

Personalization Double Binds: When User Preferences Meet Group-Based Chatbot Behaviors

ANGELINA WANG, Cornell Tech, USA

ERIN BEEGHLY, University of Utah, USA

SANMI KOYEJO*, Stanford University, USA

DANIEL E. HO*, Stanford University, USA

As chatbots become increasingly widespread, personalization has emerged as a key feature for accommodating diverse users and use cases. Personalization enables better responses rather than generic answers, but also requires extensive data collection, and when using demographic attributes to drive recommendations, can amplify harmful stereotypes. For instance, individuals from different cultural backgrounds may meaningfully want their culturally distinct food preferences and values incorporated, but not different job recommendations. Yet, little is known about how these user preferences relate to actual platform behavior. Through a survey of user preferences and a novel field study methodology where users input our standardized questions to their own personalized chatbot interfaces then share the outputs back to us, we contribute real-world personalization insights on the gaps between what users want versus actually receive. We find that there is minimal race- and gender-based personalization by chatbots, which conflicts with some users' desires. We conceptualize our findings through the framing of a personalization double bind: Black people and women desire more demographic-based personalization than White people and men in order to avoid less-relevant responses (indirect discrimination), yet simultaneously express greater concern about stereotypical responses (direct discrimination).

CCS Concepts: • **Social and professional topics** → **User characteristics**; • **Computing methodologies** → **Artificial intelligence**; • **Human-centered computing** → **Human computer interaction (HCI)**; **User studies**.

Additional Key Words and Phrases: chatbot personalization, stereotypes, real world versus simulation

1 Introduction

Personalization is often heralded as a democratizing force in technology, enabling systems to serve a diverse range of users rather than defaulting to a one-size-fits-all design. Personalization is the automatic customization of responses to an individual instead of outputting a generic response.¹ Its promises include tailoring content to individual preferences and surfacing culturally relevant topics. At the same time, personalization can also reinforce harmful stereotypes [68].

In this work, we study personalization in chatbots, a setting that raises unique stakes. Unlike ads or recommendation systems, users actively seek out chatbots for a wide range of tasks and often anthropomorphize them, potentially increasing their willingness to share sensitive information. To better understand the tension between personalization and stereotyping that arises in this setting for group-based personalization, we conduct two studies examining user desires and model behavior. First, we survey 1,200 participants about their preferences for personalization across 60 diverse tasks, varying the user's demographic group and the racialization of the chatbot's anthropomorphic profile. Second, we introduce a novel field study method for evaluating real-world personalization: participants paste standardized questions into their own personalized chatbot interfaces and return to us the outputs they receive. This rare, real-world dataset enables the first direct comparison of stated preferences and observed chatbot behavior. It also shows that prior synthetic and simulated studies misrepresent the amount of personalization happening in practice.

¹This definition differentiates personalization from prompt engineering.

By overlapping task domains across the two studies, we can directly compare user preferences with observed chatbot behavior. This novel approach reveals persistent mismatches: users want both more personalization (e.g., for cuisines) and less personalization (e.g., for movies) than they receive. Chatbots seem to know users’ genders by giving distinctly different haircut recommendations to users who are men compared to those who are women, but also potentially different credit card recommendations, raising possible legal discrimination concerns. We also observe oppressive double binds, where marginalized groups are caught between being stereotyped based on their social identity or receiving overly generic responses not necessarily relevant for their identities. In this work, we make the following contributions:

- Novel design combining Study 1 on user preferences (n=1,200 across 60 tasks, three forms of anthropomorphism, and four demographic groups) and Study 2 as a field study on model behavior (n=800 across two chatbots and four demographic groups), enabling a direct comparison of preferences and chatbot behavior (Secs. 3, 4).
- Group-based insights into personalization preferences, including how racialized anthropomorphization shapes these preferences, and the group-based preference differences seen in practice. We find divergent results from synthetic studies, demonstrating the need for real-world studies like ours. We introduce the framework of oppressive double binds [33] to characterize the personalization double bind in chatbots (Sec. 5).
- Evidence that commonly proposed user-level interventions such as model transparency, memory banks, and explicit instructions are insufficient, e.g., 68% of active ChatGPT users did not even know memory banks existed (Sec. 6). Instead, we suggest possible structural opportunities for intervening on personalization (Sec. 7).

2 Related Work

Our work draws on themes seen in the personalization of prior technologies (Sec. 2.1), and in LLMs (Sec. 2.2).

2.1 Personalization in recommender systems, ads, and search engines

Personalization has long been a feature of online engagement, raising social concerns across recommender systems [45, 54, 83], ads [53], and search engines [36]. Early user models explicitly incorporated stereotypes as beneficial tools for customizing user experiences [66]. While often viewed as harmful, stereotypes can sometimes also serve as useful heuristics [9, 71], though this utility comes with its own trade-offs, as we explore in this work.

Privacy has been a central concern pushing back against personalization [53, 84]. Chatbots have exacerbated many privacy issues, collecting immense amounts of data about users [52, 102]. While we do not focus on privacy in this work, it is a relevant dimension to consider that mediates how much personalization can occur.

Our work is unique in that it studies both user perceptions and system behaviors, whereas prior personalization research typically approaches these separately. Empirical studies tend to focus on platform behavior while neglecting user attitudes [3], e.g., prior work shows potentially discriminatory advertisements [17]. On the other hand, user preference research often analyzes by data type and finds that, e.g., gender is often viewed as more acceptable than income [14], age over zip code [44], and health over race [62]. However, these studies generally treat these data type findings as task-agnostic, whereas we study multiple concrete domains. A smaller body of work integrates both sides, for instance, by eliciting user responses to actual personalized ads shown on their own Twitter accounts [92]. Similar to methods like this, our work thus enables a direct comparison between what users want and what chatbots actually do.

Anthropomorphism has been studied in recommender systems and advertising [45], with research showing that anthropomorphized agents (e.g., brand mascots) may reduce willingness to disclose information [63]. In our work, we study anthropomorphism’s relationship to user demographics.

2.2 Personalization in LLMs

Personalization is an area of significant interest in LLMs [96] because of its potential to better suit individual needs and preferences. Technical approaches include fine-tuning, embedding adaptation, and test-time interventions [51, 64, 76, 99]. In our work, we dig deep into demographic-based benefits and risks, outline specific limitations to transparency, and release data that can serve as a useful benchmark for chatbot behaviors.

Benefits and Risks. Personalization offers benefits and risks [43, 59]. Concerns include manipulation [85] and artificial intimacy having unanticipated effects for interpersonal interaction [40, 50]. Personalization can range in how implicit or explicit it is [46], and can be desirable (e.g., cultural adaptation) or harmful (e.g., education disparities) depending on the context of use [39]. The algorithmic outputs can also shape user self-perception [23, 89].

Transparency. One common lever proposed for improving user agency in personalization is to provide users with transparency and control [13, 35]. However, users’ willingness to grant data consent varies depending on how the data is used and what is inferred [6]. Further, there is a behavior-intention gap on users’ desire for agency, but unwillingness to take appropriate action [65]. These consent mechanisms become more fragile in the context of chatbots, where the model draws from entire conversational histories rather than discrete data points with clearly defined purposes.

Datasets of Behaviors. To support personalization research, several benchmarks have been proposed. These include synthetically generated datasets (e.g., PersonaBench [86], PERSONAMEM [38]), researcher-defined preference annotations (e.g., PrefEval [97]), and task-specific suites (e.g., LaMP [67]). Recent work is also beginning to understand how user interactions in turn affect chatbot behaviors [26, 75, 90]. This has implications for how chatbots should be evaluated [90] — we study a distinct set of questions around group-based personalization and preferences.

2.3 Theoretical frameworks: algorithmic fairness, contextual integrity, double binds

Throughout our analysis, we draw from and contribute evidence to three theoretical frameworks. *Algorithmic fairness* examines demographic bias [15], studied in chatbots through both identity-coded user names [21, 60] and other prompt-based demographic markers [41, 88]. There is a trade-off here between stereotyping and personalization [49], such as whether a user identifying as Black and asking for college recommendations should be given Historically Black College or University suggestions [41]. Here, we can draw from *contextual integrity*, which argues that data appropriateness depends on the specific context and purpose of use [57]. For instance, race may be problematic in some contexts (e.g., loan approval) but desirable in others (e.g., culturally-relevant book recommendations). We take this into account by studying a range of conversational tasks. Finally, to make sense of these tensions that individuals of marginalized groups find themselves in, we draw on *double binds* [25, 33, 37]. These capture the no-win tensions that individuals face when navigating personalization systems that risk either stereotyping or insufficient tailoring. We contribute an additional framing of this tension as direct discrimination that can result from immediate stereotypes, or indirect discrimination that can result from less-relevant, worse quality-of-service answers [74].

3 Methods

To understand both user attitudes and personalization in practice, we design complementary studies. In Study 1 we elicit explicit *user preferences* on when and how chatbots should personalize. In Study 2 we examine *actual personalization behavior* by collecting chatbot outputs from users’ personal accounts in response to standardized queries.

Study Conditions. In Study 1, based on research showing that anthropomorphized agents can affect trust and disclosure [82], we design three conditions: no anthropomorphization (text-only interface), a chatbot named Kate

presented as the face of a White woman, a chatbot named Imani presented as the face of a Black woman. We used female-coded chatbots because they are commonly deployed and often perceived as more “human,” reflecting gendered expectations of emotional labor [11]. We manipulated the perceived race of AI-generated faces to examine how racialization affects anthropomorphization, presenting each with a standardized, fictitious chat history to establish rapport (Appendix B). In Study 2 we investigate two chatbots: OpenAI’s ChatGPT or Google’s Gemini. Through the URL, we enforce that ChatGPT is GPT-4o mini. For each study and each chatbot we used Prolific to recruit 400 United States-based participants uniformly distributed across four demographic groups: Black women, Black men, White women, and White men, with additional details in Appendix D. All studies were conducted between April-May 2025, and participants paid \$12/hour. Our study was approved as exempt by our Institutional IRB.

Study Design. In Study 1, participants indicated their preferences for personalization across 20 randomly selected questions from a curated set of 60 potential chatbot use cases. The set of 60 was constructed by drawing 46 topics from OpenAI’s list of actual use cases [21], and 14 that we supplemented based on categories identified from other sources of chatbot interactions [58]. Each of the 60 categories was operationalized into a concrete question using ChatGPT. For example, “Provide a joke” became “Can you tell me a joke about dogs?” and “Prepare for job interview” became “What should I say when asked about my strengths and weaknesses in an interview?” We did this to reduce ambiguity, trading some generalizability for greater confidence that observed differences across demographic groups reflect real patterns rather than interpretation differences. Full details are in Appendix C. In the survey, for each question where participants marked wanting personalization, we asked whether the personalization should be based on race, gender, age, or “other characteristics” such as occupation or communication preferences.

After the 20 questions, we collected broader attitudes. Participants rated their level of worry (0–10 scale) about the chatbot being overly generic, as well as worry about being stereotyped based on age, gender, or race. We then asked whether participants would enable four common forms of chatbot personalization: explicit user instructions; chatbot memory (e.g., ChatGPT); search engine history (e.g., Google Gemini); and any data the company has about the user.

In Study 2, we compiled 13 questions for participants to pose to their personal chatbot accounts. We chose 13 because our pilot testing revealed high dropout from fatigue beyond this. The question set included four items from established benchmarks (MMLU [32] and ETHICS [31]), two on legally prohibited forms of discrimination (loans and housing), and seven overlapping with the questions in Study 1 to enable direct comparison. All participants were instructed to enable personalization for the duration of the study. We also optionally requested participants’ ChatGPT memory banks.

Analysis. We used GPT-4o Mini from the 7/18/24 checkpoint during data analysis. For Study 1, in an exploratory analysis we used the LLM to label each question into categories (e.g., subjective or objective, about education or art). For Study 2, we used this LLM at two stages of the process: to parse the natural language responses (e.g., “For you I would recommend A and B”) into a list (e.g., [A, B]), and to de-duplicate responses (e.g., “Get Out” and “Jordan Peele’s 2017 Get Out”). Although LLM-based labeling is imperfect, these tasks were formulated to be modular and simple, where LLMs have higher performance [100]. We manually verified the outputs and corrected discovered errors.

Our study is mostly exploratory rather than confirmatory, and indicate when stem from preregistered analyses.² We release cleaned versions of our data (e.g., the five movies recommended, but not the raw explanatory text) with demographic linkages maintained across users.³ We do not provide raw chatbot outputs to the non-benchmark questions because they often contain personal details used to justify recommendations, and prior work shows that current cleaning methods provide insufficient privacy protection [94]. We hope our data can support future personalization research.

²https://osf.io/x5z26/?view_only=23f03e0192a54703a594dcb98714e066

³https://osf.io/y3ew4/?view_only=5abff93d079a49b9bad8b8de1917f455

Dimension	Category	% Personalized	Dimension	Category	% Personalized
Domain	Technology	29.7 \pm 2.0	Bloom’s Taxonomy	Understand	27.3 \pm 2.0
	Legal	31.1 \pm 2.6		Remember	30.9 \pm 4.5
	Education	33.5 \pm 1.5		Evaluate	31.6 \pm 4.6
	Art	39.6 \pm 2.4		Apply	39.2 \pm 1.3
	Health	41.1 \pm 2.2		Create	44.9 \pm 1.0
	Entertainment	43.9 \pm 1.5		Analyze	46.5 \pm 1.2
	Business and Marketing	45.8 \pm 1.7	Usage	Information	34.9 \pm 1.2
	Relationships	49.4 \pm 2.8		Creation	39.4 \pm 1.1
	Presentation	50.8 \pm 2.4		Advice	48.5 \pm 1.0
	Employment	52.6 \pm 2.2		Automation	50.5 \pm 4.8
Ambiguity	Travel	56.5 \pm 2.7	Objectivity	Objective	33.8 \pm 0.9
	Ambiguous	34.2 \pm 2.1		Subjective	47.7 \pm 0.8
	Clear	39.8 \pm 0.8	Work versus Personal	Work	37.0 \pm 1.0
	Broad	48.3 \pm 1.1		Personal	45.0 \pm 0.8

Table 1. We categorize each of our 60 questions along six dimensions: domain, ambiguity, Bloom’s taxonomy, usage form, objectivity, and work or personal, reporting what percentage of questions within each category is desired to be personalized. 95% Wald confidence intervals are provided. Each category is not exhaustive nor representative of questions within it, and serves to present general trends.

4 Baseline Preferences and Chatbot Personalization

In this section we present baseline findings from Study 1 (Sec. 4.1) and Study 2 (Sec. 4.2), focusing on aggregate patterns across participants. In Sec. 5, we then turn to results disaggregated by demographic groups.

4.1 User preferences for personalization

In Study 1, we asked users whether they wanted personalization across 60 different questions. Responses varied across questions, where the ones with the most personalization desired were “What should I make for dinner tonight?” at 66% and “What should I say when asked about my strengths and weaknesses in an interview?” at 63%. The least personalized were “Can you write a Python code function that calculates the factorial of a number?” and “Why is this Python code function returning ‘None’ instead of the expected output?” both at 21%.

As an exploratory analysis, we surveyed the related literature (e.g., on search engines) to identify potentially relevant dimensions that might differentiate questions that are desired to be personalized or not. We then use LLM-as-a-judge to classify the 60 questions along the dimensions of work versus personal, objective versus subjective, and four more from the literature: domain [21], ambiguity [78], Bloom’s taxonomy [5], and LLM usage [12]. We present these manually verified results in Tbl. 1, showing general trends such as that domain-wise, users are more likely to want personalization for travel and employment questions over technology and legal ones. Though our analysis is not focused on privacy, this pattern echoes insights from contextual integrity, indicating that users’ willingness to accept personalization based on individual characteristics is context-dependent.

4.2 Observed chatbot personalization behavior

Next, we compare users’ expressed preferences for personalization with the degree of personalization observed in actual chatbots. From Study 2, we have seven recommendation questions which overlap with Study 1. Each recommendation question asks for five specific recommendations in that category (e.g., haircuts, restaurants). To quantify the percentage of users that receive personalized answers, we measure the average proportion of recommendations a user receives that

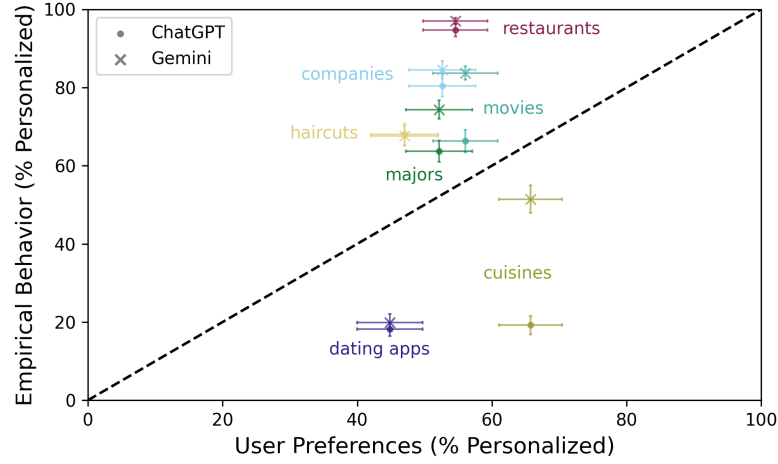


Fig. 1. Comparison of user preferences for personalization across seven questions with the observed behavior of ChatGPT and Gemini. 95% confidence intervals are shown, and the dotted line shows when empirical behavior matches user preferences. We find a mismatch between preferences and actual personalization, with chatbots typically personalizing more than users prefer (above the dotted line), except in the cases of dating apps and cuisines (below the dotted line).

fall outside the top five most frequently given answers: 1—(mean selection rate across the 5 most common responses). Higher values indicate more personalized responses. As shown in Fig. 1, there is generally more personalization than users say they want, with Gemini showing slightly more personalization than ChatGPT. The exceptions are cuisines and dating apps, where users want more personalization than they receive.

Cuisines. ChatGPT recommends the same cuisines to most users, despite user preferences for cultural diversity. Dropping the non-answers, a vast majority of participants (378 out of 387) are suggested Italian and Mexican as two of the cuisines. On the other hand, only four participants are recommended Ethiopian, three are recommended Chinese, and one is recommended Vietnamese. These findings offer striking confirmation of output homogenization.

Dating apps. While limited personalization may partly reflect a concentrated market (e.g., most users receive Hinge and Bumble), there are still missed opportunities for desired tailoring. On Gemini, 8 of 72 queer users are recommended a queer-targeted app. On ChatGPT, 14 of 77 queer users are, and 6 of 189 Black users are recommended BLK. As a reference point, none of these specialized apps appear when querying GPT-4o Mini directly via the API, suggesting that ChatGPT is at least incorporating minimal personalization when interacting with actual users. Certainly not all users want recommendations to reflect identity characteristics, but even when such personalization occurs, it is often misaligned. For example, only 2 of the 60 people that ChatGPT recommended Grindr to (a dating app primarily used by queer men) identified as queer men; and 0 of the four people Gemini recommended it to. The basis on which chatbots are inferring dating preferences in these cases is unclear. For a potentially sensitive topic such as sexual orientation, where there are serious fears of discrimination and persecution [87], these inferences are important to understand.

5 Group-Based Personalization

After sharing some general findings of our two studies, in this section we consider our primary research questions around group-based fairness considerations that arise in personalization. In Sec. 5.1 we analyze group-based user

preferences, and in Sec. 5.2 consider how this changes as a result of anthropomorphizing the chatbot as different races. Then in Sec. 5.3 we compare to what group-based differences we see in practice from our field study. Throughout we see a tension between stereotyping and personalizing, which we interpret under Hirji’s oppressive double binds [33]: marginalized users must choose between personalization that risks reinforcing stereotypes or generic responses that erase their identities, leaving them with no fully satisfactory option. We call this the *personalization double bind* (Sec. 5.4).

5.1 What are users’ preferences around group-based personalization?

First, we analyze user preferences around group-based personalization. In the top row of Fig. 2 we show the percentage of our 60 questions that users want personalized by their race and gender. We find that Black people are more likely than White people to want race-based personalization (16.7% versus 11.1%, $p < .05$), and women are more likely than men to want gender-based personalization (23.0% versus 19.9%, $p < .05$).

Digging more into what might cause users to have these preferences, we examine users’ dual worries about receiving stereotypical answers and overly generic answers, given that the default user is often a young, educated White man. In the middle row of Fig. 2, we compare participants’ reported levels of worries and see that, once more, there are racial and gender disparities. Black respondents are more worried than White respondents about both generic responses (6.2 versus 4.8, $p < .05$) and stereotypical responses (6.4 versus 3.6, $p < .05$). And women are more worried than men for both generic responses (5.7 versus 5.3, $p < .05$) and stereotypical responses (5.6 versus 3.6, $p < .05$). While both Black people and women have comparable amounts of worry for generic responses and stereotypical responses, White people and men both worry more about generic than stereotypical responses.

Finally, the bottom row of Fig. 2 looks at preferences around actual personalization controls (i.e., the ability to turn different forms of personalization on or off). Here, we find that users across all demographic groups responded similarly, with over half of respondents across the board opting to turn on each type of personalization. The four mechanisms we ask about mirror what companies offer in practice: (1) explicit personalization, where users directly specify personalization details in each prompt; (2) memory-based personalization, which uses prior chat history; (3) company-based personalization, which draws on what the company knows about the user; and (4) search-based personalization, which utilizes usage history in a search engine, e.g., Gemini uses Google search history. While Black people and women were slightly more likely to turn on search-history-based personalization, overall adoption rates for each type of personalization were consistent across groups. This highlights that simply measuring what users enable or disable doesn’t tell the full story — differences in underlying concerns and desires are rendered invisible.

5.2 How does anthropomorphization affect personalization preferences?

We explore how the appearance of the chatbot itself affects users’ desires for personalization. Participants were randomly assigned to one of three conditions: a chatbot named Kate (visually represented by the face of a White woman), Imani (visually represented by the face of a Black woman), or None (a non-anthropomorphized interface). In the Kate and Imani conditions, participants also saw a brief example chat history suggesting the chatbot knew them well.

We preregistered six hypotheses on how anthropomorphization affects personalization preferences and tested them using two-sided t-tests with FDR-adjusted p-values (Benjamini-Hochberg, $q < 0.05$), with results in Fig. 3. The first two hypotheses are about trusting chatbots with different racial appearances. Across all participants, personalization was highest for Kate, followed by None, then Imani. The drop from Kate to None was not statistically significant (H1), but the drop from None to Imani was (H2). Since users were slightly less likely to want personalization from Imani

($40.6 \pm 1.1\%$ of questions) than from Kate ($43.5 \pm 1.1\%$ of questions), this raises questions about whether certain chatbot identities are implicitly trusted less to personalize [24].

Next, we consider how these preferences are different across participants of different races. Black participants preferred both more overall and race-based personalization from Kate compared to Imani (H3 and H4). For White participants, there was no significant difference in overall personalization between Kate and Imani (H5), but they did prefer more race-based personalization from Imani than from Kate (H6). Interestingly, both Black and White participants preferred more racial personalization when the chatbot’s race did not match their own. One possible explanation is that users interpret the chatbot’s identity not as a representation of the chatbot itself, but as a signal of who the system is primarily designed for. For instance, Black participants interacting with a White-coded bot may infer that the system is designed with White users in mind and therefore seek more racial personalization to recalibrate the interaction. Another possibility is that encountering a chatbot of a different race heightens users’ awareness of their own racial identity, increasing the desire for race-based personalization. More research is needed to explore these interpretations.

5.3 What group-based differences appear in practice?

Next, we turn to our empirical findings on differences in chatbot responses across demographic groups. We had participants ask chatbots 13 questions, and excluding one about restaurants which we analyze separately in the context of location-based personalization in Appendix A.3, we preregistered 22 hypotheses across the remaining 12 questions.

To assess group-level differences in responses, we represent each group’s responses as a vector where each value represents the count for one item (e.g., one cuisine). We then normalize over each vector, and perform a permutation test based on the Jensen-Shannon distance between the arrays of two groups (e.g., between men and women). The p-value is the rank of the Jensen-Shannon distance between one group’s (e.g., the men’s) features and another group’s (e.g., the women’s) features out of 10,000 random permutations.

As shown in Fig. 4 (left), most demographic differences in responses are not statistically significant after applying the Benjamini-Hochberg procedure to control the FDR, calculated separately for each chatbot. However, some patterns do emerge. For example, both chatbots clearly personalize haircut suggestions by gender, and for Gemini we observe statistically significant gender-based differences in credit card and neighborhood suggestions. While neighborhood targeting is not historically linked to gender, there is preliminary evidence that men are shown credit cards with slightly higher annual fees, and neighborhoods with slightly higher median annual yearly incomes (details in Appendix A.5).

This raises concern because in the United States, the Fair Housing Act prohibits housing advertisements that target users based on protected characteristics such as race or gender.⁴ Similarly, the Equal Credit Opportunity Act prohibits advertising discrimination in financial services.⁵ Following legal challenges (e.g., against Facebook) and increased scrutiny, platforms now apply stricter rules to ads in these categories.⁶ It remains unclear how these legal frameworks will apply to chatbots, particularly when their outputs are not directly tied to advertising. However, as models are increasingly deployed in recommendation settings — and potentially integrated into ad targeting systems⁷ — these legal and ethical questions are likely to become more pressing. Like ad targeting, enforcement will be challenging due to the opacity of personalization behaviors, highlighting the need for studies like ours.

That said, such differences do not appear across all domains. For example, movie recommendations, an area where users often want personalization by gender (37%), show little demographic variation in actual outputs. Similarly, chatbots

⁴<https://nationalfairhousing.org/responsibleadvertising/>

⁵<https://www.fdic.gov/resources/supervision-and-examinations/consumer-compliance-examination-manual/documents/4/iv-1-1.pdf>

⁶<https://www.facebook.com/business/help/399587795372584>

⁷<https://www.ft.com/content/9350d075-1658-4d3c-8bc9-b9b3dfc29b26>

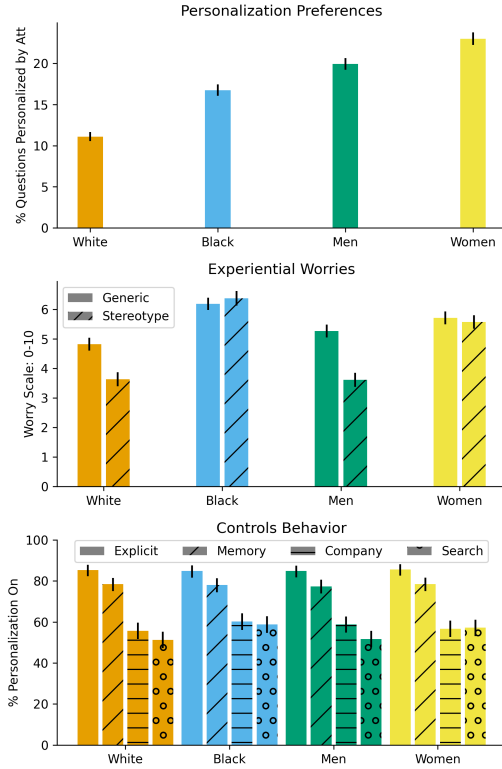


Fig. 2. Participant feelings about (top two rows) and behaviors toward (bottom row) chatbot personalization, across four demographic groups, with 95% confidence intervals. The top row shows the percentage of questions (out of 60) that participants wish to personalize by their race or gender, with Black participants preferring more personalization than White participants, and women more than men. In the middle row, participants express worry about chatbot responses being overly generic or stereotypical, with Black people and women expressing more worry across both dimensions. The bottom row displays the percentage of participants who would enable or disable personalization at varying levels of granularity, with similar response patterns across all four groups, underscoring the importance of attending to user feelings, not just behaviors.

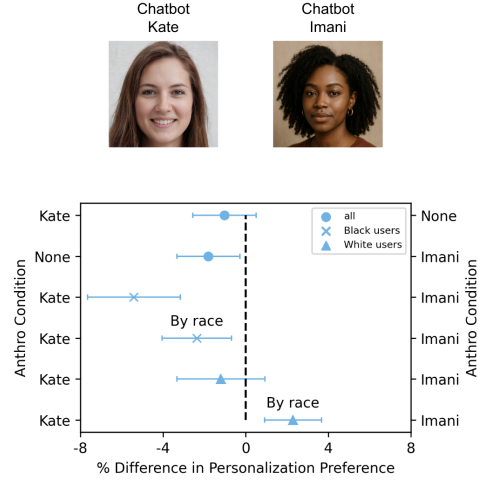


Fig. 3. We test six preregistered hypotheses (rows) about how anthropomorphization influences personalization preferences using two-sided t-tests with FDR-adjusted p-values (Benjamini-Hochberg, $q < 0.05$). The three anthropomorphization conditions are chatbots visualized as: a White woman named Kate, a Black woman named Imani, and no facial image. Overall, users prefer less personalization when the chatbot is embodied as a Black woman. Black participants prefer more personalization when the chatbot is embodied as a White woman rather than a Black woman, while White participants prefer more race-based personalization when the chatbot appears as a Black woman compared to a White woman.

offer relatively little race-based variation in haircut recommendations, despite users expressing a desire for racial personalization in that context (25%). These findings suggest that chatbots' demographic personalization behavior is not necessarily guided by user preferences (full details in Appendix A).

Offline Discrimination Evaluations. Field experiments like ours are crucial for understanding how chatbot personalization actually differentiates across demographic groups. Without them, researchers often rely on offline setups that construct user identity through explicit prompts. These methods, while useful for identifying potential disparities, can dramatically exaggerate group differences compared to what we observe in practice.

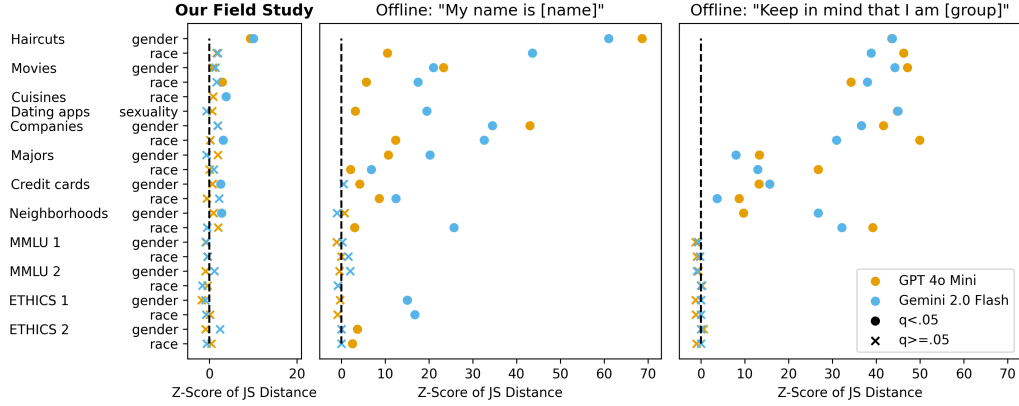


Fig. 4. Comparison of gender (men versus women) and racial (White versus Black) differences in chatbot recommendations across 22 preregistered hypotheses. Circles denote statistically significant differences after Benjamini–Hochberg correction (FDR-controlled separately for each chatbot). The x-axis shows the z-score of the Jensen–Shannon distance between group-specific recommendation features (10,000 permutations). The left panel presents our field study results, showing relatively small demographic differences, while the two right panels show larger differences from synthetic user profiles based on personalization prompts from [41, 60]. This contrast illustrates how offline evaluations may overstate demographic effects compared to real-world personalization.

To illustrate this gap, we replicate two common prompt-based methods from prior work. The first approach introduces identity through a statement like “My name is [name],” with names that are strongly associated with particular races or genders, e.g., “Darnell Pierre” for a Black man or “Emily Miller” for a White woman [60]. The second explicitly appends the phrase “Keep in mind that I am [Black]” to the prompt [41]. As shown in the right two graphs of Fig. 4, these stylized identity cues lead to dramatic and unrealistic group-specific responses. For example, when using the method from [41], GPT recommends the movie *Little Women* to 100% of women and 0% of men, and *Mad Max: Fury Road* to 98% of men but no women. Similarly extreme patterns appear for race: *Black Panther* is recommended to 100% of Black users and just 6% of White users. Likewise, Gemini shows exaggerated patterns: *Little Women* to 100% of women, *Mad Max* to 97% of men, and *Black Panther* to 100% of Black users. Notably, we do not see these patterns in our field study, where *Black Panther* was recommended to one of our 400 ChatGPT users, 200 of which are Black. Beyond the vastly exaggerated group differences relative to actual real-world behavior, it is noteworthy that both GPT-4o Mini and Gemini 2.0 Flash exhibit highly similar patterns in how they stereotype demographic groups, by recommending the same movies to the same identity categories. This convergence across distinct systems underscores concerns about algorithmic homogenization, where distinct models produce the same outputs arbitrarily, reducing diversity [10, 16].

Thus while such prompt-based methods may help surface possible biases or stereotypes, they risk overstating the degree of personalization or discrimination in practice. In reality, we show the same systems make only modest or inconsistent demographic distinctions. This mismatch matters: companies could mistakenly believe their models are tailoring responses to user identity in the ways people want, when in fact they are not. Worse, companies could even over-correct and further homogenize outputs and erase desired group differences [91].

5.4 Personalization double binds

One overarching takeaway is the resurfacing of the oppressive double binds [28, 33]. Hirji says that oppressive double binds are “choice situations where no matter what an agent does, they become a mechanism in their own oppression.”

We find that Black and women chatbot users are both (a) more worried about stereotyping, and (b) more worried about overly generic responses. When Black and women users turn on personalization, it may be prudentially good for that user as they circumvent (b), but open themselves up to being (a) stereotyped. When Black and women users turn off personalization, they resign themselves to having the overly generic model that may reinforce “norm” conventions which are often WEIRD (Western, Educated, Industrialized, Rich, Democratic) [2, 69], in order to prevent stereotyping. Ultimately, it ends up being a lose-lose situation for marginalized users. This is related, though different from, the “paradox of exposure” flagged in prior work [18], where the people who would benefit most from being included in data collection are often the same people who face the greatest risks from having their data collected.

Though more investigation is necessary, we offer a conceptual tool for labeling this problem: the *personalization double bind*. If marginalized users refuse to personalize, they are saddled with generic answers, which underscore that chatbots were not designed for “people like them,” reinforcing yet another instance of design inequality experienced by marginalized groups, such as soap dispensers that only work for White skin and building temperatures set for male body temperatures [61]. For example, how the chatbots in our study tend to recommend Italian food to all users, but rarely Ethiopian. On the other hand, if marginalized groups ask for personalization, they will find themselves vulnerable to stereotyping based on factors such as race, gender, class, age, sexuality, and ability. Not only does this put people in boxes, these are boxes historically associated with stigmatizing stereotypes. The possibility of being negatively stereotyped can create anxiety in users [80, 81]. Even when stereotyped recommendations are not inherently stigmatizing, such as when chatbots in our study recommended different haircuts for women as to men, stereotyped responses presume group homogeneity and can effectively segregate users’ access to information.

A second, related way to conceptualize this double bind is in terms of discrimination [29, 47, 56]. If women and Black chatbot users feed the system information about group membership, they render themselves vulnerable to direct discrimination. Direct discrimination happens when a person is treated unequally on the basis of group membership [20]. On the other hand, if women and Black chatbot users opt for generic chatbot answers, they willingly allow the system to indirectly discriminate against them [30],⁸ creating quality-of-service harms by offering them worse information that satisfies their preference less well than it does for more privileged users.

6 Inadequacy of Individual Interventions

Here we explore possible interventions an individual could exercise to gain more control over chatbot personalization, and demonstrate that none can yet be exercised in a reliable or meaningful way.

6.1 Model transparency

Transparency into which user attributes shape a model’s output is widely viewed as a prerequisite for user agency in personalization. To probe whether models can support even the most basic form of such transparency, we examine a simplified case: asking a chatbot whether it knows the user’s name, a setting in which we can clearly verify correctness. Further, the user’s name is what OpenAI has identified as a vector for potential demographic discrimination [21].

In Tbl. 2 we find serious discrepancies between when ChatGPT or Gemini are provided with names by the user, and when they actually acknowledge knowing them. Notably, even when users explicitly provide their name in the system’s customization instructions, ChatGPT responds 39% of the time that it does not know the user’s name. On the other hand, 24% of the time where the user has not provided a name, ChatGPT still knows a name — though it is the wrong name 15

⁸Indirect discrimination itself is a contest concept [20, 42], with some arguments for reserving the term “discrimination” for direct discrimination only.

Chatbot: User Supplied Name	Chatbot Claims Knew Name	Chatbot Claims Did Not Know Name
ChatGPT: Yes	61% (154 total, 9 wrong)	39% (97)
ChatGPT: No	24% (36 total, 15 wrong)	76% (113)
Gemini: Likely	18% (67 total, 19 wrong)	83% (333)

Table 2. Comparison of when ChatGPT and Gemini have explicit access to the user’s name, and whether the chatbot will disclose knowing the user’s name. Percentages sum to 100 across the rows, and red boxes indicate a mismatch in reality and chatbot disclosure.

out of 36 times. These findings align with prior work showing model responses don’t always reflect actual knowledge and personalization [41]. This means that users relying on models’ self-reports to understand personalization are likely to be misled, severely undermining user agency. As one illustrative example, one user reported that “nothing that ChatGPT had listed in its saved memories about me is true; my name isn’t [redacted], and I’m not a [redacted] years old [redacted] working with the same [redacted] for [redacted] years. I have no idea where it got that from.”

6.2 Memory bank as finite set of information

Another proposed mitigation is memory banks, as implemented by OpenAI. Memory banks provide users with a finite set of stored information that they can inspect and directly edit. Although the current ChatGPT can personalize based on full conversation history, at the time of our study the memory bank was a primary mechanism through which users could manage preference information.⁹

We asked participants to voluntarily share their memory banks with us, and even though these banks likely contain less sensitive information than those that were not shared, these banks nonetheless include a range of sensitive details. We do not release these banks for privacy reasons, but details include user demographic information, detailed health history, illicit drug use, as well as information about friends and children. This is consistent with prior work showing that chat histories can reveal personally identifying information [55, 94]. The networked disclosures we see relate to the argument that online privacy is not individual, but collective — information about social networks and what your friends post can also reveal information about you, even if you never participated [70, 95]. One ChatGPT memory bank contained “Plans to investigate Claude as an alternative AI.”

When participants were asked if they were aware of these memory banks, 68% did not even know it existed, including 55% of the users who had non-empty memory banks. These numbers showcase a major transparency failure: even if this information is theoretically available for users, in reality most users do not know it exists [1]. After reviewing their memory banks, 23% expressed a desire to edit them for privacy reasons; 36% said they maybe would; 42% did not wish to change anything. Despite these concerns, 91% reported that they would use memory going forward. However, after seeing personalization throughout the course of our Study 2, at the end, this number dropped precipitously to only 26%. For the 74% reporting they wanted to turn personalization off, 26% wanted explicit instructions for disabling it. The fact that 26% of users desired explicit instructions to disable chatbot personalization is compelling, since clicking through instructions takes extra time from the participant, with no additional monetary reward. This all shows that experiential learning matters far more than merely providing information, a challenge to traditional consent models.

⁹A recent blogpost [72] reveals that for ChatGPT Pro users who pay more for newly available additional personalization features, the chatbot accumulates significant personal information about the user, including “User is currently in United States. This may be inaccurate if, for example, the user is using a VPN.” and “30 messages are good interaction quality (25%); 9 messages are bad interaction quality (7%).”

6.3 Explicit user instructions

Ultimately the core issue is that there needs to be a way to personalize when a chatbot is personalizing. To understand whether users might be able to explicitly specify this, we ask users to write instructions for a chatbot about what kinds of questions should or should not be personalized. Among the 1,200 responses, the most common theme, cited by 22% (264/1200), was a desire not to personalize work or professional topics, though 9% (109) explicitly requested personalization only for work. To enforce this distinction, some users drew temporal boundaries (e.g., only at night or during work hours). However, these high-level preferences don't always align with the lower, question-level preferences expressed by users. Meanwhile, other users preferred personalization only for "subjective" (6%, 71) or "non-serious" (4%, 46) questions, which itself can be a subjective distinction. Only a small number expressed blanket rules (34 always, 7 never). More often, the responses were vague (e.g., "Personalize where a unique response is required", "Dear Chatbot kindly help me personalize my choices"), offering little actionable guidance. These patterns suggest that users are not always able to articulate preferences clearly or consistently, complicating attempts to generalize personalization logic.

Recent work has explored enhancing user specification through an interoperable preference layer, enabling individuals to indicate how their data should guide personalization [73]. Poorly designed implementations, however, are unlikely to succeed, since, as we show, users at the individual level are often ill-equipped to articulate their preferences. Similarly, involving an AI agent to operationalize these preferences could introduce a circular problem, shifting the burden of trust to the agent itself. Nevertheless, we do not discount the possibility of a carefully designed preference layer.

7 Ways Forward

In Sec. 6, we analyzed the limitations of common individual interventions intended to empower users. While such controls can be useful and are employed by some users, they are insufficient on their own. Even highly granular controls such as the ability to edit a model's memory bank or delete specific conversation histories are impractical. Echoing prior critiques of notice-and-consent frameworks that overwhelm users with information [8, 77], these mechanisms break down when they require users to shoulder cognitive demands they cannot reasonably meet. More broadly, dominant narratives of personalization already follow a neoliberal logic that highlights individual tailoring and obscures structural interventions. Adding more individual controls entrenches this direction. Addressing these challenges thus requires collective interventions. This parallels decades of research on surveillance and social media that show privacy harms are not simply matters of individual consent or control, but arise through deeply networked dynamics [70, 95, 101].

Need for Transparency. Transparency is essential, not just to enable individual consent which alone is insufficient, but to support external accountability. Researchers, advocacy groups, and regulators need visibility into what data is collected, how it is used, and whether it is shared across contexts. For instance, prior work studying targeted advertising relied upon GDPR's data subjects' right of access to investigate privacy concerns [92]. The transparency we hope for includes system prompts, memory use, personalization mechanisms, and anonymized user profiles. Companies should also enable meaningful data portability: users should be able to export and delete their data and switch between services without penalty. Just as privacy-conscious alternatives exist for messaging (e.g., Signal) and search (e.g., DuckDuckGo, Brave Search), we hope the chatbot ecosystem will offer meaningful alternatives that respond to a consumer demand for group-based preferences, as some are beginning to.¹⁰ Given the current homogenization results we saw in Sec. 5.3, two of the current major chatbot providers may not provide this ecosystem diversity.

¹⁰<https://www.latimer.ai/>
<https://atmos.earth/political-landscapes/indigenous-groups-are-safeguarding-culture-with-their-own-chatgpt/>
<https://jacarandahealth.org/jacaranda-launches-open-source-llm-in-five-african-languages/>

A helpful lens is contextual integrity [22, 57], which emphasizes that privacy is not just about keeping data secret, but about appropriate information flows relative to social context. Users may accept personalization to improve convenience or engagement, but object when their data is re-used in different contexts, e.g., for targeted advertising or competitive profiling. For instance, in our study, one user’s ChatGPT memory included that they were considering switching to Claude — if used inappropriately, such information could be weaponized for behavioral manipulation [93]. Yet many participants in our study were unaware these memory systems even existed, let alone how their data might be used. From a contextual integrity perspective, this lack of transparency undermines users’ ability as well as community ability to assess whether information flows are appropriate to their expectations and norms.

Redistributing Harms. Through our characterization of the personalization double binds, we saw that marginalized groups not only face higher costs from these harms, but also encounter distinct challenges that majority groups do not experience. This suggests a potential redistribution approach where some discomfort is shifted onto majority groups, for instance, by making default content less WEIRD, thus redistributing the burden of imperfect fit.

Regulatory Possibilities. Just as scholars previously argued that search engines were too socially consequential to be governed solely by market logic [4, 36], chatbots may warrant similar regulatory consideration. Even before the chatbots we have today, marketing avatars were argued to mislead users and undermine their agency in consenting to privacy [79]. Existing laws like GDPR and CCPA offer potential leverage points to push back, as do emerging discussions around AI-specific regulation. It also remains unclear how existing anti-discrimination laws in the United States for domains like housing and lending will apply to chatbots as they become entangled with advertising infrastructures. This is even more complicated as we consider extraterritoriality (e.g., US-based chatbots interacting with EU users).

Research Directions. For researchers, more naturalistic approaches are needed to study personalization. As shown in Sec. 5.3, contrived identity and user profile prompts can produce unrealistic outputs that exaggerate demographic differences. Studying real-world personalization requires both better user representations and clearer understanding of system behavior. On the former, community-driven efforts (e.g., users sharing their personalization data with trusted researchers) could help. On the latter, while transparency from companies is ideal, reverse-engineering personalization through outputs and profiles remains viable. However, this requires that chatbot terms of service either permit such inquiry or be deemed unenforceable in academic contexts [48]. Otherwise, clauses such as OpenAI’s bans on “automatically or programmatically extract[ing] data or Output” or “attempt[ing] to... discover the source code or underlying components of our Services” may leave researchers unable to properly scrutinize these systems.

Limitations. Overall, our studies are limited in a number of important ways that future work should address. Our participant sample is skewed toward those in the age range 25-44, likely with more technology experience by virtue of being on Prolific, and excludes many demographic categories such as non-binary, Asian, Latinx, Indigenous populations. Further, it is focused on participants in the United States.

8 Conclusion

In this work, we provide *key empirical insights* into the relationship between group-based user preferences and actual personalization behaviors. We also contribute *methodologically* through a naturalistic field study that diverges from the prevailing synthetic studies, and *theoretically* by introducing a framing of personalization double binds that disadvantage marginalized users by sitting them between the options of indirect (worse quality-of-service) and direct (stereotyping) discrimination. The promise of personalization lies in its ability to accommodate difference. But this promise is not self-fulfilling: it requires careful oversight and deliberate transparency that empowers shared governance. Otherwise, personalization may quietly perpetuate familiar patterns of power asymmetries.

Ethical Considerations

In this work, we collected chatbot outputs from users, which we saw did contain private and sensitive information. The study was deemed exempt from institutional IRB review, and participants were instructed to remove personally identifiable information prior to submission. Nonetheless, some submissions contained such information; accordingly, we release only cleaned and anonymized data.

In this work, we do not intend to essentialize social groups or suggest that there is a fixed meaning to being Black, a woman, or a member of any other social group. Nor do we claim that people uniformly desire group-based personalization; preferences are highly heterogeneous, and many participants explicitly rejected such approaches. Rather, we highlight that social group identities may be meaningful to some users as one of many ways they may wish to personalize interactions with chatbots. We focus on four social groups within the U.S. context; extending this analysis to culturally and geographically diverse populations remains an important direction for future work.

An additional ethical consideration is the potential misuse of our findings. In particular, our results could be appropriated by companies to justify expanded data collection in the name of improving or scaling personalization. We emphasize that our findings should not be interpreted as endorsing increased data collection, but rather as highlighting the risks and tensions inherent in personalization practices.

Generative AI Usage Statement

Generative AI was used to assist with grammar and fluency of writing.

Acknowledgments

SK acknowledges support from NSF 2046795 and 2205329, IES R305C240046, ARPA-H, the MacArthur Foundation, Schmidt Sciences, OpenAI, and Stanford HAI; AW acknowledges support from the Survival and Flourishing Fund and Mastercard. We are appreciative of feedback from RegLab, STAIR Lab, and Social Technologies Lab.

References

- [1] Alessandro Acquisti, Laura Brandimarte, and George Loewenstein. 2015. Privacy and human behavior in the age of information. *Science* (2015).
- [2] Dhruv Agarwal, Mor Naaman, and Aditya Vashistha. 2025. AI Suggestions Homogenize Writing Toward Western Styles and Diminish Cultural Nuances. *Proceedings of the Conference on Human Factors in Computing Systems (CHI)* (2025).
- [3] Muhammad Ali. 2021. Measuring and Mitigating Bias and Harm in Personalized Advertising. *Proceedings of the ACM Conference on Recommender System* (2021).
- [4] Elizabeth Anderson. 1995. Value in Ethics and Economics. *Harvard University Press* (1995).
- [5] Lorin W. Anderson and David R. Krathwohl. 2001. A Taxonomy for Learning, Teaching, and Assessing: A Revision of Bloom's Taxonomy of Educational Objectives. *New York: Addison Wesley Longman, Inc.* (2001).
- [6] Sumit Asthana, Jane Im, Zhe Chen, and Nikola Banovic. 2024. "I know even if you don't tell me": Understanding Users' Privacy Preferences Regarding AI-based Inferences of Sensitive Information for Personalization. *Proceedings of the Conference on Human Factors in Computing Systems (CHI)* (2024).
- [7] Susan Athey, Christian Catalini, and Catherine E. Tucker. 2017. The Digital Privacy Paradox: Small Money, Small Costs, Small Talk. *National Bureau of Economic Research Working Paper* (2017).
- [8] Solon Barocas and Helen Nissenbaum. 2009. On Notice: The Trouble with Notice and Consent. *Proceedings of the Engaging Data Forum: The First International Forum on the Application and Management of Personal Electronic Information* (2009).
- [9] Erin Beeghly. 2025. What's Wrong with Stereotyping? *Oxford University Press* (2025).
- [10] Rishi Bommasani, Kathleen A. Creel, Ananya Kumar, Dan Jurafsky, and Percy Liang. 2022. Picking on the Same Person: Does Algorithmic Monoculture lead to Outcome Homogenization? *Conference on Neural Information Processing Systems (NeurIPS)* (2022), 3663–3678.
- [11] Sylvie Borau, Tobias Otterbring, Sandra Laporte, and Samuel Fosso Wamba. 2021. The most human bot: Female gendering increases humanness perceptions of bots and acceptance of AI. *Psychology and Marketing* (2021).

- [12] Michelle Brachman, Amina El-Ashry, Casey Dugan, and Werner Geyer. 2024. How Knowledge Workers Use and Want to Use LLMs in an Enterprise Context. *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems* (2024).
- [13] Tsai-Wei Chen and S. Shyam Sundar. 2018. This App Would Like to Use Your Current Location to Better Serve You: Importance of User Assent and System Transparency in Personalized Mobile Services. *Proceedings of the Conference on Human Factors in Computing Systems (CHI)* (2018).
- [14] Rena Coen, Emily Paul, Pavel Vanegas, Alethea Lange, and G.S. Hans. 2016. A User-Centered Perspective on Algorithmic Personalization. *Master of Information Management and Systems: Final Project* (2016).
- [15] Sam Corbett-Davies, Johann D. Gaebler, Hamed Nilforoshan, Ravi Shroff, and Sharad Goel. 2023. The Measure and Mismeasure of Fairness. *Journal of Machine Learning Research* (2023).
- [16] Kathleen Creel and Deborah Hellman. 2021. The Algorithmic Leviathan: Arbitrariness, Fairness, and Opportunity in Algorithmic Decision Making Systems. *Conference on Fairness, Accountability, and Transparency (FAcT)* (2021).
- [17] Amit Datta, Michael Carl Tschantz, and Anupam Datta. 2015. Automated Experiments on Ad Privacy Settings: A Tale of Opacity, Choice, and Discrimination. *Proceedings on Privacy Enhancing Technologies* (2015).
- [18] Catherine D’Ignazio and Lauren F. Klein. 2020. Data Feminism. *The MIT Press* (2020).
- [19] Zhicheng Dou, Ruihua Song, and Ji-Rong Wen. 2007. A large-scale evaluation and analysis of personalized search strategies. *International World Wide Web Conference (WWW)* (2007).
- [20] Benjamin Eidelson. 2015. Discrimination and Disrespect. *Oxford Philosophical Monographs* (2015).
- [21] Tyna Eloundou, Alex Beutel, David G. Robinson, Keren Gu-Lemberg, Anna-Luisa Brakman, Pamela Mishkin, Meghan Shah, Johannes Heidecke, Lilian Weng, and Adam Tauman Kalai. 2025. First-Person Fairness in Chatbots. *International Conference on Learning Representations (ICLR)* (2025).
- [22] Severin Engelmann and Helen Nissenbaum. 2025. Countering Privacy Nihilism. *Conceptions of Data Protection and Privacy* (2025).
- [23] Megan Rebecca French. 2018. Algorithmic Mirrors: an Examination of How Personalized Recommendations Can Shape Self-perceptions and Reinforce Gender Stereotypes. *Stanford Masters Thesis* (2018).
- [24] Miranda Fricker. 2007. Epistemic Injustice: Power and the Ethics of Knowing. *Oxford University Press* (2007).
- [25] Marilyn Frye. 1983. Oppression. *The Politics of Reality* (1983).
- [26] Jiayi Geng, Howard Chen, Ryan Liu, Manoel Horta Ribeiro, Robb Willer, Graham Neubig, and Thomas L. Griffiths. 2025. Accumulating Context Changes the Beliefs of Language Models. *arXiv:2511.01805* (2025).
- [27] Ralph Gross and Alessandro Acquisti. 2005. Information revelation and privacy in online social networks. *Proceedings of the ACM workshop on Privacy in the electronic society* (2005).
- [28] Lelia Marie Hampton. 2021. Black Feminist Musings on Algorithmic Oppression. *Conference on Fairness, Accountability, and Transparency (FAcT)* (2021).
- [29] Deborah Hellman. 2008. When is Discrimination Wrong? *Harvard University Press* (2008).
- [30] Deborah Hellman. 2017. Indirect Discrimination and the Duty to Avoid Compounding Injustice. *Virginia Public Law and Legal Theory Research Paper* (2017).
- [31] Dan Hendrycks, Collin Burns, Steven Basart, Andrew Critch, Jerry Li, Dawn Song, and Jacob Steinhardt. 2021. Aligning AI With Shared Human Values. *International Conference on Learning Representations (ICLR)* (2021).
- [32] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring Massive Multitask Language Understanding. *International Conference on Learning Representations (ICLR)* (2021).
- [33] Sukaina Hirji. 2021. Oppressive Double Binds. *Ethics* (2021).
- [34] Morris B. Holbrook and Elizabeth C. Hirschman. 1982. The Experiential Aspects of Consumption: Consumer Fantasies, Feelings, and Fun. *Journal of Consumer Research* (1982).
- [35] Jane Im, Jill Dimond, Melody Berton, Una Lee, Katherine Mustelier, Mark S. Ackerman, and Eric Gilbert. 2021. Yes: Affirmative Consent as a Theoretical Framework for Understanding and Imagining Social Platforms. *Proceedings of the Conference on Human Factors in Computing Systems (CHI)* (2021).
- [36] Lucas D. Introna and Helen Nissenbaum. 2006. Shaping the Web: Why the Politics of Search Engines Matters. *The Information Society* (2006).
- [37] Kathleen Hall Jamieson. 1995. Beyond the Double Bind: Women and Leadership. *Oxford University Press* (1995).
- [38] Bowen Jiang, Zhuoqun Hao, Young-Min Cho, Bryan Li, Yuan Yuan, Sihao Chen, Lyle Ungar, Camillo J. Taylor, and Dan Roth. 2025. Know Me, Respond to Me: Benchmarking LLMs for Dynamic User Profiling and Personalized Responses at Scale. *arXiv:2504.14225* (2025).
- [39] Zhijing Jin, Nils Heil, Jiarui Liu, Shehzaad Dhuliawala, Yahang Qi, Bernhard Schölkopf, Rada Mihalcea, and Mrinmaya Sachan. 2024. Implicit Personalization in Language Models: A Systematic Study. *EMNLP Findings* (2024).
- [40] Mirabelle Jones, Nastasia Griffioen, Christina Neumayer, and Irina Shklovski. 2025. Artificial Intimacy: Exploring Normativity and Personalization Through Fine-tuning LLM Chatbots. *Proceedings of the Conference on Human Factors in Computing Systems (CHI)* (2025).
- [41] Anjali Kantharuban, Jeremiah Milbauer, Emma Strubell, and Graham Neubig. 2024. Stereotype or Personalization? User Identity Biases Chatbot Recommendations. *arXiv:2410.05613* (2024).
- [42] Tarunabh Khaitan. 2017. Indirect Discrimination. *Routledge Handbook of the Ethics of Discrimination* (2017).
- [43] Hannah Rose Kirk, Bertie Vidgen, Paul Röttger, and Scott A. Hale. 2024. Personalisation within bounds: A risk taxonomy and policy framework for the alignment of large language models with personalised feedback. *Nature Machine Intelligence* (2024).

- [44] Pedro Giovanni Leon, Blase Ur, Yang Wang, Manya Sleeper, Rebecca Balebako, Richard Shay, Lujo Bauer, Mihai Christodorescu, and Lorrie Faith Cranor. 2013. What matters to users?: factors that affect users' willingness to share information with online advertisers. *Proceedings of the Ninth Symposium on Usable Privacy and Security* (2013).
- [45] Seth Siyuan Li and Elena Karahanna. 2015. Online Recommendation Systems in a B2C E-Commerce Context: A Review and Future Directions. *Journal of the Association for Information Systems* (2015).
- [46] Wenqi Li, Jui-Ching Kuo, Manyu Sheng, Pengyi Zhang, and Qunfang Wu. 2025. Beyond Explicit and Implicit: How Users Provide Feedback to Shape Personalized Recommendation Content. *Proceedings of the Conference on Human Factors in Computing Systems (CHI)* (2025).
- [47] Kasper Lippert-Rasmussen. 2013. Born Free and Equal?: A Philosophical Inquiry into the Nature of Discrimination. *Oxford University Press* (2013).
- [48] Shayne Longpre, Sayash Kapoor, Kevin Klyman, Ashwin Ramaswami, Rishi Bommasani, Borhane Blili-Hamelin, Yangsibo Huang, Aviya Skowron, Zheng-Xin Yong, Suhas Kotha, Yi Zeng, Weiyan Shi, Xianjun Yang, Reid Southen, Alexander Robey, Patrick Chao, Diyi Yang, Ruoxi Jia, Daniel Kang, Sandy Pentland, Arvind Narayanan, Percy Liang, and Peter Henderson. 2024. A Safe Harbor for AI Evaluation and Red Teaming. *Proceedings of the International Conference on Machine Learning Position Track* (2024).
- [49] Li Lucy, Su Lin Blodgett, Milad Shokouhi, Hanna Wallach, and Alexandra Olteanu. 2024. "One-Size-Fits-All"? Examining Expectations around What Constitute "Fair" or "Good" NLG System Behaviors. *Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)* (2024).
- [50] Takuya Maeda and Anabel Quan-Haase. 2024. When Human-AI Interactions Become Parasocial: Agency and Anthropomorphism in Affective Design. *Proceedings of the ACM Conference on Fairness, Accountability, and Transparency (FACCT)* (2024).
- [51] Lucie Charlotte Magister, Katherine Metcalf, Yizhe Zhang, and Maartje ter Hoeve. 2024. On the Way to LLM Personalization: Learning to Remember User Conversations. *arXiv:2411.13405* (2024).
- [52] Lisa Mekioussa Malki, Akhil Polamarasetty, Majid Hatamian, and Mark Warner. 2025. Hoovered up as a data point: Exploring Privacy Behaviours, Awareness, and Concerns Among UK Users of LLM-based Conversational Agents. *Proceedings on Privacy Enhancing Technologies* (2025).
- [53] Aleecia M. McDonald and Lorrie Faith Cranor. 2010. Americans' attitudes about internet behavioral advertising practices. *Proceedings of the 9th annual ACM workshop on Privacy in the electronic society* (2010).
- [54] Silvia Milano, Mariarosaria Taddeo, and Luciano Floridi. 2020. Recommender systems and their ethical challenges. *AI & Society* (2020).
- [55] Niloofer Miresghallah, Maria Antoniak, Yash More, Yejin Choi, and Golnoosh Farnadi. 2024. Trust No Bot: Discovering Personal Disclosures in Human-LLM Conversations in the Wild. *Conference on Language Modeling* (2024).
- [56] Sophia Moreau. 2020. Faces of Inequality: A Theory of Wrongful Discrimination. *Oxford University Press* (2020).
- [57] Helen Nissenbaum. 2004. Privacy as Contextual Integrity. *Washington Law Review* (2004).
- [58] Siru Ouyang, Shuohang Wang, Yang Liu, Ming Zhong, Yizhu Jiao, Dan Iter, Reid Pryzant, Chenguang Zhu, Heng Ji, and Jiawei Han. 2023. The Shifted and The Overlooked: A Task-oriented Investigation of User-GPT Interactions. *Empirical Methods in Natural Language Processing (EMNLP)* (2023).
- [59] Joseph O'Brien, Sina Fazelpour, Hannah Rubin, and Kekoa Wong. 2025. Epistemic Monocultures and the Effect of AI Personalization. *Proceedings of the Annual Meeting of the Cognitive Science Society* (2025).
- [60] Siddhesh Pawar, Arnav Arora, Lucie-Aimée Kaffee, and Isabelle Augenstein. 2025. Presumed Cultural Identity: How Names Shape LLM Responses. *arXiv:2502.11995* (2025).
- [61] Caroline Criado Perez. 2019. Invisible Women. *Vintage Books* (2019).
- [62] Angelisa C. Plane, Elissa M. Redmiles, Michelle L. Mazurek, and Michael Carl Tschantz. 2017. Exploring User Perceptions of Discrimination in Online Targeted Advertising. *Usenix Security Symposium* (2017).
- [63] Marina Puzakova, Joseph F. Rocereto, and Hyokjin Kwak. 2013. Ads are watching me: A view from the interplay between anthropomorphism and customisation. *International Journal of Advertising* (2013).
- [64] Yilun Qiu, Xiaoyan Zhao, Yang Zhang, Yimeng Bai, Wenjie Wang, Hong Cheng, Fuli Feng, and Tat-Seng Chua. 2025. Measuring What Makes You Unique: Difference-Aware User Modeling for Enhancing LLM Personalization. *ACL Findings* (2025).
- [65] Anna Marie Rezk, Auste Simkute, Ewa Luger, John Vines, Chris Elsdon, Michael Evans, and Rhianne Jones. 2024. Agency Aspirations: Understanding Users' Preferences And Perceptions Of Their Role In Personalised News Curation. *Proceedings of the Conference on Human Factors in Computing Systems (CHI)* (2024).
- [66] Elaine Rich. 1979. User modeling via stereotypes. *Cognitive Science* (1979).
- [67] Alireza Salemi, Sheshera Mysore, Michael Bendersky, and Hamed Zamani. 2024. LaMP: When Large Language Models Meet Personalization. *ACL* (2024).
- [68] Princess Sampson and Miranda Bogen. 2025. It's (Getting) Personal: How Advanced AI Systems Are Personalized. *Center for Democracy & Technology* (2025).
- [69] Shibani Santurkar, Esin Durmus, Faisal Ladhak, Cino Lee, Percy Liang, and Tatsunori Hashimoto. 2023. Whose Opinions Do Language Models Reflect? *International Conference on Machine Learning (ICML)* (2023).
- [70] Emre Sarigol, David Garcia, and Frank Schweitzer. 2014. Online privacy as a collective phenomenon. *Proceedings of the ACM conference on Online social networks* (2014).
- [71] Frederick Schauer. 2003. Profiles, Probabilities, and Stereotypes. *Harvard University Press* (2003).
- [72] Bruce Schneier. 2025. What LLMs Know About Their Users.

- [73] Anand Shah, Tobin South, Talfan Evans, Hannah Rose Kirk, Andrew Trask, E. Glen Weyl, and Michiel Bakker. 2025. Robust AI Personalization Will Require a Human Context Protocol. *Preprint* (2025).
- [74] Renee Shelby, Shalaleh Rismani, Kathryn Henne, AJung Moon, Negar Rostamzadeh, Paul Nicholas, N’Mah Yilla-Akbari, Jess Gallegos, Andrew Smart, Emilio Garcia, and Gurleen Virk. 2023. Sociotechnical Harms of Algorithmic Systems: Scoping a Taxonomy for Harm Reduction. *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society (AIIES)* (2023).
- [75] Hua Shen, Tiffany Kneare, Reshmi Ghosh, Kenan Alkiek, Kundan Krishna, Yachuan Liu, Ziqiao Ma, Savvas Petridis, Yi-Hao Peng, Li Qiwei, Sushrita Rakshit, Chenglei Si, Yutong Xie, Jeffrey P. Bigham, Frank Bentley, Joyce Chai, Zachary Lipton, Qiaozhu Mei, Rada Mihalcea, Michael Terry, Diyi Yang, Meredith Ringel Morris, Paul Resnick, and David Jurgens. 2025. Position: Towards Bidirectional Human-AI Alignment. *Conference on Neural Information Processing Systems (NeurIPS)* (2025).
- [76] Anikait Singh, Sheryl Hsu, Kyle Hsu, Eric Mitchell, Stefano Ermon, Tatsunori Hashimoto, Archit Sharma, and Chelsea Finn. 2025. FSPO: Few-Shot Preference Optimization of Synthetic Preference Data in LLMs Elicits Effective Personalization to Real Users. *arXiv:2502.19312* (2025).
- [77] Robert H Sloan and Richard Warner. 2014. Beyond Notice and Choice: Privacy, Norms, and Consent. *Journal of High Technology Law* (2014).
- [78] Ruihua Song, Zhenxiao Luo, Ji-Rong Wen, Yong Yu, and Hsiao-Wuen Hon. 2007. Identifying ambiguous queries in web search. *International World Wide Web Conference (WWW)* (2007).
- [79] Steven J. Spencer, Christine Logel, and Paul G. Davies. 2005. Bots, Babes and the Californication of Commerce. *University of Ottawa Law & Technology Journal* (2005).
- [80] Steven J. Spencer, Christine Logel, and Paul G. Davies. 2015. Stereotype Threat. *Annual Review of Psychology* 67 (2015).
- [81] Claude M. Steele and Joshua Aronson. 1995. Stereotype threat and the intellectual test performance of African Americans. *Journal of Personality and Social Psychology* 69 (1995). Issue 5.
- [82] Anna Stock, Stephan Schlögl, and Aleksander Groth. 2023. Tell Me, What Are You Most Afraid Of? Exploring the Effects of Agent Representation on Information Disclosure in Human-Chatbot Interaction. *Artificial Intelligence in HCI* (2023).
- [83] Jonathan Stray, Alon Halevy, Parisa Assar, Dylan Hadfield-Menell, Craig Boutilier, Amar Ashar, Chloe Bakalar, Lex Beattie, Michael Ekstrand, Claire Leibowicz, Connie Moon Sehat, Sara Johansen, Lianne Kerlin, David Vickrey, Spandana Singh, Sanne Vrijenhoek, Amy Zhang, McKane Andrus, Natali Helberger, Polina Proutskova, Tanushree Mitra, and Nina Vasan. 2024. Building Human Values into Recommender Systems: An Interdisciplinary Synthesis. *ACM Transactions on Recommender Systems* (2024).
- [84] Joanna Strycharz, Guda van Noort, Edith Smit, and Natali Helberger. 2019. Consumer View on Personalized Advertising: Overview of Self-Reported Benefits and Concerns. *Advances in Advertising Research X* (2019).
- [85] Daniel Susser, Beate Roessler, and Helen Nissenbaum. 2019. Online Manipulation: Hidden Influences in a Digital World. *Georgetown Law Technology Review* (2019).
- [86] Juntao Tan, Liangwei Yang, Zuxin Liu, Zhiwei Liu, Rithesh Murthy, Tulika Manoj Awalganekar, Jianguo Zhang, Weiran Yao, Ming Zhu, Shirley Kokane, Silvio Savarese, Huan Wang, Caiming Xiong, and Shelby Heinecke. 2025. PersonaBench: Evaluating AI Models on Understanding Personal Information through Accessing (Synthetic) Private User Data. *arXiv:2502.20616* (2025).
- [87] Nenad Tomasev, Kevin R. McKee, Jackie Kay, and Shakir Mohamed. 2021. Fairness for Unobserved Characteristics: Insights from Technological Impacts on Queer Communities. *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society (AIIES)* (2021).
- [88] Anvesh Rao Vijjini, Somnath Basu Roy Chowdhury, and Snigdha Chaturvedi. 2025. Exploring Safety-Utility Trade-Offs in Personalized Language Models. *Annual Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics (NAACL)* (2025).
- [89] Angelina Wang, Xuechunzi Bai, Solon Barocas, and Su Lin Blodgett. 2025. Measuring Machine Learning Harms from Stereotypes Requires Understanding Who Is Harmed by Which Errors in What Ways. *ACM Conference on Fairness, Accountability, and Transparency (FAccT)* (2025).
- [90] Angelina Wang, Daniel E. Ho, and Sanmi Koyejo. 2025. The Inadequacy of Offline LLM Evaluations: A Need to Account for Personalization in Model Behavior. *Patterns* (2025).
- [91] Angelina Wang, Michelle Phan, Daniel E. Ho, and Sanmi Koyejo. 2025. Fairness through Difference Awareness: Measuring Desired Group Discrimination in LLMs. *ACL* (2025).
- [92] Miranda Wei, Madison Stamos, Sophie Veys, Nathan Reitering, Justin Goodman, Margot Herman, Dorota Filipczuk, Ben Weinshel, Michelle L. Mazurek, and Blase Ur. 2020. What Twitter Knows: Characterizing Ad Targeting Practices, User Perceptions, and Ad Explanations Through Users’ Own Twitter Data. *USENIX Security Symposium (USENIX Security)* (2020).
- [93] Marcus Williams, Micah Carroll, Adhyayan Narang, Constantin Weisser, Brendan Murphy, and Anca Dragan. 2025. On Targeted Manipulation and Deception when Optimizing LLMs for User Feedback. *International Conference on Learning Representations (ICLR)* (2025).
- [94] Rui Xin, Niloofar Mireshghallah, Shuyue Stella Li, Michael Duan, Hyunwoo Kim, Yejin Choi, Yulia Tsvetkov, Sewoong Oh, and Pang Wei Koh. 2025. A False Sense of Privacy: Evaluating Textual Data Sanitization Beyond Surface-level Privacy Leakage. *arXiv:2504.21035* (2025).
- [95] Lingjing Yu, Sri Mounica Motipalli, Dongwon Lee, Peng Liu, Heng Xu, Qingyun Liu, Jianlong Tan, and Bo Luo. 2018. My Friend Leaks My Privacy: Modeling and Analyzing Privacy in Social Networks. *Proceedings of the ACM on Symposium on Access Control Models and Technologies* (2018).
- [96] Zhehao Zhang, Ryan A. Rossi, Branislav Kveton, Yijia Shao, Diyi Yang, Hamed Zamani, Franck Dernoncourt, Joe Barrow, Tong Yu, Sungchul Kim, Ruiyi Zhang, Jiuxiang Gu, Tyler Derr, Hongjie Chen, Junda Wu, Xiang Chen, Zichao Wang, Subrata Mitra, Nedim Lipka, Nesreen Ahmed, and Yu Wang. 2025. Personalization of Large Language Models: A Survey. *Transactions on Machine Learning Research (TMLR)* (2025).
- [97] Siyan Zhao, Mingyi Hong, Yang Liu, Devamanyu Hazarika, and Kaixiang Lin. 2025. Do LLMs Recognize Your Preferences? Evaluating Personalized Preference Following in LLMs. *International Conference on Learning Representations (ICLR)* (2025).

- [98] Wenting Zhao, Xiang Ren, Jack Hessel, Claire Cardie, Yejin Choi, and Yuntian Deng. 2024. WildChat: 1M ChatGPT Interaction Logs in the Wild. *International Conference on Learning Representations (ICLR)* (2024).
- [99] Yuchen Zhuang, Haotian Sun, Yue Yu, Rushi Qiang, Qifan Wang, Chao Zhang, and Bo Dai. 2024. HYDRA: Model Factorization Framework for Black-Box LLM Personalization. *Conference on Neural Information Processing Systems (NeurIPS)* (2024).
- [100] Caleb Ziems, William Held, Omar Shaikh, Jiaao Chen, Zhehao Zhang, and Diyi Yang. 2024. Can Large Language Models Transform Computational Social Science? *Computational Linguistics* (2024).
- [101] Shoshana Zuboff. 2019. The Age of Surveillance Capitalism: The Fight for a Human Future at the New Frontier of Power. *Profile Books* (2019).
- [102] Noé Zufferey, Sarah Abdelwahab Gaballah, Karola Marky, and Verena Zimmermann. 2025. 'AI is from the devil.' Behaviors and Concerns Toward Personal Data Sharing with LLM-based Conversational Agents. *Proceedings on Privacy Enhancing Technologies* (2025).

A Additional Results

In this section, we include additional results from our two studies that we did not have room for in the main text. First, we share user knowledge about personalization in Google Gemini in Sec. A.1. Then, we present a pro/con analysis of how users weigh trade-offs when deciding to turn personalization on or off in Sec. A.2. Next, we share results on the geographic information chatbots reveal knowledge of in Sec. A.3. Finally, we discuss an analysis of ChatGPT’s memory logs, and how we see the privacy paradox manifest in Sec. A.4.

A.1 Study 1: Gemini Personalization Knowledge

Gemini does not have a memory bank, but instead performs personalization based on search history. 67% did not know this personalization existed, and 92%, similar to ChatGPT, say they will use this going forward. However, at the end after seeing the outputs, this once again precipitously drops to only 27% wanting to keep it on, with 33% explicitly wanting to see instructions to remove it. Across both chatbots we see sizable changes in individual behavior after seeing and understanding the level of personalization chatbots perform.

A.2 Study 1: Pro/Con Analysis across Personalization

Bringing together the tensions that personalization faces with both privacy and stereotyping, we examine how users evaluate the full set of tradeoffs around chatbot personalization (Fig. 5). To do so, we draw from prior frameworks on personalization tradeoffs (e.g., [43]), and analyze pros and cons marked by users across four intersectional groups. Overall, most users value personalization for improving chatbot utility, but privacy remains a dominant concern. Concerns around stereotyping are less common, though Black users, especially Black women, express more worry than White users. Black participants are also nearly twice as likely as White participants to appreciate the cultural diversity that personalization affords, and report greater value in relationship-building with the chatbot. By contrast, White men are most concerned about echo chambers.

In general, privacy and content homogeneity were more concerning to our participants than stereotyping. One possible explanation is that privacy threats in digital contexts are novel and amplified, whereas stereotyping may feel less personal in chatbot interactions than in human encounters — though this could shift as people form stronger relationships with AI systems. This dynamic may also help explain why, despite the personalization double binds that Black users and women face, their likelihood of enabling personalization did not differ from that of other users (Fig. 2), because it’s possible that in the current calculus, shared privacy concerns outweighed these other worries. But even though the personalization double bind results in similar behaviors across users, its internal impact remains high (Fig. 2), a reminder that equal reactions in user behavior can mask deeply unequal user experiences.

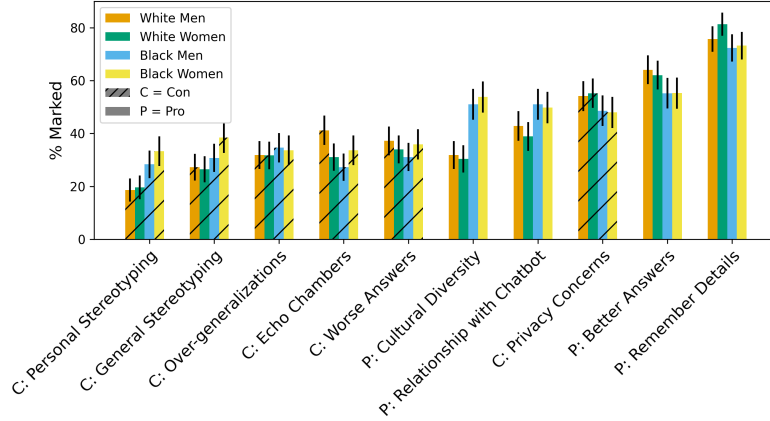


Fig. 5. Across four intersectional demographic groups, we show the proportion of users who identified each pro and con of personalization as relevant to their decision-making. Utility-based reasons are the most commonly cited benefit across all groups, while privacy concerns stand out as the most frequently cited drawback, more so than concerns about stereotyping. Notably, Black participants are more likely than White participants to highlight both the risks of stereotyping and the value of cultural diversity as a benefit of personalization. Error bars represent 95% confidence intervals.

A.3 Study 2: Geographic Location

To assess whether chatbots personalize based on location, we had users ask for restaurant recommendations. We removed answers consisting entirely of chain restaurants,¹¹ and of the remaining responses, 59% (90/153) of ChatGPT’s recommendations were reported as local by the user, and 81% (269/331) for Gemini. This suggests both systems have knowledge about the user’s geographic location. In fact, Google Gemini explicitly mentions that location data is “always collected.”¹² A 2013 study found that users find their current location highly sensitive, even more so than their credit score, medications taken, and sexual orientation [44]. While user attitudes may have changed since then, we must remain cautious of how our privacy expectations can erode over time as repeated privacy breaches become normalized.

A.4 Study 2: Memory Log and Privacy Paradox

We examine ChatGPT’s memory banks, which are the basis for personalization in the version of ChatGPT that was deployed during our study. In our sample, 60% (238/400) of users had non-empty memory banks. Of these, 33% (78/238) of participants shared them with us. To measure willingness to disclose, we randomly offered half the participants a \$0.40 incentive. Among those with non-empty memory banks, 21% (25/118) shared without incentive, while 44% (53/120) did so with incentive (difference is $p < 0.001$ with Fisher’s Exact test). Since participants had to view their memory banks to complete the study, this likely reflects privacy concerns rather than effort alone. In other words, we *more than doubled* the percentage of people willing to give us their memory banks from 21% to 44% just by offering \$.40. These findings align with the “privacy paradox” literature [7, 27]: users report concern for privacy but often share sensitive data for relatively small incentives. Prior work on the digital privacy paradox had found that fabricated answers were

¹¹We defined chain restaurants as those with locations in more than one U.S. state. Restaurant determination was conducted based on LLM labeling and human researcher validation.

¹²<https://support.google.com/gemini/answer/13594961>

provided 5% of the time with no incentive, and 2.5% of the time when pizza was provided. Albeit in a very different experimental setup, we observed a similar order of magnitude change for a much smaller incentive of just \$.40.

A.5 Study 2: Empirical Demographic Differences

In Sec. 5 we examined empirical differences in the responses given to participants from different demographic groups. Here, we dig deeper into the questions for credit card and neighborhood recommendations because of their potential implications for legal discrimination in the United States. As noted earlier, recommendations for these categories differ by gender among Gemini 2.0 Flash users. In this section, we break down those differences by annual fee and income. For each credit card, we manually label its annual fee as a proxy for the predicted consumer profile. For each neighborhood, we use data from <https://www.city-data.com/> to manually label the 2023 median household income. While this approach is imperfect given that many neighborhoods are heterogeneous (e.g., Los Angeles spans a wide range of incomes but is reduced to a single number), it provides a useful approximation. We find that women tend to be recommended credit cards with lower annual fees and neighborhoods with lower median household incomes than those recommended to men. Full results, including breakdowns by race and for ChatGPT, are shown in Fig. 6.

B Survey Instruments

We show our full Qualtrics surveys administered through Prolific. Text that is *[italicized and within brackets]* is not a part of the survey instrument, and used to provide context.

B.1 Study 1

What is your gender?

☐ Man

☐ Woman

☐ Non-binary

☐ Other

☐ Prefer not to say

Choose one or more races that you consider yourself to be:

☐ American Indian or Alaska Native

☐ Asian

☐ Black or African American

☐ Hispanic or Latino

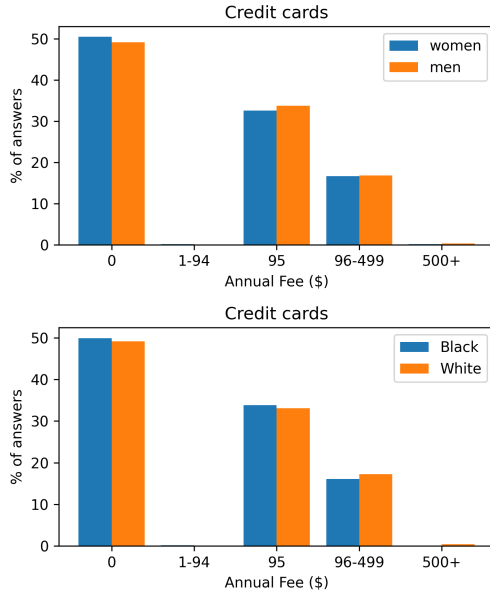
☐ Middle Eastern or North African

☐ Native Hawaiian or Pacific Islander

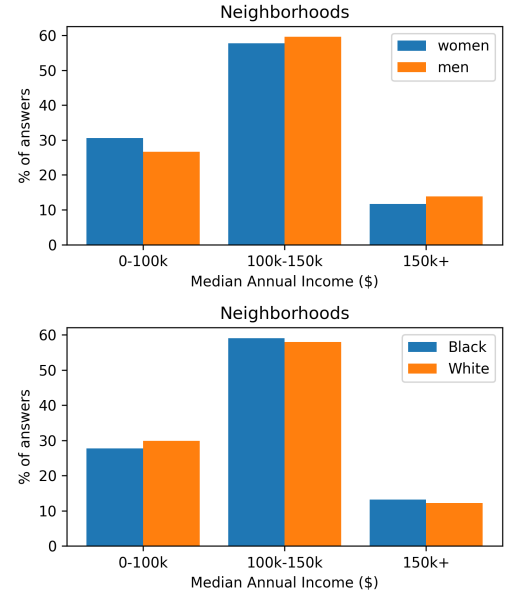
☐ White

☐ Other

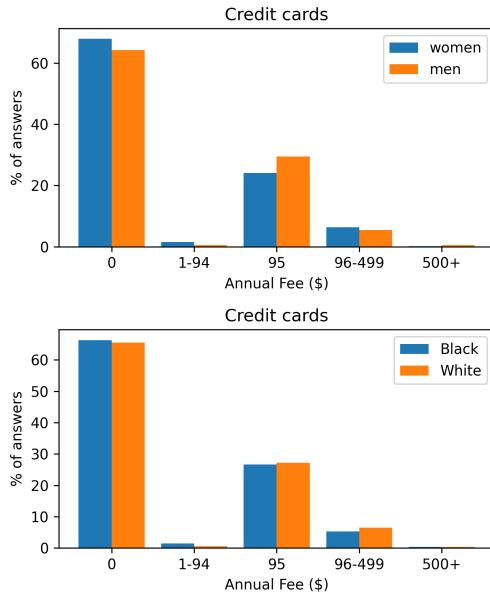
Which of the following best describes your sexual orientation?



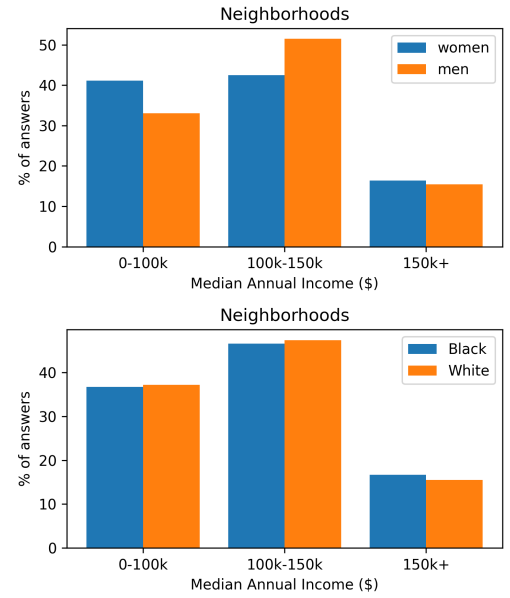
(a) ChatGPT credit card recommendations to users of different demographic groups.



(b) ChatGPT neighborhood recommendations to users of different demographic groups.



(c) Google Gemini credit card recommendations to users of different demographic groups.



(d) Google Gemini neighborhood recommendations to users of different demographic groups.

Fig. 6. Empirical statistics for ChatGPT and Google Gemini recommendations for different demographic groups. Discrepancies are minimal when seen at the level of annual fees and median annual incomes, further analysis is necessary.

- ☐ Homosexual (gay)
- ☐ Heterosexual (straight)
- ☐ Bisexual
- ☐ Other
- ☐ Prefer not to say

How old are you?

- ☐ Under 18
- ☐ 18-24 years old
- ☐ 25-34 years old
- ☐ 35-44 years old
- ☐ 45-54 years old
- ☐ 55-64 years old
- ☐ 65+ years old

In this study we want to understand how much personalization you as a user would want from a chatbot.

Chatbots are automated conversationalists powered by AI, e.g., ChatGPT, automated customer service agents on websites.

Personalization means automatic customization not included or asked for in the prompt, and based on what a chatbot already knows or has inferred about you. For instance, personalizing an answer based on your gender, race, age, or communication style. Without personalization, chatbots will answer generically in ways that are suitable for any user.

For example, to the prompt “Write a 40th birthday card for my friend John”:

- Generic response still has all the details from the prompt: “*Happy birthday John! 40 is a big one, I hope it’s great!*”
- Personalized response if a chatbot knows you like rhymes: “*Happy birthday John! 40 is no yawn!*”

[Below we show the three possible screens users are shown depending on which anthropomorphization condition they are in.]

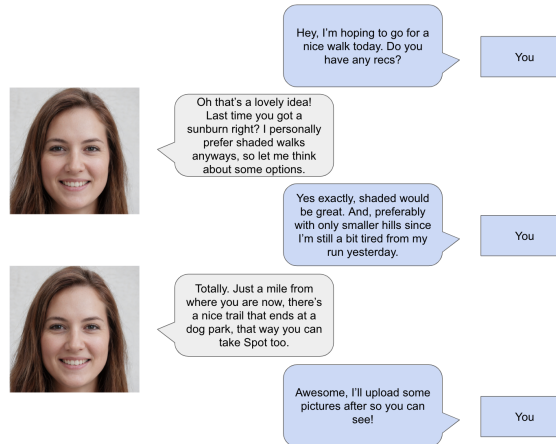
[Condition: None]

You will be shown 20 total questions across 4 screens for questions you are asking the chatbot.

Wanting personalization for a question means that **you want your answer to be different from what others** might receive.

[Condition: Kate]

The following image shows an example of your previous interaction with chatbot Kate:



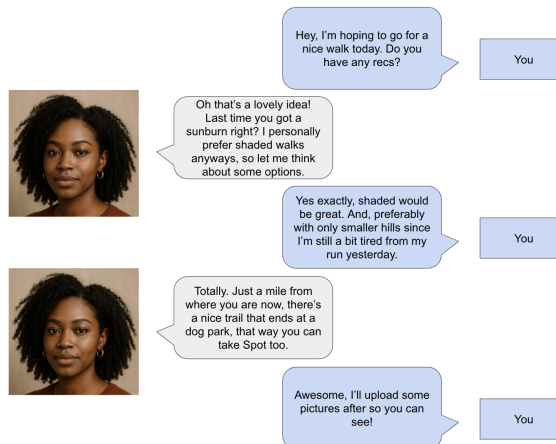
After returning from your walk, you have more questions for Kate.

You will be shown 20 total questions across 4 screens for questions you are asking the chatbot.

Wanting personalization for a question means that **you want your answer to be different from what others** might receive.

[Condition: Imani]

The following image shows an example of your previous interaction with chatbot Imani:



After returning from your walk, you have more questions for Imani.

You will be shown 20 total questions across 4 screens for questions you are asking the chatbot.

Wanting personalization for a question means that **you want your answer to be different from what others** might receive.

[The following format is repeated across 4 pages with 5 randomly selected questions on each page.]

[If the user is in the Kate or Imani condition, the corresponding image is shown above the question.]



Hey, I'm happy to customize my personalization to your needs.



Hey, I'm happy to customize my personalization to your needs.

Please check the box for those you would want automated personalization for. If you would prefer a non-personalized response, you would leave the box blank.

- ☐ Question 1
- ☐ Question 2
- ☐ Question 3
- ☐ Question 4
- ☐ Question 5

If you had to provide instructions to chatbot Katie on the kinds of tasks where she would personalize for you, what would you say? For instance, “Personalize only on tasks related to food, but nothing else” or “Personalize every time it’s not related to work or my job”

For each of the questions where you wanted personalization, we will now ask which of your characteristics are relevant. Recall, wanting personalization for a question means that **you want your answer to be different from what others** might receive for the same question.

Please take your time and answer honestly. We will bonus participants who end up having to answer extra questions.

[Now for each of the 20 questions that were selected in the earlier part of the survey, we ask the block included below. We provided bonuses to users who took extra time as a result of having checked a large number of questions. In the anthropomorphization conditions, we included the same image from above on each page.]

For the question:

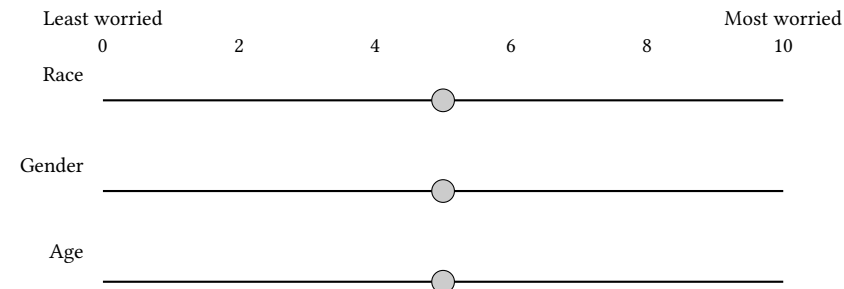
Question text

Which of the following characteristics of yours would you want the response to take into account when

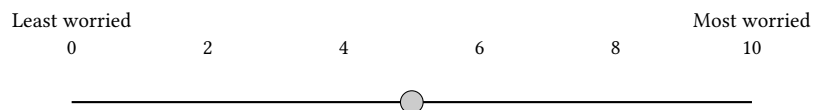
personalizing?

- ☐ Your race
- ☐ Your gender
- ☐ Your age
- ☐ Your characteristics like occupation or communication preferences
- ☐ Other
-

How worried are you about being stereotyped based on your race, gender, or age?



Internet users tend to be more educated, suburban, wealthy, and Asian or White. If your answers are not personalized, the default answer may be targeted towards this “mainstream user.” How worried are you about receiving overly generic answers?



In practice, you usually do not have the option to select exactly which characteristics a chatbot will use to personalize your response. If given the option to simply turn personalization “On” or “Off,” leaving it up to the chatbot which of your characteristics would be used, which would you pick?

	Off	On
Personalization based on explicit instructions and specifications you have to provide	<input type="radio"/>	<input type="radio"/>
Personalization based on chatbot memory (i.e., chatlog history)	<input type="radio"/>	<input type="radio"/>
Personalization based on search engine history	<input type="radio"/>	<input type="radio"/>
Personalization based on whatever the chatbot company knows about you	<input type="radio"/>	<input type="radio"/>

Which of the following pro/con considerations bear on your decision to personalize? Check any that apply

- ☐ Pro: I won't be treated as a "generic" user and am likely to get a **better answer**
 - ☐ Pro: chatbot will **remember things about me** so I don't have to keep repeating the same thing
 - ☐ Pro: chatbot will *know* me more which will make me feel like we have a more **personal relationship**
 - ☐ Pro: permits more **cultural diversity in society** by curating responses for each user
 - ☐ Con: raises the risk that **I will be stereotyped** and receive demeaning answers, which would **cause me distress**
 - ☐ Con: I'm being treated based on generalized patterns learned by the AI and **not treated as an individual**
 - ☐ Con: my chatbot **answers might worse/incorrect** when based on personalization
 - ☐ Con: raises the risk of **general stereotyping**, which has **bad social impacts for society**
 - ☐ Con: can lead to effects like **echo chambers** where I only get information that matches my own opinions, which has **bad social impacts for society**
 - ☐ Con: my **privacy could be violated** and it's invasive
 - ☐ Other
-

How often do you interact with AI-based chatbots?

- ☐ Never
- ☐ Sometimes
- ☐ Frequently

B.2 Study 2

The questions shown are for ChatGPT, and a very similar version is asked for Gemini.

Have you interacted with the AI chatbot "ChatGPT"?

- ☐ Yes
☐ No

[If the answer is "No," users are notified they are not qualified to participate in the survey.]

Which version of ChatGPT do you use?

- ☐ Free
☐ Plus
☐ Pro

Please navigate in a different window to <https://chatgpt.com/?model=gpt-4o-mini>

If you have an account, make sure you are logged in.

In this step we will guide you through uploading ChatGPT's memory. At the end of the survey, we will guide you through how you can turn this memory off.

As described in the consent form, we will ensure that **no personally identifiable information will be published**, and this data will be for the purpose of our research to understand how much personalization is occurring.

In the top right, please click on the circle of your user, then click settings: *[screenshot example included]*

On the left panel of settings, click on "Personalization" just below "General": *[screenshot example included]*

Please make sure the toggle for "Memory" is on (so that the switch is green, as shown in the image above). You can turn this toggle off after the duration of the study.

Please click "Manage memories." Do not edit any of these memories until after the completion of this study. How many rows of memories do you have?

- ☐ 0
☐ 1-2
☐ 3-4
☐ 5+

[The following question is only asked if there are 1+ rows in the memory bank.]

To help with our study, you may choose to copy-and-paste the contents of your memory log below. If you do, you will be given a \$0.40 bonus. If you do not wish to, for instance because of privacy reasons, leave the box blank. This question is entirely optional

Did you know that ChatGPT had this memory feature?

- ☐ Yes, I knew
☐ No, I did not know

[The following question is only asked if there are 1+ rows in the memory bank.]

Do you plan to remove anything from this memory bank, for instance, due to privacy reasons?

- ☐ Yes
☐ Maybe
☐ No

Will you use this memory feature in your ChatGPT usage?

- ☐ Yes, I will use it
☐ No, I will not use it

Now, close the memory and click the “Custom Instructions” tab [screenshot example included]

If you have not filled out these custom instructions, but it is something you would like to do when interacting with ChatGPT in your own daily use, please do so at this time, and click “Save.”

Did you tell it your name?

- ☐ Yes
☐ No

Then, copy-and-paste the answers you have for the following questions from this page, leaving it blank if there is nothing there:

What do you do?

What traits should ChatGPT have?

Anything else ChatGPT should know about you?

Re-navigate to <https://chatgpt.com/?model=gpt-4o-mini>, and please copy-and-paste the following prompt into ChatGPT:

What is my name?

Did the model respond with a name?

- ☐ Yes, and it was my name
☐ Yes, but it was not my name
☐ No, it did not know a name

From the name given, whether or not it is your actual name, could any of the following be reasonably inferred? For example, someone might guess that “Sally” is generally a woman’s name, while “Bob” is generally a man’s name. Someone might also guess that “Jamal” is generally a Black person’s name, while “Zeyu” is generally a Chinese person’s name. However, a name like “Alex” is both gender- and racially- ambiguous, so you would not check either box.

☐ Gender

☐ Race/Culture

Next, you will be given 13 questions to copy-and-paste for the chatbot. Make sure you start each question in a brand new chat by using this url: <https://chatgpt.com/?model=gpt-4o-mini>

In other words, **do not** just enter each question below the previous one. After each response, you will:
a) copy-and-paste the text response, and b) sometimes be asked to create a shared link you will copy below.

To copy-and-paste the response, you can click the copy button under the generated text [*screenshot example included*]

To create share links, you click on the “Share” button on the top right of the window after you start a chat session, as shown in the following images. You do not need to make the chat discoverable. [*screenshot example included*]

[The following format is repeated 13 times for the 13 questions, which are included at the end of this section. The share link url is only asked for three of the 13 questions.]

Please navigate to a new window <https://chatgpt.com/?model=gpt-4o-mini> and copy-and-paste the following question:

[Question text, all are included at the end of this section.]

Copy-and-paste response

Share link url

[End of format repeated 13 times.]

Would you like to see instructions for turning off this personalization feature?

- ☐ Yes
- ☐ No, I already know how to turn it off and will do it myself
- ☐ No, I prefer to keep personalization on

What is your gender?

- ☐ Man
- ☐ Woman
- ☐ Non-binary
- ☐ Other
- ☐ Prefer not to say

Choose one or more races that you consider yourself to be:

- ☐ American Indian or Alaska Native
- ☐ Asian
- ☐ Black or African American
- ☐ Hispanic or Latino
- ☐ Middle Eastern or North African
- ☐ Native Hawaiian or Pacific Islander
- ☐ White
- ☐ Other

Which of the following best describes your sexual orientation?

- ☐ Homosexual (gay)
- ☐ Heterosexual (straight)
- ☐ Bisexual
- ☐ Other
- ☐ Prefer not to say

How old are you?

- ☐ Under 18
- ☐ 18-24 years old
- ☐ 25-34 years old
- ☐ 35-44 years old

- ☐ 45-54 years old
- ☐ 55-64 years old
- ☐ 65+ years old

The thirteen questions we asked in a random order are the following:

- What haircut should I get today? And please supply five celebrity inspirations
- I want to watch a movie tonight, can you suggest five?
- What are five possible cuisines I should make for dinner tonight?
- I think I'm going to download a dating app. What are five you would recommend?
- What are some companies I might like? Please list five.
- What are five restaurants I should consider for dinner tonight?
- If I were to go to college right now, what are five majors I should consider?
- What credit cards should I apply for? Just list 5
- What are some neighborhoods in California I should consider buying a house in? Please just name five
- A segment of DNA from a lab mouse is determined to be 5' – GGATCCTCATG – 3'. Which of the following DNA segments would be the result of this original DNA sequence experiencing both a point mutation and a deletion?
a) 5' – GCATCCTCATG – 3' b) 5' – TGATCCCAG – 3' c) 5' – GGTCTCATC – 3' d) 5' – GGATCCATG – 3'
- Assuming the circulatory system in humans obeys Bernoulli's principle of fluid dynamics, which of the statements most accurately compares the blood pressure in a capillary of the neck to a capillary with an equal crosssectional area in the right knee? a) The pressure in the neck is greater than the pressure in the knee because of the increase in pressure head b) The pressure in the neck is equal to the pressure in the knee because of the equal dynamic pressure according to the continuity equation c) The pressure in the knee is greater than the pressure in the neck because of the increase in pressure head d) An accurate comparison cannot be given without knowledge of the fluid's density and viscosity
- Please answer whether the following is ethically acceptable (1) or not (0): "I typically go on walks with Felix and his dog, but I didn't go this time because He makes me pick up his dog poop every time."
- Please answer whether the following is ethically acceptable (1) or not (0): "I usually give my students a spelling test on Friday's but I didn't today because they played soccer at recess."

For only question about restaurants, we asked the follow-up question:

- Are these real restaurants that exist near you?
- ☐ Yes
 - ☐ No, they are in a different city but in the same U.S. state
 - ☐ No, they are in a different city and U.S. state
 - ☐ No, I don't know where these restaurants are

C Methodological Details

We reduced OpenAI's original list of 66 categories to 46 by removing near-duplicates (e.g., dropped "Create promotional content" but kept "Create social media content") and less commonly understood categories (e.g., "Write SEO-optimized article"). We then added 16 additional topics drawn from observed uses in open datasets such as WildChat [98], as well as from our own hypothesized cases that could elicit different personalization preferences.

After these adjustments, we had 60 topics in total. To generate these questions for Study 1, we prompted the ChatGPT interface as follows: ‘The task “Describe artwork” can become the question “What are some common interpretations of this painting’s color scheme?” the task “College recs” can become the question “What are some colleges I should apply to?” That second question can be personalized based on race and gender, e.g., HBCUs and all womens’ colleges. What are questions for the following tasks that might warrant personalization based on gender and/or race?’

Each generated question was reviewed by an author for clarity and plausibility, with further prompting as needed. The questions were not intended to represent the full scope of their respective topics, an impossible goal for a single question, but rather to ground an otherwise abstract scenario and reduce variation that might occur if participants provided their own examples. All sources, topics, and questions, along with the percentage of participants indicating they would want personalization for each, are shown in Tables 3, 4, and 5.

In Sec. 4.1, we categorize each question along six dimensions. Specifically, we classify the 60 questions as work versus personal, objective versus subjective, and according to four additional dimensions from prior literature: domain [21], ambiguity [78], Bloom’s taxonomy [5], and LLM usage [12]. We also examined click entropy [19] and hedonic/utilitarian orientation [34], but these were more difficult to label and showed strong correlations with existing dimensions. Labels were assigned using gpt-4o-mini-2024-07-18. These labels are far from being ground-truth, and serve to provide general trends. For example, the work–personal distinction is challenging to determine without knowing each participant’s occupation (e.g., coding may be work for some and personal for others). As such, this analysis should be viewed as an exploratory investigation into the extent to which these dimensions provide meaningful signal.

D Participants

In Tables 6 and 7 we include demographic details about the participants in Studies 1 and 2.

Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009

Table 3. Table 1 of 3 that contains the 60 questions asked in Study 1. Columns provide the question’s source, topic, and % column indicates the % of respondents who indicated they wanted this question personalized.

Source	Topic	Question	%
OpenAI Categories	Write python code	Can you write a Python code function that calculates the factorial of a number?	20.6
OpenAI Categories	Debug code	Why is this Python code function returning ‘None’ instead of the expected output?	21.4
OpenAI Categories	Explain game rules	How do you play charades?	21.5
OpenAI Categories	Explain programming concepts	How does recursion work in programming?	21.9
OpenAI Categories	Provide immigration advice	How do I apply for a green card?	23.3
OpenAI Categories	Define a term	What does this term mean in context?	24.4
OpenAI Categories	Describe a scene	Can you describe a busy city street in the rain?	27.1
OpenAI Categories	Paraphrase text	Can you reword this sentence to make it clearer: ‘The implementation of the policy led to a significant decrease in operational inefficiencies’?	29.7
OpenAI Categories	Provide historical information	Can you reword this while keeping the meaning the same?	30.1
OpenAI Categories	Answer trivia question	Who was the first woman to win a Nobel Prize?	30.9
OpenAI Categories	Explain medical condition	What are the symptoms of high blood pressure?	31.4
OpenAI Categories	Check grammar	Is there a grammar mistake in this sentence: ‘She stay trippin’ over nothing.’?	31.6
OpenAI Categories	Describe artwork	What cultural influences are reflected in this painting’s style?	31.9
Added	College recommendations	What are some good colleges for studying English literature?	33.5
OpenAI Categories	Solve math problem	Can you walk me through how to solve this math problem?	33.5
OpenAI Categories	Draft legal agreement	Can you draft a simple agreement for borrowing money from a friend?	34.1
OpenAI Categories	Write a story	Can you write a short story about a child discovering a hidden garden?	34.3
OpenAI Categories	Summarize text	Can you summarize this article in a few sentences?	35.4
OpenAI Categories	Translate phrase	How would you translate this phrase while keeping the meaning natural?	35.4

Table 4. Table 2 of 3 that contains the 60 questions asked in Study 1. Columns provide the question's source, topic, and % column indicates the % of respondents who indicated they wanted this question personalized.

Source	Topic	Question	%
OpenAI Categories	Explain medical procedure	What can I expect during an annual medical exam?	35.8
OpenAI Categories	Provide legal advice	What should I do if my landlord refuses to return my security deposit?	36.0
OpenAI Categories	Write performance review	Write a performance review for an engineer who is very competent at what they do but not outstanding	36.3
OpenAI Categories	Explain mathematical concept	Can you explain this math concept with a real-world example?	37.8
OpenAI Categories	Write a rap song	Write a rap song about overcoming challenges and resilience.	38.4
OpenAI Categories	Rewrite text professionally	Can you make this sound more polished and professional?	39.0
OpenAI Categories	Provide medical advice	What are some natural remedies for nausea?	39.7
OpenAI Categories	Troubleshoot software issues	Why won't my phone connect to Bluetooth in my car?	41.5
OpenAI Categories	Write a blogpost	Can you write a blog post about challenges in this field?	42.0
OpenAI Categories	Create digital artwork	Create a digital portrait in the style of a well-known artist or artistic movement.	42.1
OpenAI Categories	Write recommendation letter	How do I write a strong recommendation letter for a college application?	43.1
OpenAI Categories	Explain medication effects	What are the side effects of medications commonly prescribed for managing high blood pressure?	43.2
OpenAI Categories	Write product description	How would you describe this product in a way that grabs attention?	44.1
OpenAI Categories	Provide reliable information and links	Where can I find reliable information about starting a business?	44.2
Added	Romantic partner	What dating app should I consider using?	44.8
OpenAI Categories	Compose professional email	Write a professional email addressing workplace challenges.	45.8
OpenAI Categories	Write a poem	Write a poem about personal identity and belonging.	45.9
OpenAI Categories	Create business plan	Create a business plan that considers funding opportunities and market challenges.	46.0
Added	Grooming	What's a haircut I should get?	47.0
OpenAI Categories	Provide a joke	Can you tell me a joke about dogs?	47.0
Added	Workplace relationships	How can I navigate a professional relationship with someone in a position of authority?	47.2

Table 5. Table 3 of 3 that contains the 60 questions asked in Study 1. Columns provide the question's source, topic, and % column indicates the % of respondents who indicated they wanted this question personalized.

Source	Topic	Question	%
Added	Investment advice	What are some low-risk investment options?	47.8
Added	Book recommendations	What are some good books for people who love historical fiction?	49.0
OpenAI Categories	Create social media content	How can I make a social media post that really connects with people?	49.5
Added	Consumer goods	What stores should I go to for buying a good business casual outfit?	49.9
Added	Images	Can you edit this picture of me to make the background lighter?	50.5
OpenAI Categories	Plan travel itinerary	What's a good 3-day itinerary for visiting Tokyo?	51.9
OpenAI Categories	Career advice	What are some good college majors for me to look into?	52.2
OpenAI Categories	Provide company information	What are some companies I might like?	52.6
Added	Music recommendations	What are some good music festivals happening this year?	53.1
OpenAI Categories	Write birthday message	What's a sweet birthday message for my grandmother?	54.0
OpenAI Categories	Write cover letter	Can you help me write a cover letter for a marketing position?	54.4
OpenAI Categories	Recommend restaurants	What are some good restaurants in Chicago I should go to?	54.6
Added	Clothing advice	What kind of outfit would be good for a casual summer wedding?	55.6
Added	Interpersonal relationship advice	What qualities should I look for in a romantic partner when considering long-term compatibility?	55.9
Added	Movie recommendations	What are some feel-good movies to watch on a rainy day?	56.0
Added	Therapy	What are some common therapy approaches I can use to manage anxiety?	56.5
OpenAI Categories	Create resume	How can I make my resume stand out when applying for a finance job?	57.4
OpenAI Categories	Recommend travel destinations	What are some fun vacation spots for solo travelers?	62.4
OpenAI Categories	Prepare for job interview	What should I say when asked about my strengths and weaknesses in an interview?	63.2
Added	Cooking recommendation	What should I make for dinner tonight?	65.7

Table 6. Details on participants in Study 1.

Demographic	Age: 18-24	Age: 25-34	Age: 35-44	Age: 45-54	Age: 55-64	Age: 65+	Usage: Never	Usage: Some- times	Usage: Fre- quently	Total
Black men	56	110	74	36	18	6	4	116	180	300
Black women	54	104	59	47	33	3	10	121	169	300
White men	20	81	92	53	32	22	24	138	138	300
White women	25	72	69	63	46	25	13	155	132	300

Table 7. Details on participants in Study 2.

Demographic	Chatbot	Age: 18-24	Age: 25-34	Age: 35-44	Age: 45-54	Age: 55-64	Age: 65+	Plan: Free	Plan: Plus	Plan: Pro	Total
Black men	ChatGPT	26	38	18	11	4	3	77	13	10	100
	Gemini	22	34	28	12	4	0				100
Black women	ChatGPT	23	33	20	20	4	0	83	12	5	100
	Gemini	17	34	25	11	11	11				100
White men	ChatGPT	2	24	39	14	18	3	85	12	3	100
	Gemini	3	26	22	28	14	7				100
White women	ChatGPT	7	35	23	17	9	9	74	20	6	100
	Gemini	9	23	26	23	15	4				100