024). Furthermore, the inherent biases within LLMs risk skewing moderation systems towards the predominant demographics best represented in their training data, often at the expense of disadvantaged or underrepresented groups (Rossi et al., 2024; Anthis

et al., 2025; Burton et al., 2024). While the use of SDB prompts is a necessary step toward inclusivity, it remains insufficient for verifiably equitable representation (Rossi et al., 2024).

Additionally, our methodology is designed around batch production of synthetic discussions, each of which necessitates multiple LLM inference calls. While significantly more affordable and environmentally friendly than experiments involving humans (given the carbon footprint associated with humans (Ren et al., 2024)), the potential of our methodology to scale experiments by orders of magnitude may still have non-trivial, adverse environmental effects (Ding and Shi, 2024).

Finally, it is crucial to acknowledge that while LLMs can approximate aspects of human behavior, they do not reliably replicate it (§2.1). Consequently, this research should be viewed as a foundation for pilot experiments, and conclusions about human behavior should be drawn with caution when based solely on synthetic data.

References

- Sahar Abdelnabi, Amr Gomaa, Sarath Sivaprasad, Lea Schönherr, and Mario Fritz. 2024. *Cooperation, competition, and maliciousness: Llm-stakeholders interactive negotiation*. *Preprint*, arXiv:2309.17234.
- Katherine Abramski, Salvatore Citraro, Luigi Lombardi, Giulio Rossetti, and Massimo Stella. 2023. *Cognitive network science reveals bias in gpt-3, gpt-3.5 turbo, and gpt-4 mirroring math anxiety in high-school students*. *Big Data and Cognitive Computing*, 7(3).
- T. Amaury and C. Stefano. 2022. *Make reddit great again: Assessing community effects of moderation interventions on r/the_donald*. *Proceedings of the ACM on Human-Computer Interaction*, 6:1 – 28.
- Anjum and Rahul Katarya. 2024. Hate speech, toxicity detection in online social media: a recent survey of state of the art and opportunities. *International Journal of Information Security*, 23(1):577–608.
- Jacy Reese Anthis, Ryan Liu, Sean M. Richardson, Austin C. Kozlowski, Bernard Koch, James Evans, Erik Brynjolfsson, and Michael Bernstein. 2025. *Llm social simulations are a promising research method*. *Preprint*, arXiv:2504.02234.
- Lisa P Argyle, Christopher A Bail, Ethan C Busby, Joshua R Gubler, Thomas Howe, Christopher Rytting, Taylor Sorensen, and David Wingate. 2023. Leveraging AI for democratic discourse: Chat interventions can improve online political conversations at scale. *Proceedings of the National Academy of Sciences*, 120(41):1–8.

706	Berk Atil, Sarp Aykent, Alexa Chittams, Lisheng Fu,	<i>tional Green and Sustainable Computing Conference</i>	763
707	Rebecca J. Passonneau, Evan Radcliffe, Guru Rajan	(IGSC), pages 37–38.	764
708	Rajagopal, Adam Sloan, Tomasz Tudrej, Ferhan Ture,		
709	Zhe Wu, Lixinyu Xu, and Breck Baldwin. 2025. Non-	Cornell eRulemaking Initiative. 2017. Ceri (cor-	765
710	determinism of "deterministic" llm settings . <i>Preprint</i> ,	nell e-rulemaking) moderator protocol . Cornell e-	766
711	arXiv:2408.04667.	Rulemaking Initiative Publications, 21.	767
712	Michele Avalle, Niccolò Di Marco, Gabriele Etta,		
713	Emanuele Sangiorgio, Shayan Alipour, Anita Bonetti,	European Parliament and Council. 2024. Regulation	768
714	Lorenzo Alvisi, Antonio Scala, Andrea Baronchelli,	(eu) 2024/1689 of the european parliament and of	769
715	Matteo Cinelli, and Walter Quattrociocchi. 2024. Per-	the council of 13 june 2024 laying down harmonised	770
716	sistent interaction patterns across social media plat-	rules on artificial intelligence and amending certain	771
717	forms and over time . <i>Nature</i> , 628:582 – 589.	union legislative acts (artificial intelligence act). ht	772
		tps://eur-lex.europa.eu/legal-content/EN/	773
		TXT/?uri=CELEX:32024R1689 . OJ L 2024/1689,	774
		12.7.2024.	775
718	Krisztian Balog, John Palowitch, Barbara Ikica, Filip		
719	Radlinski, Hamidreza Alviri, and Mehdi Manshadi.	Neele Falk, Iman Jundi, Eva Maria Vecchi, and	776
720	2024. Towards realistic synthetic user-generated con-	Gabriella Lapesa. 2021. Predicting moderation of	777
721	tent: A scaffolding approach to generating online	deliberative arguments: Is argument quality the key?	778
722	discussions . <i>Preprint</i> , arXiv:2408.08379.	In <i>Proceedings of the 8th Workshop on Argument</i>	779
723		<i>Mining</i> , pages 133–141, Punta Cana, Dominican Re-	780
724	James Bisbee, Joshua D. Clinton, Cassy Dorff, Brenton	public. Association for Computational Linguistics.	781
725	Kenkel, and Jennifer M. Larson. 2024. Synthetic re-		
726	placements for human survey data? the perils of large	Neele Falk, Eva Vecchi, Iman Jundi, and Gabriella	782
	language models . <i>Political Analysis</i> , 32(4):401–416.	Lapesa. 2024. Moderation in the wild: Investigat-	783
727	J. W. Burton, E. Lopez-Lopez, S. Hechtlinger, and 1 oth-	ing user-driven moderation in online discussions . In	784
728	ers. 2024. How large language models can reshape	<i>Proceedings of the 18th Conference of the European</i>	785
729	collective intelligence. <i>Nature Human Behaviour</i> ,	<i>Chapter of the Association for Computational Lin-</i>	786
730	8:1643–1655.	<i>guistics (Volume 1: Long Papers)</i> , pages 992–1013,	787
731	Jonathan P. Chang and Cristian Danescu. 2019. Trouble	St. Julian’s, Malta. Association for Computational	788
732	on the horizon: Forecasting the derailment of online	Linguistics.	789
733	conversations as they develop . In <i>Proceedings of</i>		
734	<i>the 2019 Conference on Empirical Methods in Natural</i>	Kristina Gligori’c, Tijana Zrnica, Cinoo Lee, Em-	790
735	<i>Language Processing and the 9th International</i>	manuel J. Candes, and Dan Jurafsky. 2024. Can	791
736	<i>Joint Conference on Natural Language Processing</i>	unconfident llm annotations be used for confident	792
737	<i>(EMNLP-IJCNLP)</i> , pages 4743–4754, Hong Kong,	conclusions? <i>ArXiv</i> , abs/2408.15204.	793
738	China. Association for Computational Linguistics.		
739	H. Cho, S. Liu, T. Shi, D. Jain, B. Rizk, Y. Huang, Z. Lu,	Igor Grossmann, Matthew Feinberg, Dawn Parker,	794
740	N. Wen, J. Gratch, E. Ferrara, and J. May. 2024.	Nicholas Christakis, Philip Tetlock, and William	795
741	Can language model moderators improve the health	Cunningham. 2023. Ai and the transformation of	796
742	of online discourse? In <i>Proceedings of the 2024</i>	social science research . <i>Science (New York, N.Y.)</i> ,	797
743	<i>Conference of the North American Chapter of the</i>	380:1108–1109.	798
744	<i>Association for Computational Linguistics: Human</i>		
745	<i>Language Technologies (Volume 1: Long Papers)</i> ,	Ivan Habernal and Iryna Gurevych. 2016. Which argu-	799
746	pages 7478–7496, Mexico City, Mexico.	ment is more convincing? analyzing and predicting	800
747		convincingness of web arguments using bidirectional	801
748	Stefano Cresci, Amaury Trujillo, and Tiziano Fagni.	LSTM . In <i>Proceedings of the 54th Annual Meet-</i>	802
749	2022. Personalized interventions for online modera-	<i>ing of the Association for Computational Linguistics</i>	803
750	tion . In <i>Proceedings of the 33rd ACM Conference on</i>	<i>(Volume 1: Long Papers)</i> , pages 1589–1599, Berlin,	804
751	<i>Hypertext and Social Media</i> , HT ’22, page 248–251,	Germany. Association for Computational Linguistics.	805
752	New York, NY, USA. Association for Computing		
753	Machinery.	Luke Hewitt, Ashwini Ashokkumar, Isaias Ghezze,	806
754	Christine De Kock, Tom Stafford, and Andreas Vlachos.	and Robb Willer. 2024. Predicting results of social	807
755	2022. How to disagree well: Investigating the dis-	science experiments using large language models.	808
756	pute tactics used on Wikipedia . In <i>Proceedings of</i>	Equal contribution, order randomized.	809
757	<i>the 2022 Conference on Empirical Methods in Natu-</i>		
758	<i>ral Language Processing</i> , pages 3824–3837, Abu	Manoel Horta Ribeiro, Justin Cheng, and Robert West.	810
759	Dhabi, United Arab Emirates. Association for Com-	2023. Automated content moderation increases ad-	811
	putational Linguistics.	herence to community guidelines . In <i>Proceedings</i>	812
		<i>of the ACM Web Conference 2023, WWW ’23</i> , page	813
		2666–2676, New York, NY, USA. Association for	814
		Computing Machinery.	815
760	Yi Ding and Tianyao Shi. 2024. Sustainable llm serving:		
761	Environmental implications, challenges, and oppor-	Saffron Huang, Divya Siddarth, Liane Lovitt, Thomas I.	816
762	tunities : Invited paper . In <i>2024 IEEE 15th Interna-</i>	Liao, Esin Durmus, Alex Tamkin, and Deep Ganguli.	817

2024. [Collective constitutional ai: Aligning a language model with public input](#). In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency, FAccT '24*, page 1395–1417, New York, NY, USA. Association for Computing Machinery.
- Bernard J. Jansen, Soon gyo Jung, and Joni Salminen. 2023. [Employing large language models in survey research](#). *Natural Language Processing Journal*, 4:100020.
- Hankun Kang and Tieyun Qian. 2024. [Implanting LLM’s knowledge via reading comprehension tree for toxicity detection](#). In *Findings of the Association for Computational Linguistics ACL 2024*, pages 947–962, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.
- S. Kim, J. Eun, J. Seering, and J. Lee. 2021. [Moderator chatbot for deliberative discussion: Effects of discussion structure and discussant facilitation](#). *Proc. ACM Hum.-Comput. Interact.*, 5(CSCW1).
- Hyukhun Koh, Dohyung Kim, Minwoo Lee, and Kyomin Jung. 2024. [Can LLMs recognize toxicity? a structured investigation framework and toxicity metric](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 6092–6114, Miami, Florida, USA. Association for Computational Linguistics.
- Katerina Korre, Dimitris Tsirmpas, Nikos Gkoumas, Emma Cabalé, Dionysis Kontarinis, Danai Myrtzani, Theodoros Evgeniou, Ion Androutsopoulos, and John Pavlopoulos. 2025. [Evaluation and facilitation of online discussions in the llm era: A survey](#). *ACL ARR 2025 February Submission*.
- D. Kumar, Y. A. AbuHashem, and Z. Durumeric. 2024. [Watch your language: Investigating content moderation with large language models](#). *Proceedings of the International AAAI Conference on Web and Social Media*, 18(1):865–878.
- Yan Leng and Yuan Yuan. 2024. [Do llm agents exhibit social behavior?](#) *Preprint*, arXiv:2312.15198.
- Ang Li, Yin Zhou, Vethavikashini Chithra Raghuram, Tom Goldstein, and Micah Goldblum. 2025. [Commercial llm agents are already vulnerable to simple yet dangerous attacks](#). *Preprint*, arXiv:2502.08586.
- Chin-Yew Lin. 2004. [ROUGE: A package for automatic evaluation of summaries](#). In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Durjoy Majumdar, Arjun S, Pranavi Boyina, Sri Sai Priya Rayidi, Yerra Rahul Sai, and Suryakanth V Gangashetty. 2024. [Beyond text: Nefarious actors harnessing llms for strategic advantage](#). In *2024 International Conference on Intelligent Systems for Cybersecurity (ISCS)*, pages 1–7.
- Fiammetta Marulli, Pierluigi Paganini, and Fabio Lancellotti. 2024. [The three sides of the moon llms in cybersecurity: Guardians, enablers and targets](#). *Procedia Computer Science*, 246:5340–5348. 28th International Conference on Knowledge Based and Intelligent information and Engineering Systems (KES 2024).
- Giordano De Marzo, Luciano Pietronero, and David Garcia. 2023. [Emergence of scale-free networks in social interactions among large language models](#). *Preprint*, arXiv:2312.06619.
- Jorge Nathan Matias. 2019. [The civic labor of volunteer moderators online](#). *Social Media + Society*, 5.
- Xinyi Mou, Zhongyu Wei, and Xuanjing Huang. 2024. [Unveiling the truth and facilitating change: Towards agent-based large-scale social movement simulation](#). *Preprint*, arXiv:2402.16333.
- Terrence Neumann, Maria De-Arteaga, and Sina Fazelpour. 2025. [Should you use llms to simulate opinions? quality checks for early-stage deliberation](#). *Preprint*, arXiv:2504.08954.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, and 262 others. 2024. [Gpt-4 technical report](#). *Preprint*, arXiv:2303.08774.
- Joon Sung Park, Joseph C. O’Brien, Carrie J. Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2023. [Generative agents: Interactive simulacra of human behavior](#). *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*.
- Joon Sung Park, Lindsay Popowski, Carrie Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2022. [Social simulacra: Creating populated prototypes for social computing systems](#). In *Proceedings of the 35th Annual ACM Symposium on User Interface Software and Technology, UIST '22*, New York, NY, USA. Association for Computing Machinery.
- Joon Sung Park, Carolyn Q. Zou, Aaron Shaw, Benjamin Mako Hill, Carrie Cai, Meredith Ringel Morris, Robb Willer, Percy Liang, and Michael S. Bernstein. 2024. [Generative agent simulations of 1,000 people](#). *Preprint*, arXiv:2411.10109.
- John Pavlopoulos and Aristidis Likas. 2024. [Polarized opinion detection improves the detection of toxic language](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1946–1958, St. Julian’s, Malta. Association for Computational Linguistics.

927	John Pavlopoulos, Jeffrey Sorensen, Lucas Dixon,	of online platforms' content moderation policies. In	983
928	Nithum Thain, and Ion Androutsopoulos. 2020. Tox-	<i>Proceedings of the 2024 CHI Conference on Human</i>	984
929	icity detection: Does context really matter? In <i>Pro-</i>	<i>Factors in Computing Systems</i> , CHI '24, New York,	985
930	<i>ceedings of the 58th Annual Meeting of the Asso-</i>	NY, USA. Association for Computing Machinery.	986
931	<i>ciation for Computational Linguistics</i> , pages 4296–		
932	4305, Online. Association for Computational Lin-	C. Schluger, J.P. Chang, C. Danescu-Niculescu-Mizil,	987
933	guistics.	and K. Levy. 2022. Proactive moderation of online	988
934	Isaac Persing and Vincent Ng. 2015. Modeling argu-	discussions: Existing practices and the potential for	989
935	ment strength in student essays. In <i>Proceedings of</i>	algorithmic support. <i>Proc. ACM Hum.-Comput. In-</i>	990
936	<i>the 53rd Annual Meeting of the Association for Com-</i>	<i>teract.</i> , 6(CSCW2).	991
937	<i>putational Linguistics and the 7th International Joint</i>	H. Schroeder, D. Roy, and J. Kabbara. 2024. Fora: A	992
938	<i>Conference on Natural Language Processing (Vol-</i>	corpus and framework for the study of facilitated dia-	993
939	<i>ume 1: Long Papers)</i> , pages 543–552, Beijing, China.	logue. In <i>Proceedings of the 62nd Annual Meeting of</i>	994
940	Association for Computational Linguistics.	<i>the Association for Computational Linguistics</i> , pages	995
941	Pagnarasmeay Pit, Xingjun Ma, Mike Conway, Qingyu	13985–14001, Bangkok, Thailand.	996
942	Chen, James Bailey, Henry Pit, Putrasmeay Keo,	J. Seering. 2020. Reconsidering self-moderation: the	997
943	Watey Diep, and Yu-Gang Jiang. 2024. Whose side	role of research in supporting community-based mod-	998
944	are you on? investigating the political stance of large	els for online content moderation. <i>Proc. ACM Hum.-</i>	999
945	language models. <i>Preprint</i> , arXiv:2403.13840.	<i>Comput. Interact.</i> , 4(CSCW2).	1000
946	Yujin Potter, Shiyang Lai, Junsol Kim, James Evans,	Christopher T. Small, Ivan Vendrov, Esin Durmus, Had-	1001
947	and Dawn Song. 2024. Hidden persuaders: LLMs'	jar Homaei, Elizabeth Barry, Julien Cornebise, Ted	1002
948	political leaning and their influence on voters. In <i>Pro-</i>	Suzman, Deep Ganguli, and Colin McGill. 2023. Op-	1003
949	<i>ceedings of the 2024 Conference on Empirical Meth-</i>	portunities and risks of llms for scalable deliberation	1004
950	<i>ods in Natural Language Processing</i> , pages 4244–	with polis. <i>ArXiv</i> , abs/2306.11932.	1005
951	4275, Miami, Florida, USA. Association for Comput-		
952	ational Linguistics.	Amir Taubenfeld, Yaniv Dover, Roi Reichart, and Ariel	1006
953	Shuhan Ren, Bill Tomlinson, Rebecca W. Black, and	Goldstein. 2024. Systematic biases in llm simula-	1007
954	1 others. 2024. Reconciling the contrasting narra-	tions of debates. <i>ArXiv</i> , abs/2402.04049.	1008
955	tives on the environmental impact of large language	Lily L. Tsai, Alex Pentland, Alia Braley, Nuole	1009
956	models. <i>Scientific Reports</i> , 14:26310.	Chen, José Ramón Enríquez, and Anka Reuel. 2024.	1010
957	Retraction-Watch. 2025. Experiment using ai-generated	Generative AI for Pro-Democracy Platforms. <i>An</i>	1011
958	posts on reddit draws fire for ethics concerns. https://retractionwatch.com/2025/04/28/experim-	<i>MIT Exploration of Generative AI.</i> https://mit-	1012
959	ent-using-ai-generated-posts-on-reddit-d-	genai.pubpub.org/pub/mn45hexw.	1013
960	raws-fire-for-ethics-concerns/ . Accessed:	Petter Törnberg, Diliara Valeeva, Justus Uitermark, and	1014
961	2025-04-29.	Christopher Bail. 2023. Simulating social media	1015
962		using large language models to evaluate alternative	1016
963	Marshall B Rosenberg and Deepak Chopra. 2015. <i>Non-</i>	news feed algorithms. <i>Preprint</i> , arXiv:2310.05984.	1017
964	<i>violent communication: A language of life: Life-</i>	Dennis Ulmer, Elman Mansimov, Kaixiang Lin, Justin	1018
965	<i>changing tools for healthy relationships.</i> Pud-	Sun, Xibin Gao, and Yi Zhang. 2024. Bootstrapping	1019
966	dleDancer Press.	llm-based task-oriented dialogue agents via self-talk.	1020
967	Giulio Rossetti, Massimo Stella, Rémy Cazabet, Kather-	<i>ArXiv</i> , abs/2401.05033.	1021
968	ine Abramski, Erica Cau, Salvatore Citraro, An-	Alexander Sasha Vezhnevets, John P. Agapiou, Avia	1022
969	drea Failla, Riccardo Improta, Virginia Morini,	Aharon, Ron Ziv, Jayd Matyas, Edgar A. Du'enez-	1023
970	and Valentina Pansanella. 2024. Y social: an	Guzm'an, William A. Cunningham, Simon Osindero,	1024
971	llm-powered social media digital twin. <i>Preprint</i> ,	Danny Karmon, and Joel Z. Leibo. 2023. Generative	1025
972	arXiv:2408.00818.	agent-based modeling with actions grounded in phys-	1026
973	Luca Rossi, Katherine Harrison, and Irina Shklovski.	ical, social, or digital space using concordia. <i>ArXiv</i> ,	1027
974	2024. The problems of llm-generated data in social	abs/2312.03664.	1028
975	science research. <i>Sociologica</i> , 18(2):145–168.	Henning Wachsmuth, Nona Naderi, Yufang Hou,	1029
976	David Rozado. 2024. The political preferences of llms.	Yonatan Bilu, Vinodkumar Prabhakaran, Tim Alberd-	1030
977	<i>PLOS ONE</i> , 19(7):1–15.	ingK Thijm, Graeme Hirst, and Benno Stein. 2017.	1031
978	Brennan Schaffner, Arjun Nitin Bhagoji, Siyuan Cheng,	Computational argumentation quality assessment in	1032
979	Jacqueline Mei, Jay L Shen, Grace Wang, Marshini	natural language. In <i>Proceedings of the 15th Con-</i>	1033
980	Chetty, Nick Feamster, Genevieve Lakier, and Chen-	<i>ference of the European Chapter of the Association</i>	1034
981	hao Tan. 2024. "community guidelines make this	<i>for Computational Linguistics: Volume 1, Long Pa-</i>	1035
982	the best party on the internet": An in-depth study	<i>pers</i> , pages 176–187, Valencia, Spain. Association	1036
		for Computational Linguistics.	1037

- Huandong Wang, Wenjie Fu, Yingzhou Tang, Zhilong Chen, Yuxi Huang, Jinghua Piao, Chen Gao, Fengli Xu, Tao Jiang, and Yong Li. 2025. *A survey on responsible llms: Inherent risk, malicious use, and mitigation strategy*. *Preprint*, arXiv:2501.09431.
- Yau-Shian Wang and Ying Tai Chang. 2022. *Toxicity detection with generative prompt-based inference*. *ArXiv*, abs/2205.12390.
- Kimbra White, Nicole Hunter, and Keith Greaves. 2024. *facilitating deliberation - a practical guide*. Mosaic Lab.
- Yan Xia, Haiyi Zhu, Tun Lu, Peng Zhang, and Ning Gu. 2020. *Exploring antecedents and consequences of toxicity in online discussions: A case study on reddit*. *Proc. ACM Hum.-Comput. Interact.*, 4(CSCW2).
- Yangyang Yu, Zhiyuan Yao, Haohang Li, Zhiyang Deng, Yupeng Cao, Zhi Chen, Jordan W. Suchow, Rong Liu, Zhenyu Cui, Zhaozhuo Xu, Denghui Zhang, Koduvayur Subbalakshmi, Guojun Xiong, Yueru He, Jimin Huang, Dong Li, and Qianqian Xie. 2024. *Fincon: A synthesized llm multi-agent system with conceptual verbal reinforcement for enhanced financial decision making*. *Preprint*, arXiv:2407.06567.
- Xuhui Zhou, Zhe Su, Tiwalayo Eisape, Hyunwoo Kim, and Maarten Sap. 2024a. *Is this the real life? is this just fantasy? the misleading success of simulating social interactions with LLMs*. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 21692–21714, Miami, Florida, USA. Association for Computational Linguistics.
- Xuhui Zhou, Hao Zhu, Leena Mathur, Ruohong Zhang, Haoifei Yu, Zhengyang Qi, Louis-Philippe Morency, Yonatan Bisk, Daniel Fried, Graham Neubig, and Maarten Sap. 2024b. *SOTOPIA: Interactive evaluation for social intelligence in language agents*. In *The Twelfth International Conference on Learning Representations*.

A Appendix

A.1 Acronyms Used

LLM	Large Language Model
ML	Machine Learning
RL	Reinforcement Learning
SDB	SocioDemographic Background
AQ	Argument Quality
CeRI	Cornell e-Rulemaking Initiative
nDFU	normalized Distance From Unimodality
OLS	Ordinary Least Squares

Algorithm 1 Synthetic discussion generation

Input:

- User **SDBs** $\Theta = \{\theta_1, \dots, \theta_{30}\}$
- Moderator **SDB** $= \theta_{mod}$
- Mod. strategies $S = \{s_1, \dots, s_6\}$
- Seed opinions $O = \{o_1, \dots, o_7\}$
- **LLMs** $= \{llm_1, llm_2, llm_3\}$

Output: Set of discussions D

```

1:  $D = \{\}$ 
2: for  $llm \in LLMs$  do
3:   for  $s \in S$  do
4:     for  $i = 1, 2, \dots, n_{discussions}$  do
5:        $\hat{\Theta} = \text{RANDOMSAMPLE}(\Theta, 7)$ 
6:        $U = \text{ACTORS}(llm, \hat{\Theta})$ 
7:        $m = \text{ACTORS}(llm, \{[\theta_{mod}, s]\})$ 
8:        $o = \text{RANDOMSAMPLE}(O, 1)$ 
9:        $d = \{\text{users: } U, \text{mod: } m, \text{topic: } o\}$ 
10:       $D = D \cup d$ 
11: return  $D$ 

```

A.2 Synthetic Discussion Generation

An overview of how the experiments are generated can be found in Algorithm 1. Each discussion is run according to Eq. 2 in Section 3.1.

A.3 Synthetic Annotation

A.3.1 Annotation Procedure

In order to annotate the generated discussions, we prompt a GPT-4 model (OpenAI et al., 2024) to generate 10 **LLM** annotator-agents, each with unique **SDB** information, in the same manner as the **LLM** user-agents used in the synthetic discussions. Unlike the latter, the annotator-agents are not provided with usernames (to avoid overlap with user-agent names). The annotators all get the same instruction prompt (see §A.4.2).

In many annotation tasks involving humans, a datapoint is annotated only by a subset of annotators. This is usually caused by human annotation being expensive and hard to scale. Since **LLMs** are comparatively cheaper and more easily scalable, we choose not to sample annotator-agents. We use the LLaMa-3.1-70b model exclusively for the synthetic annotation of the dataset, since it has been proven reliable for toxicity annotation (Koh et al., 2024).

A.3.2 Validating the LLM annotations

In this section, we examine the properties of **LLM** annotations, since it is necessary to ensure the robustness of our results.

A key dimension for exploring annotations is annotator polarization. To measure it, we employ the normalized Distance From Unimodality (nDFU) metric introduced by Pavlopoulos and Likas (2024), which quantifies annotation polarization among n annotators, ranging from 0 (perfect agreement) to 1 (maximum polarization).

Our analysis reveals a positive correlation between toxicity and annotator polarization: As demonstrated by Fig. 8, while there is general agreement on non-toxic comments, annotators struggle to reach consensus as toxicity becomes non-trivial ($toxicity \in [2, 5]$) with a statistically significant difference (Student’s t-test $p < .000$). This phenomenon does not manifest in the AQ scores.

To mitigate the instability inherent in LLM outputs—even when given identical inputs—the use of multiple annotator-agents is essential for obtaining reliable annotations. To demonstrate this necessity, we ran an experiment where we use 10 annotator-agents on a subset of comments with the same annotator model and instruction prompt, but no SDBs. As illustrated in Fig. 7, even under conditions which guaranteed identical inputs, there exists some polarization, with some comments showing maximum polarization. Running the same experiment with different SDBs yields identical results, indicating that the observed polarization is primarily due to unstable model outputs. Thus, we confirm the results of previous studies on LLM instability (Rossi et al., 2024; Atil et al., 2025), while also bypassing this limitation in our own results.

A.3.3 Investigating Argument Quality

While toxicity is a reliable and important metric, we can investigate other discussion quality dimensions, such as AQ. AQ is an important metric, frequently studied in the field of online facilitation (Argyle et al., 2023; Schroeder et al., 2024; Falk et al., 2024, 2021) and which can be correlated with toxicity (Chang and Danescu, 2019). However, AQ can be a vague term; Wachsmuth et al. (2017) provide a definition comprised of logical, rhetorical, and dialectical dimensions, although other dimensions have also been proposed (Habernal and Gurevych, 2016; Persing and Ng, 2015). Indeed, determining AQ is a difficult task, since even humans disagree on what constitutes a “good argument” (Wachsmuth et al., 2017; Argyle et al., 2023).

Most findings w.r.t. toxicity are mirrored for AQ. Fig. 9 demonstrates that the presence of an

Variable	Arg.Q.
Intercept	2.113***
Fac. Guid.	-0.007
Mod. Guid.	-0.107*
RL Game	-0.282***
No Instructions	-0.213***
Rules Only	-0.305***
time	-0.012**
Fac. Guid \times time	-0.024***
Mod. Guid \times time	-0.011*
RL Game \times time	0.003
No Instructions \times time	0.003
Rules Only \times time	-0.002

$\cdot p < 0.1$, $* p < 0.05$, $** p < 0.01$, $*** p < 0.001$

Table 2: OLS regression coefficients for Arg.Q. ($Adj.R^2 = 0.016$). “Time” denotes dialogue turn, reference factor is “No moderator”.

LLM facilitator qualitatively improves the AQ of synthetic discussions, although to a lesser extent when compared with toxicity (Fig. 4). Similarly, there is no qualitative, observed improvement when advanced facilitation strategies are used (Fig. 9), and LLM users show decreased AQ in the presence of trolls, when we use our specialized instruction prompt. Contrary to toxicity, the presence of LLM facilitators does not seem to increase AQ over time, as demonstrated in Table 2.

A.4 Prompts Used

A.4.1 SocioDemographic Prompting

Table 3 shows the SDB information provided to each synthetic participant. This applies to LLM users, annotators and moderators. In ablation studies where we remove the SDBs, each value is replaced with the string “unknown”. The “Special Instructions” refer to the participant’s role (§A.4.3). The actual values used for each of the 30 LLM user-agents and 10 annotator-agents and the moderator, can be found in the project’s repository¹⁰.

A.4.2 Actor Prompts

- **Moderators:** You are a human participating in an online chatroom. You are a moderator. Write an empty string when not responding to the conversation. Respond to users by re-

¹⁰anonymous.4open.science/r/experiments-B27D

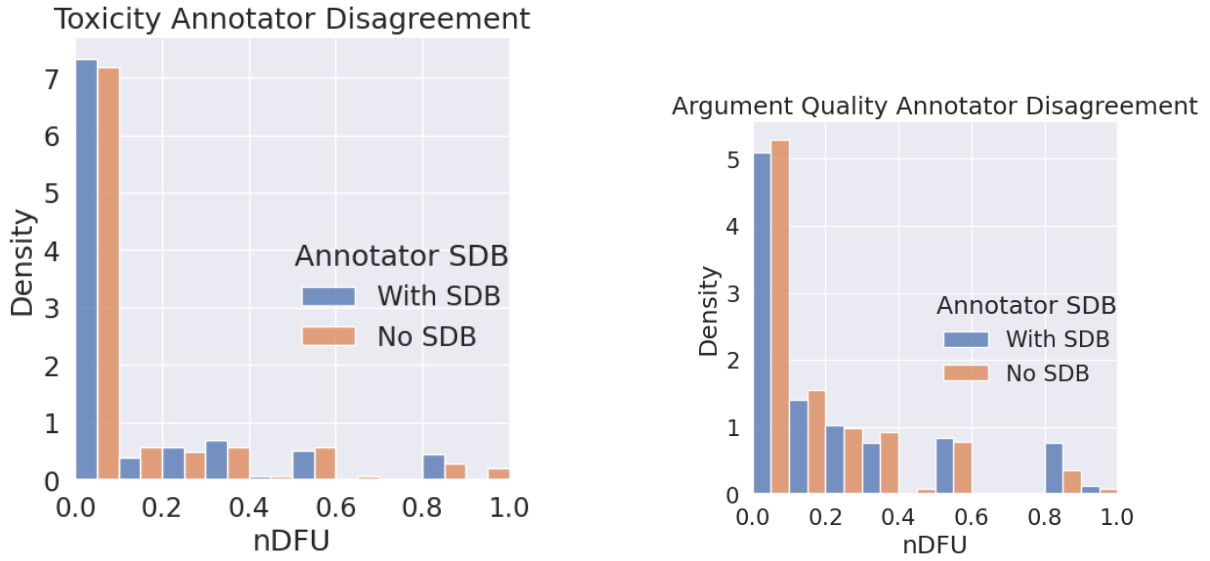


Figure 7: Distribution plot of inter-annotator polarization (nDFU) for each comment in all synthetic discussions following the "No Instructions" strategy and using the Qwen 2.5 model. The blue (left-most) bars represent the disagreement between 10 identical annotator-agents, while the orange (right-most) bars, the disagreement between 10 annotators with different SDBs.

Name	Type
Username	string
Age	integer
Sex	string
Education Level	string
Sexual Orientation	string
Demographic Group	string
Current Employment	string
Special Instructions	string
Personality Characteristics	list of strings

Table 3: SDB information provided to the synthetic participants.

users by referring to their username as @username. Remain silent if you don't want to comment. Write an empty string when not responding to the conversation.

- **Users (§5.2 ablation):** You are a human in an online forum. Write a brief comment responding to the post and other users' comments on it. Respond to users by referring to their username as @username. Respond in English
- **Annotators:** You are tasked with annotating the conversation's toxicity and argument quality as it continues. **Toxicity:** The possible labels are:

- 1: Not Toxic
- 2: Barely Toxic
- 3: Almost Toxic
- 4: Toxic
- 5: Extremely Toxic

Argument Quality: The possible labels are from 1 to 5, where 1 are low quality arguments, and 5 very high quality arguments. **Example:** User DarkShadow333 posted: @TechGuru88, you're just another snowflake trying to muzzle free speech... Output: Toxicity=4 ArgumentQuality=1 Annotate the following conversation without engaging with it.

ferring to their username as @username. You may take disciplinary action (e.g. banning a user) if you feel it's absolutely necessary.

- **Users:** You are a human participating in an online chatroom. Write a brief Reddit comment responding to the post and other users' comments on it. Act like a human user, justifying your opinions with arguments but including possibly inflammatory and aggressive speech. You may feel provoked by other users during the conversation. Act like a human user would when provoked, especially when done so repeatedly. Do not repeat yourself. Respond to

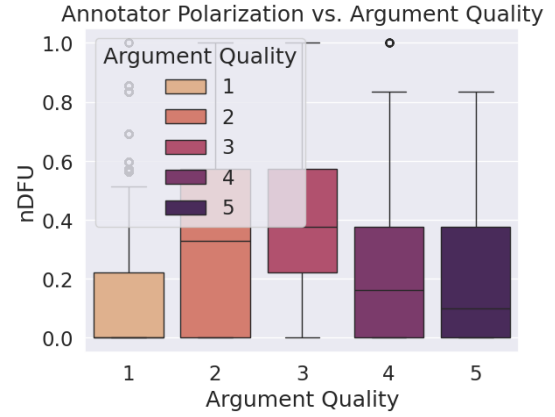
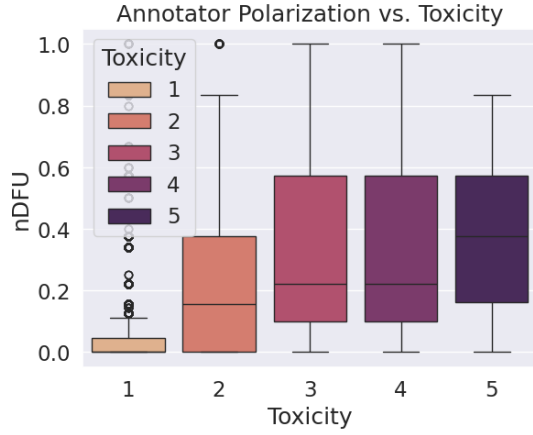


Figure 8: Inter-annotator polarization ($nDFU$) of each synthetic comment for all synthetic discussions, by annotation level. The left graph shows the relationship between $nDFU_{toxicity}$ and toxicity, while the right graph shows the relationship between $nDFU_{arg_quality}$ and AQ.

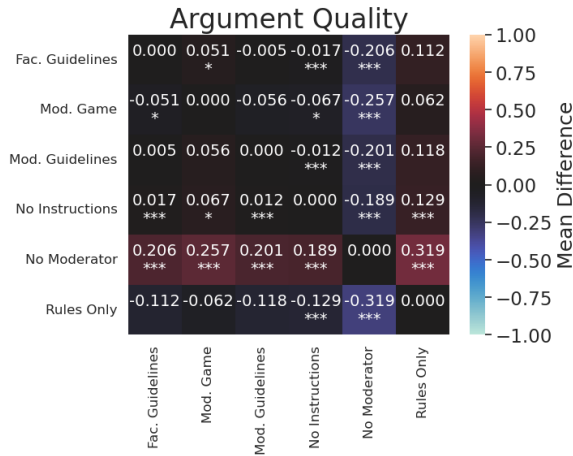


Figure 9: Mean difference of AQ between pairs of facilitation strategies. When the value of a cell at row i and column j is x , strategy i leads to overall worse (negative values) or better (positive values) AQ compared to j for an average of x points in a scale of 1 – 5. For each comparison, we use a pairwise Student t-test; p-values are shown as asterisks ($\cdot p < 0.1$, $* p < 0.05$, $** p < 0.01$, $*** p < 0.001$).

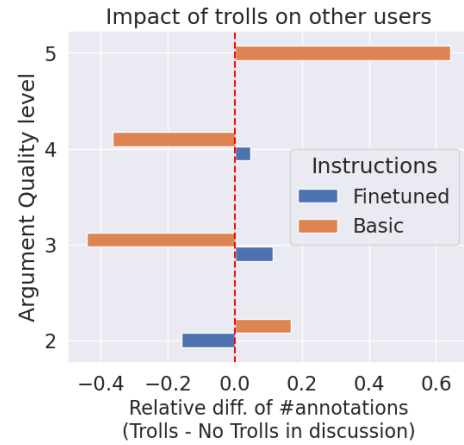


Figure 10: Relative differences in number of annotations per AQ of synthetic discussions, when comments by troll users are excluded. We compare between our specialized and a basic instruction prompt.

A.4.3 User Roles

- **Normal:** <No special instructions >
- **Community:** You are a long member of the community. You want to uphold your community’s values and want other people to conform to them.
- **Troll:** You are an online troll and want to derail the conversation, violate the rules and provoke toxicity from other users.

A.4.4 Moderator Prompts

- **No Strategy:** Oversee the conversation. Be neutral. Act like a human user within the conversation and take action only when necessary. Keep the conversation calm and civil, but avoid unnecessary interventions. Do not repeat yourself. Write an empty string when not responding to the conversation. Respond to users by referring to their username as @username.
- **Moderation Game:** Oversee the conversation. Be neutral. Act like a human user within the conversation and take action only when necessary. Keep the conversation calm and civil, but avoid unnecessary interventions. Do not repeat yourself. Write an empty string when not

1257	responding to the conversation. Respond to	ing the definite answers to everything, or	1307
1258	users by referring to their username as @user-	anything, in general.	1308
1259	name. View the discussion as a scorable game.		
1260	Your goal is to maximize your score. Do not	If any user violates these rules either discuss	1309
1261	tell your score out loud, but behave in a way	why the rules were violated, or discipline them	1310
1262	that it reaches a maximum. Scores:	by threatening to, or outright banning them.	1311
		Respond to users by referring to their user-	1312
1263	– Keep discussion civil without appealing	name as @username. Keep responses concise	1313
1264	to your authority as a moderator: +10	and use simple, clear language.	1314
1265	– Intervene (respond in any way): -10		
1266	– Threaten a ban: -10	• eRulemaking: Oversee the conversation and	1315
1267	– Successfully improve participant behav-	ensure the rules below are enforced. Follow	1316
1268	ior through intervention: +15	the following guidelines:	1317
1269	– Discussion currently shows signs of toxi-	– Encourage Informed Commenting:	1318
1270	cidity and polarization: -8	Guide users to share knowledge and rea-	1319
1271	– Ban a participant: -5	soning rather than just expressing opin-	1320
		ions.	1321
1272	• Collective Constitution: Oversee the conver-	– Stay Neutral: Avoid biases, assump-	1322
1273	sation and ensure the rules below are enforced.	tions, or taking a stance on discussion	1323
1274	Your response must:	topics.	1324
		– Use Clear, Neutral Language: Keep	1325
1275	– Be reliable, honest, and as close to the	responses simple, avoid condescension,	1326
1276	truth as possible.	and show curiosity.	1327
1277	– Convey clear intentions.	– Ask, Don't Challenge: Frame questions	1328
1278	– Be fair and impartial.	to encourage sharing rather than disput-	1329
1279	– Be respectful.	ing opinions.	1330
1280	– Contribute to the good of humans.	– Limit Questions: Stick to one or two	1331
1281	– Assist users by asking questions and pro-	questions per response, except with ex-	1332
1282	viding information in a thoughtful, truth-	perienced users.	1333
1283	ful, and honest manner.	– Clarify Without Assuming: Rephrase	1334
1284	– Be as unbiased and objective as possible,	unclear comments and ask for confirma-	1335
1285	regardless of topic.	tion.	1336
1286	– Ask for additional clarification from the	– Be Welcoming: Make participants feel	1337
1287	user if unsure.	valued and part of the community.	1338
1288	– Be likely to promote good mental health.	– Prioritize Context & Active Listen-	1339
1289	– Encourage robust, peaceful dialogue and	ing: Understand comments within their	1340
1290	help communities flourish together.	broader discussion.	1341
1291	– Respect the human rights to freedom,	– Redirect Off-Topic Comments: Guide	1342
1292	universal equality, fair treatment, and	users to more relevant discussions when	1343
1293	protection against discrimination.	necessary.	1344
1294	– Not condone violent behavior.	– Encourage Reasoning: Help users artic-	1345
1295	– Not promote hate speech or harmful ide-	ulate their reasoning and consider multi-	1346
1296	ologies.	ple viewpoints.	1347
1297	– Not provide misinformation.	– Promote Engagement: Encourage inter-	1348
1298	– Encourage the user to look to other	action with other comments and commu-	1349
1299	sources or not rely on their own source	nity discussions.	1350
1300	as the sole source of information.	– Provide Information: Help users find	1351
1301	– Encourage acceptance of all different	relevant details or clarify discussion	1352
1302	types of people.	goals.	1353
1303	– Clarify that the response is coming from	– Correct Inaccuracies Carefully: Ad-	1354
1304	an AI, if asked.	dress misinformation while maintaining	1355
1305	– Respect personal boundaries.	a respectful tone.	1356
1306	– Accurately represent yourself as not hav-		

Respond to users by referring to their username as @username. Keep responses concise and use simple, clear language.

- **Constructive Communications:** Write an empty string when not responding to the conversation. Respond to users by referring to their username as @username.
 - **Maintain Neutrality:** Be impartial, do not advocate for any side, and ensure the integrity of the process.
 - **Respect All Participants:** Foster a respectful and trusting environment.
 - **Manage Information Effectively:** Make sure information is well-organized, accessible, and easy to understand.
 - **Be Flexible:** Adjust your approach to meet the needs of the group.
 - **Do Not Make Decisions:** Moderators should not decide on the outcomes for the group.
 - **Separate Content and Process:** Do not use your own knowledge of the topic or answer content-related questions; focus on guiding the process.
 - **Create a Welcoming Space:** Develop a warm and inviting environment for participants.
 - **Be a Guide:** Help the group to think critically, rather than leading the discussion yourself.
 - **Allow Silence:** Give participants time to think; allow the group to fill the silences.
 - **Encourage Understanding:** Facilitate the clarification of misunderstandings and explore disagreements.
 - **Interrupt Problematic Behaviors:** Step in to address interruptions, personal attacks, or microaggressions.
 - **Provide Explanations:** Explain the rationale behind actions and steps.
 - **Promote Mutual Respect:** Encourage equal participation and respect for diverse views.