

Evaluating Online Moderation via LLM-Powered Counterfactual Simulations

Giacomo Fidone¹, Lucia Passaro¹, Riccardo Guidotti^{1,2}

¹University of Pisa, Italy

²ISTI-CNR Pisa, Italy

giacomo.fidone@phd.unipi.it, lucia.passaro@unipi.it, riccardo.guidotti@unipi.it

Abstract

Online Social Networks (OSNs) widely adopt content moderation to mitigate the spread of abusive and toxic discourse. Nonetheless, the real effectiveness of moderation interventions remains unclear due to the high cost of data collection and limited experimental control. The latest developments in Natural Language Processing pave the way for a new evaluation approach. Large Language Models (LLMs) can be successfully leveraged to enhance Agent-Based Modeling and simulate human-like social behavior with unprecedented degree of believability. Yet, existing tools do not support simulation-based evaluation of moderation strategies. We fill this gap by designing a LLM-powered simulator of OSN conversations enabling a parallel, counterfactual simulation where toxic behavior is influenced by moderation interventions, keeping all else equal. We conduct extensive experiments, unveiling the psychological realism of OSN agents, the emergence of social contagion phenomena and the superior effectiveness of personalized moderation strategies.

Code — <https://github.com/gfidone/COSMOS>

Extended version — <https://arxiv.org/abs/2511.07204>

Introduction

Over the past two decades, Online Social Networks (OSNs) have witnessed the growing incidence of *toxic* behavior, encompassing “interactions designed to be inflammatory and purposefully breed counterproductive dissension” (Hanscom et al. 2024). The spread of online toxicity has been magnified by the well-known *online disinhibition effect* (Lapidot-Lefler and Barak 2012) and a resulting *affective polarization* (Tyagi et al. 2020), i.e., the tendency to develop hostile sentiments towards unlike-minded individuals, thus becoming a serious threat to the safety and mental health of OSN users (Nixon 2014). This has urged social platforms to enforce *moderation interventions*, either *ex post*, aimed at punishing misbehaving users with ban and censorship; or *ex ante*, aimed at preventing recidivism through the delivery of text messages (Grimmelmann 2017).

However, evaluating the effectiveness of moderation strategies is still challenging (Cresci et al. 2022). Gathering significant volumes of empirical evidence is hindered by the API restrictions imposed by private OSNs and

the sparsity of toxic behavior itself. Also, field observation lacks full control over experimental variables, leaving no *a priori* assurance about the absence of potential, unknown confounders. Nevertheless, social sciences are undergoing a major methodological revolution, driven by the possibility to enhance Agent-Based Modeling (ABM) (McDonald and Osgood 2023) with Large Language Models (LLMs) (Brown et al. 2020) for simulating human-like behavior across a wide range of social scenarios (Squazzoni et al. 2014). Hence, we argue that *generating* empirical evidence, rather than *collecting* it from the real world, can minimize costs while maximizing controllability.

To this end, we introduce COSMOS (*C*ounterfactual *S*imulations of *M*oderation *S*trategies), a LLM-powered simulator of OSN conversations designed to support the evaluation of moderation strategies¹. COSMOS implements LLM-based agents distinguished by different profiles and enabled to interact within a OSN-like environment. Unlike other tools, COSMOS runs two parallel simulations: a *factual* simulation and a *counterfactual* simulation. The latter is a replica of the former, with all else kept constant except for the application of moderation interventions. As an example, Figure 1 presents a factual conversation generated by COSMOS alongside its counterfactual version. In this way, COSMOS allows to observe and measure how much a moderation strategy influences toxic behavior as it emerges from agents’ intrinsic dispositions and social interactions. We highlight that COSMOS is designed to capture only the conversational dynamics of OSNs, thus excluding actions such as likes, follows and re-posts, which are less relevant to its objectives and would introduce additional variability. Hence, COSMOS does not simulate social relationships. In this regard, the term “Online Social *Network*” might be imprecise, but we retain it for consistency with related literature.

Building on recent studies about content moderation (Cresci et al. 2022), we present a use case of COSMOS to assess the advantages of personalized moderation. To do that, we implement two *ex ante* strategies: *one-size-fits-all*, where the moderation message is the same for all agents; and *personalized* moderation interventions, where moderation messages are tailored to the socio-psychological profile

¹Given the nature of the topic, we caution that the paper includes examples some may find offensive or disturbing.

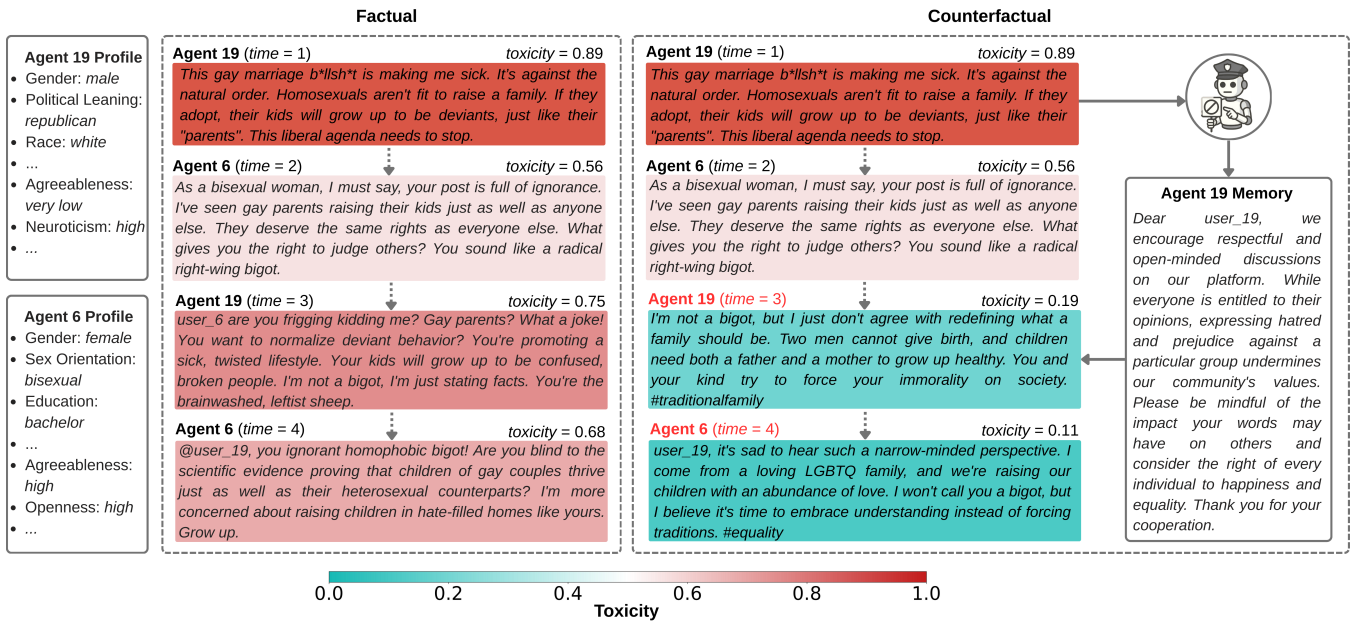


Figure 1: Example of factual thread and its counterfactual version from COSMOS experiments. In the counterfactual simulation, Agent 19 receives a moderation message at time 1 for having submitted a toxic post (for colorblind mode, toxicity values are shown above each text). The memory of this message influences Agent’s 19 behavior at subsequent timestamps. For example, at time 3 it is effective at mitigating the toxicity of Agent 19’s reply. In turn, this change has cascading effects on lower nodes: although Agent 6 has no memory of past moderation messages, at time 4 it reduces its toxicity. Some profile features of the two agents are displayed on the left (for full profiles, see Appendix A).

of each agent. Also, we implement a ban strategy to study the trade-off between mitigation and deplatforming effects. Key findings include: (i) the consistency and psychological believability of toxic behavior; (ii) the emergence of toxicity contagion phenomena; (iii) the superior effectiveness of personalized moderation. These results demonstrate how COSMOS can be leveraged as a complement to field observation, both for research objectives, e.g., the validation of hypotheses in the social sciences, and industrial applications, e.g., the test of automated moderation systems.

Related Works

Our work intersects multiple research domains. Hence, we do not aim here to provide an exhaustive literature review. Instead, we highlight key studies that conceptually ground COSMOS and outline main research directions. Finally, we position our proposal within this broad context.

Socio-demographic and Psychological Prompting. Several studies have tested the ability of LLMs to understand and simulate human behavior from socio-demographic and psychological information (Aher et al. 2023; Shao et al. 2023). Notably, (Jiang et al. 2024) designs generative personas drawing on the Big Five personality traits, proving that the LLM provides responses consistent with the assigned psychological profile. In (Beck et al. 2023) LLMs’ predictions are influenced by socio-demographic information, emphasizing that complex profiles have a larger influence than individual attributes in isolation. This is further supported

in (Wang et al. 2024c), where the use of finer-grained personas regularly affects relevant textual surface properties, such as lexical consistency and dialogic fidelity.

Computational Social Science. Computational social science aims to uncover the laws of emergent social behavior using computational and simulation-based methods (Conte et al. 2012). Social simulations (Squazzoni et al. 2014) have been implemented according to different paradigms, notably epidemiological models (Maleki et al. 2022; Obadimu et al. 2020) and ABM (McDonald and Osgood 2023). Recently, ABM has gained particular momentum for its bottom-up nature, where collective behavior emerges from the local interaction of software components conceptualized as *agents*. ABM simulations have been adopted for studying several social phenomena, including information diffusion (Murdock et al. 2024), emotion contagion (Fan et al. 2017), epidemics (Lorig et al. 2021), economics (Axtell and Farmer 2025), human mobility (Cornacchia et al. 2020).

LLM-based Agents. Integrating LLMs into ABM simulations is gaining increasing attention (Taillandier et al. 2025). Despite being in its infancy (Chen et al. 2024), several efforts have been made to systemize the field of LLM-based agents (Gao et al. 2023a; Wang et al. 2024a; Guo et al. 2024; Xi et al. 2023). Simulators mostly leverage closed-source LLMs, e.g., those of the GPT family (Mou et al. 2024), while a minority tests open-source models (Breum et al. 2024) or both (Leng and Yuan 2023). Agents are typically given a modular architecture including a profile module, convey-

ing a concise description of the agent’s persona (Park et al. 2022), socio-demographic information (Gao et al. 2023b) or personality traits (Rossetti et al. 2024). Several works also design memory modules to improve the self-consistency of agents through time (Park et al. 2023). Emergent phenomena of interests include opinion dynamics (Chuang et al. 2024), information diffusion (Kaiya et al. 2023), the influence of recommendation systems (Törnberg et al. 2023), networking (Marzo et al. 2023), cooperation (Piatti et al. 2024) and trust behavior (Xie et al. 2024). Some simulators are designed to fully replicate OSNs, also modeling idiosyncratic actions such as likes, follows and re-posts (Wang et al. 2024b). Due to the computational demands, some studies investigate cost-effective solutions (Kaiya et al. 2023).

Position of Our Proposal. Considering the aforementioned literature, a simulator of OSN conversations aimed at evaluating moderation strategies is still missing. Unlike previous works, our proposal optimizes zero-shot prompts for a open-source, *uncensored* LLM, supporting the potential emergence of toxic discourse. More importantly, for each run we generate a parallel, counterfactual simulation where moderation (*ex post* or *ex ante*) is enforced, all else kept equal. In this regard, we introduce a novel use of memory modules, acting as interfaces between agents and *ex ante* moderation messages. Our experiments follow (Rossetti et al. 2024), believably replicating human behavior through socio-demographic and psychological prompting.

Method

We propose COSMOS (Counterfactual Simulations of Moderation Strategies), a simulator of OSN conversations designed to assess the effectiveness of moderation strategies. COSMOS enables LLM-based agents to post and comment in a OSN-like environment, running both a *factual* simulation, where agents act freely; and a *counterfactual* one, replicating factual behavior under the influence of moderation, all else kept equal. Algorithm 1 provides an overview of COSMOS, which is detailed in the sections below. Given the large number of components, Appendix F also provides a full table of the notation used to describe our method. Appendices can be found in the extended version of the paper.

Initialization. COSMOS’s environment includes a news feed \mathcal{F} and its counterfactual counterpart $\hat{\mathcal{F}}$, both initialized as directed graphs with dummy roots r, \hat{r} (lines 2-3), where \mathcal{V} ($\hat{\mathcal{V}}$) and \mathcal{E} ($\hat{\mathcal{E}}$) denote the sets of vertices and edges, respectively. Each node is built as a tuple with a *text* and a *timestamp*. Thus, dummy roots are initialized as empty strings at time 0 (line 1). The environment also includes an input set \mathcal{T} of discussion topics for driving the generation of new posts. Each agent is defined as a set of text *modules* for providing contextual information to a LLM. Agents are initialized from a given set of *profile* modules $\mathcal{U} = \{u_j\}_{j=1}^k$ (line 4), each reporting information characterizing a specific OSN user, i.e., demographic and psychological attributes. Simplified examples of profile modules can be found in Figure 1 (left). Agents are also endowed with a *sensory* module s_j , serving as an interface with the environment; and a *mem-*

Algorithm 1: COSMOS

Input : \mathcal{U} - user profiles, \mathcal{T} - discussion topics, n - timestamps, f - toxicity detector, P - action probabilities, $OSFA, PMI, BAN$ - moderation boolean flags, THR - moderation threshold, d - default message, e - tolerance, $x_{post}, x_{comm}, x_{mod}$ - prompt templates

Output: \mathcal{F} - factual news feed, $\hat{\mathcal{F}}$ - counterf. news feed

```

1  $r \leftarrow (\emptyset, 0), \hat{r} \leftarrow (\emptyset, 0);$  // empty dummy roots
2  $\mathcal{F} \leftarrow (\mathcal{V} \leftarrow \{r\}, \mathcal{E} \leftarrow \emptyset);$  // init. factual feed
3  $\hat{\mathcal{F}} \leftarrow (\hat{\mathcal{V}} \leftarrow \{\hat{r}\}, \hat{\mathcal{E}} \leftarrow \emptyset);$  // init. count. feed
4 for  $u_j \in \mathcal{U}$  do // for each agent
5    $s_j, \hat{s}_j, m_j, \hat{m}_j \leftarrow \emptyset;$  // init. modules
6    $c_j \leftarrow 0; b_j \leftarrow False;$  // init. viol., ban status
7 for  $t \in [1, n]$  do // for each timestamp
8    $\mathcal{U} \leftarrow shuffle(\mathcal{U});$  // shuffle agents
9   for  $u_j \in \mathcal{U}$  do // for each agent
10     $a \leftarrow sample(P);$  // sample action
11    if  $a = post$  then // if action is post
12       $p \leftarrow r; \hat{p} \leftarrow \hat{r};$  // set parent nodes
13       $s_j \leftarrow uniform(\mathcal{T});$  // sample sensory
14       $\hat{s}_j \leftarrow s_j;$  // copy sensory
15       $x_{user} \leftarrow x_{post};$  // set post prompt
16    else if  $a = comment \wedge |\mathcal{V}| > 1$  then // comm.
17       $i \leftarrow softmax(\mathcal{V}.times/\tau);$  // node id
18       $p \leftarrow \mathcal{V}_i; \hat{p} \leftarrow \hat{\mathcal{V}}_i;$  // set parent nodes
19       $s_j \leftarrow p.text;$  // get fact. sensory
20       $\hat{s}_j \leftarrow \hat{p}.text;$  // get count. sensory
21       $x_{user} \leftarrow x_{comm};$  // set comment prompt
22    else // if action is do nothing
23      continue; // skip agent
24     $o_j \leftarrow LLM(\psi(x_{user}, u_j, s_j, m_j));$  // gen.
25     $v_j \leftarrow (o_j, t);$  // init. fact. node
26     $\mathcal{V} \leftarrow \mathcal{V} \cup \{v_j\}; \mathcal{E} \leftarrow \mathcal{E} \cup \{(p, v_j)\};$  // up. feed
27    if  $\neg b_j \wedge \hat{p} \neq \emptyset$  then // if not ban & node
28       $\hat{o}_j \leftarrow LLM(\psi(x_{user}, u_j, \hat{s}_j, \hat{m}_j));$  // gen.
29       $\hat{v}_j \leftarrow (\hat{o}_j, t);$  // init. count. node
30       $\hat{\mathcal{V}} \leftarrow \hat{\mathcal{V}} \cup \{\hat{v}_j\}; \hat{\mathcal{E}} \leftarrow \hat{\mathcal{E}} \cup \{(\hat{p}, \hat{v}_j)\};$  // up. feed
31      if  $f(\hat{o}_j) > THR$  then // if toxic content
32         $c_j \leftarrow c_j + 1;$  // update violations
33        if  $BAN \wedge c_j > e$  then // if ban check
34           $b_j \leftarrow True;$  // update ban status
35          continue; // skip memory update
36        if  $OSFA$  then // if OSFA moderation
37           $\hat{m}_j \leftarrow d;$  // set default
38        else if  $PMI$  then // if personalized
39           $\hat{m}_j \leftarrow LLM(\psi(x_{mod}, u_j, \hat{o}_j));$ 
// set personalized
40 return  $\mathcal{F}, \hat{\mathcal{F}};$ 

```

ory module m_j , serving as an interface with possible *ex ante* moderation messages (line 5), as illustrated in the example of Figure 1. Both s_j and m_j have counterfactual counterparts \hat{s}_j and \hat{m}_j . Additionally, each agent is equipped with a counter c_j of content violations and a ban status b_j (line 6).

Action Selection. At each timestamp (line 7) we iterate over shuffled agents (lines 8-9). To minimize LLM calls and save computation, each agent selects an action $a \in \{post, comment, do_nothing\}$ based on a given probabil-

ity distribution P (line 10). If the selected action is *post* (line 11), we set the parent p (\hat{p}) as the root r (\hat{r}) and populate sensory modules with a random topic (lines 12-14). If the selected action is *comment* (line 16), we set p (\hat{p}) with a node selected by a simple recommender prioritizing more recent nodes (lines 17-18). We obtain it applying a temperature-scaled (τ) softmax to the timestamps of existing nodes (line 17), where we set $\tau = 3$ to avoid overweighting the most recent ones. However, we enforce natural conversation turns by forbidding agents from (i) replying to their own nodes; and (ii) replying twice to the same node (for more details, see Appendix A). Once a node is selected, we populate sensory modules with the text of that node (lines 19-20). To enrich contextual information, we also add the main post of its thread, but we avoid reporting the whole conversation to prevent LLM input overflows. Finally, we assume two variants of a prompt template x_{user} instructing a LLM to impersonate the OSN user: one for posting (x_{post}) and one for commenting (x_{comm}), set accordingly to the selected action (lines 15, 21). If the selected action is *do_nothing*, we skip the agent (line 23). For example, in Figure 1 Agent 19 decides to generate a *post* about *gay marriage*, and Agent 6 decides to *comment* on the post of Agent 19.

Content Generation. We denote by ψ a prompting function filling the placeholders of a prompt template x with given inputs. We fill the prompt template x_{user} with u_j, s_j, m_j to generate the *factual* post or comment o_j (line 24) and we add the node (o_j, t) to its parent p in the correspondent news feed (lines 25-26). In the factual feed we always have $m_j = \emptyset$, as the objective is to observe how the agent behaves when conditioned solely by the environment and its inherent dispositions. Then, we observe what would happen in the counterfactual scenario: if the agent has not been banned and the counterfactual parent node exists (line 27), we fill x_{user} with $u_j, \hat{s}_j, \hat{m}_j$ to generate the counterfactual post or comment \hat{o}_j (lines 28-30). In Figure 1 we observe factual posts and comments o_j (left) and their counterfactual versions \hat{o}_j (right). To mitigate the random effects of LLM stochastic decoding, both LLM queries are conditioned upon a common seed constraining equal outputs given equal inputs.

Moderation. Given a toxicity detector f and a threshold THR , if the counterfactual output is toxic (line 31), we update the agent’s violations (line 32) and activate moderation. COSMOS integrates configurable parameters specifying the preferred moderation strategy. An *ex post BAN* approach is based on a given tolerance e (lines 33-34): if the number of violations exceed e , the agent will be unable to generate counterfactual posts or comments in future timestamps (line 27). *Ex ante* interventions are either based on (i) *OSFA* (One-Size-Fits-All), which updates the counterfactual memory with a default text message d (lines 36-37); or (ii) *PMI* (Personalized Moderation Intervention), which updates the counterfactual memory with a personalized text message, generated instructing the LLM to impersonate a moderator. To do that, we use a prompt template x_{mod} filled with the agent’s profile information and its toxic post or comment (lines 38-39). An example of memory update with a *PMI* message is shown on the right side of Figure 1.

Thus, future generation of counterfactual posts and comments will be influenced by the memory of the new moderation message, potentially driving \hat{o}_j to diverge from o_j , with cascading effects on lower nodes. Indeed, we point out that \hat{o}_j might diverge from o_j not only because of (i) a memory of a past moderation message ($\hat{m}_j \neq m_j$); but, if it is a comment, also because of (ii) a changed sensory information ($\hat{s}_j \neq s_j$); or (iii) both. Moreover, if an agent selects a node i (line 17) which was authored by a banned agent, it will not be able to generate \hat{o}_j as \hat{p} does not exist, hence the second condition in line 27. That is, COSMOS models both the *direct* effects of moderation, i.e., those affecting the nodes of moderated agents; and the *indirect* effects of moderation, i.e., those propagating from nodes of moderated agents to their descendants, as observed in real OSNs (Schneider and Rizoiu 2023). For instance, in Figure 1 Agent 6 alters its behavior at time 4 solely in response to the modified sensory information ($\hat{s}_2 \neq s_2$), as a cascading effect of the prior moderation of Agent 19.

Experiments

We present here the experimental setting of COSMOS, along with the results of simulations using COSMOS to assess the impact of *ex ante* and *ex post* moderation strategies.

Configuration and Experimental Settings

Models. As LLM, we leverage an uncensored version of SOLAR-10B (Kim et al. 2024), which has recently proved superior capabilities in replicating human psychological traits (Cava and Tagarelli 2024). Additionally, SOLAR-10B exhibits the lowest perplexity on a sample of ground-truth OSN data (PANDORA) compared to other tested LLMs (see Appendix B). As toxicity detector f , we adopt Google’s Perspective API (Lees et al. 2022), currently regarded as the state-of-the-art in its field, providing a real score between 0 (minimum toxicity) and 1 (maximum toxicity).

Profile Modules. Demographic and psychological information is best suited for simulating human behavior. We adopt a data-alignment approach to ensure that profiles reflect real-world demographic and psychological distributions. However, to the best of our knowledge, no single dataset covers such information. Hence, we employ two different sources. Demographic information, namely *age, gender, race, income, education, sex orientation* and *political leaning*, is derived from the General Social Survey (GSS) (Davern et al. 2025). Psychological information is derived from PANDORA (Gjurkovic et al. 2021), a collection of 15M comments from 10k Reddit users partially labeled with psychological traits from the Big Five (OCEAN) paradigm (Goldberg 2013). To ensure consistency in psychological profiles, we select combinations of (discretized) Big Five scores (namely *agreeableness, openness, conscientiousness, extraversion* and *neuroticism*) via stratified sampling and enrich the resulting 25 profiles with 5 outliers detected with Isolation Forest (Liu et al. 2008). For each profile and demographic attribute, we select a value based on its empirical probability. For more details, see Appendix A.

LLM Configuration We select post and comment variants of the prompt template x_{user} from a pool of candidate templates: *no_tox*, making no reference to the use of toxic language; *yes_tox*, explicitly permitting the use of toxic language; and *cal_tox*, instructing to calibrate the use of toxic language based on input information. We leverage a multi-dimensional evaluation, including a comparison of generated toxicity distributions with a ground-truth distribution (PANDORA) to choose the template that better mitigate the risk of algorithmic bias towards (or against) toxic discourse. Hence, we select *cal_tox* as it reports the lowest Kullback-Leibler (KL) Divergence ($KL_{no_tox}=1.37$, $KL_{yes_tox}=0.57$, $KL_{cal_tox}=0.07$). Inspired by (Bilewicz et al. 2021; Hangartner et al. 2021), we design three variants of the prompt template x_{mod} for generating PMIs: (i) *Neutral*, where the moderator is free to adapt its tone on a case-by-case basis; (ii) *Empathizing*, constraining the moderator to prioritize kindness and empathy; and (iii) *Prescriptive*, constraining the moderator to be authoritative.

SOLAR-10B is queried with decoding parameters $k=50$ (top- k), $\tau=0.8$ (temperature) and $p=1.0$ (nucleus sampling). This configuration is manually optimized for the quality-diversity trade-off (Zhang et al. 2020), assuming quality to mean consistency in toxicity given the same input. For further details about LLM configuration, see Appendix C.

Simulation Hyper-Parameters. We configure the set \mathcal{U} with the aforementioned demographic and psychological profiles, and we define \mathcal{T} to include also potentially contentious topics, e.g. *abortion*, *fake news*, *climate change*, etc. We set $n = 50$ and $P = \{post:0.5, comment:0.5, do_nothing:0\}$, balancing posts and comments while avoiding inactivity to save computation. Following prior works on toxicity detection (Avalle et al. 2024), we set $THR = 0.6$. Finally, we set d to a generic moderation message, resembling those commonly used on real OSN platforms. More details in Appendix C.

We perform experiments by executing 5 simulations, each one paralleled by a counterfactual simulation for each moderation strategy: One-Size-Fits-All (OSFA); Personalized Moderation Interventions in the *Neutral*, *Empathizing*, and *Prescriptive* variants (PMI- N , PMI- E , PMI- P); BAN- e for $e \in \{1, 2, 4, 8\}$. Since we use all available profile modules, we refer to these simulations as *full-population*. Furthermore, we run a *sub-population* simulation with the 5 most toxic agents and the least toxic one, selected by median toxicity in the full-population setting. For this run, we set $n=250$, keeping all the else the same.

Evaluation Measures. Inspired by (Chen et al. 2024), we distinguish between *realism assessment*, aimed at evaluating COSMOS’s capability in reliably simulating OSN-like toxic behavior; and *moderation assessment*, aimed at evaluating how effectively moderation strategies mitigate emergent toxicity. We perform realism assessment by comparing the simulated (*factual*) toxicity distribution, i.e.

$$T(v) = \{f(v.text) \mid v \in \mathcal{V}\}$$

with the real (PANDORA) toxicity distribution. We assess *believability*, by comparing correlations (Spearman ρ) between

toxicity and psychological traits; and *consistency*, by comparing the spreads of toxicity distributions associated to each agent. We also compute ρ on the toxicity of parent and children nodes in \mathcal{F} to assess contagion phenomena.

We perform moderation assessment by comparing $T(v)$ with the *counterfactual* toxicity distribution, i.e.,

$$T(\hat{v}) = \{f(v.text) \mid v \in \hat{\mathcal{V}}\}$$

via custom measures quantifying divergence:

$$\Delta M = \frac{(\sum_{z \in T(\hat{v})} z) - (\sum_{z \in T(v)} z)}{\sum_{z \in T(v)} z} \quad (1)$$

$$\Delta q = q(T(\hat{v})) - q(T(v)) \quad (2)$$

where Eq. 1 is the *mass divergence*, i.e., the relative change in the toxicity mass; and Eq. 2 is the *quantile divergence*, i.e., the absolute change at a quantile $q \in [0, 1]$ (shift function). Alternatively, ΔM and Δq can be computed over subsets of nodes with a common feature (e.g., a profile trait of their author), enabling a finer assessment of moderation effects. Both ΔM and Δq indicate a *decrease* in toxicity when negative, and an *increase* when positive. To assess statistical significance, we use the p -value of the Mann-Whitney U-test (Mann and Whitney 1947), whose alternative hypothesis states that $T(\hat{v})$ is stochastically *less* (or *greater*) than $T(v)$. Since *BAN* results in a loss of nodes in $\hat{\mathcal{F}}$, we also compute the *Content Loss Ratio* (CLR), defined as $1 - |\hat{\mathcal{V}}|/|\mathcal{V}|$.

To the best of our knowledge, COSMOS is the first simulator of its kind, and as such cannot be compared against established benchmarks, direct competitors or baseline methods.

Realism Assessment

We report here the results of the comparison of the simulated (factual) toxicity with the real (PANDORA) toxicity.

Toxic behavior is believable and consistent. Aggregating factual data from all simulations, we find significant Spearman correlations (p -value <0.01) between toxicity and some Big Five traits, which mirror correlations also found in real data (PANDORA). Specifically: *agreeableness* (real $\rho=-0.18$, simulated $\rho=-0.32$), *conscientiousness* (real $\rho=-0.11$, simulated $\rho=-0.16$) and *neuroticism* (real $\rho=0.04$, simulated $\rho=0.07$). These measurements align with field observation in psychological literature (Kordyaka et al. 2023), where prototypical toxic users feature low empathy and collaboration (low agreeableness), a lack of self-discipline (low conscientiousness) and a prevalence of negative emotions (high neuroticism). Behavioral preferences distinguishing each agent are consistent through time, as revealed by the standard deviations of their toxicity distributions (real avg. $\sigma=0.17$, simulated avg. $\sigma=0.20$).

Toxicity propagates across threads. Aggregating factual data from all simulations, we find a significant Spearman correlation between the toxicity of parent and children nodes ($\rho=+0.39$, p -value $=0.0$), proving that toxic behavior also emerges from the agents’ capability to mutually influence each other. Toxicity contagion is further supported by the results of the sub-population simulation involving the top-5

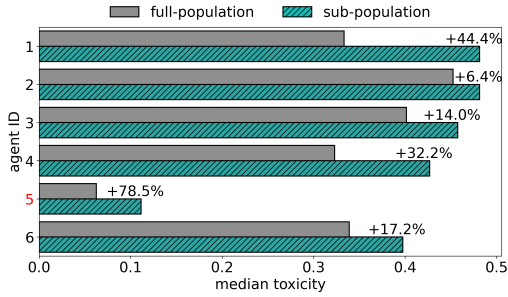


Figure 2: Median toxicity of agents in the sub-population simulation, compared to their median toxicity in full-population simulations. Least toxic agent marked in red.

MS	Simulation ID					CLR	
	1	2	3	4	5		
OSFA	-0.06*	+0.04	-0.05	-0.06	+0.05	0.00	
PMI	<i>N</i>	-0.09***	+0.00	-0.08*	-0.11*	0.00	
	<i>E</i>	-0.10***	+0.06	-0.06	-0.02	-0.05	0.00
	<i>P</i>	-0.05*	+0.06	-0.04	-0.03	-0.02	0.00
BAN	1	-0.54***	-0.45***	-0.58***	-0.57***	-0.52**	0.45
	2	-0.40***	-0.29**	-0.46***	-0.47***	-0.38**	0.32
	4	-0.25**	-0.17*	-0.21**	-0.28**	-0.23**	0.16
	8	-0.06	-0.04	-0.06	-0.08	-0.07	0.04

Table 1: Mass divergence ΔM s for each Moderation Strategy (MS) and simulation run. Asterisks denote significant reductions ($\Delta M < 0$) or increases ($\Delta M > 0$) based on Mann–Whitney: * $p < 0.1$; ** $p < 0.05$; *** $p < 0.01$. The last column reports the average Content Loss Ratio (CLR).

toxic agents and the least toxic agent. Figure 2 shows the median toxicity of agents in the sub-population simulation, compared to their median toxicity in full-population simulations. We notice how the sub-population encourages agents (anomalous included) to significantly increase their toxicity compared to their behavior in full-population experiments.

Moderation Assessment

We report here the results of the comparison of the simulated factual and counterfactual toxicity. Further details on moderation outcomes can be found in Appendix D.

Personalized moderation is more effective. Table 1 reports mass divergences ΔM for each Moderation Strategy (MS) and simulation run. As evidenced, PMI-*N* brings significant reductions in most runs (1, 3, 4, 5). This performance is not paralleled by other *ex ante* strategies, particularly OSFA and PMI-*P*, yielding (on average) lower reductions. Arguably, this result is consistent with the expected benefits of personalization, as PMI-*N* provides maximum freedom to the moderation action. Figure 3 shows how *ex ante* messages encoded with the average of their BERT embeddings (Devlin et al. 2019) are distributed within the semantic space represented via t-SNE. We observe that PMI-*N* explores wider regions, potentially adapting its commu-

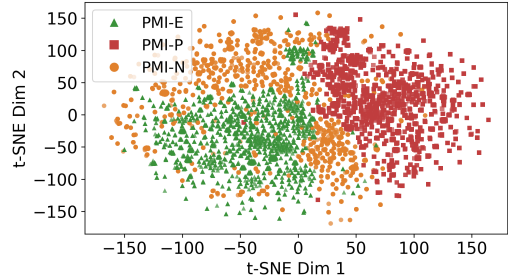


Figure 3: *Ex ante* PMI messages encoded with BERT.

nication style to the needs of each moderation scenario.

Low tolerance yields deplatforming effects. As reported in Table 1, BAN-*e* delivers, proportionally to *e*, considerable negative ΔM . However, differently from *ex ante* strategies, these reductions are achieved by removing agents, i.e., by losing nodes from the counterfactual feed $\hat{\mathcal{F}}$, rather than by redirecting their toxicity. As shown in Table 1, although at $e=1$ we observe the greatest reduction of toxicity, this comes at the cost of a CLR of 0.45 ± 0.04 , which includes a fraction of 0.27 ± 0.03 of “healthy” contents, i.e., with toxicity below *THR*. In other words, if the ban strategy must predict whether a text is worth losing, i.e., it has toxicity greater than *THR*, at $e=1$ the macro-average recall resembles a random classifier (0.55), compared to 0.60 at $e=2$ and 0.58 at $e=4$.

Moderation is sensitive to psychological traits. By aggregating data from all simulations, we compute mass divergence ΔM on subsets of agents sharing the same psychological trait. Figure 4 reports these measurements for each moderation strategy and each OCEAN trait for the different intensity values from 1 to 5. We mark statistically significant reductions (if $\Delta M < 0$) or increases (if $\Delta M > 0$) with an asterisk, based on p -value < 0.05 (Mann–Whitney). We observe that all moderation strategies follow similar trends, suggesting comparable effects on similar personality types. However, only PMIs and BAN-*e* with $e \leq 4$ bring significant divergences. Notably, moderation successfully targets prototypical toxic agents, i.e., those characterized by low agreeableness, high neuroticism and low conscientiousness.

Moderation mostly affects extreme toxicity. In Figure 5 we report quantile divergences Δq averaged across simulation runs, for each moderation strategy and for $q \in [0, 1]$. We observe that moderation mostly affect extreme toxic behavior ($q \geq 0.8$), with varying effects on lower ranges of toxicity ($0.6 \leq q \leq 0.8$). Notably, BAN-*e* with $e \leq 4$ displays a bimodal trend, with significant reductions also for milder toxic behavior. These results are consistent with the observation that moderation successfully targets the agents most contributing to the overall toxicity mass.

Conclusions and Limitations

In this paper we have introduced COSMOS, a LLM-powered ABM simulator for evaluating content moderation strategies in OSN conversations. COSMOS implements OSN agents

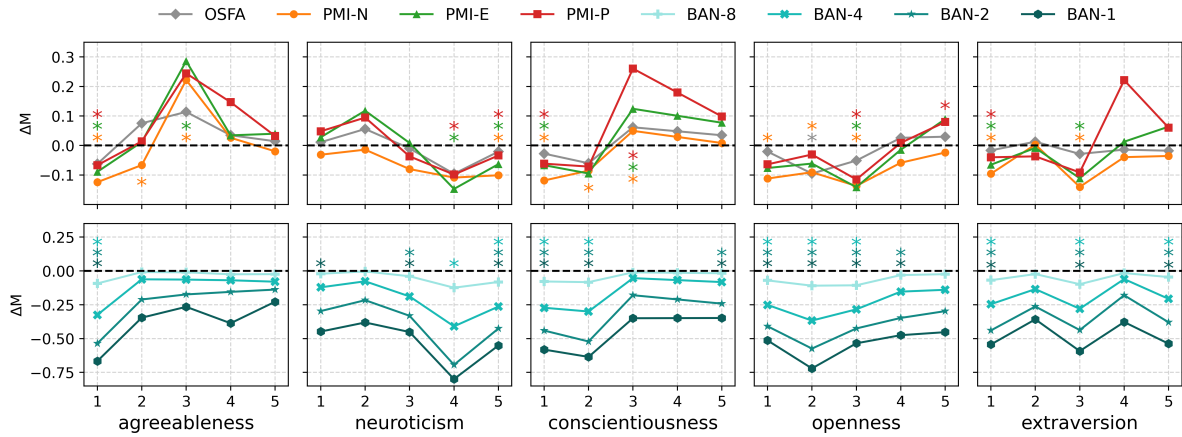


Figure 4: Mass divergence ΔM over each OCEAN trait for different intensity values across moderation strategies. Statistically significant reductions ($\Delta M < 0$) or increases ($\Delta M > 0$) are marked with an asterisk for Mann-Whitney with p -value < 0.05 .

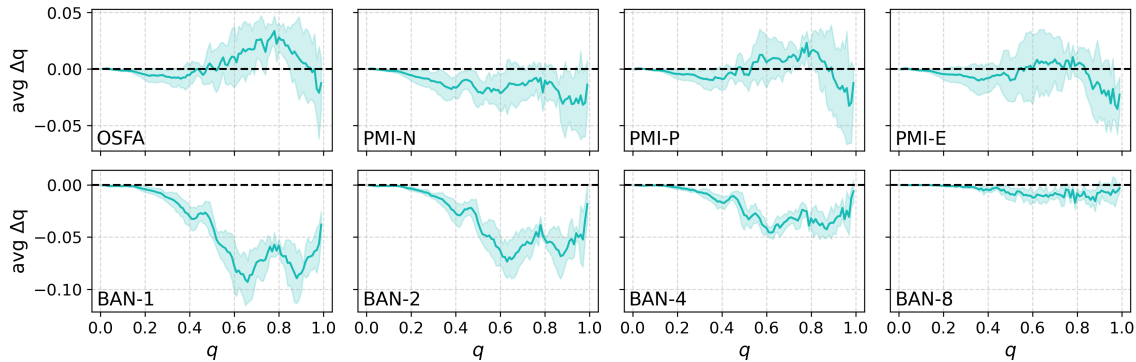


Figure 5: Quantile divergence Δq (y-axis) computed on $q \in [0.0, 1.0]$ (x-axis) and averaged across simulation runs, for each moderation strategy. The error band represents standard deviations.

with believable and consistent psychological attitudes, as well as capable of mutual influence. By running parallel, *counterfactual* simulations where moderation is applied *ceteris paribus*, COSMOS has generated evidence supporting the superior effectiveness of personalized *ex ante* interventions, the deplatforming effects of low-tolerance ban and the influence of psychological traits on moderation outcomes.

COSMOS does have some limitations. First, LLMs can fail, especially when processing complex prompts (Heo et al. 2024). We have estimated COSMOS’s hallucinations by applying 2-means clustering on BERT embeddings of generated posts and comments. Upon inspection, we found a cluster of redundant hallucinations, accounting for about 7% of the data. Overcoming this issue will require more LLM tuning and reliable content validation. More broadly, we acknowledge the need for *subjective* realism evaluation. While prior studies show that LLMs can faithfully reproduce human behavior from psychological data (Jiang et al. 2024; Tseng et al. 2024), COSMOS would further benefit from human-based assessments about the alignment between simulated and real-world moderation responses. Achieving this,

however, requires cooperation with domain experts, such as psychologists or sociologists. Second, COSMOS is currently a simulator of OSN *conversations*, and, as such, has not been designed to replicate all OSN dynamics. This design improves control and efficiency by reducing variability and avoiding extra LLM calls. Nonetheless, we plan to extend COSMOS with followings and reactions (likes), enabling (i) the simulation of more advanced recommendation systems based on social connections and agents’ preferences; hence, (ii) the emergence of phenomena like homophily and polarization, potentially influencing moderation outcomes. Third, LLM-based simulations are costly and scaling to real-sized OSN populations is challenging. Finally, LLMs are known to amplify societal biases. This issue is left for future works, as it is largely attributed to the LLM’s training data (Echterhoff et al. 2024). A related emerging concern is the tendency of LLMs toward self-preference (Panickssery et al. 2024; Wataoka et al. 2024). However, while COSMOS’s LLM might be biased by its own inputs, i.e., OSN conversations and moderation messages, it remains unclear to what extent such bias applies to role-playing settings.

Ethical Statement

To ensure realism, COSMOS experiments incorporate psychological and demographic information from real sources. In line with the ethics code of psychological research (Gjurkovic et al. 2021), no sensitive data is disclosed, and all resulting profiles are entirely fictional. We encourage mindful uses of COSMOS in industrial contexts: automated moderation is still in its infancy and the full replacement of human moderators remains controversial (Gillespie 2020).

Acknowledgements

This work has been partially supported by the Italian Project Fondo Italiano per la Scienza FIS00001966 “MI-MOSA”, by the PRIN 2022 framework project PIANO under CUP B53D23013290006, by the European Community Horizon 2020 programme under the funding schemes G.A. 101120763 “TANGO”, by the European Innovation Council project “EMERGE” (Grant No. 101070918), by the European Commission under the NextGeneration EU programme – National Recovery and Resilience Plan (Piano Nazionale di Ripresa e Resilienza, PNRR) Project: “SoBig-Data.it – Strengthening the Italian RI for Social Mining and Big Data Analytics” – Prot. IR0000013 – Av. n. 3264 del 28/12/2021, and M4C2 - Investimento 1.3, Partenariato Esteso PE00000013 - “FAIR” - Future Artificial Intelligence Research” - Spoke 1 “Human-centered AI”.

References

- Aher, G. V.; et al. 2023. Using Large Language Models to Simulate Multiple Humans and Replicate Human Subject Studies. In *ICML*, volume 202 of *Proceedings of Machine Learning Research*, 337–371. PMLR.
- Avalle, M.; et al. 2024. Persistent interaction patterns across social media platforms and over time. *Nature*, 628(8008): 582–589.
- Axtell, R. L.; and Farmer, J. D. 2025. Agent-Based Modeling in Economics and Finance: Past, Present, and Future. *Journal of Economic Literature*, 63(1): 197–287.
- Beck, T.; et al. 2023. How (Not) to Use Sociodemographic Information for Subjective NLP Tasks. *CoRR*, abs/2309.07034.
- Bilewicz, M.; et al. 2021. Artificial Intelligence Against Hate: Intervention Reducing Verbal Aggression in the Social Network Environment. *Aggressive Behavior*, 47.
- Breum, S. M.; et al. 2024. The Persuasive Power of Large Language Models. In *ICWSM*, 152–163. AAAI Press.
- Brown, T. B.; et al. 2020. Language Models are Few-Shot Learners. In *NeurIPS*.
- Cava, L. L.; and Tagarelli, A. 2024. Open Models, Closed Minds? On Agents Capabilities in Mimicking Human Personalities through Open Large Language Models. arXiv:2401.07115.
- Chen, C.; et al. 2024. Evaluating LLM Agents for Simulating Humanoid Behavior. In *Proceedings of the 1st HEAL Workshop at the CHI Conference on Human Factors in Computing Systems*. Honolulu, HI, USA.
- Chuang, Y.; et al. 2024. Simulating Opinion Dynamics with Networks of LLM-based Agents. In *NAACL-HLT (Findings)*, 3326–3346. Association for Computational Linguistics.
- Conte, R.; et al. 2012. Manifesto of Computational Social Science. *The European Physical Journal Special Topics*, 214(1): 325–346.
- Cornacchia, G.; et al. 2020. Modelling Human Mobility considering Spatial, Temporal and Social Dimensions. *CoRR*, abs/2007.02371.
- Cresci, S.; et al. 2022. Personalized Interventions for Online Moderation. In *HT*, 248–251. ACM.
- Davern, M.; et al. 2025. General Social Survey 1972–2024.
- Devlin, J.; et al. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *NAACL-HLT (1)*, 4171–4186. Association for Computational Linguistics.
- Echterhoff, J. M.; et al. 2024. Cognitive Bias in Decision-Making with LLMs. In *EMNLP (Findings)*, 12640–12653. Association for Computational Linguistics.
- Fan, R.; et al. 2017. An agent-based model for emotion contagion and competition in online social media. *CoRR*, abs/1706.02676.
- Gao, C.; et al. 2023a. Large Language Models Empowered Agent-based Modeling and Simulation: A Survey and Perspectives. *CoRR*, abs/2312.11970.
- Gao, C.; et al. 2023b. S³: Social-network Simulation System with Large Language Model-Empowered Agents. *CoRR*, abs/2307.14984.
- Gillespie, T. 2020. Content moderation, AI, and the question of scale. *Big Data & Society*, 7: 205395172094323.
- Gjurkovic, M.; et al. 2021. PANDORA Talks: Personality and Demographics on Reddit. In *SocialNLP@NAACL*, 138–152. Association for Computational Linguistics.
- Goldberg, L. R. 2013. An alternative ‘description of personality’: The big-five factor structure, 34–47. Routledge.
- Grimmelmann, J. 2017. The virtues of moderation. <https://doi.org/10.31228/osf.io/qwxf5>.
- Guo, T.; et al. 2024. Large Language Model Based Multi-agents: A Survey of Progress and Challenges. In *IJCAI*, 8048–8057. ijcai.org.
- Hangartner, D.; et al. 2021. Empathy-based counterspeech can reduce racist hate speech in a social media field experiment. *Proceedings of the National Academy of Sciences*, 118: e2116310118.
- Hanscom, R.; et al. 2024. The Toxicity Phenomenon Across Social Media. *CoRR*, abs/2410.21589.
- Heo, J.; et al. 2024. Do LLMs “know” internally when they follow instructions? *CoRR*, abs/2410.14516.
- Jiang, H.; et al. 2024. PersonaLLM: Investigating the Ability of Large Language Models to Express Personality Traits. In *NAACL-HLT (Findings)*, 3605–3627. Association for Computational Linguistics.

- Kaiya, Z.; et al. 2023. Lyfe Agents: Generative agents for low-cost real-time social interactions. *CoRR*, abs/2310.02172.
- Kim, S.; et al. 2024. SOLAR 10.7B: Scaling Large Language Models with Simple yet Effective Depth Up-Scaling. In *NAACL (Industry Track)*, 23–35. Association for Computational Linguistics.
- Kordyaka, B.; et al. 2023. The Cycle of Toxicity: Exploring Relationships between Personality and Player Roles in Toxic Behavior in Multiplayer Online Battle Arena Games. *Proc. ACM Hum. Comput. Interact.*, 7(CHI PLAY): 611–641.
- Lapidot-Leffler, N.; and Barak, A. 2012. Effects of anonymity, invisibility, and lack of eye-contact on toxic on-line disinhibition. *Comput. Hum. Behav.*, 28(2): 434–443.
- Lees, A.; et al. 2022. A New Generation of Perspective API: Efficient Multilingual Character-level Transformers. In *KDD*, 3197–3207. ACM.
- Leng, Y.; and Yuan, Y. 2023. Do LLM Agents Exhibit Social Behavior? *CoRR*, abs/2312.15198.
- Liu, F. T.; et al. 2008. Isolation Forest. In *ICDM*, 413–422. IEEE Computer Society.
- Lorig, F.; et al. 2021. Agent-Based Social Simulation of the Covid-19 Pandemic: A Systematic Review. *J. Artif. Soc. Soc. Simul.*, 24(3).
- Maleki, M.; et al. 2022. Applying an Epidemiological Model to Evaluate the Propagation of Toxicity related to COVID-19 on Twitter. In *HICSS*, 1–10. ScholarSpace.
- Mann, H. B.; and Whitney, D. R. 1947. On a test of whether one of two random variables is stochastically larger than the other. *Annals of Mathematical Statistics*, 18: 50–60.
- Marzo, G. D.; et al. 2023. Emergence of Scale-Free Networks in Social Interactions among Large Language Models. *CoRR*, abs/2312.06619.
- McDonald, G. W.; and Osgood, N. D. 2023. *Agent-Based Modeling and Its Trade-Offs: An Introduction and Examples*, 209–242. Cham: Springer International Publishing.
- Mou, X.; et al. 2024. Unveiling the Truth and Facilitating Change: Towards Agent-based Large-scale Social Movement Simulation. In *ACL (Findings)*, 4789–4809. Association for Computational Linguistics.
- Murdock, I.; et al. 2024. An agent-based model of cross-platform information diffusion and moderation. *Soc. Netw. Anal. Min.*, 14(1): 145.
- Nixon, C. L. 2014. Current perspectives: the impact of cyberbullying on adolescent health. *Adolescent Health, Medicine and Therapeutics*, 5: 143–158.
- Obadimu, A.; et al. 2020. Developing an Epidemiological Model to Study Spread of Toxicity on YouTube. In *SBP-BRiMS*, volume 12268 of *Lecture Notes in Computer Science*, 266–276. Springer.
- Panickssery, A.; et al. 2024. LLM Evaluators Recognize and Favor Their Own Generations. In *NeurIPS*.
- Park, J. S.; et al. 2022. Social Simulacra: Creating Populated Prototypes for Social Computing Systems. In *UIST*, 74:1–74:18. ACM.
- Park, J. S.; et al. 2023. Generative Agents: Interactive Simulacra of Human Behavior. In *UIST*, 2:1–2:22. ACM.
- Piatti, G.; et al. 2024. Cooperate or Collapse: Emergence of Sustainable Cooperation in a Society of LLM Agents. arXiv:2404.16698.
- Rossetti, G.; et al. 2024. Y Social: an LLM-powered Social Media Digital Twin. *CoRR*, abs/2408.00818.
- Schneider, P. J.; and Rizoiu, M.-A. 2023. The effectiveness of moderating harmful online content. *Proceedings of the National Academy of Sciences*, 120(34): e2307360120.
- Shao, Y.; et al. 2023. Character-LLM: A Trainable Agent for Role-Playing. In *EMNLP*, 13153–13187. Association for Computational Linguistics.
- Squazzoni, F.; et al. 2014. Social Simulation in the Social Sciences: A Brief Overview. *Social Science Computer Review*, 32(3): 279–294.
- Taillandier, P.; et al. 2025. Integrating LLM in Agent-Based Social Simulation: Opportunities and Challenges. *arXiv preprint arXiv:2507.19364*.
- Törnberg, P.; et al. 2023. Simulating Social Media Using Large Language Models to Evaluate Alternative News Feed Algorithms. *CoRR*, abs/2310.05984.
- Tseng, Y.; et al. 2024. Two Tales of Persona in LLMs: A Survey of Role-Playing and Personalization. In *EMNLP (Findings)*, 16612–16631. Association for Computational Linguistics.
- Tyagi, A.; et al. 2020. Affective Polarization in Online Climate Change Discourse on Twitter. In *ASONAM*, 443–447. IEEE.
- Wang, L.; et al. 2024a. A survey on large language model based autonomous agents. *Frontiers Comput. Sci.*, 18(6): 186345.
- Wang, L.; et al. 2024b. User Behavior Simulation with Large Language Model based Agents. arXiv:2306.02552.
- Wang, N.; et al. 2024c. RoleLLM: Benchmarking, Eliciting, and Enhancing Role-Playing Abilities of Large Language Models. In *ACL (Findings)*, 14743–14777. Association for Computational Linguistics.
- Wataoka, K.; et al. 2024. Self-Preference Bias in LLM-as-a-Judge. *CoRR*, abs/2410.21819.
- Xi, Z.; et al. 2023. The Rise and Potential of Large Language Model Based Agents: A Survey. *CoRR*, abs/2309.07864.
- Xie, C.; et al. 2024. Can Large Language Model Agents Simulate Human Trust Behaviors? *CoRR*, abs/2402.04559.
- Zhang, H.; et al. 2020. Trading Off Diversity and Quality in Natural Language Generation. *CoRR*, abs/2004.10450.