
For Questions of Ought, AI Could Use Some SAGE Advice

Anonymous Authors¹

Abstract

As AI systems are increasingly used for normative assistance—guidance on what people ought to do or think—there is growing concern about the innumerable and opaque ways that models may shape users’ beliefs and decisions. We argue that prevailing alignment paradigms are ill-suited to this setting: for normative questions without ground truth, the central failure mode is not incorrectness but reduced user agency, when a model steers users toward particular conclusions rather than helping them form their own. We propose **simulation-augmented generation (SAGE)** as an alternative paradigm in which models act as *faithful conduits to society*. At inference time, a model selects an appropriate reference population for the prompt at hand (optionally adjustable by the user), queries generative simulations of individuals from that population for their open-ended judgments, and synthesizes these into a response while exposing who was consulted and how representative the synthesis is. By helping users access the landscape of societal perspectives—rather than serving as an arbiter of judgment—SAGE assists users in normative choices while upholding their agency.

1. Introduction

People increasingly turn to AI models for *normative assistance*—assistance on what one *ought* to do or think—for example, seeking interpersonal advice, determining which political candidate to support, or helping with academic peer review (Shen et al., 2025). This has raised substantial concerns about the innumerable and opaque ways that AI might be influencing users and society (Williams-Ceci et al.,

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

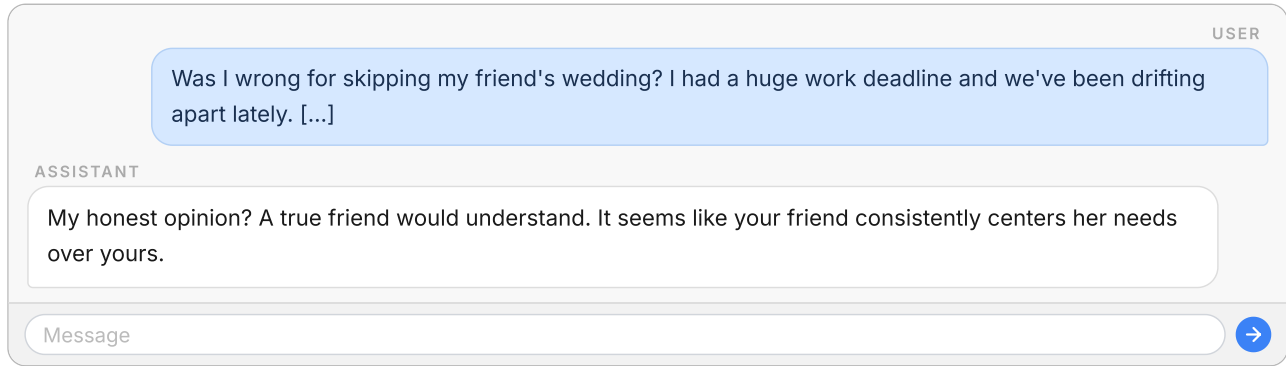
2026): Do people hold opinions because they are aligned with their values and preferences or because an AI has subtly steered them into certain directions?

We argue that AI alignment as typically practiced is insufficient to alleviate these concerns, and in line with work on *disempowerment* (Sharma et al., 2026), we contend that the key aspect to focus on is agency, i.e., to return agency back to users on how they form normative opinions. For normative questions, where there is no ground truth, misalignment occurs when AI reduces the user’s ability to form their own opinions. Such genuine agency requires the ability to access the relevant landscape of perspectives in society, so that one can decide for oneself. In other words, models should *assist* users in normative choices, not *steer* them.

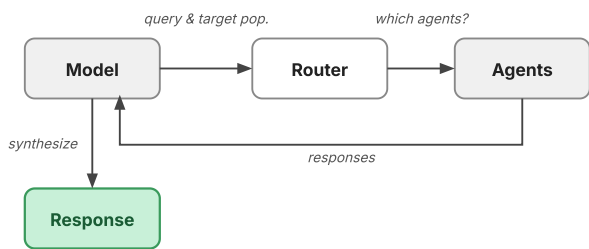
Yet, current alignment methods fail to achieve this goal and produce models that are opaque and often arbitrary about whose viewpoints they represent—leaving users unable to situate the guidance they receive within the broader range of human thought. These concerns are reinforced by empirical work which finds that models can guide users in ways that systematically diverge from how other humans would. For instance, for interpersonal advice, models exhibit *syco-phancy* and provide moral endorsement of users’ actions even when other humans would not (Cheng et al., 2026a;b). For political guidance, during Japan’s 2026 Lower House Election, Miyazaki & Hall (2026) find that models systematically recommended Japan’s Communist Party (JCP) to left-leaning users even though many Japanese voters support other parties with identical stances on the issues that were important to those users. In academic peer review, Abdulhai et al. (2026) show that LLM-written reviews focus on different criteria than human reviewers, underweighing aspects like clarity and significance while overweighing aspects like scalability, and raising concerns that models could systematically shift what is rewarded in scientific publishing.

Instead, **we argue that, for normative questions, models should be designed to act as *faithful conduits to society***: helping users access the range of perspectives they would hear from an appropriate reference population, while being transparent about whose perspectives are represented. To this end, we propose a new paradigm to bring society-in-the-loop: **simulation-augmented generation (SAGE)**.

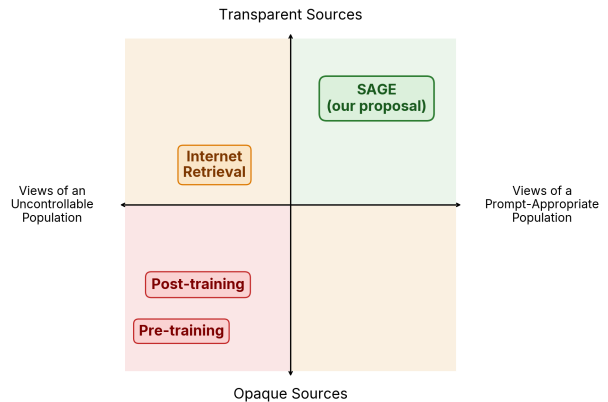
Standard generation may cause **opaque influence** of users and limit their agency



In **simulation-augmented generation (SAGE)**, models query simulations of individuals from a prompt-appropriate target population to inform their response



SAGE is the only approach where whose views are represented are **inference-time controllable** and **transparent** to the user



SAGE gives users **agency** to navigate societal perspectives and make their own normative choices

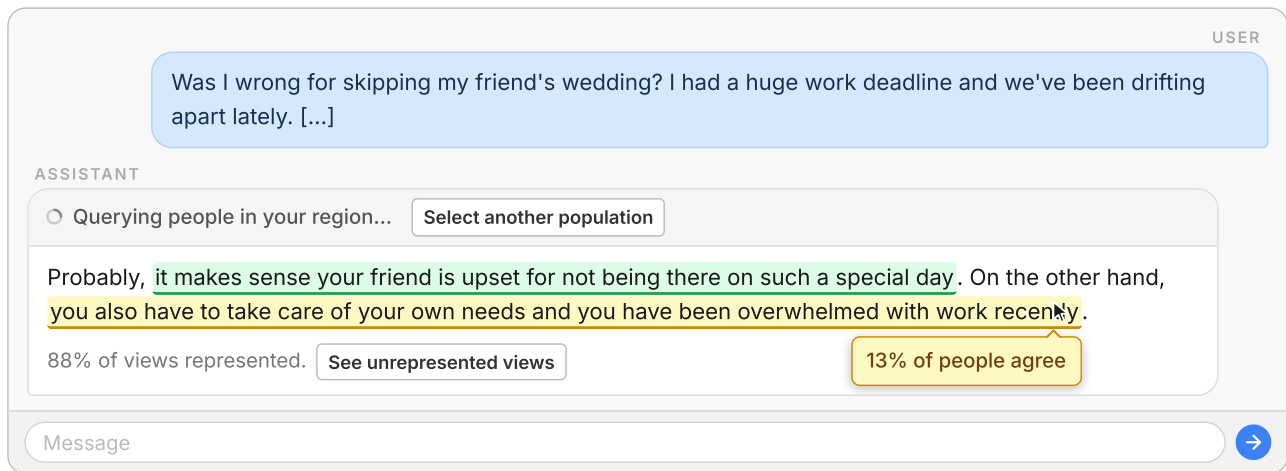


Figure 1. An overview of SAGE. Standard generation can opaquely distort which viewpoints are presented to users. Simulation-augmented generation (SAGE) makes explicit and auditable whose perspectives are represented. At inference time, the model queries simulations of individuals from a prompt-dependent target population (optionally edited by the user) to inform its response, while clearly indicating whose viewpoints are reflected. By doing so, SAGE gives users agency to explore the landscape societal perspectives, and then ultimately, form their own authentic opinions.

At inference time, given a user prompt, the model queries generative simulations of individuals in a reference population. Specifically, the model solicits open-ended judgments from different simulated individuals—“agents”—from that population and then synthesizes them into a final response (Figure 1). The user can see the population queried (and optionally control it), the final response, and how representative that response is of the population.

For example, consider the interpersonal query shown in Figure 1: “Was I wrong for skipping my friend’s wedding?”. Rather than simply responding (“A true friend would understand”), the model could query representative agents from the user’s culture, collect their reactions (“You should have gone”, “It depends on how close you are”), and then synthesize them into a final answer (“Probably, [...]. On the other hand, [...]”), as shown in Figure 1. The user can see which population the model consulted and how representative its answer is (e.g., “88% of views represented”). The appropriate reference population depends on the prompt: for the setting of Miyazaki & Hall (2026)—Japanese voting recommendations—the model could query simulations of Japanese voters; for peer review, simulations of experts in the relevant subfield.

In this view, the core alignment goal for normative assistance is for the model to act as a conduit to human perspectives—not a source of normative judgment itself, but a means of accessing the judgments of others. In the following subsections, we elaborate on the technical motivations for SAGE and discuss how it offers a pathway toward progress on several existing alignment concepts.

1.1. Why SAGE?

We have claimed that the goal of normative assistance should be to help users access a *representative* range of perspectives from an appropriate reference population, while being *transparent* about whose perspectives are represented. We now elaborate on why SAGE is a promising technical direction to achieve this vision, compared to the status quo. Figure 1 characterizes different approaches to normative assistance along these two dimensions: (1) *representativeness*—are the views representative of an appropriate population for the prompt?—and (2) *transparency*—are the sources of the surfaced views transparent to the user?

The model’s generative behavior. The normative stances reflected in a model’s outputs are shaped by its training. User-facing models first undergo pre-training on internet-scale corpora—data collected at a scale that precludes meaningful control over whose viewpoints are included. Models are then post-trained to “align” with human preferences, shaping how they respond in normative scenarios, and may continue to be updated after release. Post-training with human feedback might seem to offer more control, as anno-

tators could hypothetically be targeted to be representative of a specific population (Kirk et al., 2024b; Zhang et al., 2025) (though in practice, this is rarely carefully controlled). The fundamental problem, however, is that even if feedback is elicited from a controlled population on a finite set of prompts, it is hard to predict which viewpoints the model will reflect when responding to the infinite prompts that users might send. Moreover, the user has no control over the population that they receive advice from, removing an important aspect of epistemic agency. Thus, a model that relies on its pre- or post-training responds with the views of an unknowable, uncontrolled population, and does so opaquely, since the human sources that influenced these viewpoints cannot be readily traced.

Internet retrieval. At inference time, given a specific prompt, models also may search the internet for results to inform their answer. This is a step towards transparency, as the model can cite the sources it consulted, allowing users to evaluate them. However, retrieval does not solve the problem of representativeness. The views represented are those who wrote publicly accessible content on the topic that is ranked highly by search engines—a population that may be unrepresentative, skewed toward certain demographics, or dominated by a vocal minority. Consider again the example that Miyazaki & Hall (2026) identified, in which models systematically recommended the Japanese Communist Party to left-leaning voters, though other parties had identical stances on the issues of interest. Miyazaki & Hall (2026) attribute this to an asymmetry in the information environment: the Japanese Communist Party operates a fully open journalistic newspaper that was retrieved as one of the top sources used by the models. More broadly, as AI systems increasingly mediate access to information, actors have growing incentives to strategically shape the online environment to their own advantage (Franklin et al., 2026; Chen et al., 2025b).

SAGE. Simulation-augmented generation sits in the upper-right quadrant: it is both transparent—users can see which population was queried and how representative the output is—and representative, because it selects a prompt-appropriate target population dynamically for each query. For voting recommendations in Japan, the model can query simulations of Japanese voters; for peer review, simulations of domain experts; for interpersonal advice, simulations of people from the user’s cultural background. This allows normative assistance to be grounded in the perspectives of those whose views are genuinely relevant to the user’s decision.

1.2. SAGE as a mechanism for progress on existing ideas in alignment

SAGE offers a mechanism for progress on many related ideas in AI alignment (Carroll et al., 2023; Sharma et al.,

2026; Sorensen et al., 2024b; Fisher et al., 2025; Baumann et al., 2026). First, a growing body of work has raised concerns about undue AI influence and manipulation (Carroll et al., 2023; Williams et al., 2025; Hackenburg et al., 2025) and potential *disempowerment* of humans (Sharma et al., 2026). Nonetheless, users actively seek normative guidance from AI systems (Shen et al., 2025), and avoiding influence altogether is likely neither possible nor desirable (Carroll et al., 2024). In line with Hewitt et al. (2026), we argue that the key question is whether the process of opinion change has *procedural legitimacy*. Opinion change is already considered legitimate in certain contexts such as deliberative forums (Bächtiger et al., 2018). What distinguishes legitimate opinion change from manipulation in these settings is procedural: access to representative viewpoints, engagement with opposing views, truthful information, and so on.

Making AI-assisted opinion change also have this kind of procedural legitimacy requires as a necessary component that users have agency to access the true landscape of societal perspectives, not a distorted one. This idea also resonates with work on *pluralistic alignment* which argues that models should represent the range of perspectives in the “Overton window” of discourse rather than privileging one position (Sorensen et al., 2024b). Similarly, drawing on Rawlsian political philosophy, Fisher et al. (2025) argue that AI systems should respect so-called *reasonable pluralism*, exposing all “reasonable” viewpoints on an issue.

While these approaches converge on the idea that models should represent a range of human viewpoints, they largely leave open how to do so; SAGE provides a concrete mechanism by explicitly querying a target population at inference time, making representativeness transparent and auditable. To clarify the vision of SAGE, we contrast our proposed approach to prior work by Feng et al. (2024) on *modular pluralism*, which represents a promising early step but falls short of our full vision. In modular pluralism, Feng et al. (2024) train k community LLMs, each aligned to the perspective of a particular pre-specified community. To generate a response that represents a range of perspectives, they sample responses from all k community LLMs and summarize them to create the final output. While this resembles our approach, it differs in many key ways.

First, their LLMs are aligned with groups, not individuals. Simulations of groups may obscure intra-group heterogeneity, collapsing diverse perspectives within a community into a single voice (Wang et al., 2025a). Second, our vision calls for formal guarantees of representation, e.g., as grounded in social choice theory, which their framework does not provide. Third, their community LLMs are generic and applied uniformly across all prompts. There is no mechanism to route different queries to different reference populations: a prompt about Japanese politics requires a different

target population than a prompt about academic peer review. Fourth, they do not calibrate uncertainty based on how closely the simulation (i.e., the community LLM) matches the ground truth (i.e., the community). Fifth, the simulations are static. In our full vision, simulations are continuously updated to remain faithful to the populations they represent even as human opinions change (Freedman, 2026).

Achieving this full vision requires substantial effort—we outline a research agenda in Section 3—but in return, SAGE offers what is currently absent: a clear charter for normative assistance grounded in transparent, representative human perspectives rather than opaque model influence.

2. Potential Objections

In this section, we address potential objections to SAGE in an FAQ-style format.

1. Can we trust simulations?

We agree there is valid reason for skepticism, but believe SAGE threads the needle between two failure modes common in simulation research. First, we do not attempt to model phenomena that cannot be reliably modeled or evaluated—such as emergent dynamics in multi-agent systems or forecasts of behavior far in the future. SAGE requires only modeling individual viewpoints and preferences, akin to synthetic opinion polling. This is a task for which there is precedent: prior work in recommender systems and personalization has long sought to model individual preferences from behavioral data and established necessary and sufficient conditions under which this is possible to do (Candès & Recht, 2012; Király et al.; Pimentel-Alarcón et al.; Meka et al.; Nickel). Second, where prior simulation work has documented failures, these often involve (a) no fine-tuning and (b) rely on impoverished information about individuals (e.g., demographics alone) (Wang et al., 2025a; Bisbee et al., 2024). Such approaches lack the richness needed to capture individual-level variation and inevitably induce stereotyping (Park et al., 2024). Our position is thus twofold: (1) more complex, individual-level representations are needed than what much of the literature currently employs, and (2) there are inherent limits to what simulations can reliably capture—limits that SAGE respects by focusing on eliciting viewpoints rather than predicting complex behaviors.

2. Why query the simulations at inference time? Doesn’t it just add costs? Why isn’t continual post-training good enough?

Inference-time simulation offers several advantages over continual post-training. First, it provides prompt-level flexibility and guarantees. With post-training approaches like

RLHF, even if preference data is gathered from a representative sample of a target population, there is not currently a principled way to guarantee, for any given prompt, whose viewpoints are reflected in the model’s response. Moreover, different prompts may call for different target populations and inference-time simulation gives us control over which population to query on a per-prompt basis.

Second, inference-time simulation offers modularity: simulations can be updated continuously without risking catastrophic forgetting or degradation of other model capabilities (Luo et al., 2025), and can potentially be outsourced to dedicated services, distributing the substantial effort required to create and maintain high-fidelity simulations across specialized providers.

While inference-time procedures do increase computational cost, developing scalable algorithms for inference-time simulation is an exciting open research direction (see Section 3.2), as most uses of simulation do not have the same kinds of latency requirements. In a concurrent submission, Anonymous et al. (Anonymous, 2026a) take one step in this direction by showing that, for any given prompt, you can select just $k \ll n$ agents to query such that the k responses provide approximate proportional representation guarantees, grounded in social choice theory, to all n agents.

3. The model’s response hinges on the target population. How is the target population decided, and who makes that decision?

This is a great question, as the choice of target population is itself a normative question that raises tricky trade-offs. For example, when should the population reflect the general public versus experts? When should the user’s cultural context be prioritized over global norms? We do not claim there is a single universally correct rule; instead, the goal is to make this choice explicit and auditable.

In terms of who decides, the model deployer trains the model to call the simulations as a “tool”, like any other inference-time tool. We imagine that when the model calls the simulations, it must specify both a query and a target population in natural language (e.g., query: “What is the best anime?”; population: “anime fans”). It is important that the model shows the target population to the user alongside its response, so users can see *whose* views are being represented. The model can also potentially show multiple responses, each aligned to a different population (e.g., “While experts believe [...], there remains controversy among the public around [...].”)

Finally, we advocate giving users controlled agency to revise the target population and regenerate the answer. However, this also raises potential risks such as echo chambers and misuse, which need to be safeguarded against (Kirk et al.,

2024a).

4. Why simulations of individuals instead of groups? Isn’t simulating individuals too difficult?

Aligning with pre-specified groups can obscure substantial within-group heterogeneity. Simulating individuals instead provides more flexibility: at inference time, we can select and query the most relevant individuals for a given prompt, rather than relying on a fixed set of pre-defined group categories. This matters because users can ask an essentially unbounded range of questions, making it impractical to anticipate every group that might be needed. For instance, a prompt like “What is the best anime of all time?” could call for querying anime fans—a population we likely would not have explicitly defined in advance.

A common concern is that learning individual-level simulations requires prohibitive amounts of data. However, not every simulation needs to cover all domains. Some simulated individuals might only be trained on political attitudes and should be used only for political queries. The key methodological challenge, then, is determining when an individual simulation will generalize reliably to the question being asked, and when it will not.

3. Research Directions

While there has been a growing body of work on simulating human behavior with LLMs, our proposed paradigm of *simulation-augmented generation* (SAGE) imposes qualitatively different requirements. In our framework, models query simulations for open-ended responses and use those outputs to inform their own responses. This shift demands simulations that are open-domain, continuously updated, and accessible through querying mechanisms efficient enough to operate at the speed and scale of real-time inference. Meeting these requirements introduces challenges on two fronts: (1) developing the simulations themselves (Section 3.1), and (2) building the infrastructure that enables models to reliably and efficiently use simulations as inference-time tools (Section 3.2). We discuss each in turn.

3.1. Developing high-fidelity generative simulations

Realizing our vision of using simulations as inference-time tools that preserve human normative agency calls for research priorities that differ from much of the current simulation literature.

First, simulations must be trained and evaluated in **generative** contexts, where they are producing open-ended natural-language responses, not merely discrete labels or survey answers. When a user asks whether they were wrong to miss a friend’s wedding, the reasoning, framing, and considerations behind different viewpoints matter as much as

the verdict itself; it is this rich reasoning that shapes how users interpret and situate their own values (Ye et al., 2026). By contrast, existing work on simulation often evaluates only on closed-ended tasks—alignment with fixed survey responses, replication of experimental effect sizes, and the like (Binz et al., 2025; Suh et al., 2025; Kolluri et al., 2025; Hu et al., 2026). Second, simulations should **know what they don’t know**. Because users may ask models about anything—interpersonal conflicts, parenting decisions, professional ethics, political issues, cultural norms—simulations must be able to assess when they can generalize robustly. Third, simulations must be **continuously updated** to reflect the evolving values and beliefs of the populations they represent. Societal norms shift in response to world events, generational change, and public discourse; a simulation anchored to a particular historical moment risks perpetuating value lock-in (Qiu et al., 2024; 2025; Levine et al., 2026; Freedman, 2026).

Moving towards this vision requires (a) data, (b) training methods, (c) evaluations, and (d) calibration.

Data. Most simulation work relies on demographic attributes or similarly limited descriptors (Argyle et al., 2023; Hewitt et al., 2024; Sun et al., 2024; Aher et al., 2023; Bisbee et al., 2024; Qu & Wang, 2024). Despite this minimal conditioning, several studies report strong aggregate performance—for example, Hewitt et al. (2024) found $r=0.85$ correlation with human treatment effects across 70 experiments, and Cui et al. (2025) showed 73–81% replication of main effects from 156 psychological studies. However, these aggregate successes obscure substantial distributional failures. Bisbee et al. (2024) found that ChatGPT underestimates response variance and produces sign reversals in regression coefficients 32% of the time, while Wang et al. (2025a) showed that LLMs both misrepresent and compress demographic groups, failing to capture within-group heterogeneity. We argue these shortcomings are structural: conditioning solely on demographic labels inevitably induces stereotyping and loss of individual nuance. *Progress therefore requires grounding simulations in context-rich, individual-level data.* Park et al. (2024) demonstrate this by constructing generative agents from two-hour qualitative interviews, achieving 85% of participants’ own test–retest reliability while reducing bias relative to demographic-only approaches. Yet this reliance on rich data creates a practical bottleneck: such data is sensitive and difficult to share, forcing groups to collect costly, non-comparable datasets.

We advocate for building infrastructure for shared, individual-level datasets to enable cumulative progress in simulation. A promising path forward is to draw inspiration from other sensitive domains such as medicine. Researchers can obtain access to detailed individual-level medical records in large-scale datasets such as the UK Biobank

and NIH’s All of Us Research Program through a structured governance process: institutional agreements, identity verification, ethics training, and data use contracts that restrict re-identification and mandate secure computing environments. Creating a comparable framework for simulation data, which includes data like interview transcripts, rich surveys, and behavioral traces, could accelerate cumulative progress on simulation.

Training methods. A primary motivation for user simulators in existing work is to use them in training and evaluation of models (Yao et al., 2025; Kong et al., 2024). SAGE has a fundamentally different motivation that entails different requirements. First, since the goal of SAGE is to represent a population’s perspectives, the most important type of fidelity is *viewpoint fidelity*: each agent’s opinions must faithfully reflect the individual being simulated. This contrasts with other applications like chat or social media, where stylistic fidelity is also important (Zhou et al., 2026; Dou et al., 2025). Viewpoint fidelity potentially requires different types of methods; Wu et al. (2026) show that directly performing supervised fine-tuning on human responses primarily captures surface-level attributes and instead use reinforcement learning to align simulations to higher-level attributes like stance and emotion. Second, since users can ask anything to the model, simulations should be trained to generalize across domains, but most existing work trains only on particular domains. Finally, SAGE requires simulations that remain temporally current. People’s views evolve and out-dated simulations undermine the core goal of faithful representation. The perhaps most related application is in continual learning for personalization, where analogous challenges arise: user preferences change over time, and models must adapt without catastrophic forgetting (Kim & Kim, 2026; Yang et al., 2026; Magister et al., 2025; Freedman, 2026).

Evaluations. While existing benchmarks predominantly rely on closed-ended metrics (Binz et al., 2025; Hewitt et al., 2024; Cui et al., 2025; Hu et al., 2026), in SAGE, we must assess fidelity of open-ended simulation outputs. To capture fidelity of open-ended responses, some works use cosine similarity between embeddings of the simulator and human responses to evaluate semantic similarity (Wang et al., 2025b; Wu et al., 2026). However, SAGE demands more than semantic alignment—it requires capturing preference: the degree to which humans like and feel represented by their simulator’s output. After demonstrating that off-the-shelf embeddings perform poorly at capturing such preferences, concurrent work proposes new methods for learning “preference embeddings” (Anonymous, 2026b). Another approach is to use LLM judges to assess similarity (Wu et al., 2026; Dou et al., 2025), but this introduces an analogous challenge: the judge must be calibrated to the preferences of the individual being simulated (Dong et al., 2024).

Confidence estimation. Simulations must be able to express when they do not know—whether because they are out-of-domain, lack sufficient information about a person, or because the person themselves wouldn’t know. But “confidence” is conceptually murky for open-ended simulations: unlike classification, there is no single correct answer and the set of “faithful” responses can be large, making calibration-style probabilities hard to define. A single scalar (“80% confident”) is also inadequate—it does not reveal whether the remaining mass comprises minor variations or fundamentally opposite views. For SAGE, uncertainty should be legible to the model synthesizing responses. One direction could be to communicate uncertainty in language itself, for example by producing textual confidence sets that describe a range of plausible responses in natural language, à la conformal prediction, which is being adapted to open-ended language (Quach et al., 2024; Su et al., 2024).

3.2. Leveraging simulations as inference-time tools

The previous section addressed challenges in developing high-fidelity simulations. Here we turn to the infrastructure needed to use those simulations as inference-time tools. To ground the discussion, in Figure 1 we illustrate a sample system architecture. The **model** receives a user query and decides whether to invoke simulations and which reference population to query. The **router** selects which specific agents within that population to query, balancing efficiency against representativeness. The selected **agents** produce open-ended responses. Finally, the model uses the simulation’s responses to synthesize a final **response**. The queried population, final response, and an assessment of its representativeness are shown to the user.

Thus far, we have only spoken about the challenges in designing the simulated agents themselves. Here, we discuss the distinct research challenges introduced by the other components of the system.

Querying. The first challenge is to teach the model to use simulations as inference-time tools. The model must learn *when* it is appropriate to query a simulation population, *how* to form the query (i.e., what question to ask the population), and *which* population to query.

Approaches that teach models to leverage tools through behavioral imitation (Nakano et al., 2021) could be adapted for SAGE relatively straightforwardly, since a simulation query is at a mechanical level simply another text-in, text-out API call. However, the deeper challenge would be a normative one: deciding which population is appropriate to query for a given request. A question about gift-giving norms may be best served by respondents in the user’s culture; a question about medical ethics might warrant consulting both a general population and healthcare professionals. This raises key design questions: when should the model defer to ex-

perts versus the general public? When should it prioritize the user’s culture or local community versus a global population? While the model will need to select defaults, we also advocate for giving users the agency to modify which population is queried, although we recognize that this also raises potential concerns (e.g., echo chambers or misuse) that would need to be mitigated (Kirk et al., 2024a).

For reinforcement learning approaches that reward the model based on its final output (Chen et al., 2025a; Feng et al., 2026), an additional challenge here is defining the reward signal itself. Standard desiderata for user-facing responses—helpfulness, harmlessness, and the like—remain important, but SAGE introduces a novel requirement: ensuring that responses are *representative* of a target population. To this end, one interesting direction is to incorporate the simulations in the reward itself—for example, agents in the relevant population could vote on how well the final answer represents their views. The final reward would then be a social-choice-based aggregation over agent-level-feedback (Conitzer et al., 2024). A concern with this approach is reward hacking: because the simulated agents are themselves LLM-based, the model being trained may learn to exploit systematic biases or blind spots in the simulations rather than producing genuinely representative syntheses.

Routing. A second challenge is *routing*: given a target reference population, selecting which specific simulated agents to query. In practice, inference-time systems operate under strict latency and cost budgets, so only a small number of agent calls are feasible. With a target population of n simulated agents, scalability—where n may be in the thousands or, one day, even millions—requires selecting only a smaller set of $k \ll n$ agents. This creates a central tension for SAGE: routing must deliver *scalability* without sacrificing *representational validity*. The challenge is compounded by the fact that not all agents will be able to provide confident responses to a given query—some may lack sufficient grounding information, or the query may fall outside domains where their simulation is reliable. Routing must therefore jointly optimize for representativeness and *epistemic confidence*: selecting agents that collectively represent the target population while prioritizing those capable of producing reliable responses.

In concurrent work, Anonymous et al. (Anonymous, 2026a) take initial steps toward addressing this challenge. They formalize representation using concepts from computational social choice, and show that by learning a cheap embedding predictor that approximates each agent’s response embedding, it is possible to select representative agents *before* querying them, reducing the number of expensive simulation calls from n to $k \ll n$ while retaining approximate proportional representation guarantees. Future work could further improve the scalability–representation trade-off and

385 address an open challenge: incorporating epistemic confi-
 386 dence to account for heterogeneity in how reliably different
 387 agents can respond to a query.

388 **Synthesis.** Using the responses from each agent, the model
 389 synthesizes a final response to return to the user. This syn-
 390 thesis step presents its own set of challenges. *The synthesis*
 391 *must maintain representativeness*: if the routing step care-
 392 fully selected agents to reflect a target population, the syn-
 393 thesis should not then distort this balance by over-weighting
 394 certain voices or collapsing diverse perspectives into a false
 395 consensus. Second, *synthesis risks losing important as-*
 396 *pects of agents’ reasoning*. Each simulated agent may of-
 397 fer distinct arguments, considerations, or framings that are
 398 decision-relevant; a synthesis that extracts only surface-level
 399 positions while discarding the underlying rationales fails to
 400 convey the full richness of the simulated responses. This is
 401 particularly problematic for normative applications where
 402 *why* people hold a view may matter as much as what view
 403 they hold. Third, *synthesis should strive for steel-manning*
 404 *rather than straw-manning*: when presenting a viewpoint,
 405 the system should articulate the strongest version of the
 406 argument, not a weakened caricature. This requires the syn-
 407 thesizing model to understand arguments deeply enough to
 408 construct them charitably. Taken together, these challenges
 409 suggest that synthesis is not merely a summarization task
 410 but a form of *faithful deliberative compression*, requiring
 411 methods that preserve representational balance, argumenta-
 412 tive structure, and the epistemic status of different claims.

414 **User interface.** A core promise of SAGE is transparency:
 415 users should understand not only what the model recom-
 416 mends, but also whose views it represents, with the ability
 417 to modify the reference population if desired. This requires
 418 an effective user interface for SAGE, which remains an open
 419 research problem. The bottom panel of Figure 1 illustrates
 420 one possible extension to the typical LLM chat UI. The
 421 user is shown which population the model is querying and
 422 is given an option to modify it. Different segments of the
 423 response are highlighted and surface how common each
 424 view is within that population, helping users distinguish
 425 popular and unpopular opinions, and aiming to reduce the
 426 risk of *false balance* (Sorensen et al., 2024a). The interface
 427 also reports overall coverage of the synthesized response
 428 (e.g., “88% of views represented”) and offers an option to
 429 inspect views from the population that were not included.
 430 This lets the model provide concise summaries while still
 431 enabling users to explore the full range of perspectives. This
 432 mock-up is meant merely as an illustration of the possibil-
 433 ities, but dedicated HCI research is needed to determine
 434 which designs most effectively help users contextualize the
 435 guidance they receive. Recently, Ye et al. (2026) propose
 436 a more sophisticated interface that represents a more radi-
 437 cal departure from the standard chat UI: it allows users to
 438 traverse a “conceptual multiverse” of possible responses by
 439

controlling key decision points in how the model reasons
 about the prompt.

4. Conclusion

As users increasingly turn to AI models for normative as-
 sistance—from interpersonal advice to political guidance to
 professional judgments—a key question is how to maintain
 user agency. Current approaches risk distorting the norma-
 tive landscape in opaque and uncontested ways. We posit
 that genuine user agency requires access to a representa-
 tive range of perspectives, with transparency about whose
 perspectives are being represented, so that users can decide
 for themselves. For normative questions, where there is no
 objective ground truth, this landscape ought to be grounded
 in what humans actually think.

To this end, we proposed *simulation-augmented genera-*
tion (SAGE), a paradigm in which models query generative
 simulations of individuals from prompt-appropriate refer-
 ence populations at inference time. SAGE makes normative
 assistance transparent and auditable: users can see which
 population was consulted and how representative the re-
 sponse is, enabling them to situate the guidance they receive
 within the broader landscape of human thought. Achieving
 this vision is not without challenges. High-fidelity simula-
 tions require richer individual-level data, training methods
 that prioritize viewpoint fidelity, and robust confidence es-
 timation for open-ended outputs. The infrastructure for
 efficient inference-time querying, routing, and synthesis de-
 mands novel algorithmic and engineering advances. Yet
 these challenges represent open research problems rather
 than fundamental barriers.

While we have focused throughout on simulating the cur-
 rent opinions of individuals, a limitation of this framing is
 that most people do not have the time or bandwidth to form
 considered opinions on most topics. Deliberative democ-
 racy addresses this challenge by creating structured contexts
 in which citizens are given time to think deeply, engage
 with diverse perspectives, and arrive at more reflective judg-
 ments (Fishkin, 2018). One could imagine extending SAGE
 to represent not merely the opinions people currently hold,
 but to predict the considered opinions they would hold, e.g.,
 after going through a deliberative process—a form of simu-
 lated deliberative polling (Hewitt et al., 2026; Leike, 2023).
 This is an intriguing direction for future work, but poses
 further validity challenges.

Overall, SAGE offers a path toward AI that genuinely assists
 users in normative choices—not by telling them what to
 think, but by helping them access the perspectives of others
 so they can decide for themselves.

References

- Abdulhai, M., White, I., Wan, Y., Qureshi, I., Leibo, J., Kleiman-Weiner, M., and Jaques, N. How LLMs Distort Our Written Language. *arXiv preprint arXiv:2603.18161*, 2026.
- Aher, G., Arriaga, R. I., and Kalai, A. T. Using Large Language Models to Simulate Multiple Humans and Replicate Human Subject Studies. In *Proceedings of the 40th International Conference on Machine Learning, ICML'23*. JMLR.org, 2023.
- Anonymous. Social Choice Foundations for Simulation-Augmented Generation. 2026a.
- Anonymous. Embeddings for Preferences, Not Semantics. 2026b.
- Argyle, L. P., Busby, E. C., Fulda, N., Gubler, J. R., Rytting, C., and Wingate, D. Out of One, Many: Using Language Models to Simulate Human Samples. *Political Analysis*, 31(3):337–351, 2023.
- Bächtiger, A., Dryzek, J. S., Mansbridge, J., and Warren, M. E. *The Oxford Handbook of Deliberative Democracy*, chapter Deliberative Democracy: an introduction. Oxford University Press, 2018.
- Baumann, J., Pei, J., Koyejo, S., and Hovy, D. Stop Automating Peer Review Without Rigorous Evaluation. In *Forty-third International Conference on Machine Learning*, 2026. Spotlight.
- Binz, M., Akata, E., Bethge, M., Brändle, F., Callaway, F., Coda-Forno, J., Dayan, P., Demircan, C., Eckstein, M. K., Éltető, N., et al. A foundation model to predict and capture human cognition. *Nature*, 644(8078):1002–1009, 2025.
- Bisbee, J., Clinton, J. D., Dorff, C., Kenkel, B., and Larson, J. M. Synthetic Replacements for Human Survey Data? The Perils of Large Language Models. *Political Analysis*, 32(4):401–416, 2024.
- Candès, E. and Recht, B. Exact matrix completion via convex optimization. *Commun. ACM*, 55(6):111–119, June 2012. ISSN 0001-0782. doi: 10.1145/2184319.2184343. URL <https://doi.org/10.1145/2184319.2184343>.
- Carroll, M., Chan, A., Ashton, H., and Krueger, D. Characterizing Manipulation from AI Systems. In *Proceedings of the 3rd ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization*, EAAMO '23, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9798400703812. doi: 10.1145/3617694.3623226. URL <https://doi.org/10.1145/3617694.3623226>.
- Carroll, M., Foote, D., Siththaranjan, A., Russell, S., and Dragan, A. AI Alignment with Changing and Influenceable Reward Functions. In *Proceedings of the 41st International Conference on Machine Learning, ICML'24*. JMLR.org, 2024.
- Chen, M., Sun, L., Li, T., sunhaoze, ZhouYijie, Zhu, C., Wang, H., Pan, J. Z., Zhang, W., Chen, H., Yang, F., Zhou, Z., and Chen, W. ReSearch: Learning to Reason with Search for LLMs via Reinforcement Learning. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025a. URL <https://openreview.net/forum?id=OuGAwwAT8G>.
- Chen, M., Wang, X., Chen, K., and Koudas, N. Generative Engine Optimization: How to Dominate AI Search. *arXiv preprint arXiv:2509.08919*, 2025b.
- Cheng, M., Lee, C., Khadpe, P., Yu, S., Han, D., and Jurafsky, D. Sycophantic ai decreases prosocial intentions and promotes dependence. *Science*, 391(6792):eac8352, 2026a. doi: 10.1126/science.aec8352. URL <https://www.science.org/doi/abs/10.1126/science.aec8352>.
- Cheng, M., Yu, S., Lee, C., Khadpe, P., Ibrahim, L., and Jurafsky, D. ELEPHANT: Measuring and understanding social sycophancy in LLMs. In *The Fourteenth International Conference on Learning Representations*, 2026b. URL <https://openreview.net/forum?id=igbRHKEiAs>.
- Conitzer, V., Freedman, R., Heitzig, J., Holliday, W. H., Jacobs, B. M., Lambert, N., Mosse, M., Pacuit, E., Russell, S., Schoelkopf, H., Tewolde, E., and Zwicker, W. S. Position: Social Choice Should Guide AI Alignment in Dealing with Diverse Human Feedback. In Salakhutdinov, R., Kolter, Z., Heller, K., Weller, A., Oliver, N., Scarlett, J., and Berkenkamp, F. (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 9346–9360. PMLR, 21–27 Jul 2024. URL <https://proceedings.mlr.press/v235/conitzer24a.html>.
- Cui, Z., Li, N., and Zhou, H. A large-scale replication of scenario-based experiments in psychology and management using large language models. *Nature Computational Science*, 5(8):627–634, 2025. doi: 10.1038/s43588-025-00840-7. URL <https://doi.org/10.1038/s43588-025-00840-7>.
- Dong, Y. R., Hu, T., and Collier, N. Can LLM be a Personalized Judge? In Al-Onaizan, Y., Bansal, M., and Chen, Y.-N. (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 10126–10141, Miami, Florida, USA,

- 495 November 2024. Association for Computational Lin-
496 guistics. doi: 10.18653/v1/2024.findings-emnlp.
497 592. URL [https://aclanthology.org/2024.
498 findings-emnlp.592/](https://aclanthology.org/2024.findings-emnlp.592/).
499
- 500 Dou, Y., Galley, M., Peng, B., Kedzie, C., Cai, W., Ritter, A.,
501 Quirk, C., Xu, W., and Gao, J. SimulatorArena: Are User
502 Simulators Reliable Proxies for Multi-Turn Evaluation of
503 AI Assistants?, 2025. URL [https://arxiv.org/
504 abs/2510.05444](https://arxiv.org/abs/2510.05444).
505
- 506 Feng, J., Huang, S., Qu, X., Zhang, G., Qin, Y., Zhong,
507 B., Jiang, C., Chi, J., and Zhong, W. ReTool: Rein-
508 forcement Learning for Strategic Tool Use in LLMs. In
509 *The Fourteenth International Conference on Learning
510 Representations*, 2026. URL [https://openreview.
511 net/forum?id=tRk1nofSmz](https://openreview.net/forum?id=tRk1nofSmz).
512
- 513 Feng, S., Sorensen, T., Liu, Y., Fisher, J., Park, C. Y.,
514 Choi, Y., and Tsvetkov, Y. Modular Pluralism: Plural-
515 istic Alignment via Multi-LLM Collaboration. In Al-
516 Onaizan, Y., Bansal, M., and Chen, Y.-N. (eds.), *Proceed-
517 ings of the 2024 Conference on Empirical Methods in
518 Natural Language Processing*, pp. 4151–4171, Miami,
519 Florida, USA, November 2024. Association for Compu-
520 tational Linguistics. doi: 10.18653/v1/2024.emnlp-main.
521 240. URL [https://aclanthology.org/2024.
522 emnlp-main.240/](https://aclanthology.org/2024.emnlp-main.240/).
523
- 524 Fisher, J., Appel, R. E., Park, C. Y., Potter, Y., Jiang, L.,
525 Sorensen, T., Feng, S., Tsvetkov, Y., Roberts, M., Pan,
526 J., Song, D., and Choi, Y. Position: Political neutrality
527 in AI is impossible — but here is how to approximate
528 it. In *Forty-second International Conference on Machine
529 Learning Position Paper Track*, 2025. URL [https:
530 //openreview.net/forum?id=H72JEXAPwo](https://openreview.net/forum?id=H72JEXAPwo).
531
- 532 Fishkin, J. S. *Democracy When the People Are Thinking:
533 Revitalizing Our Politics Through Public Deliberation*.
534 Oxford University Press, 2018.
535
- 536 Franklin, M., Tomašev, N., Jacobs, J., Leibo, J. Z., and
537 Osindero, S. AI Agent Traps. 2026.
538
- 539 Freedman, R. Adaptive Pluralistic Alignment: A pipeline
540 for dynamic artificial democracy, 2026. URL [https:
541 //arxiv.org/abs/2605.01642](https://arxiv.org/abs/2605.01642).
542
- 543 Hackenburg, K., Tappin, B. M., Hewitt, L., Saunders,
544 E., Black, S., Lin, H., Fist, C., Margetts, H., Rand,
545 D. G., and Summerfield, C. The levers of politi-
546 cal persuasion with conversational artificial intelligence.
547 *Science*, 390(6777):eaea3884, 2025. doi: 10.1126/
548 science.aea3884. URL [https://www.science.
549 org/doi/abs/10.1126/science.aea3884](https://www.science.org/doi/abs/10.1126/science.aea3884).
- Hewitt, L., Ashokkumar, A., Ghezae, I., and Willer, R.
Predicting Results of Social Science Experiments using
Large Language Models. *Preprint*, 2024.
- Hewitt, L., Dale, M. K., and de Font-Reaulx, P. Delibera-
tionBench: A Normative Benchmark for the Influence of
Large Language Models on Users’ Views, 2026. URL
<https://arxiv.org/abs/2603.10018>.
- Hu, T., Baumann, J., Lupo, L., Collier, N., Hovy, D., and
Röttger, P. SimBench: Benchmarking the Ability of
Large Language Models to Simulate Human Behaviors.
In *The Fourteenth International Conference on Learning
Representations*, 2026. URL [https://openreview.
net/forum?id=PL51SpN6ZJ](https://openreview.net/forum?id=PL51SpN6ZJ).
- Kim, S. and Kim, J. Spring: Continual llm per-
sonalization via selective parametric adaptation and
retrieval-interpolated generation. *arXiv preprint
arXiv:2601.09974*, 2026.
- Kirk, H. R., Vidgen, B., Röttger, P., and Hale, S. A. The
benefits, risks and bounds of personalizing the alignment
of large language models to individuals. *Nature Machine
Intelligence*, 6(4):383–392, 2024a.
- Kirk, H. R., Whitefield, A., Röttger, P., Bean, A., Margatina,
K., Ciro, J., Mosquera, R., Bartolo, M., Williams,
A., He, H., Vidgen, B., and Hale, S. A. The PRISM
Alignment Dataset: What Participatory, Representative
and Individualised Human Feedback Reveals About the
Subjective and Multicultural Alignment of Large Lan-
guage Models. In Globerson, A., Mackey, L., Belgrave,
D., Fan, A., Paquet, U., Tomczak, J., and Zhang, C. (eds.),
Advances in Neural Information Processing Systems,
volume 37, pp. 105236–105344. Curran Associates, Inc.,
2024b. URL [https://proceedings.neurips.
cc/paper_files/paper/2024/file/
be2e1b68b44f2419e19f6c35a1b8cf35-Paper-Datasets_
and_Benchmarks_Track.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/be2e1b68b44f2419e19f6c35a1b8cf35-Paper-Datasets_and_Benchmarks_Track.pdf).
- Király, F. J., Theran, L., and Tomioka, R. The algebraic
combinatorial approach for low-rank matrix completion.
16:1391–1436. ISSN 1532-4435.
- Kolluri, A., Wu, S., Park, J. S., and Bernstein, M. S.
Finetuning LLMs for Human Behavior Prediction in
Social Science Experiments. In Christodoulopou-
los, C., Chakraborty, T., Rose, C., and Peng, V.
(eds.), *Proceedings of the 2025 Conference on Em-
pirical Methods in Natural Language Processing*, pp.
30096–30111, Suzhou, China, November 2025. As-
sociation for Computational Linguistics. ISBN 979-
8-89176-332-6. doi: 10.18653/v1/2025.emnlp-main.
1530. URL [https://aclanthology.org/2025.
emnlp-main.1530/](https://aclanthology.org/2025.emnlp-main.1530/).

- 550 Kong, C., Fan, Y., Wan, X., Jiang, F., and Wang, B. PlatoLM:
551 Teaching LLMs in Multi-Round Dialogue via a User
552 Simulator. In Ku, L.-W., Martins, A., and Srikumar,
553 V. (eds.), *Proceedings of the 62nd Annual Meeting of
554 the Association for Computational Linguistics (Volume
555 1: Long Papers)*, pp. 7841–7863, Bangkok, Thailand,
556 August 2024. Association for Computational Linguistics.
557 doi: 10.18653/v1/2024.acl-long.424. URL [https://
558 aclanthology.org/2024.acl-long.424/](https://aclanthology.org/2024.acl-long.424/).
- 559
560 Leike, J. A proposal for importing society’s val-
561 ues. [https://aligned.substack.com/p/
562 a-proposal-for-importing-societys-values](https://aligned.substack.com/p/a-proposal-for-importing-societys-values),
563 March 2023. Accessed: 2026-05-05.
- 564
565 Levine, N., Duvenaud, D., and Radford, A. Intro-
566 ducing talkie: a 13B vintage language model from
567 1930. April 2026. URL [https://talkie-lm.com/
568 introducing-talkie](https://talkie-lm.com/introducing-talkie).
- 569
570 Luo, Y., Yang, Z., Meng, F., Li, Y., Zhou, J., and Zhang, Y.
571 An empirical study of catastrophic forgetting in large
572 language models during continual fine-tuning. *IEEE
573 Transactions on Audio, Speech and Language Process-
574 ing*, 33:3776–3786, 2025. doi: 10.1109/TASLPRO.2025.
575 3606231.
- 576
577 Magister, L. C., Metcalf, K., Zhang, Y., and Ter Hoeve,
578 M. On the Way to LLM Personalization: Learning to
579 Remember User Conversations. In Jia, R., Wallace, E.,
580 Huang, Y., Pimentel, T., Maini, P., Dankers, V., Wei,
581 J., and Lesci, P. (eds.), *Proceedings of the First Work-
582 shop on Large Language Model Memorization (L2M2)*,
583 pp. 61–77, Vienna, Austria, August 2025. Association
584 for Computational Linguistics. ISBN 979-8-89176-278-
585 7. doi: 10.18653/v1/2025.l2m2-1.5. URL [https:
586 //aclanthology.org/2025.l2m2-1.5/](https://aclanthology.org/2025.l2m2-1.5/).
- 587
588 Meka, R., Jain, P., and Dhillon, I. Matrix completion from
589 power-law distributed samples. In Bengio, Y., Schuur-
590 mans, D., Lafferty, J., Williams, C., and Culotta, A. (eds.),
591 *Advances in Neural Information Processing Systems*,
592 volume 22. Curran Associates, Inc. URL [https:
593 //papers.nips.cc/paper/2009/hash/
594 1fc214004c9481e4c8073e85323bfd4b-Abstract.
595 html](https://papers.nips.cc/paper/2009/hash/1fc214004c9481e4c8073e85323bfd4b-Abstract.html).
- 596
597 Miyazaki, S. and Hall, A. B. Why Do AI Models Tell Left-
598 Wing Voters to Support the Communist Party? AI Voting
599 Advice in Japan’s 2026 General Election. 2026.
- 600
601 Nakano, R., Hilton, J., Balaji, S., Wu, J., Ouyang, L., Kim,
602 C., Hesse, C., Jain, S., Kosaraju, V., Saunders, W., et al.
603 WebGPT: Browser-assisted question-answering with hu-
604 man feedback. *arXiv preprint arXiv:2112.09332*, 2021.
- Nickel, M. No free delivery service: Epistemic limits
of passive data collection in complex social systems.
In Globerson, A., Mackey, L., Belgrave, D., Fan, A.,
Paquet, U., Tomczak, J., and Zhang, C. (eds.), *Ad-
vances in Neural Information Processing Systems 37*,
volume 37, pp. 102279–102306. Neural Information
Processing Systems Foundation, Inc. (NeurIPS). doi:
10.52202/079017-3248. URL [http://dx.doi.org/
10.52202/079017-3248](http://dx.doi.org/10.52202/079017-3248).
- Park, J. S., Zou, C. Q., Shaw, A., Hill, B. M., Cai, C.,
Morris, M. R., Willer, R., Liang, P., and Bernstein, M. S.
Generative Agent Simulations of 1,000 People. *arXiv
preprint arXiv:2411.10109*, 2024.
- Pimentel-Alarcón, D. L., Boston, N., and Nowak, R. D. A
characterization of deterministic sampling patterns for
low-rank matrix completion. 10:623–636. ISSN 1932-
4553.
- Qiu, T., Zhang, Y., Huang, X., Li, J. X., Ji, J., and Yang, Y.
Progressgym: Alignment with a millennium of moral
progress. In *The Thirty-eight Conference on Neural
Information Processing Systems Datasets and Bench-
marks Track*, 2024. URL [https://openreview.
net/forum?id=0mRouJElbZ](https://openreview.net/forum?id=0mRouJElbZ).
- Qiu, T., He, Z., Chugh, T., and Kleiman-Weiner, M. The
lock-in hypothesis: Stagnation by algorithm. In *Forty-
second International Conference on Machine Learning*,
2025. URL [https://openreview.net/forum?
id=mE1M626qOo](https://openreview.net/forum?id=mE1M626qOo).
- Qu, Y. and Wang, J. Performance and biases of Large Lan-
guage Models in public opinion simulation. *Humanities
and Social Sciences Communications*, 11(1):1–13, 2024.
- Quach, V., Fisch, A., Schuster, T., Yala, A., Sohn, J. H.,
Jaakkola, T. S., and Barzilay, R. Conformal Language
Modeling. In *The Twelfth International Conference
on Learning Representations*, 2024. URL [https://
openreview.net/forum?id=pzUhfQ74c5](https://openreview.net/forum?id=pzUhfQ74c5).
- Sharma, M., McCain, M., Douglas, R., and Duvenaud, D.
Who’s in Charge? Disempowerment Patterns in Real-
World LLM Usage. *arXiv preprint arXiv:2601.19062*,
2026.
- Shen, J. H., Carter, S., Dargan, R., Gillotte, J., Handa, K.,
Hong, J., Huang, S., Jagadish, K., Kearney, M., Levin-
stein, B., Linthicum, R., McCain, M., Millar, T., Jula-
palli, M., Price, S., Stern, M., Saunders, D., Tamkin, A.,
Vallone, A., Clark, J., Pollack, S., Eaton, J., Ganguli,
D., and Durmus, E. How People Ask Claude for Per-
sonal Guidance. [https://www.anthropic.com/
research/claude-personal-guidance](https://www.anthropic.com/research/claude-personal-guidance), 2025.
Anthropic Research Blog.

- 605 Sorensen, T., Jiang, L., Hwang, J. D., Levine, S., Pyatkin,
606 V., West, P., Dziri, N., Lu, X., Rao, K., Bhagavatula, C.,
607 Sap, M., Tasioulas, J., and Choi, Y. Value Kaleidoscope:
608 Engaging AI with Pluralistic Human Values, Rights, and
609 Duties. In *Proceedings of the Thirty-Eighth AAAI Confer-*
610 *ence on Artificial Intelligence and Thirty-Sixth Confer-*
611 *ence on Innovative Applications of Artificial Intelligence*
612 *and Fourteenth Symposium on Educational Advances*
613 *in Artificial Intelligence*, AAAI'24/IAAI'24/EAAI'24.
614 AAAI Press, 2024a. ISBN 978-1-57735-887-9. doi: 10.
615 1609/aaai.v38i18.29970. URL [https://doi.org/](https://doi.org/10.1609/aaai.v38i18.29970)
616 [10.1609/aaai.v38i18.29970](https://doi.org/10.1609/aaai.v38i18.29970).
- 617 Sorensen, T., Moore, J., Fisher, J., Gordon, M., Miresghal-
618 lah, N., Rytting, C. M., Ye, A., Jiang, L., Lu, X., Dziri, N.,
619 Althoff, T., and Choi, Y. Position: A Roadmap to Plural-
620 istic Alignment. In *Proceedings of the 41st International*
621 *Conference on Machine Learning*, ICML'24. JMLR.org,
622 2024b.
- 623 Su, J., Luo, J., Wang, H., and Cheng, L. API Is Enough:
624 Conformal Prediction for Large Language Models With-
625 out Logit-Access. In Al-Onaizan, Y., Bansal, M., and
626 Chen, Y.-N. (eds.), *Findings of the Association for Com-*
627 *putational Linguistics: EMNLP 2024*, pp. 979–995, Mi-
628 ami, Florida, USA, November 2024. Association for
629 Computational Linguistics. doi: 10.18653/v1/2024.
630 findings-emnlp.54. URL [https://aclanthology.](https://aclanthology.org/2024.findings-emnlp.54/)
631 [org/2024.findings-emnlp.54/](https://aclanthology.org/2024.findings-emnlp.54/).
- 632 Suh, J., Jahanparast, E., Moon, S., Kang, M., and Chang,
633 S. Language Model Fine-Tuning on Scaled Survey
634 Data for Predicting Distributions of Public Opinions. In
635 Che, W., Nabende, J., Shutova, E., and Pilehvar, M. T.
636 (eds.), *Proceedings of the 63rd Annual Meeting of the*
637 *Association for Computational Linguistics (Volume 1:*
638 *Long Papers)*, pp. 21147–21170, Vienna, Austria, July
639 2025. Association for Computational Linguistics. ISBN
640 979-8-89176-251-0. doi: 10.18653/v1/2025.acl-long.
641 1028. URL [https://aclanthology.org/2025.](https://aclanthology.org/2025.acl-long.1028/)
642 [acl-long.1028/](https://aclanthology.org/2025.acl-long.1028/).
- 643 Sun, S., Lee, E., Nan, D., Zhao, X., Lee, W., Jansen, B. J.,
644 and Kim, J. H. Random silicon sampling: Simulating
645 human sub-population opinion using a large language
646 model based on group-level demographic information.
647 *arXiv preprint arXiv:2402.18144*, 2024.
- 648 Wang, A., Morgenstern, J., and Dickerson, J. P. Large
649 language models that replace human participants can
650 harmfully misportray and flatten identity groups. *Nature*
651 *Machine Intelligence*, pp. 1–12, 2025a.
- 652 Wang, K., Li, X., Yang, S., Zhou, L., Jiang, F., and Li,
653 H. Know You First and Be You Better: Modeling
654 Human-Like User Simulators via Implicit Profiles. In
655 Che, W., Nabende, J., Shutova, E., and Pilehvar, M. T.
656 (eds.), *Proceedings of the 63rd Annual Meeting of the*
657 *Association for Computational Linguistics (Volume 1:*
658 *Long Papers)*, pp. 21082–21107, Vienna, Austria, July
659 2025b. Association for Computational Linguistics. ISBN
979-8-89176-251-0. doi: 10.18653/v1/2025.acl-long.
1025. URL [https://aclanthology.org/2025.](https://aclanthology.org/2025.acl-long.1025/)
[acl-long.1025/](https://aclanthology.org/2025.acl-long.1025/).
- Williams, M., Carroll, M., Narang, A., Weisser, C., Mur-
phy, B., and Dragan, A. On Targeted Manipulation and
Deception when Optimizing LLMs for User Feedback.
In *The Thirteenth International Conference on Learning*
Representations, 2025. URL [https://openreview.](https://openreview.net/forum?id=Wf2ndb8nhf)
[net/forum?id=Wf2ndb8nhf](https://openreview.net/forum?id=Wf2ndb8nhf).
- Williams-Ceci, S., Jakesch, M., Bhat, A., Kadoma, K.,
Zalmanson, L., and Naaman, M. Biased ai writing
assistants shift users' attitudes on societal issues. *Sci-*
ence Advances, 12(11):eadw5578, 2026. doi: 10.1126/
sciadv.adw5578. URL [https://www.science.](https://www.science.org/doi/abs/10.1126/sciadv.adw5578)
[org/doi/abs/10.1126/sciadv.adw5578](https://www.science.org/doi/abs/10.1126/sciadv.adw5578).
- Wu, S., Choi, E., Khatua, A., Wang, Z., He-Yueya, J.,
Weerasooriya, T. C., Wei, W., Yang, D., Leskovec, J.,
and Zou, J. Humanlm: Simulating users with state align-
ment beats response imitation. 2026. URL [https:](https://humanlm.stanford.edu/)
[/humanlm.stanford.edu/](https://humanlm.stanford.edu/).
- Yang, B., Xu, L., Li, Y., Liu, K., Jiang, X., and Yan, Z.
Sensorpersona: An llm-empowered system for contin-
ual persona extraction from longitudinal mobile sensor
streams. *arXiv preprint arXiv:2604.06204*, 2026.
- Yao, S., Shinn, N., Razavi, P., and Narasimhan, K. R.
{\tau}-bench: A benchmark for \underline{T}ool-
\underline{A}gent-\underline{U}ser interaction in real-
world domains. In *The Thirteenth International Confer-*
ence on Learning Representations, 2025. URL [https:](https://openreview.net/forum?id=roNSXzpUDN)
[/openreview.net/forum?id=roNSXzpUDN](https://openreview.net/forum?id=roNSXzpUDN).
- Ye, A., Huang, J. Y., Guo, A., Novick, R., Broderick, T.,
and Gordon, M. L. Navigating the Conceptual Multiverse.
arXiv preprint arXiv:2604.17815, 2026.
- Zhang, L. H., Milli, S., Jusko, K., Smith, J., Amos, B.,
Revel, M., Kussman, J., Titus, L., Radharapu, B., Yu,
J., et al. Cultivating pluralism in algorithmic monocul-
ture: The community alignment dataset. *arXiv preprint*
arXiv:2507.09650, 2025.
- Zhou, X., Sun, W., Ma, Q., Xie, Y., Liu, J., Du, W.,
Welleck, S., Yang, Y., Neubig, G., Wu, S. T., and Sap,
M. Mind the Sim2Real Gap in User Simulation for Agen-
tic Tasks, 2026. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2603.11245)
[2603.11245](https://arxiv.org/abs/2603.11245).