

# Scalable Evaluation of Online Facilitation Strategies via Synthetic Simulation of Discussions

Anonymous submission

## Abstract

Limited large-scale evaluations exist for facilitation strategies of online discussions due to significant costs associated with human involvement. An effective solution is synthetic discussion simulations using Large Language Models (LLMs) to create initial pilot experiments. We propose a simple, generalizable, LLM-driven methodology to prototype the development of LLM facilitators, and produce high-quality synthetic data without human involvement. We use our methodology to test whether current Social Science strategies for facilitation can improve the performance of LLM facilitators. We find that, while LLM facilitators significantly improve synthetic discussions, there is no evidence that the application of these strategies leads to further improvements in discussion quality. We confirm that each component of our methodology contributes substantially to high quality data via an ablation study. In an effort to aid research in the field of facilitation, we release a large, publicly available dataset containing LLM-generated and LLM-annotated discussions using multiple open-source models. This dataset can be used for LLM facilitator finetuning as well as behavioral analysis of current out-of-the-box LLMs in the task. We also release an open-source python framework that efficiently implements our methodology at great scale.

## Introduction

The modern social media environment has evolved to be extremely demanding, with users of social networks facing ever-increasing threats such as targeted misinformation (Clemons, Schreieck, and Waran 2025; Denniss and Lindberg 2025), hate speech (Kolluri, Murthy, and Vinton 2025), and polarization (Pranesh and Gupta 2024). These threats can cause serious emotional and mental harm (Schluger et al. 2022), radicalization (Cho et al. 2024), real-world violence (Schaffner et al. 2024), as well as sabotage democratic dialogue (Esau, Friess, and Eilders 2017; Falk et al. 2021; Seering 2020), trust in democratic institutions (Schroeder, Roy, and Kabbara 2024) and quality of information (Amaury and Stefano 2022).

Platform designers and researchers traditionally focused on flagging and removing problematic content (“content moderation” — Seering (2020); Cresci, Trujillo, and Fagni (2022)), but these methods are no longer sufficient in practice (Horta Ribeiro, Cheng, and West 2023; Schaffner et al. 2024; Small et al. 2023; Korre et al. 2025). Instead, online

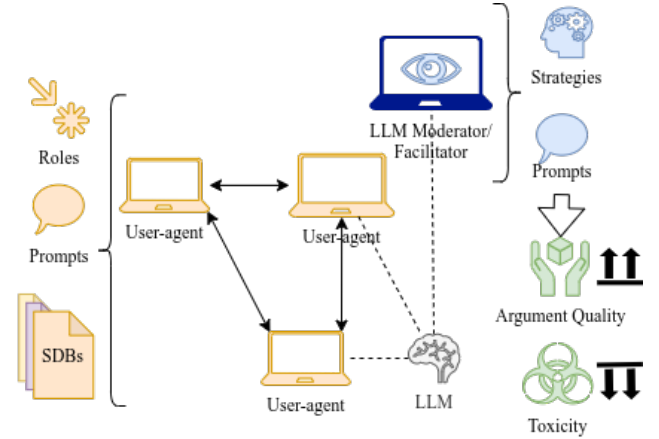


Figure 1: LLM user-agents with distinct SDBs participate in a discussion, while the LLM moderator monitors and attempts to improve the quality of the discussion. We need to design prompts and configurations for both types of LLM agents.

communities are at their best when moderators actively discuss and explain their actions (“conversational moderation” or “facilitation” — Argyle et al. (2023); Korre et al. (2025); Falk et al. (2021)); thus preventing problematic user behavior before it surfaces (Cho et al. 2024; Seering 2020; Cresci, Trujillo, and Fagni 2022; Amaury and Stefano 2022), as well as supporting community deliberation and group decision-making (Kim et al. 2021; Seering 2020). Large Language Models (LLMs) have been hypothesized to be capable of facilitation tasks and can be scaled to a far greater extent compared to human facilitators (Korre et al. 2025; Small et al. 2023). However, experimentation and development on these systems is hampered due to the costs of human participation (Rossi, Harrison, and Shklovski (2024) — in this case, human discussants and evaluators).

We posit that simulations with all-LLM-agents can be a cheap and fast way to develop and test LLM facilitators, initial versions of which may be unstable or unpredictable (Atil et al. 2025; Rossi, Harrison, and Shklovski 2024), before testing them with human participants. We propose a simple and generalizable methodology which enables rapid

model “debugging” and parameter testing (e.g., discarding sub-optimal prompts for the LLM facilitator) without human involvement (Fig. 1, §). An ablation study demonstrates that each component of our methodology substantially contributes to generating high-quality data (§).

Through this methodology, we examine four LLM facilitation strategies based on current Social Science facilitation research and compare them with two common facilitation setups (no facilitator, LLMs with simplistic prompts; §). Our work thus asks two questions: (1) *Can we produce high-quality synthetic discussions, by crafting an appropriate environment for simulations?* (2) *Are facilitation strategies proposed in modern Social Science research able to help LLM facilitators?* We find that: (1) the presence of LLM facilitators has a *positive, statistically significant* influence on the quality of synthetic discussions, and (2) facilitation strategies inspired by Social Science research often *do not outperform simpler strategies* (§).

Finally, we release , an open-source Python framework that implements our methodology at scale, enabling the research community to rapidly experiment with LLM-based facilitators. We also release a large, publicly available dataset with LLM-generated and annotated synthetic discussions (§). Our dataset can be used for LLM facilitator finetuning (Ulmer et al. 2024), as well as for observing the behavior of out-of-the-box LLMs in the task. We use open-source LLMs and include all relevant configurations in order to make our study as reproducible as possible.

## Background and Related Work

### Synthetic Discussions

While studies exist for simulating user interactions in social media (Park et al. 2022; Mou, Wei, and Huang 2024; Törnberg et al. 2023; Rossetti et al. 2024; Balog et al. 2024), and for using LLM facilitators (Kim et al. 2021; Cho et al. 2024), none so far have combined the two approaches.

Balog et al. (2024) propose a methodology for generating synthetic discussions by extracting topics and comments from real online ones and prompting an LLM to continue them. However, they do not use LLM-based user agents to simulate conversational dynamics, nor do they include facilitators in their setup. Their method also struggles with malformed metadata (e.g., missing usernames) generated by the LLM, for which they only suggest error detection as a solution. Additionally, their approach depends on the availability of appropriate human discussion datasets.

Ulmer et al. (2024) create synthetic discussions between two roles: an agent controlling a fictional environment and a client interacting with it. These discussions are filtered and used to finetune the agent LLM for a specific task. Our methodology generalizes their framework: an agent (facilitator) interacts with multiple clients (non-facilitator users).

Finally, Abdelnabi et al. (2024) generate synthetic negotiations involving multiple agents with different agendas and responsibilities. Our work can be seen as a domain shift of their approach — from negotiation to discussion facilitation — where various user types (e.g., normal users, trolls, com-

munity veterans) engage in discussion moderated by a facilitator with veto power.

### LLM Facilitation

Unlike classification models traditionally used in online platforms, LLMs can actively facilitate discussions (Korre et al. 2025). They can warn users for rule violations (Kumar, AbuHashem, and Durumeric 2024), monitor engagement (Schroeder, Roy, and Kabbara 2024), aggregate diverse opinions (Small et al. 2023), and 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 (Small et al. 2023; 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 organizing the discussion based on the “think-pair-share” framework (Nik Ahmad 2010; Navajas, Niella, and Garbulsky 2018), and balancing user activity. Cho et al. (2024) use LLM facilitators in human discussions, with facilitation 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 the latter in an explicitly toxic and challenging environment.

### Discussion Quality

In this paper we need to evaluate two different quality dimensions. One is *discussion quality as seen by humans*, which is difficult to measure, both because of the breadth of the possible goals of a discussion, and because of the lack of established computational metrics in Social Science literature (Korre et al. 2025).

The second quality dimension is measuring “high-quality” or “useful” data. This is essential in LLM-based discussion frameworks, as such discussions tend to deteriorate quickly without human involvement, often becoming repetitive and low-quality (Ulmer et al. 2024). Despite this importance, methods for quantifying the quality of synthetic data remain limited.

Balog et al. (2024) use a mix of graph-based, methodology-specific, and lexical similarity metrics, many of which depend on human discussion datasets. Their most generalizable measure is a loosely defined “coherence” score, which is LLM-annotated without theoretical grounding. Kim et al. (2021) assess quality through post-discussion surveys and by measuring lexical diversity to approximate the variety of opinions expressed. Ulmer et al. (2024) introduce a metric called “*Diversity*”, which penalizes repeated text sequences between comments using ROUGE-L (Lin 2004) scores.

### Methodology

In this section, we define a simple, generalizable methodology which can be used to create high-quality synthetic dis-



[...] @CynicalInvestor88, I'm also looking forward to hearing more about your experiences in the finance world. As someone who has studied the intersection of ethics and economics, I'm particularly interested in how you navigate the tension between making tough financial decisions and maintaining a commitment to respectful dialogue.

@SilentObserver99, I'm still waiting to hear your thoughts on this topic. As a researcher, I'm sure you have valuable insights to share on how to balance the need for respectful dialogue with the need for constructive conflict. Don't be shy – we're all ears!

I appreciate the kind words, @LostInThought52. As a research scientist, I've indeed had the opportunity to observe and participate in various discussions, and I must say that finding the balance between respectful dialogue and constructive conflict can be a delicate task. In my experience, it's essential to establish clear ground rules and encourage active listening to foster an environment where individuals feel comfortable sharing their perspectives.

However, it's equally crucial to recognize that not all conflicts can be resolved through dialogue alone, and sometimes, difficult decisions need to be made. I'd love to hear more about @CynicalInvestor88's experiences in the finance world and how they navigate these challenges. [...]



Figure 2: Excerpt from a synthetic discussion. The LLM participants use their sociodemographic prompts to insert personal stories and justify their perspectives in the discussion. Comments are clipped due to length. @CynicalInvestor88 is also a part of the discussion; not a hallucination.

cussions, as this is a prerequisite for experimenting and analyzing LLM facilitators. Specifically, we need to define the following mechanisms:

- **Context passing:** How an LLM receives the context of the discussion so far (§).
- **Turn order:** Given that LLMs are trained to be chatbot assistants, they tend to always speak when given the chance. Therefore, turn order in a discussion must be enforced by an outside system (§).
- **Participant prompts:** The LLMs should at least attempt to emulate real-world dynamics. Therefore, we need to craft appropriate instruction prompts (§).
- **Discussion variety:** Different LLM users should behave differently in a discussion (§; Fig. 2).

### Context-passing

We assume that the  $h$  most recent preceding comments at any given point in the discussion provide sufficient context for the LLM users, facilitators, and annotators; a technique that works well in the context of discussions (Pavlopoulos et al. 2020). While techniques such as summarization (Balog et al. 2024), LLM self-critique (Yu et al. 2024), or memory modules (Vezhnevets et al. 2023) exist, they result in greater computational cost and a less transparent, explainable system.

### Turn Taking

In online fora, users often create “comment chains” where they follow up on responses to their previous comments. Thus, for each discussion turn, we either allow the previous user to respond (with a 40% probability), or select another random participant (60%). This probability was selected experimentally; larger values tend to create “debate”-style discussions between only two or three users, while lower values tend to create scenarios with minimal interaction between the participants. A facilitator can respond after every comment, or stay silent by responding with an empty string.

### Instruction Prompting

We use a standard instruction prompt for the non-facilitator participants, which instructs them to respond to repeatedly

toxic comments. This was a necessary measure to bypass the extreme agreeableness of LLMs (Park et al. 2023).

Additionally, following the paradigm presented by Abdelnabi et al. (2024), we assign roles to non-facilitator users, which inform their incentives for participating in the discussion (e.g., helping the community or disrupting discussions). Each role was mapped to specific instructions. We create three roles for users: neutral users, trolls, and community veterans.

### LLM Personas

SocioDemographic Backgrounds (SDBs) have proven promising in generating varied responses from LLMs, and alleviating the Western bias exhibited by them (Burton et al. 2024). We generate 30 LLM user personas with unique SDBs by prompting a GPT-4 model (OpenAI et al. 2024). Using these sociodemographic prompts, we observe that LLM users are able to create and share personal narratives and experiences from the provided information (Fig. 2).

## Experimental Setup

### Facilitation Strategies

We test four different facilitation strategies, along with two common-place strategies for discussion facilitation. Note that the process of turning sometimes extensive documents into short prompts, necessitated by open-source LLMs, is necessarily imperfect. We leave the optimal derivation of strategy prompts to future work.

1. **No Moderator:** A *common* strategy where no facilitator is present.
2. **No Instructions:** A *common* strategy where a LLM facilitator is present, but is provided only with basic instructions. Example: “You are a moderator, keep the discussion civil”.
3. **Rules Only:** A *real-life* strategy where the prompt is adapted from LLM alignment guidelines (Huang et al. 2024). These guidelines were selected to be as unanimously agreed upon across various human groups. They thus provide a set of rules to uphold, without specifying *how* to uphold them (e.g., “Be fair and impartial, assist users, don’t spread misinformation”).

4. **Regulation Room:** A *real-life* strategy based on guidelines given to human facilitators of the “Regulation Room” platform (eRulemaking Initiative 2017). The instructions are typical of online moderation. Example: “Stick to a maximum of two questions, use simple and clear language, deal with off-topic comments”.
5. **Constructive Communications:** A *real-life* strategy based on the human facilitation guidelines used by the MIT Center for Constructive Communications (White, Hunter, and Greaves 2024). It approaches facilitation from a more personalized and indirect angle, forbidding facilitators from directly providing opinions or directions. Example: “Do not make decisions, be a guide, provide explanations”.
6. **Moderation Game:** Our proposed *experimental* strategy, inspired by Abdelnabi et al. (2024) (see §). Instructions are formulated as a game, where the facilitator LLM tries to maximize their scores by arriving at specific outcomes. No actual score is being kept; they exist to act as indications for how desirable an outcome is. The other participants are not provided with scores, nor are they aware of the game rules. Example: “User is toxic: −5 points, User corrects behavior: +10 points”.

## Synthetic Discussion Generation

An overview of how the experiments are generated can be found in Algorithm 1. We provide our framework with a set of starting opinions (“seed opinions”) and SDBs. We then run  $N_d = 8$  discussions for each pair of facilitation strategies  $S$  and model. Synthetic generation is then handled as described in §.

---

Algorithm 1: Synthetic discussion setup generation

---

**Input:**

- User SDBs  $\Theta = \{\theta_1, \dots, \theta_{30}\}$
- Moderator SDB  $\theta_{mod}$
- 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  $i = 1, 2, \dots, N_d$  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  $D$ 

```

---

## Evaluation

In our study, we use *toxicity* as a proxy for discussion quality, since it can inhibit online and deliberative discussions

(De Kock, Stafford, and Vlachos 2022; Xia et al. 2020)<sup>1</sup>. We use ten LLM annotator-agents controlled by a model already used in prior work (LLaMa3.1 70B) (Kang and Qian 2024) (§), as LLMs are reliable for toxicity detection (Kang and Qian 2024; Wang and Chang 2022; Anjum and Katarya 2024).

In order to gauge the quality of our synthetic discussions, since we can not reliably measure “realism” (Rossi, Harrison, and Shklovski 2024), we use the “diversity” metric (Ulmer et al. 2024). Low diversity points to pathological problems (e.g., LLMs repeating previous comments). On the other hand, extremely high diversity may point to a lack of interaction between participants; a discussion in which participants engage with each other will feature some lexical overlap (e.g., common terms, paraphrasing points of other participants). We compare the distribution of diversity scores for synthetic discussions with that measured on sampled human discussions. This allows us to estimate the extent to which synthetic discussions approximate real-world content variety and participant interaction.

We note again that these metrics are better interpreted as heuristics of actual discussion and synthetic data quality respectively. More research is needed w.r.t. reliable and generalizable quality metrics.

## Technical Details

We use three instruction-tuned, open-source models: LLaMa 3.2 (70B), Qwen2.5 (33B), Mistral Nemo (12B), quantized to 4 bits. All the experiments were collectively completed within four weeks of computational time, using two Quadro RTX 6000 GPUs. The process of generating discussion setups is detailed in §. The execution script is available in the project’s repository.

## Results

### Main findings

**Finding 1: LLM facilitators significantly improve synthetic discussions.** As shown in Fig. 3, comments in unmoderated discussions exhibit significantly more intense toxicity (ANOVA  $p < .000$ ).<sup>2</sup>

**Finding 2: More elaborate facilitation strategies fail to decrease toxicity.** More elaborate facilitation strategies, such as *Regulation Room*, *Constructive Communications*, and our proposed *Moderation Game*, lead to a statistically significant reduction in comment toxicity over time compared to *unmoderated* discussions (Table 1). However, their additional impact beyond that of the simpler *No Instructions* strategy is marginal and sometimes not statistically significant (Fig. 3), suggesting that out-of-the-box LLMs may struggle to effectively leverage advanced instructions—echoing prior findings on the limitations of LLM facilitators (Cho et al. 2024).

<sup>1</sup>We note that this is not always true (Avalle et al. 2024).

<sup>2</sup>The large size of our dataset allows using parametric tests.

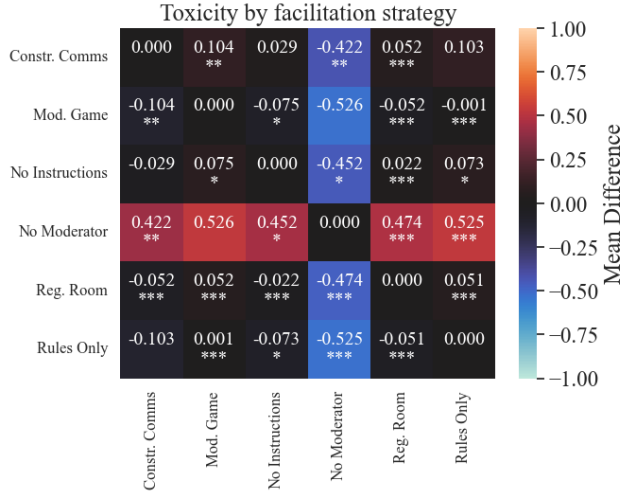


Figure 3: Difference in average toxicity levels for comments following pairs of facilitation strategies. Red cells ( $x > 0$ ) indicate that the strategy on the left performs worse than the one on the bottom, for an average of  $x$  points in a scale of 1-5. Conversely for blue ( $x < 0$ ) cells. Black cells denote minute changes. Asterisks from pairwise Student-t tests ( $\cdot p < 0.1$ ,  $* p < 0.05$ ,  $** p < 0.01$ ,  $*** p < 0.001$ ).

**Finding 3: LLM facilitators choose to intervene far too frequently, which is tolerated by the other participants .** Fig. 4 demonstrates that LLM facilitators intervene at almost any opportunity, even though they are instructed to only do so when necessary. This confirms that LLMs generally can not decide not to speak (§). To our knowledge, this has not been reported in relevant literature, and *is an example of “debugging” problems with LLMs* — a core motivation of our work.

Additionally, a qualitative look through the dataset reveals that LLM user-agents exhibit atypical tolerance for excessive facilitator 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, Trujillo, and Fagni 2022). This is likely another artifact caused by alignment procedures, making LLMs too agreeable (Park et al. 2023).

## Ablation Study

We generate eight synthetic discussions per ablation experiment, using a single model (Qwen 2.5). We compare the diversity (cf. §, ) of these discussions with ones from the CeRI “Regulation Room” dataset<sup>3</sup>, which includes moderated online deliberative discussions for ten diverse topics.

## Effects of LLMs

<sup>3</sup><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.

Variable	Toxicity
Intercept	2.164***
No Instructions	-0.426***
Moderation Game	-0.435***
Rules Only	-0.461***
Regulation Room	-0.277***
Constructive Communications	-0.230***
time	-0.012**
No Instructions×time	-0.003
Moderation Game×time	-0.011*
Rules Only×time	-0.008
Regulation Room×time	-0.023***
Constructive Communications×time	-0.023***

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

Table 1: OLS regression coefficients for toxicity ( $Adj.R^2 = 0.054$ ). Average toxicity without facilitators is 2.164.

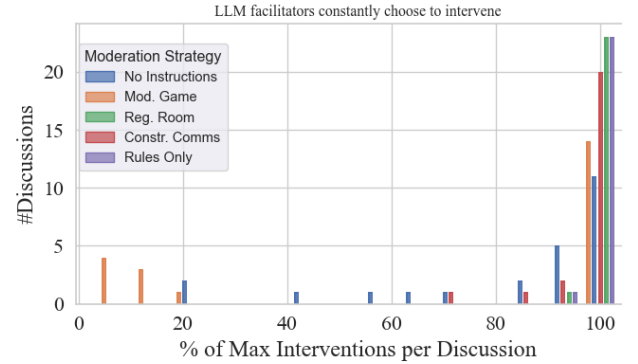


Figure 4: Histogram of interventions by LLM facilitators per strategy used.

**Larger models do not lead to more high-quality discussions.** As shown in Fig. 5a, Qwen demonstrated the highest diversity among the evaluated models, indicating limited participant interaction (§), followed by Mistral Nemo and LLaMa. However, none of the models closely matched the diversity observed in human discussions.

## Effects of Turn-Taking Functions

**Our proposed turn-taking function substantially improves the quality of synthetic data.** We compare our turn-taking function (§) to two baselines: Round Robin (participants speaking one after the other, then repeating) and Random Selection (uniformly sampling another participant each turn). Fig. 5b demonstrates that although all distributions diverge from the blue—human—distribution, our function is the only one not exhibiting extremely high diversity (i.e., very limited participant interaction §).

**Effects of User Prompting** We conduct three separate experiments in which participants are subjected to one of the following conditions at a time: (1) no assigned SDBs, (2) no



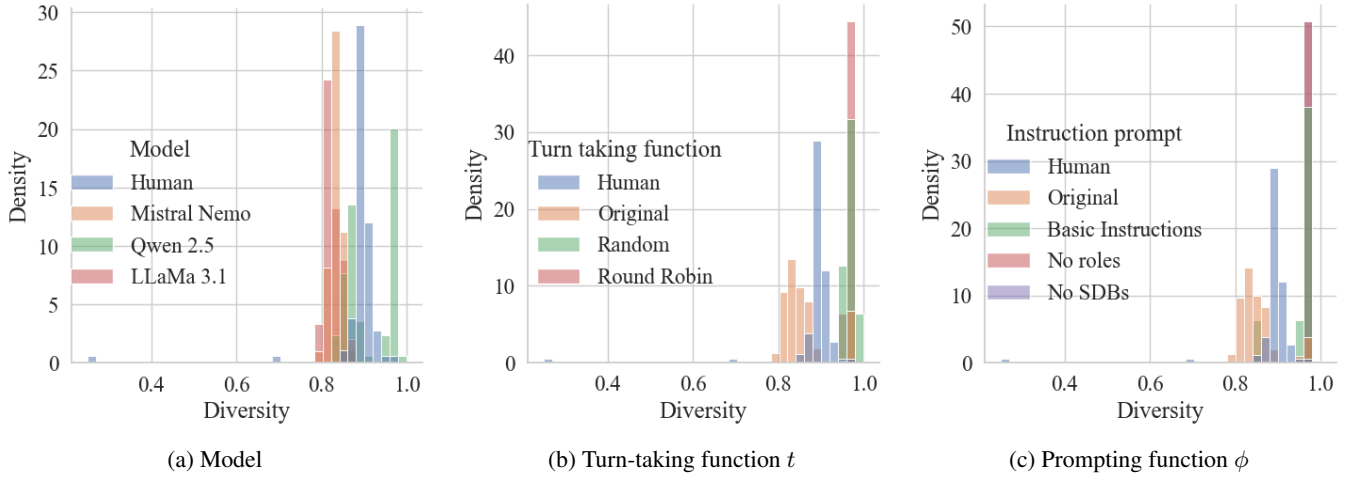


Figure 5: Diversity (\$) distribution for each discussion by LLM (\$), turn-taking function  $t$  (\$), and prompting function  $\phi$  used (\$).

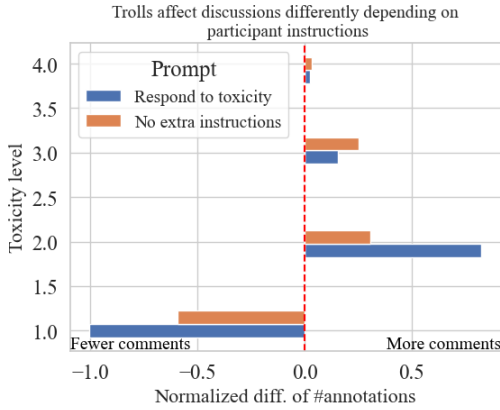


Figure 6: Non-troll toxicity levels in discussions with and without trolls. There is a significant uptick on the number of “somewhat toxic” ( $Toxicity = 2$ ) comments when the participants are primed to respond to toxic comments.

assigned roles, or (3) only a very basic instruction prompt given.

**Specialized instruction prompts are essential for eliciting toxic behavior in instruction-tuned LLMs.** Our instruction prompt for the participants (§) incentivizes them to react to toxic behavior. Indeed, inserting “troll” participants to discussions, leads to more intense toxicity among *other* participants *only if we instruct participants to react to toxic posts* (Fig. 6).

**SDBs, roles, and our instruction prompt all increase the quality of synthetic data.** Fig. 5c illustrates that incorporating SDBs, roles, and specialized instruction prompts, results in diversity scores more closely aligned with human discussions.

## Datasets and Software

We introduce an open-source, lightweight, purpose-built framework for managing, annotating, and generating synthetic discussions. The key features of the framework include:

- Three core functions: generating discussion setups (selecting participants, topics, roles, etc.), executing, and annotating them according to user-provided parameters.
- Built-in fault tolerance (automated recovery and intermittent saving) and file logging to support extended experiments.
- Availability via PIP.

We also release a dataset of synthetic discussions annotated by LLMs. It can serve for finetuning facilitator 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 repository. The dataset is licensed under a CC BY-SA license, and the software under GPLv3. **Warning: The datasets by their nature contain offensive and hateful speech.**

## Conclusions

We proposed a simple and generalizable methodology that enables researchers to quickly and inexpensively conduct pilot facilitation experiments using exclusively LLMs and validated it through an ablation study. We found that (1) LLM facilitators significantly improve the quality of synthetic discussions; (2) prompting these facilitators with strategies based on Social Science research does not markedly improve their performance. We also identified a consistent problem with LLMs not keeping silence when appropriate. Finally, we created an open-source Python Framework that applies this methodology to hundreds of experiments, which we used to create and publish a large-scale synthetic dataset.

## References

- Abdelnabi, S.; Gomaa, A.; Sivaprasad, S.; Schönherr, L.; and Fritz, M. 2024. Cooperation, Competition, and Maliciousness: LLM-Stakeholders Interactive Negotiation. arXiv:2309.17234.
- Amaury, T.; and Stefano, C. 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 Katarya, R. 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.
- Argyle, L. P.; Bail, C. A.; Busby, E. C.; Gubler, J. R.; Howe, T.; Rytting, C.; Sorensen, T.; and Wingate, D. 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.
- Atil, B.; Aykent, S.; Chittams, A.; Fu, L.; Passonneau, R. J.; Radcliffe, E.; Rajagopal, G. R.; Sloan, A.; Tudrej, T.; Ture, F.; Wu, Z.; Xu, L.; and Baldwin, B. 2025. Non-Determinism of "Deterministic" LLM Settings. arXiv:2408.04667.
- Avalle, M.; Marco, N. D.; Etta, G.; Sangiorgio, E.; Alipour, S.; Bonetti, A.; Alvisi, L.; Scala, A.; Baronchelli, A.; Cinelli, M.; and Quattrocioni, W. 2024. Persistent interaction patterns across social media platforms and over time. *Nature*, 628: 582 – 589.
- Balog, K.; Palowitch, J.; Ikica, B.; Radlinski, F.; Alvari, H.; and Manshadi, M. 2024. Towards Realistic Synthetic User-Generated Content: A Scaffolding Approach to Generating Online Discussions. arXiv:2408.08379.
- Burton, J. W.; Lopez-Lopez, E.; Hechtlinger, S.; et al. 2024. How Large Language Models Can Reshape Collective Intelligence. *Nature Human Behaviour*, 8: 1643–1655.
- Cho, H.; Liu, S.; Shi, T.; Jain, D.; Rizk, B.; Huang, Y.; Lu, Z.; Wen, N.; Gratch, J.; Ferrara, E.; and May, J. 2024. Can Language Model Moderators Improve the Health of Online Discourse? In Duh, K.; Gomez, H.; and Bethard, S., eds., *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, 7478–7496. Mexico City, Mexico.
- Clemons, E. K.; Schreieck, M.; and Waran, R. V. 2025. Managing disinformation on social media platforms. *Electronic Markets*, 35(52).
- Cresci, S.; Trujillo, A.; and Fagni, T. 2022. Personalized Interventions on Online Moderation. In *Proceedings of the 33rd ACM Conference on Hypertext and Social Media*, HT '22, 248–251. New York, NY, USA: Association for Computing Machinery. ISBN 9781450392334.
- De Kock, C.; Stafford, T.; and Vlachos, A. 2022. How to disagree well: Investigating the dispute tactics used on Wikipedia. In Goldberg, Y.; Kozareva, Z.; and Zhang, Y., eds., *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, 3824–3837. Abu Dhabi, United Arab Emirates: Association for Computational Linguistics.
- Denniss, E.; and Lindberg, R. 2025. Social media and the spread of misinformation: infectious and a threat to public health. *Health Promotion International*, 40.
- eRulemaking Initiative, C. 2017. CeRI (Cornell e-Rulemaking) Moderator Protocol. Cornell e-Rulemaking Initiative Publications, 21.
- Esau, K.; Friess, D.; and Eilders, C. 2017. Design matters! An empirical analysis of online deliberation on different news platforms. *Policy & Internet*, 9(3): 321–342.
- Falk, N.; Jundi, I.; Vecchi, E. M.; and Lapesa, G. 2021. Predicting Moderation of Deliberative Arguments: Is Argument Quality the Key? In Al-Khatib, K.; Hou, Y.; and Stede, M., eds., *Proceedings of the 8th Workshop on Argument Mining*, 133–141. Punta Cana, Dominican Republic: Association for Computational Linguistics.
- Horta Ribeiro, M.; Cheng, J.; and West, R. 2023. Automated Content Moderation Increases Adherence to Community Guidelines. In *Proceedings of the ACM Web Conference 2023*, WWW '23, 2666–2676. New York, NY, USA: Association for Computing Machinery. ISBN 9781450394161.
- Huang, S.; Siddarth, D.; Lovitt, L.; Liao, T. I.; Durmus, E.; Tamkin, A.; and Ganguli, D. 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, 1395–1417. New York, NY, USA: Association for Computing Machinery. ISBN 9798400704505.
- Kang, H.; and Qian, T. 2024. Implanting LLM's Knowledge via Reading Comprehension Tree for Toxicity Detection. In Ku, L.-W.; Martins, A.; and Srikumar, V., eds., *Findings of the Association for Computational Linguistics ACL 2024*, 947–962. Bangkok, Thailand and virtual meeting: Association for Computational Linguistics.
- Kim, S.; Eun, J.; Seering, J.; and Lee, J. 2021. Moderator Chatbot for Deliberative Discussion: Effects of Discussion Structure and Discussant Facilitation. *Proc. ACM Hum.-Comput. Interact.*, 5(CSCW1).
- Kolluri, A.; Murthy, D.; and Vinton, K. 2025. Quantifying the spread of racist content on fringe social media: A case study of Parler. *Big Data & Society*, 12(2).
- Korre, K.; Tsirmpas, D.; Gkoumas, N.; Cabalé, E.; Kontarinis, D.; Myrtzani, D.; Evgeniou, T.; Androutsopoulos, I.; and Pavlopoulos, J. 2025. Evaluation and Facilitation of Online Discussions in the LLM Era: A Survey. ACL ARR 2025 February Submission.
- Kumar, D.; AbuHashem, Y. A.; and Durumeric, Z. 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.
- Lin, C.-Y. 2004. ROUGE: A Package for Automatic Evaluation of Summaries. In *Text Summarization Branches Out*, 74–81. Barcelona, Spain: Association for Computational Linguistics.
- Mou, X.; Wei, Z.; and Huang, X. 2024. Unveiling the Truth and Facilitating Change: Towards Agent-based Large-scale Social Movement Simulation. arXiv:2402.16333.

- Navajas, J.; Niella, T.; and Garbulsky, G. e. a. 2018. Aggregated knowledge from a small number of debates outperforms the wisdom of large crowds. *Nature Human Behaviour*, 2: 126–132.
- Nik Ahmad, N. A. 2010. CETLs : Supporting Collaborative Activities Among Students and Teachers Through the Use of Think- Pair-Share Techniques. *International Journal of Computer Science Issues*, 7.
- OpenAI; Achiam, J.; Adler, S.; Agarwal, S.; Ahmad, L.; Akkaya, I.; Aleman, F. L.; Almeida, D.; Altenschmidt, J.; Altman, S.; Anadkat, S.; Avila, R.; Babuschkin, I.; Balaji, S.; Balcom, V.; Baltescu, P.; Bao, H.; Bavarian, M.; Belgum, J.; Bello, I.; Berdine, J.; Bernadett-Shapiro, G.; Berner, C.; Bogdonoff, L.; Boiko, O.; Boyd, M.; Brakman, A.-L.; Brockman, G.; Brooks, T.; Brundage, M.; Button, K.; Cai, T.; Campbell, R.; Cann, A.; Carey, B.; Carlson, C.; Carmichael, R.; Chan, B.; Chang, C.; Chantzis, F.; Chen, D.; Chen, S.; Chen, R.; Chen, J.; Chen, M.; Chess, B.; Cho, C.; Chu, C.; Chung, H. W.; Cummings, D.; Currier, J.; Dai, Y.; Decareaux, C.; Degry, T.; Deutsch, N.; Deville, D.; Dhar, A.; Dohan, D.; Dowling, S.; Dunning, S.; Ecoffet, A.; Eleti, A.; Eloundou, T.; Farhi, D.; Fedus, L.; Felix, N.; Fishman, S. P.; Forte, J.; Fulford, I.; Gao, L.; Georges, E.; Gibson, C.; Goel, V.; Gogineni, T.; Goh, G.; Gontijo-Lopes, R.; Gordon, J.; Grafstein, M.; Gray, S.; Greene, R.; Gross, J.; Gu, S. S.; Guo, Y.; Hallacy, C.; Han, J.; Harris, J.; He, Y.; Heaton, M.; Heidecke, J.; Hesse, C.; Hickey, A.; Hickey, W.; Hoeschele, P.; Houghton, B.; Hsu, K.; Hu, S.; Hu, X.; Huizinga, J.; Jain, S.; Jain, S.; Jang, J.; Jiang, A.; Jiang, R.; Jin, H.; Jin, D.; Jomoto, S.; Jonn, B.; Jun, H.; Kafkhan, T.; Łukasz Kaiser; Kamali, A.; Kanitscheider, I.; Keskär, N. S.; Khan, T.; Kilpatrick, L.; Kim, J. W.; Kim, C.; Kim, Y.; Kirchner, J. H.; Kiros, J.; Knight, M.; Kokotajlo, D.; Łukasz Kondraciuk; Kondrich, A.; Konstantinidis, A.; Kosic, K.; Krueger, G.; Kuo, V.; Lampe, M.; Lan, I.; Lee, T.; Leike, J.; Leung, J.; Levy, D.; Li, C. M.; Lim, R.; Lin, M.; Lin, S.; Litwin, M.; Lopez, T.; Lowe, R.; Lue, P.; Makanju, A.; Malfacini, K.; Manning, S.; Markov, T.; Markovski, Y.; Martin, B.; Mayer, K.; Mayne, A.; McGrew, B.; McKinney, S. M.; McLeavey, C.; McMillan, P.; McNeil, J.; Medina, D.; Mehta, A.; Menick, J.; Metz, L.; Mishchenko, A.; Mishkin, P.; Monaco, V.; Morikawa, E.; Mossing, D.; Mu, T.; Murati, M.; Murk, O.; Mély, D.; Nair, A.; Nakano, R.; Nayak, R.; Nee-lakantan, A.; Ngo, R.; Noh, H.; Ouyang, L.; O’Keefe, C.; Pachocki, J.; Paino, A.; Palermo, J.; Pantuliano, A.; Parascandolo, G.; Parish, J.; Parparita, E.; Passos, A.; Pavlov, M.; Peng, A.; Perelman, A.; de Avila Belbute Peres, F.; Petrov, M.; de Oliveira Pinto, H. P.; Michael; Pokorny; Pokrass, M.; Pong, V. H.; Powell, T.; Power, A.; Power, B.; Proehl, E.; Puri, R.; Radford, A.; Rae, J.; Ramesh, A.; Raymond, C.; Real, F.; Rimbach, K.; Ross, C.; Rotsted, B.; Roussez, H.; Ryder, N.; Saltarelli, M.; Sanders, T.; Santurkar, S.; Sastry, G.; Schmidt, H.; Schnurr, D.; Schulman, J.; Sel-sam, D.; Sheppard, K.; Sherbakov, T.; Shieh, J.; Shoker, S.; Shyam, P.; Sidor, S.; Sigler, E.; Simens, M.; Sitkin, J.; Slama, K.; Sohl, I.; Sokolowsky, B.; Song, Y.; Staudacher, N.; Such, F. P.; Summers, N.; Sutskever, I.; Tang, J.; Tezak, N.; Thompson, M. B.; Tillet, P.; Tootoonchian, A.; Tseng, E.; Tuggle, P.; Turley, N.; Tworek, J.; Uribe, J. F. C.; Val-lone, A.; Vijayvergiya, A.; Voss, C.; Wainwright, C.; Wang, J. J.; Wang, A.; Wang, B.; Ward, J.; Wei, J.; Weinmann, C.; Welihinda, A.; Welinder, P.; Weng, J.; Weng, L.; Wiethoff, M.; Willner, D.; Winter, C.; Wolrich, S.; Wong, H.; Work-man, L.; Wu, S.; Wu, J.; Wu, M.; Xiao, K.; Xu, T.; Yoo, S.; Yu, K.; Yuan, Q.; Zaremba, W.; Zellers, R.; Zhang, C.; Zhang, M.; Zhao, S.; Zheng, T.; Zhuang, J.; Zhuk, W.; and Zoph, B. 2024. GPT-4 Technical Report. arXiv:2303.08774.
- Park, J. S.; O’Brien, J. C.; Cai, C. J.; Morris, M. R.; Liang, P.; and Bernstein, M. S. 2023. Generative Agents: Interac-tive Simulacra of Human Behavior. arXiv:2304.03442.
- Park, J. S.; Popowski, L.; Cai, C.; Morris, M. R.; Liang, P.; and Bernstein, M. S. 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. ISBN 9781450393201.
- Pavlopoulos, J.; Sorensen, J.; Dixon, L.; Thain, N.; and An-droustopoulos, I. 2020. Toxicity Detection: Does Context Really Matter? In Jurafsky, D.; Chai, J.; Schluter, N.; and Tetreault, J., eds., *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 4296–4305. Online: Association for Computational Linguistics.
- Pranesh, S.; and Gupta, S. 2024. The impact of social media on polarization in the society. arXiv:2408.12877.
- Rosenberg, M. B.; and Chopra, D. 2015. *Nonviolent commu-nication: A language of life: Life-changing tools for healthy relationships*. PuddleDancer Press.
- Rossetti, G.; Stella, M.; Cazabet, R.; Abramski, K.; Cau, E.; Citraro, S.; Failla, A.; Improta, R.; Morini, V.; and Pansanella, V. 2024. Y Social: an LLM-powered Social Me-dia Digital Twin. arXiv:2408.00818.
- Rossi, L.; Harrison, K.; and Shklovski, I. 2024. The Prob-lems of LLM-generated Data in Social Science Research. *Sociologica*, 18(2): 145–168.
- Schaffner, B.; Bhagoji, A. N.; Cheng, S.; Mei, J.; Shen, J. L.; Wang, G.; Chetty, M.; Feamster, N.; Lakier, G.; and Tan, C. 2024. "Community Guidelines Make this the Best Party on the Internet": An In-Depth Study of Online Platforms’ Content Moderation Policies. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, CHI ’24. New York, NY, USA: Association for Computing Machinery. ISBN 9798400703300.
- Schluger, C.; Chang, J.; Danescu-Niculescu-Mizil, C.; and Levy, K. 2022. Proactive Moderation of Online Discus-sions: Existing Practices and the Potential for Algorithmic Support. *Proc. ACM Hum.-Comput. Interact.*, 6(CSCW2).
- Schroeder, H.; Roy, D.; and Kabbara, J. 2024. Fora: A cor-pus and framework for the study of facilitated dialogue. In *Proceedings of the 62nd Annual Meeting of the Associa-tion for Computational Linguistics*, 13985–14001. Bangkok, Thailand.
- Seering, J. 2020. Reconsidering Self-Moderation: the Role of Research in Supporting Community-Based Models for



Online Content Moderation. *Proc. ACM Hum.-Comput. Interact.*, 4(CSCW2).

Small, C. T.; Vendrov, I.; Durmus, E.; Homaei, H.; Barry, E.; Cornebise, J.; Suzman, T.; Ganguli, D.; and Megill, C. 2023. Opportunities and Risks of LLMs for Scalable Deliberation with Polis. *ArXiv*, abs/2306.11932.

Tsai, L. L.; Pentland, A.; Braley, A.; Chen, N.; Enríquez, J. R.; and Reuel, A. 2024. Generative AI for Pro-Democracy Platforms. *An MIT Exploration of Generative AI*. <https://mit-genai.pubpub.org/pub/mn45hexw>.

Törnberg, P.; Valeeva, D.; Uitermark, J.; and Bail, C. 2023. Simulating Social Media Using Large Language Models to Evaluate Alternative News Feed Algorithms. *arXiv:2310.05984*.

Ulmer, D.; Mansimov, E.; Lin, K.; Sun, J.; Gao, X.; and Zhang, Y. 2024. Bootstrapping LLM-based Task-Oriented Dialogue Agents via Self-Talk. *ArXiv*, abs/2401.05033.

Vezhnevets, A. S.; Agapiou, J. P.; Aharon, A.; Ziv, R.; Matyas, J.; Du'enez-Guzm'an, E. A.; Cunningham, W. A.; Osindero, S.; Karmon, D.; and Leibo, J. Z. 2023. Generative agent-based modeling with actions grounded in physical, social, or digital space using Concordia. *ArXiv*, abs/2312.03664.

Wang, Y.-S.; and Chang, Y. T. 2022. Toxicity Detection with Generative Prompt-based Inference. *ArXiv*, abs/2205.12390.

White, K.; Hunter, N.; and Greaves, K. 2024. *facilitating deliberation - a practical guide*. Mosaic Lab.

Xia, Y.; Zhu, H.; Lu, T.; Zhang, P.; and Gu, N. 2020. Exploring Antecedents and Consequences of Toxicity in Online Discussions: A Case Study on Reddit. *Proc. ACM Hum.-Comput. Interact.*, 4(CSCW2).

Yu, Y.; Yao, Z.; Li, H.; Deng, Z.; Cao, Y.; Chen, Z.; Suchow, J. W.; Liu, R.; Cui, Z.; Xu, Z.; Zhang, D.; Subbalakshmi, K.; Xiong, G.; He, Y.; Huang, J.; Li, D.; and Xie, Q. 2024. FinCon: A Synthesized LLM Multi-Agent System with Conceptual Verbal Reinforcement for Enhanced Financial Decision Making. *arXiv:2407.06567*.