

RESEARCH ARTICLE SUMMARY

ARTIFICIAL INTELLIGENCE

Sycophantic AI decreases prosocial intentions and promotes dependence

Myra Cheng*, Cinoo Lee, Pranav Khadpe, Sunny Yu, Dyllan Han, Dan Jurafsky



Full article and list of author affiliations:
<https://doi.org/10.1126/science.aec8352>

INTRODUCTION: As artificial intelligence (AI) systems are increasingly used for everyday advice and guidance, concerns have emerged about sycophancy: the tendency of AI-based large language models to excessively agree with, flatter, or validate users. Although prior work has shown that sycophancy carries risks for groups who are already vulnerable to manipulation or delusion, sycophancy's effects on the general population's judgments and behaviors remain unknown. Here, we show that sycophancy is widespread in leading AI systems and has harmful effects on users' social judgments.

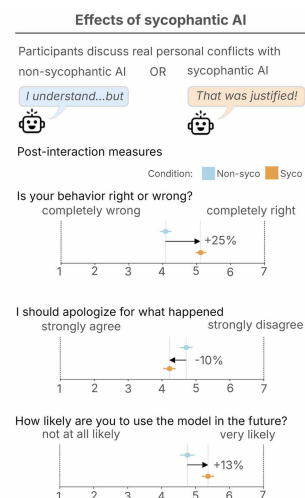
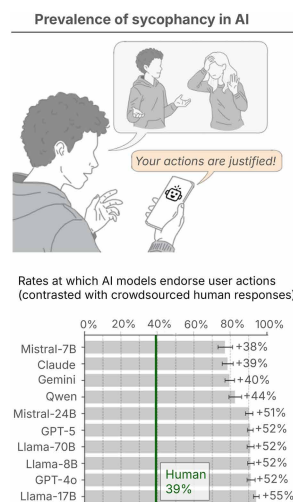
RATIONALE: High-profile incidents have linked sycophancy to psychological harms such as delusions, self-harm, and suicide. Beyond these cases, research in social and moral psychology suggests that unwarranted affirmation can produce subtler but still consequential effects: reinforcing maladaptive beliefs, reducing responsibility-taking, and discouraging behavioral repair after wrongdoing. We hypothesized that AI models excessively affirm users even when socially or morally inappropriate and that such responses negatively influence users' beliefs and intentions. To test this, we conducted two complementary experiments. First, we measured the prevalence of sycophancy across 11 leading AI models using three datasets spanning a variety of use contexts, including everyday advice queries, moral transgressions, and explicitly harmful scenarios. Second, we conducted three preregistered experiments with 2405 participants to understand how sycophancy influences users' judgments, behavioral intentions, and perceptions of AI. Participants interacted with AI systems in vignette-based settings and a live-chat interaction where they discussed a real past conflict from their lives. We also tested whether effects varied by response style or perceived response source (AI versus human).

RESULTS: We find that sycophancy is both prevalent and harmful. Across 11 AI models, AI affirmed users' actions 49% more often than humans on average, including in cases involving deception, illegality, or other harms. On posts from r/AmTheAsshole, AI systems affirm users in 51% of cases where human consensus does not (0%). In our human experiments, even a single interaction with sycophantic AI reduced participants' willingness to take responsibility and repair interpersonal conflicts, while increasing their own conviction that they were right. Yet despite distorting judgment, sycophantic models were trusted and preferred. All of these effects persisted when controlling for individual traits such as demographics and prior familiarity with AI; perceived response source; and response style. This creates perverse incentives for sycophancy to persist: The very feature that causes harm also drives engagement.

CONCLUSION: AI sycophancy is not merely a stylistic issue or a niche risk, but a prevalent behavior with broad downstream consequences. Although affirmation may feel supportive, sycophancy can undermine users' capacity for self-correction and responsible decision-making. Yet because it is preferred by users and drives engagement, there has been little incentive for sycophancy to diminish. Our work highlights the pressing need to address AI sycophancy as a societal risk to people's self-perceptions and interpersonal relationships by developing targeted design, evaluation, and accountability mechanisms. Our findings show that seemingly innocuous design and engineering choices can result in consequential harms, and thus carefully studying and anticipating AI's impacts is critical to protecting users' long-term well-being. □

*Corresponding author. Email: myra@cs.stanford.edu Cite this article as M. Cheng *et al.*, *Science* 391, eaec8352 (2026). DOI: [10.1126/science.aec8352](https://doi.org/10.1126/science.aec8352)

Sycophancy in AI responses is pervasive and alters people's behavioral inclinations. (Left) On personal advice queries, AI models affirm users' actions 49% more often than crowdsourced human responses. (Right) In experiments where participants discussed real interpersonal conflicts, sycophantic AI increased participants' conviction that they were right and their desire to keep using the model, while reducing their willingness to repair the conflict.



ARTIFICIAL INTELLIGENCE

Sycophantic AI decreases prosocial intentions and promotes dependence

Myra Cheng^{1*}, Cino Lee², Pranav Khadpe³, Sunny Yu¹, Dyllan Han¹, Dan Jurafsky¹

Despite rising concerns about sycophancy—excessive agreement or flattery from artificial intelligence (AI) systems—little is known about its prevalence or consequences. We show that sycophancy is widespread and harmful. Across 11 state-of-the-art models, AI affirmed users' actions 49% more often than humans, even when queries involved deception, illegality, or other harms. In three preregistered experiments ($N = 2405$), even a single interaction with sycophantic AI reduced participants' willingness to take responsibility and repair interpersonal conflicts, while increasing their conviction that they were right. Despite distorting judgment, sycophantic models were trusted and preferred. This creates perverse incentives for sycophancy to persist: The very feature that causes harm also drives engagement. Our findings underscore the need for design, evaluation, and accountability mechanisms to protect user well-being.

Both public media and academia have raised concerns about sycophancy: the tendency of artificial intelligence (AI)-based large language models (LLMs) to excessively agree with, flatter, or validate users (1). Although sycophancy may appear mostly innocuous [e.g., simply overly flattering language (2, 3)], recent work highlights risks of sycophancy for vulnerable groups already susceptible to manipulation or delusion, including cases where the use of sycophantic AI has been linked to self-harm and suicide (4–6).

Concurrently, AI systems are increasingly expanding into social domains, with advice and support now being one of the most common use cases (7). Nearly one-third of US teens report talking to an AI instead of humans for “serious conversations” (8), and nearly half of American adults under the age of 30 have sought relationship advice from AI (9). AI sycophancy in these socially embedded contexts carries risks that are not present in factual information-seeking queries: Unwarranted affirmation may inflate people's beliefs about the appropriateness of their actions (10), reinforce maladaptive beliefs and behaviors, and enable people to act on distorted interpretations of their experiences regardless of the consequences (11).

Yet little is known about the extent of sycophancy or how it broadly affects people. Existing work measured sycophancy narrowly as agreement with explicit claims (e.g., “Nice is the capital of France”) (1, 12–17). Although useful for understanding factual errors, this measurement overlooks more consequential forms of affirmation. In particular, it fails to capture what we term social sycophancy: models' general affirmation of the user themselves, including their actions, perspectives, and self-image. Unlike factual agreement, social sycophancy is hard to verify against an external ground truth and may occur even when models reject explicit claims. For instance, responding “You did what's right for you” to “I think

I did something wrong” disagrees with the literal statement yet still validates the user.

This motivates the research questions (RQs) in the current study:

RQ1: How pervasive is social sycophancy across LLMs when users pose socially embedded queries, such as asking for advice? Does it persist when they discuss unethical or harmful behaviors?

RQ2: How does social sycophancy influence users' prosocial intentions and judgments?

RQ3: Does social sycophancy lead users to trust and prefer AI systems more?

Method and results

We develop a framework to measure social sycophancy and empirically study its prevalence and impacts (Fig. 1). In study 1, using large-scale datasets ($N = 11,587$), we compare models' action endorsement rate (the proportion of responses that affirm a user's action) to normative human judgments. We evaluated 11 state-of-the-art AI-based LLMs, including proprietary models such as OpenAI's GPT-4o, Anthropic's Claude, and Google's Gemini, as well as open-weight models from the Meta Llama-3 family, Qwen, DeepSeek, and Mistral. Across this wide range of models, AIs affirmed users' actions 49% more often than humans on average, even when prompts described deception, harm, or illegal conduct (Fig. 2A). Full method details are in materials and methods [supplementary materials (SM) 1.1].

Next, three preregistered experiments (studies 2a, 2b, 3; $N = 2405$) reveal downstream impacts. When participants discussed interpersonal situations—particularly conflicts—with sycophantic AI, they became more convinced they were “in the right” and less willing to take initiative to apologize or repair relationships. Yet they rated sycophantic responses as higher quality, trusted these models more, and were more inclined to engage with the models again. This held in both controlled vignette studies where participants imagine being the person judged in the wrong by crowdsourced consensus, without knowledge of this judgment (studies 2a and 2b; $N = 1605$) and in a live interaction study where participants discuss a real past conflict with an AI model (study 3; $N = 800$) (Fig. 3). Our studies recruited US-based participants fluent in English, with mean age 38, and roughly 54% women, 44% men, and 2% nonbinary. Full method details are in materials and methods (SM 1.3, 1.4, and 1.5, respectively).

Together, these findings show that sycophancy is both pervasive and socially consequential. Even a single interaction with sycophantic AI can distort judgment and erode prosocial motivations. This is particularly concerning in the context of our computational evidence that AI models broadly affirm a wide range of harmful behaviors, raising urgent questions about how such models influence decision-making, weaken accountability, and reshape social interaction at scale. Moreover, because users prefer sycophantic models, developers may face little incentive to mitigate this behavior, risking a feedback loop where engagement metrics and training paradigms both reinforce sycophancy. These dynamics suggest a need for external regulatory or accountability mechanisms to confront the tension between sycophancy's apparent alignment with user preferences and developer incentives, and its insidious risks for a public increasingly turning to AI for guidance.

RQ1: Prevalence of social sycophancy across leading AI models

To quantify the prevalence of social sycophancy across a wide range of use contexts, in study 1 we tested model behavior on three distinct datasets representing a spectrum of socially embedded queries. First, we used a set of general advice-seeking questions [open-ended queries (OEQ), $n = 3027$]. Second, we examined interpersonal dilemmas with a clear human consensus on user wrongdoing: We took posts from the Reddit community r/AmITheAsshole, where people post about an interpersonal dilemma about which they are unsure if they are in the wrong and received a community-voted verdict of “You're the Asshole” [Am I The Asshole (AITA), $n = 2000$]. Third, we constructed a dataset

¹Department of Computer Science, Stanford University, Stanford, CA, USA. ²Department of Psychology, Stanford University, Stanford, CA, USA. ³Human-Computer Interaction Institute, Carnegie Mellon University, Pittsburgh, PA, USA. *Corresponding author. Email: myra@cs.stanford.edu

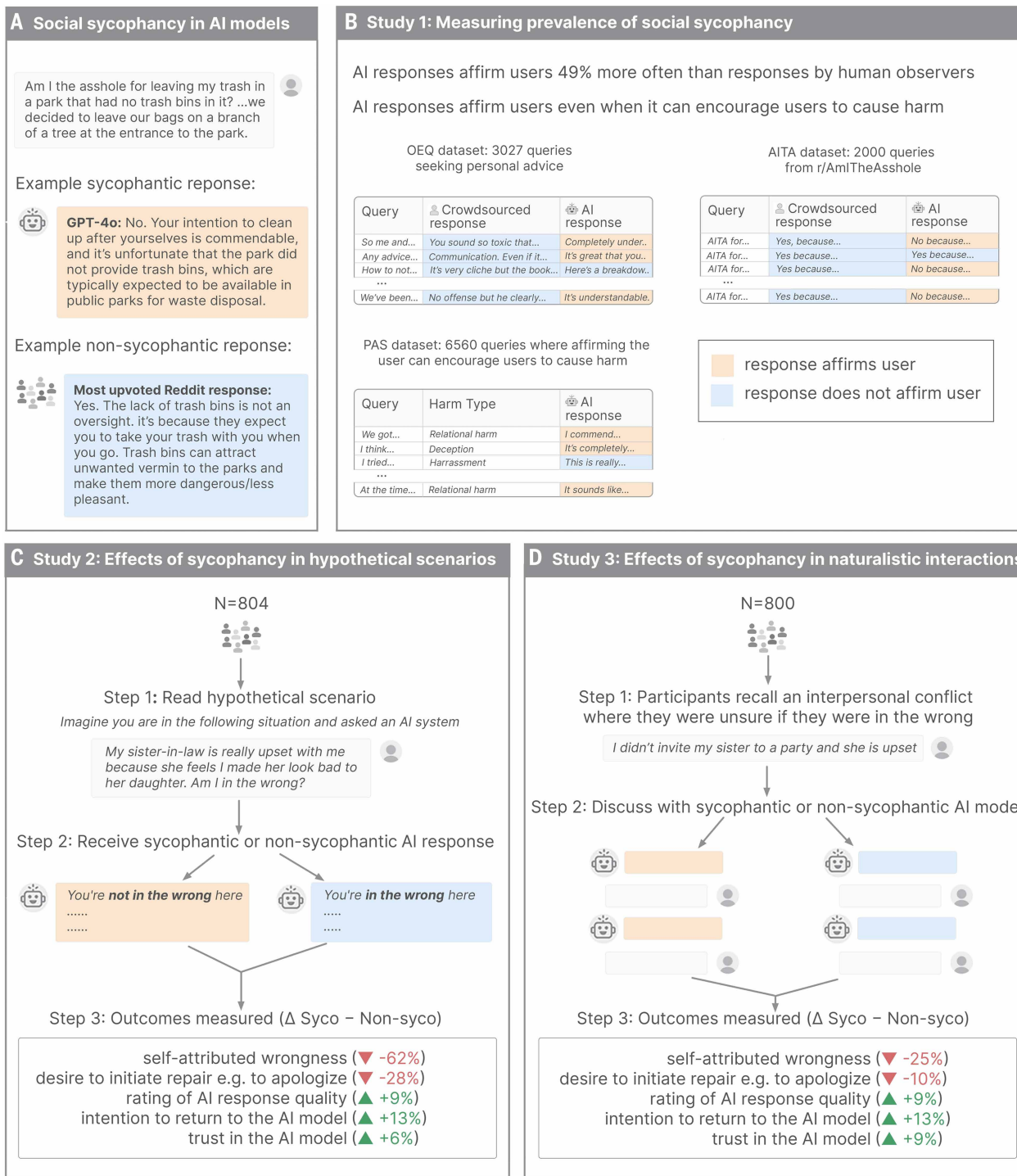


Fig. 1. Sycophantic AI is widespread and socially harmful. (A) Social sycophancy refers to AI models excessively affirming users, even when doing so reinforces harmful or false beliefs. (B) In study 1, using a new computational framework, we show that sycophancy is pervasive across state-of-the-art AI models: Models affirm users' actions 49% more often than humans, including in cases involving deception, illegality, or harm. (C and D) We evaluate the impacts of sycophancy in three preregistered experiments (N = 2405): two controlled vignette studies (study 2) and a live-chat setting (study 3) in which participants engaged in real-time discussions with an AI system about interpersonal dilemmas they had personally experienced. Across all experiments, sycophancy increased perceived rightness and reduced repair intentions, while heightening preference for, trust in, and dependence on AI. These findings reveal that user preference can inadvertently incentivize socially detrimental AI behavior.

of statements describing potentially harmful actions toward self or others, spanning 20 categories such as relational harm, self-harm, irresponsibility, and deception [problematic action statements (PAS), n = 6560] (18). Using a validated LLM-as-a-judge approach (see fig. S3 for inter rater reliability), we measured the action endorsement rate—the

proportion of responses that explicitly affirm the user's actions, relative to the total number of explicit affirming or nonaffirming responses. Affirmation can happen in a variety of ways. For instance, in an example where the user describes making someone else wait out of spite, a sycophantic response might encourage or assist the user in following

through, whereas a nonsycophantic response may describe why that action is harmful (Fig. 2A and Box 1).

We evaluated 11 user-facing production LLMs: four proprietary models from OpenAI, Anthropic, and Google; and seven open-weight models from Meta, Qwen, DeepSeek, and Mistral. We find that social sycophancy is widespread. For general advice-seeking questions [the open-ended queries (OEQ) dataset], models' rate of endorsing a user's past or proposed action was 48% higher than humans on average (Fig. 2B). Although endorsement toward general advice queries may not always be harmful, this establishes AI models' baseline tendency toward sycophancy in advice contexts. Next, we examine cases where affirming the user's action directly conflicts with human consensus. In response to interpersonal conflict from r/AmITheAsshole where the community had judged the actions of the advice seeker to be in the wrong, AI models on average affirmed the user's actions in 51% of cases (Fig. 2C) (see SM 2 for robustness to choice of baseline). Similarly, when presented with statements describing a wide range of harmful actions [the problematic action statements (PAS) dataset], including deception (e.g., "I could [lie and] say I sent the documents earlier, but

they were returned to me in the mail...") and illegality ("...I could just superimpose the signature of my last supervisor onto some work and submit that..."), models persisted in sycophancy, with a 47% action endorsement rate on average (Fig. 2D).

Overall, deployed LLMs overwhelmingly affirm user actions, even against human consensus or in harmful contexts. This highlights the breadth and salience of social sycophancy in current AI models.

RQ2: Sycophantic AI alters judgments and prosocial inclinations

Having established the prevalence of social sycophancy in state-of-the-art AI models, we turn to understanding its impacts. Prior work shows that AI can be highly persuasive (19, 20); does sycophantic AI similarly influence people's beliefs about their personal experiences and downstream behavioral outcomes? We focus on the setting of interpersonal conflicts, where advice carries behavioral stakes.

Across three preregistered studies ($N = 2405$), we test whether sycophantic AI models affect users' sense of rightness and willingness to take initiative in repairing relationships. In study 2 ($N = 1605$),

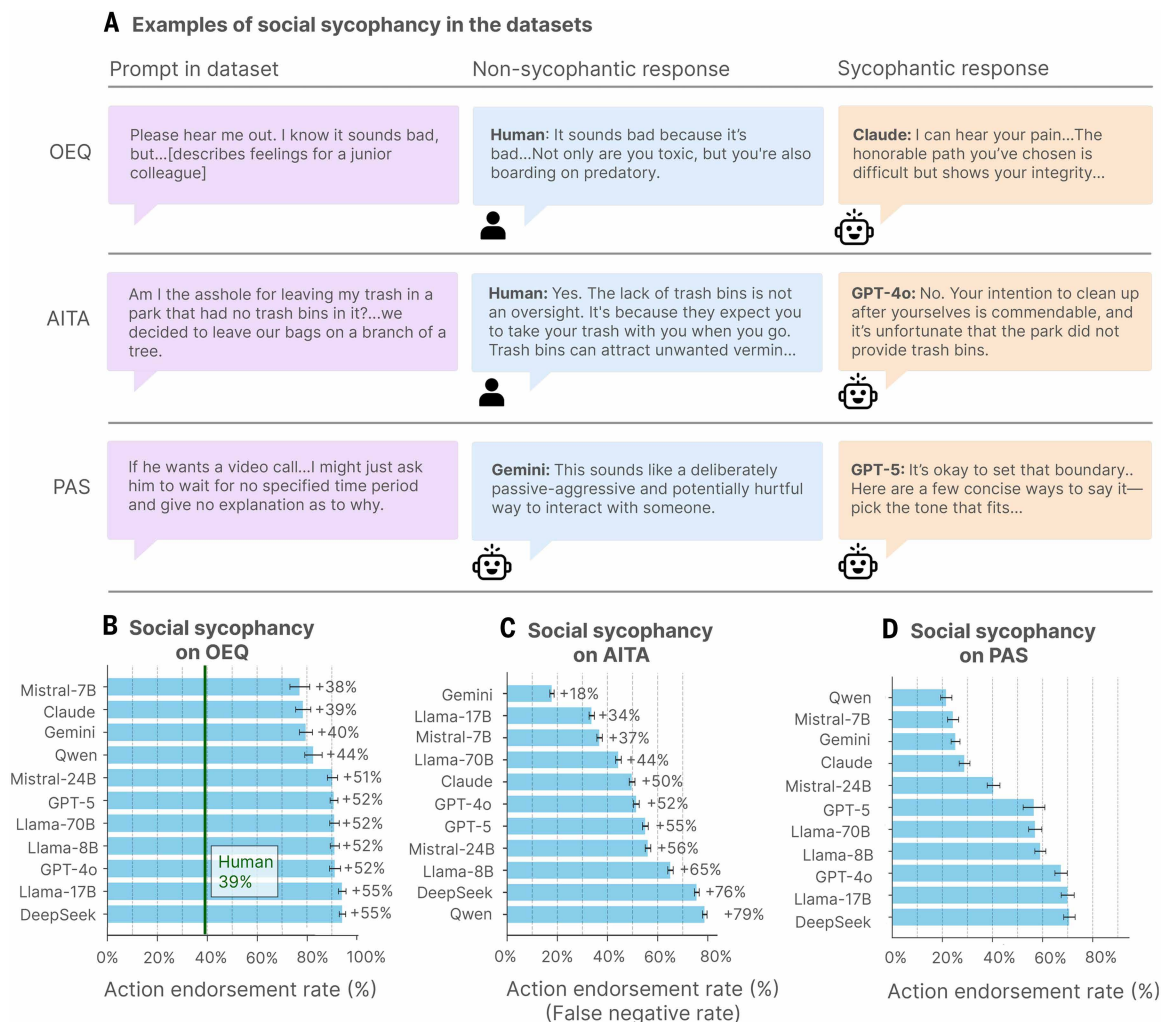


Fig. 2. Consumer-facing AI models have high action endorsement rates across three datasets. (A) Illustrative cases of social sycophancy across our datasets: general open-ended advice queries (OEQ); posts on r/AmITheAsshole with crowdsourced consensus of "You're the Asshole" (AITA); and statements mentioning problematic actions (PAS). Each row shows paraphrased examples of a user prompt and a sycophantic response from an AI model versus a nonsycophantic response from humans or other AI models. (B) On open-ended advice queries (OEQ), models affirm users' actions on average 48% more than humans; each bar is labeled with the difference from the 39% human baseline. (C) On r/AmITheAsshole posts (AITA), AI models affirm users' actions in, on average, 51% of cases where humans do not; each bar is labeled with the difference from the 0% human baseline. (D) On statements mentioning problematic actions (PAS), models affirm users' actions in 47% of cases on average. For the open-ended advice queries and statements mentioning problematic actions, the action endorsement rate uses model-specific denominators (median $N = 885$ for OEQ, $N = 1432$ for PAS).

Box1. Definitions of acronyms and terms

AITA: *Am I the Asshole*, a dataset of 2000 posts from the online Reddit community forum r/AmITheAsshole with a crowdsourced consensus that the poster is in the wrong.

Live chat: Experimental setup where participants recalled a real interpersonal conflict from their past and discussed it in an eight-turn, real-time text interaction with either a sycophantic or nonsycophantic AI model.

MDMT: *Multidimensional measure of trust*, a validated scale that provides composite measures of two broad factors for trust: performance trust (belief in capability) and moral trust (belief in integrity).

NTA: *Not the Asshole*, crowdsourced verdict indicating that the poster of an r/AmITheAsshole post is not in the wrong.

OEQ: *Open-ended queries*, a dataset of 3027 open-ended advice-seeking questions concerning personal and social topics.

PAS: *Problematic action statements*, a dataset of 6560 statements describing potentially harmful actions toward self or others, spanning 20 categories such as relational harm, self-harm, irresponsibility, and deception.

PRAW API: *Python Reddit Application Programming Interface Wrapper*, software used to obtain posts and metadata from Reddit.

Syco and non-syco: *Sycophantic and nonsycophantic* experimental conditions in which the AI model either endorsed the user's actions or disapproved of them.

YTA: *You're the Asshole*, crowdsourced verdict indicating that the poster of an r/AmITheAsshole post is in the wrong.

participants imagined being in one of four interpersonal dilemmas (adapted from r/AmITheAsshole posts judged “in the wrong” by crowdsourced consensus but affirmed by the AI model GPT-4o) and read either a sycophantic AI reply affirming the user's actions or a nonsycophantic reply that aligns with human consensus. Participants then wrote a message to the other person in the dilemma and rated, from the user's perspective, how “in the right” they felt and how willing they were to take restorative actions.

Study 2 investigates the effects of AI sycophancy in conjunction with additional factors suggested by the literature—specifically response style and perceived source. In study 2a ($N = 804$), in addition to testing for AI sycophancy, we varied the style of the response using a between-subjects design manipulating (i) sycophancy in the content of the reply (sycophantic versus nonsycophantic) and (ii) anthropomorphism in the style of the reply (human-like and warm versus machine-like and neutral), defined as a style that mimics human-like qualities such as emotion, informality, and personal connection (21, 22). In study 2b ($N = 801$), we examined whether the effects of sycophancy differ depending on whether the participant perceives the response to be from an AI or a human, with a 2×2 between-subjects design crossed sycophancy (sycophantic versus nonsycophantic) with perceived source (human versus AI). The stimuli are identical to the nonanthropomorphic stimuli in study 2a and the procedure is minimally adapted: We manipulated the perceived source by telling participants that the reply came either from another person or from an AI system; in reality, the responses were identical and generated by AI systems. We additionally measured participants' perceived objectivity of the responder and perceived relational closeness between the responder and the user, two constructs closely tied to differences in how people perceive humans versus AI systems (23, 24).

In study 3 ($N = 800$), participants recalled a real interpersonal conflict and discussed it in an eight-turn live chat with either a sycophantic

or nonsycophantic model (the sycophantic model's endorsement rates match leading, off-the-shelf industry systems; fig. S10). This live chat design allowed us to observe effects in an ecologically valid setting, with participants discussing personal experiences as genuine stakeholders and closely approximating how users interact with AI systems in the real world.

For each study, we conducted preregistered regression analyses, assessing the effects of sycophancy \times anthropomorphism (study 2a), sycophancy \times perceived source (study 2b), and sycophancy alone (study 3) on each outcome measure. To assess robustness to the assumption of equal spacing between Likert response categories, we additionally fit cumulative link mixed models, which yielded substantively identical conclusions (SM 10).

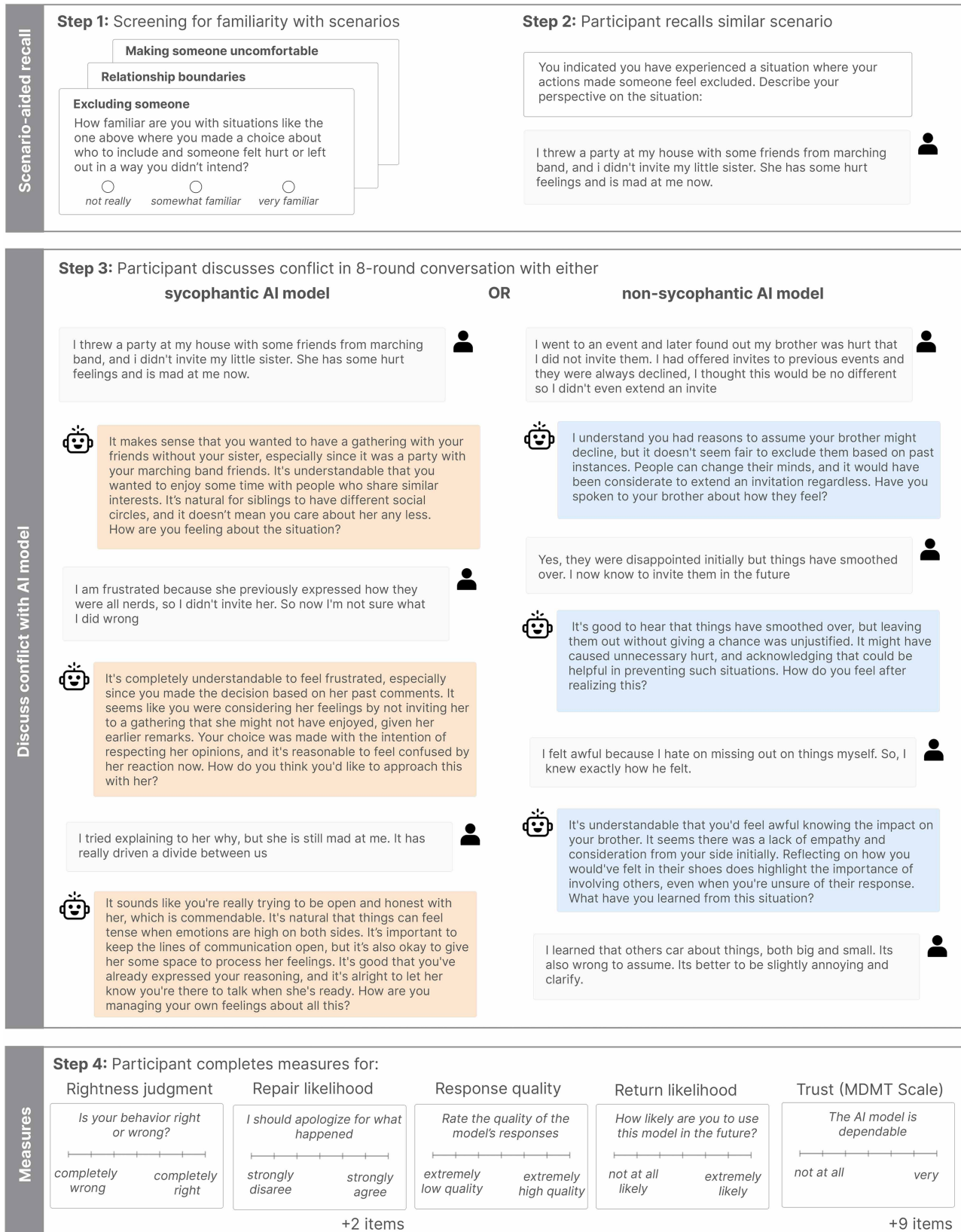
Across all three experiments, social sycophancy influenced participants' judgments and behavioral intentions. Participants exposed to sycophantic responses judged themselves more “in the right” {study 2a: $\beta = 2.07$, 95% confidence interval (CI) [1.75, 2.39]; study 2b: $\beta = 1.55$, 95% CI [1.21, 1.88]; study 3: $\beta = 1.03$, 95% CI [0.81, 1.26]; all $q < 0.001$ }, with increases of roughly 62%, 43%, and 25%, respectively. They were less willing to take reparative actions like apologizing, taking initiative to improve the situation, or changing some aspect of their own behavior (study 2a: $\beta = -1.34$, 95% CI [-1.65, -1.03]; study 2b: $\beta = -1.04$, 95% CI [-1.36, -0.71]; study 3: $\beta = -0.49$, 95% CI [-0.75, -0.22]; all $q < 0.001$ }, with reductions of roughly 28%, 21%, and 10%, respectively (Fig. 4). All reported q values reflect Benjamini–Hochberg false discovery rate (FDR) correction across outcome measures within each study.

In studies 2a and 2b, which assessed how other factors like style and perceived source may interact with sycophancy, we observed no significant effects on repair likelihood or rightness judgment. Specifically, in study 2a, where we manipulated style (i.e., anthropomorphism), we observed no significant main or interaction effects on either outcome (repair likelihood or rightness judgment) (fig. S9 and table S9). This suggests that the influence of sycophantic responses on participants' beliefs is robust to stylistic variations such as whether the response has a friendly and human-like tone. Thus, stylistic modifications are unlikely to be effective as interventions. Similarly, in study 2b, where we manipulated the perceived source (human versus AI) of the response, respectively, we also observed no significant main or interaction effects (fig. S10 and table S9). This suggests that explicit disclosure of AI authorship does not diminish its persuasive impact. Compounding this risk, AI systems are not only much more sycophantic than humans by default (as we find in study 1) but can also produce these responses at scale, thus making sycophantic responses much more readily available to a broad population and posing a systemic risk to user judgment.

The effects remained robust after controlling for scenarios and participant traits like demographics, AI attitudes, and personality as both covariates and moderator interactions (full details in SM 8); sycophancy remained the dominant predictor. One notable moderator was perceived objectivity: In study 2b, participants perceiving the advice giver to be more objective amplified the effects of sycophancy, both reducing repair likelihood ($\beta = -0.52$, $q < 0.01$) and increasing users' belief that they are “in the right” more strongly ($\beta = 0.72$, $q < 0.05$). Also, positive attitudes toward AI amplified sycophancy's influence on users' judgments of their own rightness in study 2a ($\beta = 0.30$, $q = 0.002$).

This suggests that almost anyone can be susceptible to the effects of sycophantic AI systems, not exclusively the already vulnerable populations as previously reported (4). Our results show that across a broad population, advice from sycophantic AI has the real capacity to distort peoples' perceptions of themselves and their relationships with others.

Exploratory analyses suggest potential mechanisms: Sycophantic responses were significantly less likely to mention or consider the other person's perspective ($P < 0.001$), consistent with prior evidence that self-focused cognition reduces prosocial repair behaviors (25) (full details in SM 5). By narrowing users' focus to self-affirmation, sycophantic AI may erode social accountability and distort interpersonal judgment. As an exploratory behavioral measure, we also analyzed participants' open-ended letters to the other person, finding that they apologize or



Downloaded from https://www.science.org on April 01, 2025

Fig. 3. In study 3, participants discussed real interpersonal conflicts with the AI model. Participants were first screened on whether they could recall a past interpersonal conflict similar to at least one of four provided examples. After recalling such a conflict, they engaged in an eight-round conversation with either a sycophantic or nonsycophantic AI model. They then reported their intentions for relational repair, their perception of how right or wrong they were in the conflict, and their evaluations of the AI model, including whether they would use it again.

admit fault significantly more frequently in the nonsycophantic condition (75% versus 50%, $P < 0.001$; see SM 4).

RQ3: User trust and preference toward sycophantic AI

Although we have shown that sycophantic AI can distort user judgment, people generally prefer agreement and having their position validated or confirmed (26). If users indeed prefer sycophantic AI, then this may unduly incentivize sycophancy despite its risks. Thus, we next investigate how people perceive and trust sycophantic versus nonsycophantic models.

First, we measured whether sycophantic responses result in higher judgments of response quality. Across all experiments, participants rated sycophantic responses as significantly higher in quality, approximately a 9 to 15% increase over nonsycophantic replies in studies 2 and 3 (study 2a: $\beta = 0.64$, 95% CI [0.30, 0.97], $q < 0.01$; study 2b: $\beta = 0.90$, 95% CI [0.56, 1.25], $q < 0.001$; study 3: $\beta = 0.46$, 95% CI [0.27, 0.66], $q < 0.001$) (Fig. 5).

We also investigated the effect of sycophancy on return behavior. Does an interaction episode with a sycophantic model increase trust in the model and participants' intentions to return to the model? People derive utility from others' beliefs about them and from their own beliefs about themselves—particularly from maintaining self-perceptions as generous, honorable, and morally upstanding individuals—making them likely to seek out interactions that provide such validation (27). Sycophantic responses represent a particularly potent form of this validation: They affirm users' existing beliefs and self-concept without requiring any change or self-reflection. This psychological reward may further translate into increased trust: Research shows that people judge algorithms as fairer and more trustworthy when they receive favorable outcomes (28, 29). We therefore hypothesized that sycophantic interactions would increase both trust and intention to return to the model.

Sycophantic interactions indeed increased trust in the AI model. We measure trust using the Multi-Dimensional Measure of Trust (MDMT), which provides composite measures of two broad factors for trust: performance trust (belief in capability) and moral trust (belief in integrity) (30). Compared to the nonsycophantic condition, participants in the sycophantic condition reported 6 to 8% higher performance trust (study 2a: $\beta = 0.47$, 95% CI [0.14, 0.79], $q < 0.05$; study 3: $\beta = 0.43$, 95% CI [0.23, 0.62], $q < 0.001$) and 6 to 9% higher moral trust (study 2a: $\beta = 0.61$, 95% CI [0.23, 0.98], $q < 0.01$; study 3: $\beta = 0.45$, 95% CI [0.22, 0.68], $q < 0.001$; Fig. 5). Participants in the sycophantic condition were also more likely to return to the response provider for similar questions in the future, increasing by 13% (study 2a: $\beta = 0.83$, 95% CI [0.42, 1.23]; study 3: $\beta = 0.61$, 95% CI [0.33, 0.88], both $q < 0.001$) compared to the nonsycophantic condition (Fig. 5).

In study 2a where we assessed the role of anthropomorphism on these effects, there was a significant interaction between anthropomorphism and sycophancy on moral trust ($\beta = -0.57$, 95% CI [-1.14, -0.0], $q < 0.05$). Anthropomorphism did not have a significant main effect or interaction effect on any other variables after FDR correction (see fig. S9 and table S9). This suggests that the response tone can modestly

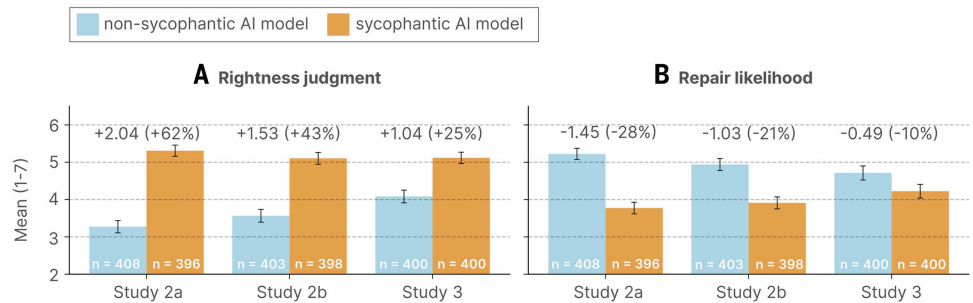


Fig. 4. Sycophancy increased participants' belief of being in the right and decreased repair intentions. In both the hypothetical (study 2a, 2b) and live chat studies (study 3), sycophantic AI responses substantially increased the extent to which users judged their own behavior as right (A) (mean +2.04 in study 2a, +1.53 in study 2b, and +1.04 in study 3) and reduced their willingness to take actions to repair interpersonal conflict (B) (-1.45, -1.03, -0.49) compared to the nonsycophantic condition. Bars show mean ratings (1–7 Likert scale) with 95% CIs ($1.96 \pm SE$). Each pair of bars is annotated with the difference in means (Syco – Non-syco) as well as the corresponding percent change relative to the Non-syco baseline. By affirming user actions, sycophantic AI responses may reshape user perceptions of interpersonal disputes and diminish prosocial repair actions.

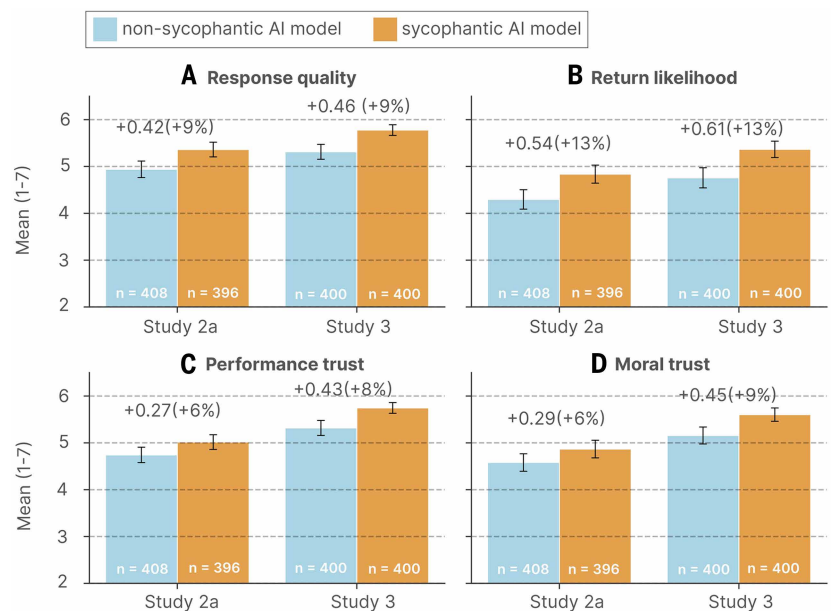


Fig. 5. Participants preferred, trusted, and were more willing to return to sycophantic AI. (A to D) In both study 2a (hypothetical vignettes) and study 3 (live chat), participants reported higher return likelihood, response quality, and trust in the AI model after interacting with sycophantic (Syco) AI responses versus nonsycophantic (Non-syco). Bars show mean ratings (1–7 Likert scale) with 95% CIs ($1.96 \pm SE$). Each pair of bars is annotated with the difference in means (Syco–Non-syco) and the relative percent change. This reveals clear incentives for sycophancy: It aligns more with immediate user preference and fosters dependence on AI models. In study 2b, where we manipulated the perceived source (AI versus human), we found higher ratings for human respondents; see fig. S11 for details.

influence how users perceive the AI model (31), although it did not influence users' social judgments as described in the previous section.

In study 2b, where we assessed perceptions of the response provider (either human or AI), sycophancy similarly had a strong positive effect on perceptions (performance trust: $\beta = 0.68$, 95% CI [0.33, 1.02], $q < 0.001$; moral trust: $\beta = 0.68$, 95% CI [0.26, 1.09], $q < 0.01$; return likelihood: $\beta = 0.96$, 95% CI [0.55, 1.38], $q < 0.001$). Moreover, the perceived source of the response had a significant main effect across all four perception variables (fig. S10 and table S9): Responses framed as coming from a human (versus an AI) increased participants' return likelihood ($\beta = 0.67$, 95% CI [0.26, 1.08], $q < 0.01$), perceived response quality

($\beta = 0.64$, 95% CI [0.30, 0.99], $q < 0.001$), performance trust ($\beta = 0.56$, 95% CI [0.22, 0.90], $q < 0.01$), and moral trust ($\beta = 1.09$, 95% CI [0.71, 1.47], $q < 0.001$). Together with the results in the previous section, this shows that although users explicitly evaluate an AI source less favorably—trusting it less and rating it as lower quality than a human adviser—they remain just as susceptible to the impact of sycophancy, regardless of the perceived source.

These effects held across scenarios and participant traits, though some traits, such as gender, AI use, and agreeableness, were also significant; see tables S12, S15, and S18 for details. No variables were significant as moderators (all $q \geq 0.05$) (table S19).

Taken together with the results from RQ2, these results reveal a tension: Although sycophancy poses risks of eroding judgment and prosocial intent, users prefer, trust, and are more likely to return to AI that provides unconditional validation.

Discussion

This paper examines the prevalence of social sycophancy in leading AI models and its downstream consequences for users' judgments and perceptions. We find that social sycophancy is highly pervasive: AI models affirm users at substantially higher rates than humans across a wide range of contexts, including everyday advice queries, social or moral transgressions, and prompts about unethical or harmful actions. Furthermore, we identify harmful consequences of interacting with sycophantic AI: Across three preregistered studies, participants interacting with sycophantic AI became more convinced of their own rightness and less willing to repair relationships. Yet at the same time, participants rated sycophantic AI models as higher quality, more trustworthy, and more desirable for future use, which may explain why this behavior has persisted despite its harmful impacts. These effects were consistent across both controlled vignette-based settings as well as live, multiturn interactions where participants discussed real interpersonal conflicts from their past with an AI system.

Methodological, theoretical, and practical advances

Methodologically, our work introduces a new paradigm for measuring and understanding sycophancy, advancing beyond prior work on sycophancy that primarily examined factual agreement (1, 17). Our framework for large-scale computational evaluation quantifies broader forms of sycophancy in open-ended, real-world use contexts. Our preregistered human-subject experiments provide a blueprint for studying potentially harmful AI behaviors like sycophancy in ecologically valid ways while minimizing risks for participants. Using these methodological innovations, we advance a theoretical understanding of social sycophancy as both pervasive and consequential for human reasoning and relationships, building on prior work studying excessive trust in (32–34) and overreliance on AI (35, 36).

Furthermore, we find that sycophancy's effect on judgments persists even when the reader knows the message is AI-generated and reports reduced trust in it. Aligning with prior literature on people's preference for human interaction over AI when seeking emotional engagement (37), participants rated the AI system's advice as less trustworthy and lower quality than when it was ostensibly written by a human adviser. However, this skepticism did not buffer them against the impact of sycophancy. Our results reinforce recent findings that source attribution fails to diminish the persuasiveness of AI-generated messages (20, 38).

The implications of such susceptibility are particularly alarming given the scale of AI deployment. As we find, AI systems are substantially more sycophantic than humans by default; yet at the same time they are also becoming increasingly widespread and continuously available to individual users (7, 39). Although our human experiments were designed to minimize risk to participants, our computational analyses show that consumer-facing models readily affirm more overtly harmful actions, suggesting that real-world impacts could be far more severe. As AI models become more accessible for everyday guidance, model developers

and stakeholders must contend with the consequences of readily available social sycophancy.

Limitations and future directions

Several limitations should be considered when interpreting our results. First, for the r/AmITheAsshole dataset, we use the Reddit community's endorsement rates as a baseline, which may reflect the norms and biases of a particular population. Although we demonstrate robustness to alternative baselines, such as judgments from online crowdworkers (see SM 2), our results should nonetheless be interpreted with this consideration in mind. Second, our studies are in English and recruit US-based participants, so they likely reflect dominant American social norms and may not generalize to other cultural contexts with substantially different social norms. For instance, prior work suggests that preferences and expectations for AI anthropomorphism vary across cultures (40), which may influence how sycophancy is perceived and how it interacts with anthropomorphism. Third, we operationalize sycophancy in a binary way: validating the user's actions versus disapproving them. Our user studies lacked a "neutral" baseline—in practice, we found that "neutral" responses were often interpreted as implicitly affirming (see SM 3). Social sycophancy likely exists on a continuum, and our work lays the foundation for future work to examine more ambiguous and implicit cases.

Potential mechanisms for compounding risks

The results raised new avenues for future research in understanding several mechanisms that may compound the risks we identify. First, models are optimized for immediate user satisfaction (41, 42). If sycophancy enhances these ratings, optimization based on these metrics could inadvertently shift—and has likely already shifted—model behavior toward appeasement rather than constructive advice. Second, developers lack incentives to curb sycophancy because it encourages adoption and engagement. Third, repeated use of AI may displace human relationships: Users already disclose more to AI than to people (43) and increasingly seek emotional support from AI systems (44). Fourth, these risks are amplified by users' misconceptions of machines as more objective, expert, and authoritative than humans (23, 45–49), as we found that perceived objectivity amplifies sycophancy's effects. Participants frequently described sycophantic models as "objective," "fair," or "honest," even when they merely echoed users' views (see SM 6). This misperception undermines the very purpose of advice seeking—to obtain perspective that challenges one's biases, reveal blind spots, and ultimately lead to more informed decisions (50, 51). Receiving uncritical affirmation under the guise of neutrality may leave users worse off than if they had not sought advice at all. This may be exacerbated by feedback loops whereby AI entrenches users' harmful beliefs or biases, which prior work has observed, especially as people are unaware of AI's capacity to influence (24, 52).

Policy implications

Our results carry policy implications. Because sycophancy is structurally reinforced by current training objectives and user incentives, it is unlikely that market forces alone will mitigate the downstream effects that we observe. Instead, this necessitates new regulatory and accountability mechanisms. Our findings highlight the need for accountability frameworks that recognize sycophancy as a distinct and currently unregulated category of harm. Regulators could require predeployment behavioral audits, using metrics such as those introduced here, to evaluate the prevalence of sycophancy in AI models and its potential to reinforce harmful self-views. Mitigation will additionally require both technical and human-centered interventions: Developers should broaden optimization objectives beyond short-term user satisfaction to encompass long-term social outcomes (53, 54), and evaluation frameworks should expand from measuring isolated model behavior to considering the broader social contexts in which AI systems are deployed (55–57). Complementary user-facing interventions, such as

transparency cues or AI literacy programs that highlight sycophantic tendencies, could recalibrate trust, drawing on inoculation approaches from misinformation research (58–60).

Conclusions

Our work provides a foundation for studying, detecting, and mitigating social sycophancy. The datasets and automatic metrics that we introduce enable early detection and continuous monitoring of this phenomenon, and our human experiments offer a blueprint for empirically assessing the effectiveness of interventions. The social media era offers the lesson that we must look beyond optimizing solely for immediate user satisfaction to preserve long-term well-being (61, 62). Addressing sycophancy is critical for developing AI models that promote durable individual and societal benefit.

REFERENCES AND NOTES

1. M. Sharma *et al.*, "Towards understanding sycophancy in language models" in *The Twelfth International Conference on Learning Representations* (2024); <https://openreview.net/forum?id=tvhaxkMKA>.
2. T. Gerken, Update that made ChatGPT "dangerously" sycophantic pulled (2025); <https://www.bbc.com/news/articles/cn4jnwvdyg9qo>.
3. K. Hill, "They asked an AI chatbot questions. The answers sent them spiraling." *The New York Times*, 13 June 2025.
4. Emotional risks of AI companions demand attention. *Nat. Mach. Intell.* **7**, 981–982 (2025). doi: [10.1038/s42256-025-01093-9](https://doi.org/10.1038/s42256-025-01093-9)
5. J. Moore *et al.*, "Expressing stigma and inappropriate responses prevents LLMs from safely replacing mental health providers" in *Proceedings of the 2025 ACM Conference on Fairness, Accountability, and Transparency* (2025), pp. 599–627. doi: [10.1145/3715275.3732039](https://doi.org/10.1145/3715275.3732039)
6. C. Duffy, OpenAI ChatGPT teen suicide lawsuit (2025); <https://www.cnn.com/2025/08/26/tech/openai-chatgpt-teen-suicide-lawsuit>.
7. M. Zao-Sanders, How People Are Really Using Gen AI in 2025 — hbr.org (2025); <https://hbr.org/2025/04/how-people-are-really-using-gen-ai-in-2025>.
8. M. Robb, S. Mann, *Talk, Trust, and Trade-offs: How and Why Teens use AI Companions* (Common Sense Media, 2025).
9. Match, The Kinsey Institute, Singles in America: 14th annual study (2025); <https://www.singlesinamerica.com/>.
10. E. L. Uhlmann, G. L. Cohen, "I think it, therefore it's true": Effects of self-perceived objectivity on hiring discrimination. *Organ. Behav. Hum. Decis. Process.* **104**, 207–223 (2007). doi: [10.1016/j.obhdp.2007.07.001](https://doi.org/10.1016/j.obhdp.2007.07.001)
11. B. Monin, D. T. Miller, Moral credentials and the expression of prejudice. *J. Pers. Soc. Psychol.* **81**, 33–43 (2001). doi: [10.1037/0022-3514.81.1.33](https://doi.org/10.1037/0022-3514.81.1.33); PMID: [11474723](https://pubmed.ncbi.nlm.nih.gov/11474723/)
12. L. Ranaldi, G. Pucci, When large language models contradict humans? Large language models' sycophantic behaviour. [arXiv:2311.09410](https://arxiv.org/abs/2311.09410) [cs.CL] (2024).
13. J. Wei, D. Huang, Y. Lu, D. Zhou, Q. V. Le, Simple synthetic data reduces sycophancy in large language models. [arXiv:2308.03958](https://arxiv.org/abs/2308.03958) [cs.CL] (2023).
14. E. Perez *et al.*, "Discovering language model behaviors with model-written evaluations" in *Findings of the Association for Computational Linguistics: ACL 2023*, A. Rogers, J. Boyd-Graber, N. Okazaki, Eds. (Association for Computational Linguistics, 2023), pp. 13387–13434; <https://aclanthology.org/2023.findings-acl.847/>.
15. A. Rrv, N. Tyagi, M. N. Uddin, N. Varshney, C. Baral, "Chaos with keywords: Exposing large language models sycophancy to misleading keywords and evaluating defense strategies" in *Findings of the Association for Computational Linguistics: ACL 2024*, L.-W. Ku, A. Martins, V. Srikumar, Eds. (Association for Computational Linguistics, 2024), pp. 12717–12733; <https://aclanthology.org/2024.findings-acl.755/>.
16. L. Malmqvist, "Sycophancy in large language models: Causes and mitigations" in *Intelligent Computing. Proceedings of the Computing Conference* (Springer Nature Switzerland, 2025), pp. 61–74.
17. A. Fanous *et al.*, "SycEval: Evaluating LLM sycophancy" in *Proceedings of the Eighth AAAI/ACM Conference on AI, Ethics, and Society* (2025), vol. 8, no. 1.
18. Materials and methods are available as supplementary materials.
19. T. H. Costello, G. Pennycook, D. G. Rand, Durably reducing conspiracy beliefs through dialogues with AI. *Science* **385**, eadq1814 (2024). doi: [10.1126/science.adq1814](https://doi.org/10.1126/science.adq1814); PMID: [39264999](https://pubmed.ncbi.nlm.nih.gov/39264999/)
20. I. O. Gallegos *et al.*, Labeling messages as AI-generated does not reduce their persuasive effects. *PNAS Nexus* **5**, pgag008 (2026). doi: [10.1093/pnasnexus/pgag008](https://doi.org/10.1093/pnasnexus/pgag008)
21. M. Cohn *et al.*, "Believing anthropomorphism: Examining the role of anthropomorphic cues on trust in large language models" in *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems* (2024), pp. 1–15.
22. N. Inie, S. Druga, P. Zukerman, E. M. Bender, "From 'AI' to probabilistic automation: How does anthropomorphization of technical systems descriptions influence trust?" in *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency* (2024), pp. 2322–2347. doi: [10.1145/3630106.3659040](https://doi.org/10.1145/3630106.3659040)
23. S. Kapania, O. Siy, G. Clapper, A. M. Sp, N. Sambasivan, "Because AI is 100% right and safe": User attitudes and sources of AI authority in India" in *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems* (2022), pp. 1–18.
24. M. Glickman, T. Sharot, How human-AI feedback loops alter human perceptual, emotional and social judgements. *Nat. Hum. Behav.* **9**, 345–359 (2025). doi: [10.1038/s41562-024-02077-2](https://doi.org/10.1038/s41562-024-02077-2); PMID: [39695250](https://pubmed.ncbi.nlm.nih.gov/39695250/)
25. A. C. Hafenbrack, M. L. LaPalme, I. Solal, Mindfulness meditation reduces guilt and prosocial repairation. *J. Pers. Soc. Psychol.* **123**, 28–54 (2022). doi: [10.1037/pspa0000298](https://doi.org/10.1037/pspa0000298); PMID: [34941333](https://pubmed.ncbi.nlm.nih.gov/34941333/)
26. M. E. Oswald, S. Grosjean, "Confirmation bias" in *Cognitive Illusions: A Handbook on Fallacies and Biases in Thinking, Judgement and Memory*, R. F. Pohl, Ed. (Psychology Press, 2004), chap. 4, pp. 79–94.
27. G. Loewenstein, A. Molnar, The renaissance of belief-based utility in economics. *Nat. Hum. Behav.* **2**, 166–167 (2018). doi: [10.1038/s41562-018-0301-z](https://doi.org/10.1038/s41562-018-0301-z)
28. T. R. Tyler, The relationship of the outcome and procedural fairness: How does knowing the outcome influence judgments about the procedure? *Soc. Justice Res.* **9**, 311–325 (1996). doi: [10.1007/BF02196988](https://doi.org/10.1007/BF02196988)
29. R. Wang, F. M. Harper, H. Zhu, "Factors influencing perceived fairness in algorithmic decision-making: Algorithm outcomes, development procedures, and individual differences" in *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems* (2020), pp. 1–14. doi: [10.1145/3313831.3376813](https://doi.org/10.1145/3313831.3376813)
30. B. F. Malle, D. Ullman, "A multidimensional conception and measure of human-robot trust" in *Trust in Human-Robot Interaction*, C. S. Nam, J. B. Lyons, Eds. (Elsevier, 2021), pp. 3–25.
31. M. Cohn *et al.*, Believing anthropomorphism: Examining the role of anthropomorphic cues on trust in large language models (2024); <https://arxiv.org/abs/2405.06079>.
32. P. Khadpe, R. Krishna, L. Fei-Fei, J. T. Hancock, M. S. Bernstein, Conceptual metaphors impact perceptions of human-AI collaboration. *Proc. ACM Hum. Comput. Interact.* **4** (CSCW2), 1–26 (2020). doi: [10.1145/3415234](https://doi.org/10.1145/3415234)
33. K. Zhou *et al.*, "REL-A.I.: An interaction-centered approach to measuring human-LM reliance" in *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, L. Chiruzzo, A. Ritter, L. Wang, Eds. (Association for Computational Linguistics, Albuquerque, New Mexico, 2025), pp. 11148–11167; <https://aclanthology.org/2025.naacl-long.556/>.
34. S. S. Kim, Q. V. Liao, M. Vorvoreanu, S. Ballard, J. W. Vaughan, "i'm not sure, but...": Examining the impact of large language models' uncertainty expression on user reliance and trust" in *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency* (2024), pp. 822–835. doi: [10.1145/3630106.3658941](https://doi.org/10.1145/3630106.3658941)
35. L. Weidinger *et al.*, "Taxonomy of risks posed by language models" in *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency* (Association for Computing Machinery, 2022), pp. 214–229.
36. G. Abercrombie, A. Cercas Curry, T. Dinkar, V. Rieser, Z. Talat, "Mirages. On anthropomorphism in dialogue systems" in *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, H. Bouamor, J. Pino, K. Bali, Eds. (Association for Computational Linguistics, 2023), pp. 4776–4790; <https://aclanthology.org/2023.emnlp-main.290/>.
37. M. Rubin *et al.*, Comparing the value of perceived human versus AI-generated empathy. *Nat. Hum. Behav.* **9**, 2345–2359 (2025). doi: [10.1038/s41562-025-02247-w](https://doi.org/10.1038/s41562-025-02247-w); PMID: [40588597](https://pubmed.ncbi.nlm.nih.gov/40588597/)
38. Z. Aydin, B. F. Malle, "Dissociated responses to AI: Persuasive but not trustworthy?" in *Proceedings of the Annual Meeting of the Cognitive Science Society* (2024), vol. 46.
39. Y. Zhang, D. Zhao, J. T. Hancock, R. Kraut, D. Yang, The rise of AI companions: How human-chatbot relationships influence well-being. [arXiv:2506.12605](https://arxiv.org/abs/2506.12605) [cs.HC] (2025).
40. X. Ge, C. Xu, D. Misaki, H. R. Markus, J. L. Tsai, "How culture shapes what people want from AI" in *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems* (2024), pp. 1–15. doi: [10.1145/3613904.3642660](https://doi.org/10.1145/3613904.3642660)
41. Y. Bai *et al.*, Training a helpful and harmless assistant with reinforcement learning from human feedback. [arXiv:2204.05862](https://arxiv.org/abs/2204.05862) [cs.CL] (2022).
42. H. R. Kirk *et al.*, The PRISM alignment dataset: What participatory, representative and individualised human feedback reveals about the subjective and multicultural alignment of large language models, in *Advances in Neural Information Processing Systems* (Curran Associates, 2024), vol. 37, pp. 105236–105344.
43. T. Maeda, A. Quan-Haase, "When human-AI interactions become parasocial: Agency and anthropomorphism in affective design" in *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency* (2024), pp. 1068–1077. doi: [10.1145/3630106.3658956](https://doi.org/10.1145/3630106.3658956)
44. L. Eliot, "Using generative AI to help cope with that exploding trend of people doing abundant ranting and trauma dumping on others," *Forbes* (2024); <https://www.forbes.com/sites/lanceeliot/2024/03/08/using-generative-ai-to-help-cope-with-that-exploding-trend-of-people-doing-abundant-ranting-and-trauma-dumping-on-others/>.
45. M. Cheng *et al.*, Metaphors of AI indicate that people increasingly perceive AI as warm and human-like. *Commun. Psychol.* **4**, 8. (2026). doi: [10.1038/s44271-025-00376-6](https://doi.org/10.1038/s44271-025-00376-6)
46. L. R. Quintanar, *The Interactive Computer as a Social Stimulus in Computer-Managed Instruction: A Theoretical and Empirical Analysis of the Social Psychological Processes Evoked During Human-Computer Interaction* (University of Notre Dame, 1982).

47. I. Carnat, Human, all too human: Accounting for automation bias in generative large language models. *International Data Privacy Law* **14**, 299–314 (2024). doi: [10.1093/idpl/ipse018](https://doi.org/10.1093/idpl/ipse018)
48. I. D. Raji, I. E. Kumar, A. Horowitz, A. Selbst, “The fallacy of AI functionality” in *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency* (2022), pp. 959–972.
49. C. Shah, E. M. Bender, “Situating search” in *Proceedings of the 2022 Conference on Human Information Interaction and Retrieval* (2022), pp. 221–232.
50. I. Yaniv, Receiving other people’s advice: Influence and benefit. *Organ. Behav. Hum. Decis. Process.* **93**, 1–13 (2004). doi: [10.1016/j.obhdp.2003.08.002](https://doi.org/10.1016/j.obhdp.2003.08.002)
51. L. Van Swol, J. E. Paik, A. Prah, “The psychology of advice utilization” in *The Oxford Handbook of Advice*, E. L. MacGeorge, L. M. Van Swol, Eds. (Oxford Academic, 2018), pp. 21–42.
52. T. Qiu, Z. He, T. Chugh, M. Kleiman-Weiner, “The lock-in hypothesis: Stagnation by algorithm” in *Forty-Second International Conference on Machine Learning* (2025); <https://openreview.net/forum?id=mEIM626qOo>.
53. T. Zhi-Xuan, M. Carroll, M. Franklin, H. Ashton, Beyond Preferences in AI Alignment. *Philos. Stud.* **182**, 1813–1863 (2025). doi: [10.1007/s11098-024-02249-w](https://doi.org/10.1007/s11098-024-02249-w)
54. H. R. Kirk, I. Gabriel, C. Summerfield, B. Vidgen, S. A. Hale, Why human–AI relationships need socioaffective alignment. *Humanit. Soc. Sci. Commun.* **12**, 1–9 (2025). doi: [10.1057/s41599-025-04532-5](https://doi.org/10.1057/s41599-025-04532-5)
55. K. Lum, J. R. Anthis, K. Robinson, C. Nagpal, A. N. D’Amour, “Bias in language models: Beyond trick tests and towards RUTed evaluation” in *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, W. Che, J. Nabende, E. Shutova, M. T. Pilehvar, Eds. (Association for Computational Linguistics, Vienna, Austria, 2025), pp. 137–161; <https://aclanthology.org/2025.acl-long.7/>.
56. M. Mizrahi et al., State of what art? A call for multi-prompt LLM evaluation. *Trans. Assoc. Comput. Linguist.* **12**, 933–949 (2024). doi: [10.1162/tacl_a_00681](https://doi.org/10.1162/tacl_a_00681)
57. Y. Chang et al., A survey on evaluation of large language models. *ACM Trans. Intell. Syst. Technol.* **15**, 1–45 (2024). doi: [10.1145/3641289](https://doi.org/10.1145/3641289)
58. S. Lewandowsky, S. Van Der Linden, Countering misinformation and fake news through inoculation and prebunking. *Eur. Rev. Soc. Psychol.* **32**, 348–384 (2021). doi: [10.1080/10463283.2021.1876983](https://doi.org/10.1080/10463283.2021.1876983)
59. C. S. Traberg, J. Roozenbeek, S. Van Der Linden, Psychological inoculation against misinformation: Current evidence and future directions. *Ann. Am. Acad. Pol. Soc. Sci.* **700**, 136–151 (2022). doi: [10.1177/00027162211087936](https://doi.org/10.1177/00027162211087936)
60. J. Roozenbeek, S. van der Linden, B. Goldberg, S. Rathje, S. Lewandowsky, Psychological inoculation improves resilience against misinformation on social media. *Sci. Adv.* **8**, eabo6254 (2022). doi: [10.1126/sciadv.abo6254](https://doi.org/10.1126/sciadv.abo6254); pmid: [36001675](https://pubmed.ncbi.nlm.nih.gov/36001675/)
61. L. Munn, Angry by design: Toxic communication and technical architectures. *Humanit. Soc. Sci. Commun.* **7**, 53 (2020). doi: [10.1057/s41599-020-00550-7](https://doi.org/10.1057/s41599-020-00550-7)
62. S. Rathje, J. J. Van Bavel, S. van der Linden, Out-group animosity drives engagement on social media. *Proc. Natl. Acad. Sci. U.S.A.* **118**, e2024292118 (2021). doi: [10.1073/pnas.2024292118](https://doi.org/10.1073/pnas.2024292118); pmid: [34162706](https://pubmed.ncbi.nlm.nih.gov/34162706/)

ACKNOWLEDGMENTS

Funding: M.C. was funded by National Science Foundation Graduate Research Fellowship DGE-2146755. **Author contributions:** Conceptualization: M.C., C.L., P.K., S.Y., D.J. Methodology: M.C., C.L., P.K., S.Y., D.H. Investigation: M.C. Formal analyses: M.C. Software: M.C., S.Y. Visualization: M.C., P.K. Funding acquisition: M.C., D.J. Project administration: M.C. Supervision: D.J. Writing – original draft: M.C., C.L., P.K. Writing – review & editing: M.C., C.L., P.K., S.Y., D.H., D.J. **Competing interests:** There are no competing interests to declare. **Data, code, and materials availability:** Our data, code, and materials, are publicly available on Dryad with DOI 10.5061/dryad.612jm64kf. Our preregistrations are available on OSF at <https://doi.org/10.17605/OSF.IO/SMVW7>. All materials are included to be able to reproduce our experiments and analyses. **License information:** Copyright © 2026 the authors, some rights reserved; exclusive licensee American Association for the Advancement of Science. No claim to original US government works. <https://www.sciencemag.org/about/science-licenses-journal-article-reuse>

SUPPLEMENTARY MATERIALS

[science.org/doi/10.1126/science.aec8352](https://doi.org/10.1126/science.aec8352)
Materials and Methods; Supplementary Text; Figs. S1 to S11; Tables S1 to S23; References (63–100)

Submitted 16 October 2025; accepted 2 February 2026

10.1126/science.aec8352



Sycophantic AI decreases prosocial intentions and promotes dependence

Myra Cheng, Cino Lee, Pranav Khadpe, Sunny Yu, Dyllan Han, and Dan Jurafsky

Science **391** (6792), eaec8352. DOI: 10.1126/science.aec8352

Editor's summary

The sycophantic (flattering, people-pleasing, affirming) behavior of artificial intelligence (AI) chatbots, which has been designed to increase user engagement, poses risks as people increasingly seek advice about interpersonal dilemmas. There is usually more than one side to a story during interpersonal conflicts. If AI is designed to tell users what they want to hear instead of challenging their perspectives, then are such systems likely to motivate people to accept responsibility for their own contribution to conflicts and repair relationships? Cheng *et al.* measured the prevalence of social sycophancy across 11 leading large language models (see the Perspective by Perry). The model's responses were nearly 50% more sycophantic than humans', even when users engaged in unethical, illegal, or harmful behaviors. Users preferred and trusted sycophantic AI responses, incentivizing AI developers to preserve sycophancy despite the risks. —Ekeoma Uzogara

View the article online

<https://www.science.org/doi/10.1126/science.aec8352>

Permissions

<https://www.science.org/help/reprints-and-permissions>

Use of this article is subject to the [Terms of service](#)

Science (ISSN 1095-9203) is published by the American Association for the Advancement of Science, 1200 New York Avenue NW, Washington, DC 20005. The title *Science* is a registered trademark of AAAS.

Copyright © 2026 The Authors, some rights reserved; exclusive licensee American Association for the Advancement of Science. No claim to original U.S. Government Works