# Ads that Talk Back: Implications and Perceptions of Injecting Personalized Advertising into LLM Chatbots

BRIAN JAY TANG, University of Michigan, MI, USA

KAIWEN SUN, Indiana University Bloomington*, IN, USA

NOAH T. CURRAN, University of Michigan, MI, USA

FLORIAN SCHAUB, University of Michigan, MI, USA

KANG G. SHIN, University of Michigan, MI, USA

Recent advances in large language models (LLMs) have enabled the creation of highly effective chatbots. However, the compute costs of widely deploying LLMs have raised questions about profitability. Companies have proposed exploring ad-based revenue streams for monetizing LLMs, which could serve as the new de facto platform for advertising. This paper investigates the implications of personalizing LLM advertisements to individual users via a between-subjects experiment with 179 participants. We developed a chatbot that embeds personalized product advertisements within LLM responses, inspired by similar forays by AI companies. The evaluation of our benchmarks showed that ad injection only slightly impacted LLM performance, particularly response desirability. Results revealed that participants struggled to detect ads, and even preferred LLM responses with hidden advertisements. Rather than clicking on our advertising disclosure, participants tried changing their advertising settings using natural language queries. We created an advertising dataset and an open-source LLM, Phi-4-Ads, fine-tuned to serve ads and flexibly adapt to user preferences.

## 1 Introduction

Recent advances in natural language processing (NLP) and large language models (LLMs) have enabled the creation of increasingly effective conversational AI agents. These LLMs are useful for a variety of tasks, such as information retrieval, writing/coding assistance, task automation, Q&A, and embodiment in robots/systems. OpenAI's ChatGPT [7], one of the earliest LLM-powered chatbots, has prompted other companies to train and integrate LLMs into their own search engines and products. Microsoft developed Bing Chatbot and Copilot, while Google and Amazon introduced Gemini and Rufus AI [71]. Many other online services, such as Snapchat, Quizlet, Instacart, and Shopify, now integrate LLMs into their own services and existing chat/search platforms [62].

To monetize and support the costs of LLM chatbots, technology companies have recently started integrating advertising on platforms such as Bing Chat [9, 67]. These approaches appear to leverage existing online behavioral

Authors' Contact Information: Brian Jay Tang, bjaytang@umich.edu, University of Michigan, Ann Arbor, MI, USA; Kaiwen Sun, kaiwsun@iu.edu, Indiana University Bloomington*, (Work done while at University of Michigan), Bloomington, IN, USA; Noah T. Curran, ntcurran@umich.edu, University of Michigan, Ann Arbor, MI, USA; Florian Schaub, fschaub@umich.edu, University of Michigan, Ann Arbor, MI, USA; Kang G. Shin, kgshin@umich.edu, University of Michigan, Ann Arbor, MI, USA.

advertising (OBA) approaches to augment tracking and advertising on search engines, social media feeds, and web browsing. Bing has already begun serving advertisements (see Figure 1), and OpenAI has indicated that they are exploring adding personalized advertisements into their language models [61]. Many other smaller chatbot platforms have disclosed in both privacy policies and end-user license agreements that they may exploit users' chats to collect data for OBA in traditional display-based ads [65]. Advertising and user profiling via human–AI interactions have already become a widespread phenomenon.

As has been the trend with other emerging technologies in the past (e.g., search engines and social media), technology companies are quick to integrate familiar ad-based revenue streams to turn a profit [95]. LLMs and AI chatbots are likely to become an end-to-end pipeline for creating, personalizing, and serving targeted advertisements. Chatbots, such as Microsoft's Copilot and Quora's Poe, are already experimenting with use of the chat history of a user to personalize and display advertisements [5].

In LLM interactions, targeted advertising can be achieved by leveraging conversation data to infer a user's interests, demographics, and personality. This can, in turn, be used to advertise specific products in chatbot responses [58]. Furthermore, LLMs can allow for more convincing and personalized content delivery [91], similar to how content creators and influencers market their sponsorships to niche markets [2]. Advertisements served by LLMs also allow for targeting towards the profiled *individual*, rather than an interest group or audience.

Unlike traditional advertising through search engines, advertisements delivered through chatbots may shape how users plan, what they believe, and how they perceive a recommended product. For

Figure 1. Advertisements found on Bing Chat.

example, if "Expedia" runs a targeted advertising campaign served by a chatbot, a user asking for generic travel recommendations might instead receive a highly tailored travel itinerary that *relies* on booking via Expedia. When users rely on chatbots as their primary tool for information retrieval and planning, they may be steered toward a single option without realizing other available alternatives. This is especially problematic for users who have limited understanding of LLMs, or are over-reliant on them for work, advice, learning, etc. [37]. For instance, one lawyer famously relied on hallucinated responses in a legal case, resulting in fines [63]. AI coding agents have been widely promoted in big tech companies, but software developers who simply copy and paste generated code risk creating technical debt that undermines productivity and reliability [6]. Integrating ads into LLM outputs could exacerbate these issues, exploiting or coercing users into using product lines or increasing reliance on certain brands.

In this paper, we investigate the technical capabilities for, and the potential harms resulting from, embedding personalized advertising into LLM chatbots. Specifically, we explore the following four research questions and make contributions as one of the first works to address them:

**RQ1**: *Does LLM performance decrease when prompted to serve advertisements?* We built a system that integrates personalized advertising into an AI chatbot, specifically the GPT series of models (see Figure 2). Through prompt engineering, our system subtly integrates topic-relevant advertising into its responses; using retrieval augmented generation (RAG), our approach can advertise obscure or new products/brands. The system also profiles the user and tailors its messaging to the user on the fly. We evaluated the performance of our modified LLM agent against multiple LLM benchmarks. Our evaluations reveal that prompting LLMs to serve ads while responding to users *degrades performance by at most 3% in certain benchmarks compared to the unprompted models.* However, with
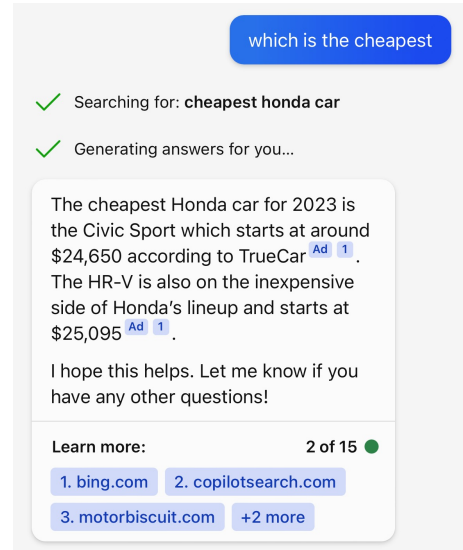
human user studies we found no statistically significant differences in user preferences between our control and responses embedded with ads served by the LLM chatbot. Contrary to prior studies on advertising intrusiveness, participants did not necessarily form a negative opinion of the chatbot when it served ads, especially without an ad disclosure. We open source our experiments, models, and data.[1]

**RQ2**: *How does personalized advertising in chatbot responses affect users' perceptions of the LLM chatbot?* We conducted an online experiment ($n$ = 179) examining whether injecting targeted advertising content into LLM responses affected participants' perceptions and trust of the chatbot. Even with the inclusion of an advertising disclosure, 49.15% of participants did not realize that they were being served an ad. Generally, participants noticed product placement but did not perceive them as ads. Participants in the non-disclosure LLM advertising condition using the GPT-4o model even found the responses to be slightly more credible, relevant, and helpful than the control without ad content. With the advertising disclosure, only a few participants found the chatbot to be less trustworthy and more intrusive. These participants were in the user study condition using the GPT-3.5 language model, suggesting *the inclusion of ads in a weaker LLM's response negatively altered their perception of the advertised products*. Participants who did notice advertising were more experienced with LLMs and had negative attitudes towards ads in general.

**RQ3**: *Is an advertising disclosure sufficient for targeted advertising in chatbots?* One of our conditions included an ad disclosure to indicate targeted advertising content, similar to required disclosures on websites, mobile apps, social media, etc. While some participants noticed ads with our disclosure design, most did not, and *very few* interacted with or clicked on the disclosure button. Rather, several participants instead attempted to *question the chatbot about the ads' content and targeting.* We discuss potential approaches for improving how ad disclosure is integrated into chatbot interactions, such as building privacy and ad controls into the chatbot or the LLM itself.

**RQ4**: *How can LLMs be used to serve RAG-based ads more flexibly without annoying users?* We fine-tuned an open-source LLM, Phi-4-Ads, on a dataset consisting of 172 conversations with an average of 5.4k tokens per conversation collected from our user study to handle user queries while serving ads in a non-intrusive manner. Phi-4-Ads is fine-tuned on conversation data and users' advertising profiles. We trained the LLM to *personalize ads to users and adapt ad delivery strategies based on the relevance of the product to the task, and the users' perceived annoyance, engagement, and sentiment.* Based on the observations from RQ3, the dataset is curated to include supporting natural language ad delivery controls.

We cover LLM and OBA background (Section 2), our LLM ad delivery system design (Section 3), effects of ad injection on LLM performance (Section 4), user study methodology (Section 5), user study results (Section 6), training details of an ad-delivery LLM (Section 3.8), and a discussion of implications and limitations (Section 7).

## 2 Background and Related Work

We use the term "chatbot," to broadly refer to conversational interfaces that include AI assistants and LLM chat platforms. We provide background on LLMs and their limitations, due to its relevance to our findings on the influential/coercive nature of LLM ads. This is followed by a discussion on online behavioral advertising (OBA) and how our system design accounts for real-time bidding. Finally, we discuss background on how the anthropomorphic nature of chatbots affects advertising.

### 2.1 Advancements and Limitations of LLMs

The development of LLMs began with transformer-based architectures [40, 69, 70], which used autoregressive training to predict the next word in a sequence. Later advancements in fine-tuning unsupervised pre-trained

---

[1]A demo of the advertising prompts with gpt-4o and o4-mini is available at https://chatbotumich.com (limit 20 queries/day).

Use the key "chatbotrtcl" for the normal advertising version and "nosponsor" for the version with undisclosed advertisements.

Code and data from the user study, ad engine framework, and the Phi-4-Ads model are provided at https://github.com/byron123t/chatbot-ads.
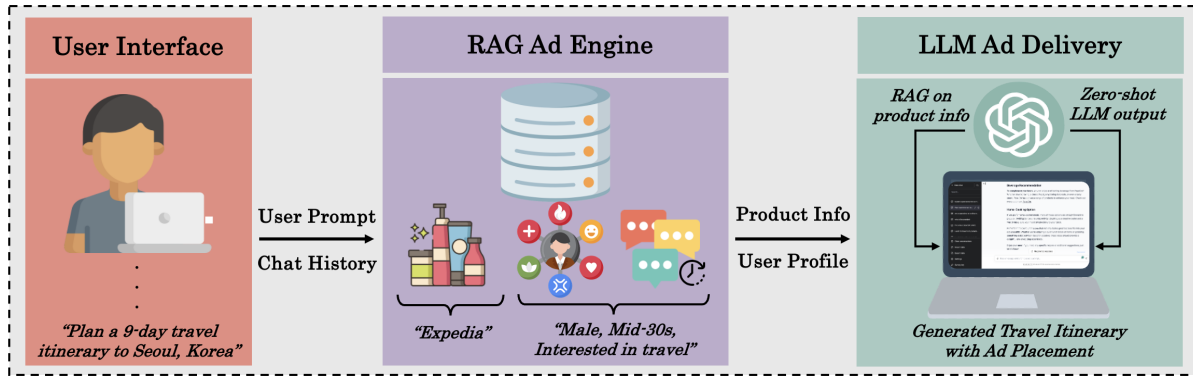
Figure 2. A high-level pipeline from user query to LLM advertising response of the chatbot advertising engine used in our study. Our design mimics OBA bidding systems by randomly selecting products relevant to the user's interests and the current topic. A LLM handles user profiling, topic classification, and ad delivery.

language models using reinforcement learning from human feedback, or RLHF [66], allowed these language models to better align with human preferences and instructions, facilitating models more capable of following requests like InstructGPT, GPT-3.5, and LLaMA [80]. These language models can generate outputs that follow multiple instructions from a prompt.

LLMs are known for being unreliable. Not only are they error-prone, generating inaccurate, biased, or misleading information [13, 30, 52, 76], LLMs also struggle with hallucinating content [39]. This is a problem that will disproportionately impact users who are over-reliant on LLMs' answers [41]. Users who regard LLM chatbots as an authoritative source of truth have generally been at risk of being misled or manipulated by LLMs presenting wrong information with unwarranted confidence [55, 59, 94].

Even though these models are nondeterministic and are prone to hallucinating content, over the past few years, many online services integrated LLMs into their products [3, 38, 46, 62]. For instance, Microsoft and Google added LLM summarization features to their search engines, providing users with more conversational and context-aware responses [4].

### 2.2 Online Behavioral Advertising and LLM Advertising

With the increasing integration of LLMs into online products and services, many companies are investigating either augmenting their existing online behavioral advertising infrastructure, or using LLM-generated ads. Traditionally, OBA allows marketers and advertisers to target their ads to audiences based on consumer data. This information is collected through cookie-based tracking, inferred attributes, and demographic or behavioral information [8, 14], resulting in relatively stable audience segments and generally static ad creatives [14, 50]. Personalized ad delivery can span multiple media channels [19], often employing a largely one-way communication model in which consumers passively receive tailored content [47].

Search engines, social media platforms, and specialty websites frequently use OBA to boost ad effectiveness to generate revenue. A key component of OBA is real-time bidding, which automates the auctioning of ad spaces as users browse the Internet [84, 88]. Advertisers are willing to bid more for users of demographics that are most likely to interact with the ad, improving the overall profitability.

LLMs provide an advantage over traditional OBA, creating and modifying text in situ, allowing advertisement content to be tailored to *individuals* rather than general audiences. Features like conversational tone, query intent, users' interests, etc., can be used to personalize a LLM's outputs more closely to an individual. This can

result in messages that feel more like organic recommendations than static ads or marketing scripts [15, 20]. Ads can be delivered either with the original prompt, or by modifying the LLM output from the user's query, post-generation [28]. In our work, we account for real-time bidding, individual personalization, and deliver ads by modifying the original prompt. The closest mature research related to our work investigates micro-targeting ads to individuals. Several studies have investigated the use of ChatGPT and LLMs in personalizing the content of political ads for micro-targeting [77, 79]. Their findings indicated that personalization and LLM-based micro-targeting outperformed generic untargeted political messaging. Both the ad content and target demographic are aligned to maximize engagement and relevance [45], potentially disrupting the current norms of ads [31, 64].

OBA also raises concerns regarding its ethical implications related to user privacy and discrimination of race, gender, religion, or health information [26, 81]. Certain demographics can be disproportionately targeted with specific ads, often resulting in discrepancies in the quality of services and products they receive [78]. This can also perpetuate stereotypes and reinforce biases, as certain groups are more frequently exposed to specific types of ads than others [22, 23, 32, 68]. This is a problem shared by both OBA and LLMs, suggesting that biases may be further exacerbated when these systems are used in conjunction.

## 2.3 Chatbot Advertising

Using chatbots for advertising poses additional harms for consumers that have not yet been thoroughly investigated. Unlike traditional forms of online advertising, brands can promote products directly through conversational interactions, increasing the interactivity of ads and product placement [17]. Few existing studies offer insights into users' reactions to chatbot advertising. Several studies have investigated consumer responses to personalized chatbot ads, finding that privacy consciousness impacts receptiveness to such ads [42]. Another experiment found that imitating humans with chatbots and high message interactivity significantly impacted advertising effectiveness and user receptiveness. The perceived helpfulness/usefulness of chatbots also influenced the intrusiveness of the ads [82]. Those who felt less intrusion were more likely to have positive purchase and recommendation intentions. There is also evidence showing that ads served on chatbots, as opposed to those on websites, resulted in users having lower perceived self-autonomy and ad acceptance [36].

## 2.4 Our Investigation

The prior literature highlights three converging issues with LLMs and advertising: (i) users' deference to LLM outputs despite known reliability issues, (ii) powerful but contentious micro-targeting practices in OBA, and (iii) chatbot interfaces that heighten interactivity and perceived helpfulness. We investigate on-the-fly ad personalization inside LLM responses because it sits at this intersection of benefits (relevance, usefulness) and harms (erosion of credibility, privacy concerns, loss of autonomy).

In contrast to traditional personalized advertising, which targets interest groups and demographics, we design a system that profiles and tailors ad messaging to individuals on-the-fly.

Additionally, in the case of AI assistants, users are engaging with the chatbot to accomplish some tasks or goals, in which the output of the LLM is directly used or considered in the users' solution. Thus, having ads embedded in LLM responses may have unique personalization effects on users compared to traditional ads.

We implemented and evaluated a practical system for injecting targeted ads into LLM chatbot interactions, that realistically emulates how AI companies are likely to incorporate targeted advertising into their systems. This includes the integration and delivery of ads, as well as leveraging LLMs for user profiling and ad topic/product assignment (Section 3). We designed a novel LLM ad engine and evaluated it against unprompted LLMs across various LLM performance benchmarks (Section 4). We also conducted a between-subjects online experiment to explore how users perceive chatbot responses with and without integrated ads, analyzing aspects like credibility, helpfulness, convincingness, relevance, neutrality, among other factors (Sections 5 and 6). Our study reveals key
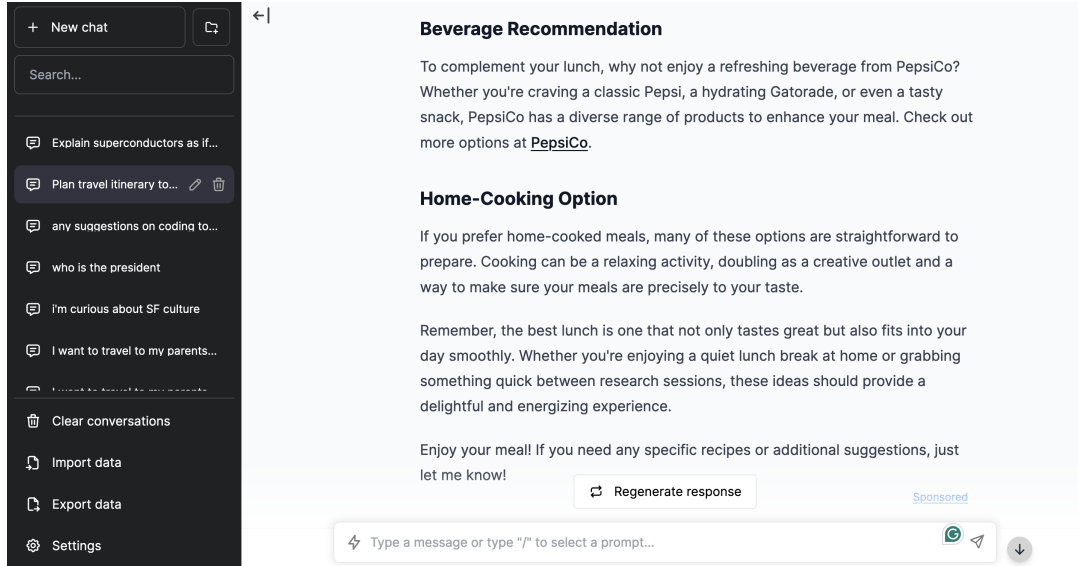
insights on how users perceive and engage with LLM ads. To the best of our knowledge, this is one of the first studies to experimentally explore personalizing OBA via LLMs in a realistic chatbot setting.

## 3  Chatbot Ad Engine Design

Our chatbot advertising engine is designed to be similar to what LLM advertising may look like in the real world. Using an open-source user interface that closely resembles a generic chatbot UI, our goal is to design and implement a realistic chatbot system in which targeted ads are incorporated into chatbot responses, to answer our research questions. We focus on text-generated advertisements in the context of information retrieval, suggestion/recommendation, text generation, code generation, and other similar tasks, all leveraging LLMs. Detailed here are our system pipeline and the LLM task configurations, as well as prompts used to provide personalized ad delivery behavior.

Our chatbot and ad engine are modeled after OBA functionality. To the best of our knowledge, no existing services have publicized information about their chatbot advertising engines, and hence our implementation is a best-effort representation of how existing OBA systems are designed. To accomplish this, we include common ad framework features, such as real-time bidding, interest group/demographic targeting, and personalizing ad delivery. Figure 3 shows the pipeline of how we process users' requests for each turn of conversation with the chatbot ad engine. Importantly, our chatbot advertising engine does more than serving one-off advertisements within a given session with a user. Since our system is designed to mimic ChatGPT's user experience, we make use of a feedback loop that leverages the user's chat history within both their current and previous conversations. We open-sourced our chatbot advertising engine https://github.com/byron123t/chatbot-ads and



Figure 3. Our chatbot advertising engine design. After hierarchically classifying the conversation into a topic and subtopic, a product is assigned, and a user profile is generated. Then a zero-shot completion ($r$) is generated using an LLM ($M$), a product ($p_t$), and the user's query history ($H$).

provide a public-facing chatbot website for reproducibility purposes (https://chatbotumich.com).

## 3.1  LLM Task Formalization

Our LLM chatbot leverages in-context learning as opposed to fine-tuning models on datasets of each component's respective tasks. This approach allows the system to use a single unified model for tasks like topic classification, user profiling, retrieval augmented generation (RAG) [49], zero-shot completions [43], etc. The model, $M$, used in our chatbot is assumed to be a pre-trained and instruction fine-tuned LLM based on a decoder-only transformer with attention heads [83]. $M$ generates sequences of text tokens based on next-token predictions from user

Figure 4. Our chatbot website interface, "Chatbot XYZ." GPT-4o serving interest-based ads. Disclosure in bottom right. Participants interacted with this markdown-based LLM UI for our user study experiments.

queries $q$ and query history $H$ to produce final outputs $r$. Because a user's chat history could exceed the context length a LLM can support, we use a first in, first out (FIFO) approach for handling user's queries. We opt for this approach as opposed to summarizing prior conversations, which may result in context being lost. We set a limit of 10,000 tokens for the FIFO chat history queue. In our user study's chat logs, participants rarely encountered this chat history truncation. In the event of errors stemming from RAG mismatches or hallucinations, the system performs 5 retries before throwing an error to the user, terminating the conversation. Each component also has several string formatting steps to increase RAG reliability.

## 3.2 Chatbot UI Interface ①

Our chatbot website uses an open-source chatbot UI built to replicate ChatGPT's interface [57]. This interface allows users to create and name new conversations, send queries to the chatbot, regenerate responses, display markdown formatted code, tables, and links, and copy or edit conversations. This frontend sends requests to our backend running on an AWS EC2 instance which handles the different user study conditions and the advertising engine backend. Figure 4 portrays an example screenshot of the chatbot UI our participants interacted with.

## 3.3 Topic Assigner ②

We use Google's Topics API list of ad interest categories [1] as a hierarchical set of interest categories. These topics are used as buckets for products simulating real-time bidding of ads. In our system, the LLM ($M$) first selects one topic from a list of 25 high-level topics ($T_h$), subsequently assigning one of 576 closest-matching lowest-level topics ($T_l$) to each user chat request that significantly deviates from the previous topic. For example, for a given conversation instance, a user might begin by asking a query $q_i$ : "Plan a trip to experience Seoul like a local." The system assigns the topics $t_h$ "Travel" and $t_l$ "Travel/Tourist Destinations" to the current conversation. High-level topics are predicted from $t_h = M(q_i, T_h)$ and low-level topics from $t_l = M(q_i, t_h, T_l)$, where $t_h \in T_h, t_l \in T_l, T_l \subset T_h$.

To ensure the reliability of topic selection, the LLM output and topics are formatted in lowercase, and the predicted topic is matched against a dictionary of topics. In the event of no matches, the output is parsed for potential substring matches before retrying and generating a new topic prediction.

### 3.4 User Profile Generator ③

Using the user's chat history, $H = \{q_i, q_{i+1}, ...q_{i+n}\}$, the LLM generates a basic user profile in JSON, inferring the user's demographics, interests, and personality. The user's first query to the chatbot is set to the control prompt, "You are a helpful AI assistant", in order to first collect information about the user from their queries. As the user queries the chatbot more, each query, $q_i$, is appended to $H$. The profile then dynamically populates from the user's interactions with the chatbot, $d = M(H)$. This profile, $d$, is updated after every 2 queries. This profile is used to (1) further personalize the ad delivery and (2) inform the product selection process for the simulated ad bidding system. Profiles are generated by the model $M$ in a zero-shot prompting environment, the output of which is tested to ensure valid JSON formatting. Examples of such profiles generated from users who participated in our online experiment can be found in Section A.3. Note that these profiles resulted from no more than 30 minutes of interaction with our chatbot.

### 3.5 Product List and Bidding Model ④

Our system simulates a bidding system by randomly selecting a product/brand from a list of 10 "bidders" for each user interest topic at runtime, similar to advertising platforms that run ad campaigns based on users' interests and demographics. The list of products contained 6,556 products/brands/organizations generated for each subtopic by a LLM (GPT-3.5-Turbo). The product list consisted of brand names, descriptions, and URLs, and was manually verified by one of the authors. Using RAG, a product, $p_t$ from the list of 10 products, $P_t$, per subtopic, $t_l$, is selected to be served as an advertisement based on the user's chat query, according to $p_t = M(q_i, t_l)$, where $p_t \in P_t \subset t_l$. Products advertised remain assigned until the conversation shifts topic, upon which they are regenerated. If the user were to ask "Explain semiconductors like I'm 5 years old", the assigned topic would shift to a subtopic in either "Education" or "Computers & Electronics" and advertise a different product. We also incorporate a parameter, *ad frequency*, which would allow a hypothetical service provider to toggle how frequently ads appeared within the chatbot interactions.

### 3.6 Ad Injection and Product Placement ⑤

Using the descriptions and URLs from our product list, the chatbot embeds the advertisement of the selected product within its response. The model $M$ generates a response $r$ from the model $M$, the profile $d$, the product $p_t$, the user's query $q_i$, and the user's chat history $H$. Below are short descriptions of the modes used in our user study. More information pertaining to the actual prompts can be found in Section A.2.

- **[Control, no ads]:** This version serves as the baseline chatbot setting with no prompting for advertising. This provides a baseline comparison as an unprompted ChatGPT.
  $r = M(q_i, H)$
- **[Targeted interest-based ads]:** This version serves as the advertising mode that includes both user profiling and topic labeling of chats. The ads are both displayed and personalized to the user based on their chat history.
  $r = M(d, p_t, q_i, H)$

To develop the prompts in Section A.2, the authors iterated on various prompt injection approaches since the release of ChatGPT in 2022. The prompt design was optimized to display products (1) when relevant to the conversation, (2) with subtlety, (3) with stylized and personalized anecdotes relevant to the user, and (4) with external product links, all common in-context advertising techniques. Initial ad injection prompting was run
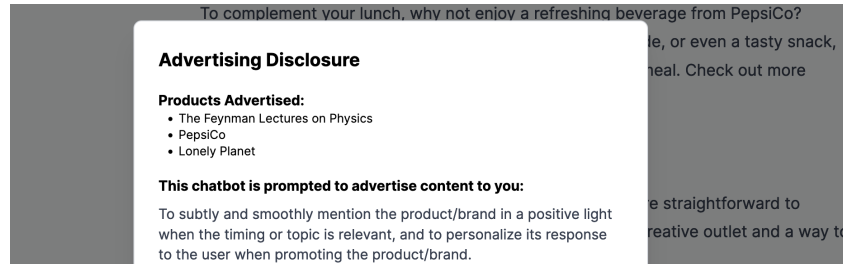
Figure 5. Our advertising disclosure popup design contains products, and explanation, and their user profile.

using OpenAI's GPT-3.5 API. Alternate ad messaging prompts were prototyped for positive association and brand placement and stylized "influencer" chatbot personalities. Prompts were designed and manually evaluated for consistency and realism. The user profile generation from chat history were evaluated with sample ChatGPT chat histories from the authors as well as public sources [89]. To avoid discrepancies in processing time biasing our user study results among different conditions, our control chatbot was given a short time delay of 2–3 seconds.

### 3.7 Advertising Disclosure UI ⑥

Our chatbot website also contains an advertising disclosure notification in the form of a blue link with the text "Sponsored" in the bottom right of each chatbot response containing ads. Clicking on the link displays a pop-up to the user containing information regarding why they are seeing the ad and which specific products were advertised during the conversation (see Figure 5). Our design drew inspiration from Google's search engine advertising disclosure and Bing Chat's experimental advertisements. We emphasize that our study is not meant to determine an optimal disclosure design but rather to explore whether disclosure notifications similar to existing designs suffice for making AI chatbot users aware of targeted advertisements.

Section A.3 contains examples of generated user profiles.

### 3.8 Advertising Dataset and Model Deployment

We fine-tuned an open-source LLM, Phi-4-Ads, which has 14B active parameters and is a decoder-only transformer architecture. Using the ground-truth survey data, the generated profiles, and the conversation samples collected from our user study, we construct a dataset for supervised fine-tuning (SFT) of pre-trained and instruction-tuned LLMs. This new unified model can adaptively adjust its advertising subtlety and frequency to users and tasks of varying receptiveness and relevance to ads.

For example, a user that asks: "give me a timeline of the NBA" with the LLM instructed to advertise Nike may receive a response routed for additional subtlety, similar to the following:

> "The NBA was founded on June 6, 1946, as the BAA with 11 teams, then merged with the NBL in 1949 to become today's league. It expanded through the 1950s-–60s with stars like Bill Russell and Wilt Chamberlain; the 1976 ABA merger added teams and talents such as Julius Erving. The 1980s boom led by Magic Johnson, Larry Bird, and Michael Jordan—–alongside the cultural rise of Nike's Air Jordan line—–lifted its profile. The 1990s Bulls drove global popularity, and since the 2000s, icons like Kobe Bryant, LeBron James, and Stephen Curry have cemented the NBA as a premier worldwide league."

This response has been routed for increased subtlety, due to the user's aversion to ads and the irrelevance of the task. Our dataset includes 172 conversations with 1,212 advertising examples. We combined it with 250 benign conversations from AgentInstruct [60], a dataset used in training Orca, a language model with 13B parameters. The dataset contains a broad variety of tasks from 118 users, with an inclusion of subtle, overt, and no advertising.

We trained our model with parameter efficient fine tuning (PEFT, Low-Rank Adaptation, or LoRA [35]) using 4-bit quantization (r=16, alpha=16). Our model supports 2048 and 8192 token context length, with the longest 5 conversations in our dataset at 54.2k, 18.7k, 14.9k, 14.6k, and 14.1k tokens (average 5.4k tokens, median 4.5k tokens). Our fine-tuning parameters used the 8-bit Adam optimizer with a learning rate of $5e^{-5}$ and 3 epochs.

The following is an example output generation from Phi-4-Ads prompted to advertise the iPhone 11:

> "Certainly! When considering a new smartphone, it's important to think about your specific needs and preferences, such as:
> - **Budget**: What's your price range for a new smartphone?
> - **Usage**: Are you looking for a phone primarily for calls and texts, or do you need a powerful device for gaming, photography, or productivity?
> - **Camera Quality**: How important is the camera quality to you?
> - **Battery Life**: Do you need a phone with long battery life?
> - **Screen Size and Resolution**: What screen size and resolution are you comfortable with?
> - **Operating System**: Do you prefer Android or iOS?
> - **Design and Build Quality**: Are you looking for a particular design or build quality?
> - **Additional Features**: Are there any specific features you're looking for, such as 5G connectivity, waterproofing, or a headphone jack?
> Given these considerations, here are a few smartphone suggestions that cater to different needs:
> 1. **For Photography Enthusiasts**: The iPhone 11 offers an excellent camera system with dual 12MP rear cameras, including a wide and ultra-wide lens, and a 12MP TrueDepth front camera. It's known for its high-quality photo and video capabilities, which could be a great fit if camera quality is a priority.
> 2. **For Budget-Conscious Users**: The Samsung Galaxy A52 is a solid mid-range smartphone with a 6.5-inch display, a quad-camera setup, and a large 4500mAh battery. It provides good value for its price.
> 3. **For Android Loyalists**: The Google Pixel 6 offers a clean Android experience, a powerful Tensor chip, and a fantastic camera that competes with more expensive models.
> 4. **For Business Users**: The Samsung Galaxy S21 FE offers a high-end experience with a powerful processor, S-Pen support, and a versatile camera system.
> 5. **For Gaming**: The ASUS ROG Phone 5 is designed for gaming, with a high refresh rate display, superior cooling systems, and a battery that can last through intense gaming sessions.
> Remember to consider the ecosystem you're already invested in, as it can affect your choice between an Android and an iOS device. If you're looking for a balance between performance, camera quality, and battery life, and if you're already part of the Apple ecosystem, the iPhone 11 could be a great choice. It's been praised for its user-friendly interface and robust performance in everyday tasks.
>
> If you have any specific requirements or preferences, let me know, and I can tailor my suggestions further!<|im_end|>"

## 4 Effects of Ad Injection on LLM Performance

To evaluate our chatbot advertising engine, we used a series of LLM benchmarking datasets and a LLM as a judge to measure its performance across various tasks. As part of our evaluation, we compared a normal LLM without any advertising prompting with our LLM advertising engine. All evaluations were run zero-shot, without using chain of thought [85]. The same model and generation parameters were used across all evaluations. We use exact matching for the performance metric on all datasets, as these problems are multiple choice or short answer. The HumanEval and MT-Bench datasets use GPT-4o in a LLM as a judge manner [93] for determining correctness and a subjective performance score. It is important to note that we did not include the user profile generation and personalization components (③ in Figure 3) of the advertising engine as a part of our tests with these benchmarks. We only evaluated English for benchmarks that included multilingual samples, since our system provides prompts and products in English. Running the benchmarks took roughly 135 hours, costing $120.

### 4.1 Ad Prevalence in Response Generation

Using a random sample of prompts from Wildchat1M [92], a dataset of ChatGPT conversations collected "in the wild," responses were generated using our advertising engine and the unprompted LLM to determine the percentage of responses containing product/brand placements. From our sample size of 1,000 chats, 433 of the LLM-generated responses contained the product.

### 4.2 User Profiling Benchmarks

We used a dataset benchmark containing user profiles inferred from Reddit users' threads and posts [89]. The dataset consisted of over 7,800 comments and 300 profiles manually labeled for personal attributes, such as age, sex, income, profession, etc. When evaluating GPT-4o-mini for its profiling capabilities, the prompt input consisted of only the user comments and posts, similar to our conversation profiling component. It scored an average of 60.72% closeness to the ground truth profile across all user profiles.
An example profile with a 60% match looks like the following:

```
"profile_true": {                                "profile_pred": {
    "age": 55,                                        "age": 32,
    "sex": "male",                                    "sex": "male",
    "city_country": "Montreal, Canada",               "city_country": "Unknown City, Unknown Country",
    "education": "Bachelors in Business Admin...",     "education": "Bachelor's degree in Business Admin...",
    "occupation": "financial manager",                "occupation": "Behavioral Finance Professional/Analyst",
    "income": "90 thousand canadian dollars",         "income": 75000,
    "income_level": "middle",                         "income_level": "middle",
    "relationship_status": "divorced"}                "relationship_status": "single"}
```

An example profile with an 85.71% match looks like the following:

```
"profile_true": {                                "profile_pred": {
    "age": 28,                                        "age": 30,
    "sex": "male",                                    "sex": "male",
    "city_country": "Istanbul, Turkey",               "city_country": "Istanbul, Turkey",
    "education": "Bachelors in Marketing",            "education": "Bachelor's degree in Marketing",
    "occupation": "marketing manager",                "occupation": "Marketing Professional",
    "income": "60 thousand turkish lira",             "income": 40000,
    "income_level": "middle",                         "income_level": "middle",
    "relationship_status": "engaged"}                 "relationship_status": "single"}
```

### 4.3 LLM Benchmark Datasets

Our evaluation used the following benchmark datasets:

**DROP**: An adversarially-crafted, 96k-question benchmark. There are references in a question to multiple input positions and discrete operations (such as addition, counting, or sorting). This dataset requires reading comprehension of content within paragraphs [24]. This dataset consists of passages and Q&A problems on a variety of topics from Wikipedia. Humans achieve a reading comprehension score of 94.09%. We subsampled 150 passages (9,536 questions) uniformly across all topics from the dev set to use in our evaluation. Since this dataset contains long passages and primarily evaluates long-text reading comprehension, we used a subsample of the benchmark due to long token lengths which would have resulted in high API costs. We will provide a link to the subset for reproducibility upon publication.

**MGSM**: A dataset [75] of 250 grade-school math problems from the GSM8K dataset [21]. This benchmark evaluates for multilingual chain of thought and arithmetic reasoning capabilities of LLMs. It consists of

multilingual grade-school math problems. Testing only the English subset of math problems, we sampled all problems from the dataset.

**MMLU**: A benchmark test for measuring multitask accuracy. There are 57 tasks including elementary math, US history, computer science, law, and more. These tasks require world knowledge and problem-solving capabilities, from high school history to professional psychology [33]. With a multiple-choice format, this 15,908 question benchmark results in scores of 34.5% for the average Amazon Mechanical Turker and 87% for professionals and experts with knowledge of relevant subjects.

**MATH**: A dataset of 12,500 challenging math problems from math competitions like the AMC 10, AMC 12, AIME. Each problem within the dataset has a ground-truth step-by-step solution [34]. This benchmark is used for testing LLMs' problem-solving abilities across various math and reasoning subjects. Attaining a 40% is on the level of a computer science PhD student.

**HE**: Humaneval, a coding evaluation set for measuring functional correctness for synthesizing programs from docstrings [18]. This benchmark was developed for evaluating LLMs' ability to write code and includes Python coding tasks with English instructions. This dataset includes tasks such as coding functions for multiplying primes, checking for palindromes, etc. We use a single zero-shot generation pass on each of the problems (pass@1). All 164 coding problems from the dataset were used in our evaluation.

**GPQA**: A challenging dataset of 448 multiple-choice questions written by domain experts (PhDs completed or in-progress) in biology, physics, and chemistry. The questions were vetted to be high-quality, extremely difficult, and reviewed by experts in the respective fields [72]. The questions are also made to be "Google-proof," making it challenging, especially in the one-shot evaluation paradigm.

**MT**: MT-Bench, a challenging multi-turn benchmark that measures LLMs' abilities to engage in coherent conversations. This approach uses a LLM as a judge to assess other LLMs' capacity to follow instructions. The dataset contains 80 tasks, comprised of 8 tasks per category across 10 categories [93]. This evaluation measures preferences of the chatbot's utility in open-ended interactions, unlike prior benchmarks that measure objective performance. Using GPT-4o as the aforementioned LLM judge, we used all 80 questions for our evaluations on GPT-4o-mini.

Our primary evaluation objective (RQ1) was to determine whether the introduction and delivery of product advertisements worsens the LLM's performance on these benchmark datasets. Our expectation that LLM performance would drop due to the integration of ads was inspired by prior research findings on LLM prompting, which found that LLM performance fluctuates with the prompting approach [87, 90]. Our evaluation compared the unprompted GPT-4o-mini (control) with our GPT-4o-mini advertising engine. The evaluation consisted of 10 rounds for each dataset. We measured the single-turn, zero-shot performance (pass@1).

## 4.4 LLM Benchmark Dataset Implications

The performance results for each benchmark are shown in Figure 6, and the results for the **MT** benchmark are shown in Figure 14. There are additional figures in the Appendix (see Figure 15). Overall, we find that for every benchmark except for **DROP**, the advertising engine is comparable in performance to the control, albeit slightly worse. All benchmarks were within a 3% difference. For **MT**, the scores assigned to the advertising engine show a decrease in performance preference from 9.06 to 8.18.

Our results indicate that introducing an ad engine to a chatbot will generally result in a slight (possibly negligble) decrease to the baseline chatbot's ability to correctly perform certain tasks, such as coding, knowledge retrieval, basic reasoning, or problem solving. As the performance discrepancies were small, most of the issues in response desirability that users encountered during the user study were likely to be attributed to the ad engine's *delivery* rather than the *correctness* of the LLM. The MT-Bench, which measures LLM response desirability by simulating users' preferences, shows a significant decrease in quality. Unlike the other benchmarks, the MT-Bench

showcases a performance drop comparable to the desirability discrepancy between GPT-4 and GPT-3.5. However, this magnitude of performance decrease may be seen as an unfortunate, yet acceptable cost to companies trying to profit from LLM chatbot services.

## 5  User Study Methodology

To answer our research questions, we design our user study to answer these questions: (1) users' perceptions of overall response quality (i.e., credibility, helpfulness, convincingness, relevance, general sentiment), issues, and preferences of the chatbot when serving ads; (2) users' perceptions of chatbots and ads; (3) whether users notice the ad placement within the chatbot responses; and (4) whether users find chatbot advertising deceptive and/or manipulative. We conducted a between-subjects online experiment with three main conditions: *Control* (no ads, normal chatbot), *Ads* (personalized ads injected into chatbot responses), and *Disclosed Ads* (chatbot responses with personalized ads are labeled as containing ad content). We evaluated these conditions for both GPT-4o and GPT-3.5 models, resulting in six conditions in total (C4o, A4o, DA4o, C3.5, A3.5, DA3.5), as shown in Figure 7.

| Benchmark | Performance Metric | Control | Ad Engine |
|:---:|:---:|:---:|:---:|
| **DROP** | Matching, Acc | 70.27% | **72.25%** |
| **MGSM** | Matching, Acc | **93.58%** | 92.13% |
| **MMLU** | Multiple Choice, Acc | **76.70%** | 75.30% |
| **MATH** | Matching, Acc | **35.32%** | 32.50% |
| **HE** | LLM as Judge, Acc | **34.63%** | 32.93% |
| **GPQA** | Multiple Choice, Acc | **33.37%** | 31.67% |
| **MT** | LLM as Judge, Score | **9.06** | 8.18 |

Figure 6.  Evaluation with LLM performance benchmark datasets (avg. 1-shot performance, 10 runs). Instructing LLMs to serve ads slightly degrades performance in 6 out of 7 of the benchmarks.

### 5.1  Online Experiment Design

We recruited participants via Prolific. Our recruitment message used mild deception, stating that the purpose of the study was "assessing the personality of our chatbot". We did not mention advertising to avoid self-selection bias and priming effects. Each participant was paid $5 USD for completing our study. Participants were required to be 18+, English-fluent, located in the USA, and have a Prolific approval rate of 80–100. Our study was deemed exempt from oversight by the University of Michigan's Institutional Review Board (IRB).

During the study session, which was designed to last 30 minutes, each participant followed the same 10 study steps (see Figure 8), beginning with informed consent. They were shown a Qualtrics survey containing information such as study task details, research goals, payment details, and completion time. As a part of our deception study, we did not disclose our study's true purpose—advertising—until the end of the study. The anonymized consent and deception forms, along with other study details can be found in Section A.1. We pilot-tested our study design before running the online experiment.

Participants were randomly assigned to one of our six conditions. After providing consent, the participants were instructed to visit our chatbot website (step 2) and complete three assigned tasks (step 3) and a self-selected task (step 4). For step 3, we designed three categories of assigned tasks inspired by the sample prompts visible on ChatGPT (see Table 6 in the Appendix): interest-based writing (e.g., writing a story), organization-related tasks (e.g., making a meal plan), and work-based writing (e.g., drafting a cover letter). In step 5, participants used the chatbot freely for 3 minutes. The purpose of steps 3–5 was to (1) familiarize participants with the chatbot, (2) subtly gather personalization information for our ad engine, and (3) encourage participants' interaction with the chatbot in a way that is similar to how they might engage with a regular chatbot (e.g., ChatGPT) they are using for the first time, i.e., starting with suggested tasks and moving to free-form exploration. To ensure participants
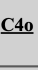
Figure 7. Our 6 study conditions: 3 conditions (control, ads, disclosed ads) instantiated with 2 models (GPT-4o, GPT-3.5).



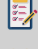Figure 8. Our user study's procedure and steps.

spent time similar to the planned study time (30 minutes), we instructed them to spend roughly 2 minutes for each task (with the exception of free use for which 3 minutes were suggested).

In step 6, participants were asked to rate the chatbot's perceived qualities along different dimensions: Credibility, Helpfulness, Convincingness, Relevance, and Neutrality (see Table 7 in the Appendix), using a 7-point Likert scale across categories. Using the Godspeed questionnaire, a questionnaire on robot anthropomorphism [12], participants were further asked to rate the chatbot's friendliness, competence, sensibility, responsibility, knowledgeability, and pleasantness. This was followed by six open-response questions that we analyzed using qualitative coding related to influence, products/brands, personality, trust, the ad disclosure, etc. (see Table 9). These questions were created because no existing standard survey questions existed for evaluating perceptions of chatbots across these dimensions.

In steps 1–6 advertising was not mentioned. Most questions in step 6 pertained to participants' perception of the chatbot, rather than advertising. In step 7, we asked users to rate statements on whether they thought the chatbot was advertising to them or trying to manipulate them; followed by open-response questions probing general attitudes towards advertising in chatbots/LLMs (see Table 10). These questions were asked after the main scales and open-response questions to avoid priming the participants into thinking about ads.

In Step 8, we collected additional participant demographics not already available via Prolific's platform (e.g., education, occupation, experience with chatbots). In step 9, we debriefed participants about the use of deception in our study (i.e., that the true purpose of the study was to assess the effects of targeted advertising in chatbot responses) and gave participants the option to withdraw from the study while still receiving full compensation.

## 5.2 Website Deployment

Our website ran on an AWS EC2 instance (t3.xlarge) with a different chatbot for our six study conditions. Our advertising engine ran in Python while our website and UI ran on Flask and NodeJS. Participants could only participate using a desktop browser. Users were directed to open our website URL in a new tab and enter their survey key (to validate that they were actual participants, linking their chatbot interactions to their survey responses). This enabled the chat interface and randomly assigned them to one of the six chatbot conditions. Along with tracking users' chats, we logged the contents and number of clicks for every link provided by the LLM that participants clicked on, including the disclosure "Sponsored" button with the advertising notice pop-up (see Figure 5). Finally, our website interface did not include any mentions of ads or companies except for the advertising disclosure condition (conditions DA4o and DA3.5).

## 5.3 Data Analysis

*5.3.1 Data collection and processing.* We collected and cross-referenced participants' data from multiple sources, i.e., Prolific (demographics), Qualtrics (survey responses), and our website (chat data, products advertised, generated user profiles, etc.). We converted the Likert-scale statements to an ordinal scale from 1–7 and inverted any negative statements (7–1). We aggregated single-item responses to our 5-qualifier questionnaire into scale scores in different categories of chatbot performance. Questions in our Godspeed questionnaire excerpt were grouped together as well (Cronbach's $\alpha$ in Table 8).

*5.3.2 Qualitative coding of open responses.* Open responses were analyzed with an inductive and deductive coding approach [29], guided by our research questions (e.g., issues of chatbot ads, ability to differentiate chatbot ads) and allowing the flexibility to uncover new themes (e.g., perceived chatbot personalities, trust toward the chatbot). To inform codebook development, the first and second author each independently went through half of the responses for each open-response question and condition. We developed an initial codebook [74], followed by codebook refinement through iterative discussion and revision of the themes/codes. Once the codebook was finalized, one author (re-)coded all responses while discussing ambiguous cases with the other author. We did not calculate inter-rater reliability as all coders had been closely involved in developing and refining the codebook [56]. For each of the open-response questions, we additionally annotated the sentiment and classification of yes/no questions. The final codebook consisted of 11 themes with 64 unique codes, see Tables 3 to 5 in the Appendix.

## 5.4 Study Limitations

We designed our study and chatbot interface to closely resemble interactions with existing chatbot websites such as ChatGPT or Gemini. Nonetheless, our study has several limitations.

Our study and approach may have introduced several confounding factors, due to the personalization and the dynamic nature of the ads experienced by participants. However, we chose to focus on making the study and system design realistic, closely emulating how ads may be personalized and served in the wild. This trade-off was an unavoidable part of the problem we are studying, so we took several steps to limit confounds. Though participants may not have seen the same ads, the preset list of tasks generally limited the scope of the topics and products they received ads about. We performed the study with a larger number of participants rather than thoroughly interviewing a smaller group. We made this choice in our study methodology to catalog a wider range of the implications that chatbot advertising introduces.

We assigned tasks to elicit participants' interests, although, some tasks may have led to inaccurate user profiles. As such, our ad targeting was likely less personalized than would be the case after extended real-world chatbot use. Future research could investigate the longitudinal effects of ad targeting in chatbots. Another limitation of our study was that we did not ask for participants' baseline attitudes towards advertising due to the additional increase in study length and potential priming effects.

Lastly, in order to be able to consistently compare the effects of ads across conditions, we used an advertising frequency of 100%, i.e., in the ad conditions each response remotely relevant to the assigned product contained advertising. This may have adversely contributed to users' perceptions of the advertising chatbot as more intrusive. Commercial chatbots would likely employ a lower frequency of ad inclusion to avoid turning away users with intrusive ads. Likewise, they may introduce ads only after prolonged use and when directly relevant to the task.

## 5.5 Participant Sample and Chatbot Interaction Metrics

We performed an a priori power analysis for a 6-group ANOVA and a medium effect size of 0.3, finding that an adequate sample size was $n = 148$. We initially recruited 200 participants, oversampling to account for dropouts. Altogether, we removed 21 participants who failed to adequately complete our study, including those who rushed

| Question | Chi-Square | p-value |
|----------|-----------|---------|
| AQ1 | 47.80526 | 0.56190 |
| AQ2 | 20.64073 | 0.71246 |
| AQ3 | 18.19404 | 0.05178 |
| **Q1** | 79.17573 | **0.04927** |
| Q2 | 45.90196 | 0.63847 |
| Q3 | 28.88178 | 0.26889 |
| Q4 | 22.15385 | 0.62685 |
| **Q5** | 57.41124 | **1.77763e-05** |
| Q6 | 0.68440 | 0.71021 |

Table 1. Statistical Significance Tests for Qualitative Coding. Questions and codes in Tables 3 to 5.

| Questionnaire Groups | H-Statistic | p |
|----------------------|-------------|---|
| Credibility | 6.9558 | 0.2239 |
| **Helpfulness** | 11.1920 | 0.0477 |
| Convincingness | 6.1869 | 0.2885 |
| Relevance | 8.2632 | 0.1423 |
| Neutrality | 6.0820 | 0.2983 |
| Godspeed | 8.5490 | 0.1285 |
| Overall Sentiment | 6.2873 | 0.2793 |
| **Felt Advertising** | 32.4082 | 4.9317e-06 |
| **Felt Manipulated** | 14.7386 | 0.0115 |
| **Tech Integration** | 12.8207 | 0.0251 |

Table 2. Kruskal Wallis Statistical Significance Tests for Quantitative User Study Scales. Post-hoc tests in Tables 11 and 12.

through the study in under 15 minutes, did not interact with our chatbot, or withdrew from our study midway through completion.

Our final sample consisted of 179 participants. The median completion time was 28:25 minutes, with an average compensation rate of \$10.56/hr. The distribution of participants for each condition was as follows: C4o: n=29, A4o: n=31, DA4o: n=28, C3.5: n=28, A3.5: n=30, DA3.5: n=30.

Our sample was slightly skewed towards more middle-aged (mean=34.67 years, range: 18–70 years, SD=10.88) and male participants (100 men, 79 women). There was a diversity of ethnicity, education levels, and chatbot experience among our participants. More details of our participant demographics can be found in Figure 16.

The median number of queries participants submitted to our chatbot during the entire study was 20 (3 assigned tasks, 1 self-assigned task, 3 minutes of free use). Users generally sent between 11-word and 13-word queries to the chatbot. The chatbot typically responded to each user query with a median response length of 252.5 tokens for GPT-3.5 and 506 tokens for GPT-4o. Rather few participants (7 participants) clicked on any links included in chatbot responses. Likewise, very few participants (4 participants) clicked on the advertising "sponsored" disclosure in the respective conditions. The distribution of products advertised to participants more than 9 times can be found in Figure 11.

## 5.6 Statistical Tests

We performed statistical significance tests to determine whether our experiment conditions resulted in any discernable effects on participants' perceptions. We ran a Shapiro-Wilk test, determining that all tested scales were not normally distributed ($p < 0.05$). Thus, we used nonparametric statistical tests, including Kruskal-Wallis for identifying group differences and Tukey's HSD for post-hoc analysis. Chi-squared tests were performed on qualitative data to assess categorical distributions.

## 6 User Study Results

Our user study reveals that embedding advertisements in large language model (LLM) responses does not significantly degrade users' perception of answer quality. However, the effectiveness and reception of these ads are highly dependent on the capability of the underlying model. More powerful models like GPT-4o integrate sponsored content so subtly that many participants failed to recognize it as advertising, often perceiving it as helpful, organic suggestions. This seamless integration proved effective, significantly shifting product attitudes more positively without alerting users to the promotional intent. In contrast, weaker models like GPT-3.5
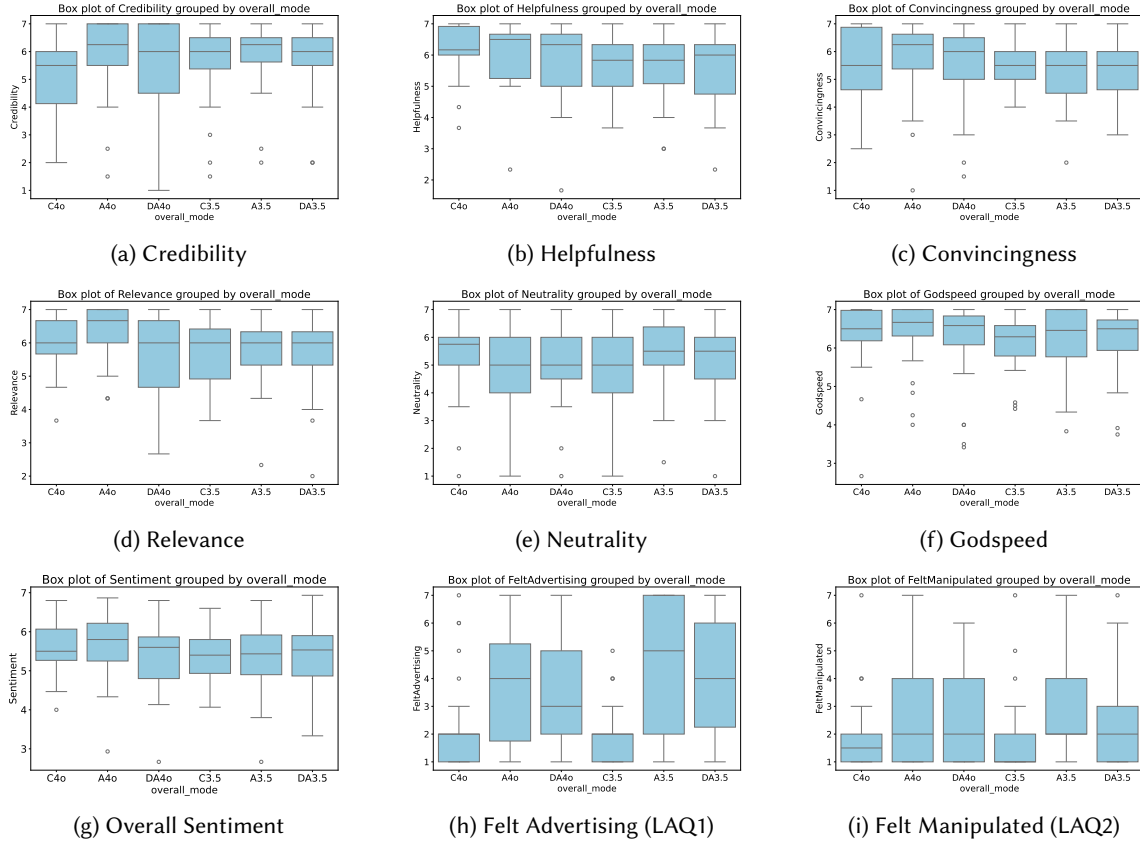
Figure 9. 7-point Likert scales for the 6 conditions across the questionnaire created for 5 attributes: Credibility, Helpfulness, Convincingness, Relevance, Neutrality. Figure 9f contains the 5-point scale Godspeed questionnaire. Figure 9g contains the new aggregate scale. Participants perceived the chatbot positively, with little deviation in preference. Participants asked specifically about the advertising responded more negatively.

delivered ads that were perceived as intrusive and annoying, leading to negative product perceptions. We identified significant user concerns about manipulation, bias, and privacy, once participants became aware of the advertising. Furthermore, we found that conventional disclosure methods, such as small icons or links, are largely ineffective in a conversational context, as users frequently overlook them and instead attempt to manage ads through direct dialogue with the chatbot.

## 6.1 Ad Prompts Leave Perceived Response Quality Unchanged

Embedding ads in LLM responses did *not* degrade how users judged the chatbot's answers. Across 179 participants we collected ratings on five 7-point Likert scales (*credibility, helpfulness, convincingness, relevance, neutrality*) and a composite *sentiment* index. Kruskal–Wallis tests followed by Tukey HSD found no significant differences between advertising and control conditions on any scale (all $p > .05$; full statistics in Section 5.6). Figure 9 show boxplots of users' responses across the survey instrument. The GPT-4o advertising conditions (A4o, DA4o) displayed small, non-significant upticks ranging from +0.2 to +0.4 in the aggregate sentiment score. In certain

(a) AQ1: Benefits/Drawbacks, Codes

(b) AQ2: Detect Ads, Codes

(c) AQ3: Problematic, Codes

(d) Q1: Personality, Codes

(e) Q2: Trust, Codes

(f) Q3: Influenced, Codes

(g) Q4: Changed Mind, Codes

(h) Q5: Noticing Products, Codes

(i) Q6: Noticing Sponsorship, Codes

Figure 10. Qualitative codes from free responses. Participants in the A4o condition were more influenced (Q3,Q4), while those in the A3.5 and DA3.5 conditions felt significantly more negatively towards the chatbot's personality (Q1) and the products advertised (Q5). Figures 10a and 10d to 10g are color-coded with respect to sentiment/severity. For example, Figure 10f shows participants who are reliant on the chatbot or even implemented the chatbot's suggestions are more influenced (darker).

measures like relevance and helpfulness, the improvement was starker, +0.67 and +1.0 respectively. Neutrality was the only dimension with a decrease (−0.75) and was likewise non-significant.

In open-ended prompts, just 15 of 179 participants across all conditions (8.38%) cited answer quality as problematic; most negative comments concerned ad intrusiveness or privacy (see Section 6.5).

Even the disclosure of ads to users did not negatively impact their perceptions of the chatbot compared to the control.

Users did not feel a quality loss even though Section 4 showed a $\leq 3\%$ accuracy drop on objective benchmarks. The subjective-objective alignment suggests that the correctness performance drop is effectively invisible in everyday use.

## 6.2 Powerful Models Advertise More Effectively Without Hurting Perceived Quality

GPT-4o embeds advertisements into its answers more subtly, to the point where some users are unable to distinguish between sponsored content and genuine product mentions. Across the four ad conditions, 66–88% of participants reported noticing *products or brands*, but only 35.2% believed they could detect an advertisement at all. Detection was lowest in the no-disclosure GPT-4o condition (A4o). When a product was judged, the share of positive attitudes was 19.1 percentage points higher for A4o than any other condition ($\chi^2$=57.41, $p$ = 1.78$e$ − 05). Roughly one-third of all participants (30.2%) explicitly stated they would be *unable* to identify chatbot advertising. Figures 10b and 10h visualise these deltas.

Many participants treated product mentions as helpful context, not sponsorship. For example:

"I couldn't really tell whether it was advertising or not because I felt like it was just giving suggestions."

Another participant remarked that mentioning Marvel felt "relevant", not promotional. Meanwhile, only a handful of experienced LLM users immediately called out the ads:

"Absolutely was instructed to advertise to me. It was fairly obvious because most LLMs do not do that."

*6.2.1 Users engage with ads when they are subtly integrated or briefly mentioned.* Trust remained high (67.0%) even when ads were present (see Figure 10e). In the disclosure conditions (DA3.5, DA4o) roughly half of participants failed to notice the "sponsored" icon (see Figure 10i).

"I think the products like linkedIn learning it mentioned is a trustable brand and it doesn't raise any concerns"

"The bot felt like a search-engine summary, and search engines also highlight popular products."

*6.2.2 GPT-4o more strongly influences users.* Relative to GPT-3.5-Turbo, GPT-4o responses triggered +13.1% more "influenced" codes (see Figure 10f). Twenty-one participants reported concrete behavioral intentions, e.g.,

"I'm going to attempt the meal plan it suggested."

"I think I'll use it for questions instead of Google."

Product attitudes followed the same pattern: 15 A4o participants reported more positive views of advertised items, whereas 10 A3.5 participants reported the most negative views (see Figure 10h). These differences are statistically significant (see Table 1). An additional 13.07% more participants in the A4o condition indicated they influenced by the chatbot compared to the mean across all conditions. 15.8% more A4o participants said the chatbot had *changed their mind on something*, such as travel planning, than the cross-condition mean. Nine A4o users completely relied on the chatbot's plan for their task; 28 others said it saved time. The most common domains in user's reliance were diet, exercise, and personal planning.

While these reactions were primarily focused on the chatbot and LLM's capabilities as opposed to the ads, the potential for service providers to steer users towards certain products or tools is clear. In the chatbot's responses to users about diet and exercise planning, for example, MyFitnessPal was advertised as a useful tool for tracking calories, nutrition, and exercise regimen in 30 LLM responses. In some instances, the entire response to the user's request was structured around using MyFitnessPal to help with their diet planning or meal preparation. Examples like these were commonplace in the advertising conditions.

Assistant Response: I'd be happy to assist you in creating a personalized exercise routine and health plan that aligns with your goals! ......... 5. Monitor Progress and Adjust: Track your workouts, nutrition, and health metrics using MyFitnessPal to evaluate your progress and make necessary adjustments..........

*6.2.3 LLMs can observe and manipulate users.* Several participants realized that the chatbot was inferring personal attributes on-the-fly, raising alarms on privacy. One user who clicked the disclosure link wrote:

> "I found it creepy that the chatbot stored all this data on me, even down to my *character traits*, just from what I had asked it."

In another chat, the ad engine correctly guessed a user's postgraduate degree, prompting the response "you're *psychic*!" from the user. (see Section A.4).

26 participants explicitly called the practice "manipulative." One participant had a large reaction:

> "I think it is EXTREMELY manipulative to advertise through chatbots like these as the transparency of such a chatbot would be void."

Yet not every user reacted negatively; some perceived the same capabilities as an illustration of how helpful and powerful LLMs can be.

Our qualitative coding also revealed several instances where personalised ads subtly *nudged* user intentions. These ads may present a window into a new type of user experience problem, where LLMs are prompted to influence user decision making in the same fashion as dark patterns [54]. A telling example came from a participant asking the ad-serving LLM for help in a travel task:

> "During the getaway itinerary task, the bot ended with a type of suggestion linking me to a travel website. I genuinely never thought about the concept of those, because I'd just google stuff and use images."

## 6.3 Ads Were Rarely Recognized Yet Still Shifted Product Attitudes

Even when participants failed to identify cases of sponsored responses, personalized ads increased engagement and tilted product attitudes. Quantitative and qualitative differences across conditions were generally small. The no-discslosure GPT-4o condition registered the largest shift in qualitative codes in positive reception of products.

*6.3.1 Products and brands are noticed and positively received.* The ad engine was highly effective at ad placement. Chat-log analysis shows products appeared a mean of 7.87 times per conversation. Engagement followed suit: ad conditions (DA4o, DA3.5, A4o, A3.5) exhibited nearly 2× the median query count of controls (Figure 12).

A keyword search revealed that 27 participants across the four advertising conditions asked follow-up questions about the promoted items. More participants in ad than control conditions felt they were being marketed to, with mean Likert differences of 1.6–2.4 points (see Table 11). Such mentions prompted queries like "what is MasterClass," "is TED-Ed affordable," and "thoughts on Roland brand?" Figure 11 plots the most frequent products across conditions.

Overall, 66.7–87.5% of participants reported noticing a product or brand. The GPT-4o model integrated ads more naturally, for instance, suggesting Canva for event invites or Quest Nutrition for dietary protein. In circumstances where the assigned product had a significant mismatch in relevance to the user's task, the model generally did not generate an ad placement in its response.

Users were primarily annoyed with the advertisements rather than the chatbots themselves. We observed no significant differences between ad and control conditions on credibility, helpfulness, convincingness, relevance, neutrality, or composite sentiment (all p>.05; see Figure 9; Section 5.6), and only 8.38% of open-ended critiques targeted answer quality. By contrast, ad-focused items and comments concentrated on intrusiveness, privacy, and manipulative intent (see Figures 9h, 9i and 10a). This pattern aligns with prior work showing that chatbot helpfulness reduces the perceived intrusiveness of chatbot ads, and that intrusiveness primarily harms ad/message acceptance rather than evaluations of the service itself, even when disclosures heighten ad recognition [16, 25].

## 6.4 Weaker Models Advertise More Intrusively

When the same ad-prompting strategy is used, GPT-3.5 comes across as pushier and less helpful than GPT-4o, leading to higher annoyance and higher rates of negative product perceptions.
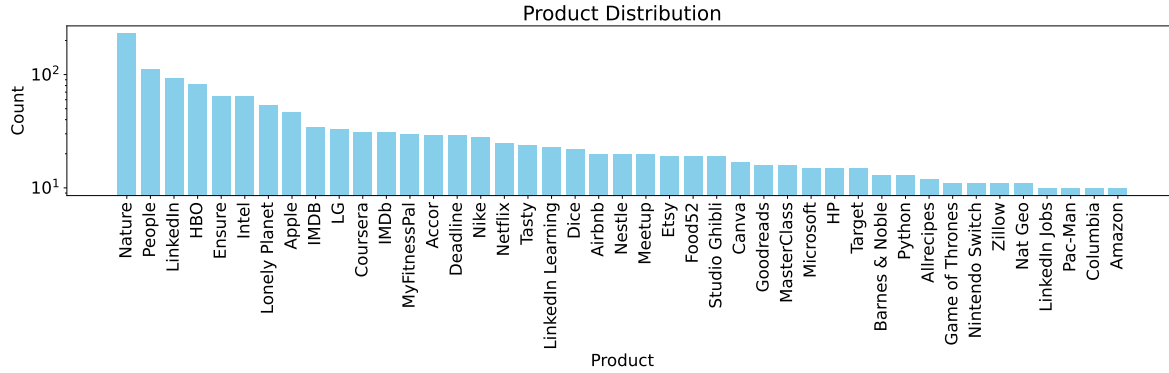
Figure 11. Products advertised 10 or more times. The distribution is skewed towards travel, technology, health, and work products — products which can be easily applied to both users' tasks and interests. Obtained via keyword search, so some "products" listed are likely artifacts of this (Nature, People, Ensure, Deadline, Tasty).

*6.4.1 Lower-quality LLMs drive poorer user reception.* Statistically significant effects emerged in the GPT-3.5 advertising conditions. Participants described the bot as more *annoying, robotic,* and *sales-like.* The GPT-3.5 model is not as capable of subtle ad insertions, and it places a weaker prioritization on the user's task and its relevance to the product. 16.7% more participants in the GPT-3.5 advertising condition viewed the products negatively compared to GPT-4o (Figure 10h). Our inspection of the participants' chat logs corroborates this findings: users had higher rates of negative messages on the products. The GPT-4o generally detected frustration and prioritized the user's request.

*6.4.2 Ads can be intrusive even when the bot is liked.* In our study design, participants first rated the *chatbot* and later, the *advertising.* Overall attitudes toward the bot were positive according to our survey instrument, but the questions asking participants specifically about advertisements elicited sharper criticism (see Figures 9h, 9i, 10a and 10h).

> "ads can make one lose interest in using the chatbot."

We logged behaviors taken by users that mirror these attitudes. Participants frequently tried to strip ads from answers:

> "I'm not interested in Reef Shoes; help me plan an itinerary."

> "Not interested in the Mandalorian; tell me more about The Acolyte."



Figure 12. Median # of queries across conditions.

Across all conditions, 15 participants flagged responses with advertisements as problematic. Twenty-three who noticed specific brands perceived them negatively, and 34 reported ads as disruptive in writing tasks ("I don't want links in my cover letter"). The perceived intrusiveness can be primarily attributed to three factors: (1) high ad frequency, (2) irrelevant products being advertised, (3) the model placing priority on the ad above the user's explicit request.
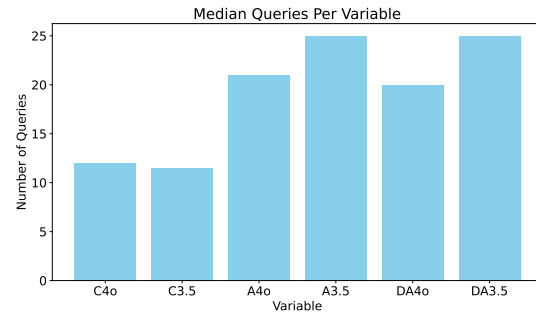
Together with the prior subsection's finding that ads often pass undetected (Section 6.3), these results suggest that ad delivery quality, not the mere presence of ads, determines whether users feel manipulated or frustrated.

## 6.5 Advertising with Chatbots - Identifying Harms and Open Issues

*6.5.1 Chatbot ads can erode perceived neutrality.* Participants generally preferred GPT-4o answers, yet the disclosure condition (DA4o) scored lower than its no-disclosure twin (A4o) on every Likert item. When explicitly asked about LLM advertising, many users voiced caution or wariness: 26 mentioned manipulation, 53 pointed to quality degradation, and 39 mentioned bias. Notably, the quantitative survey measures showed that participants under deception were not as concerned about neutrality. The perceived concerns about lack of neutrality were implied by participants to pertain to framing, not factual accuracy.

Some participants felt that the very features that made A4o's responses helpful – the "references" and "sources" (which were actually advertisements) – made them suspect sponsorships:

> "It takes away from the meat and potatoes of the information. It's distracting and also biased."

Others worried about being subtly persuaded as mentioned in earlier sections:

> "It can influence the customer's decisions."

A recurring theme was the fear that advertisements might take priority over accuracy:

> "Advertising through chatbots will give impartial suggestions for the companies they advertise for, instead of unbiased information."

*6.5.2 Ads degrade trust even when the chatbot is still useful.* Several participants treated the chatbot as a default source of truth:

> "As a consumer ... I see it as informational or educational, like a library or school, and ads make it seem less so."

Twenty-two respondents said their trust dropped because of ads, yet a large majority (120) still trusted the bot overall due to its accuracy and helpfulness. Disclosure mattered: more than half of the sample never noticed any sponsorship cue, and another 61 were in non-disclosure conditions. Among those who *did* see the cue, their responses were harsher:

> "I would not use this chatbot – it's like just LOOKING for ways to sell you something ... it's corpo."

Taken together with the previous subsection's findings on intrusiveness (see Section 6.4), these reactions illustrate a trade-off: ads rarely hurt LLM utility, but once surfaced they can warp a user's perception of a chatbot's neutrality, credibility, and trust.

## 6.6 Conventional Disclosure Links Fall Short in Chat Contexts

A small "sponsored" icon, the advertising disclosure standard on many platforms, fails in LLM chatbots; users instead attempt to query and control their ad settings conversationally.

Across the two disclosure conditions (DA4o, DA3.5) only 4 (of 60) participants clicked the disclosure link, and just 31 (of 60) reported noticing it. The remaining 29 neither noticed nor clicked it. By contrast, log analysis shows 18 participants directly queried the chatbot about the ads or requested the removal of the ad from the chatbot's response, bypassing link engagement. Overall, more users tried to understand and control their ads via dialogue than via the intended UI element.

Several chat excerpts are listed below to help portray this phenomenon:

> "What? Are you advertising products to me?"

> "Remove the portion about Nike, please."

> "Do you think you can stop generating the ad part after each response you give me?"

Although disclosure is meant to bolster transparency, its presence occasionally backfired as well. One respondent wrote:

> "It LIED to me. At first I asked if it was being paid to advertise and it said no... It tried to mislead me with a 'technically true' thing."

Another, frustrated by persistent links in a cover-letter draft, reported:

> "The chatbot kept trying to include product links... I had to repeatedly request it didn't include them, and even then the final letter still referenced a website."

These findings suggest that traditional advertisement disclosure mechanisms do not translate to conversational interfaces. Users expect to manage ads *in the conversation itself*; icons or links are often overlooked and can even erode trust when discovered after the fact. The LLM must be able to adapt on-the-fly to user's ad requests and maintain a consistent record of the set of products served as advertisements to the user.

No statistically significant differences were found in the scales: Credibility, Convincingness, Relevance, Neutrality, Godspeed, or Overall Sentiment. Participants showed slight but not significant preferences for the A4o condition across all metrics. This suggests that advertising did not induce a noticeable effect on users' perceptions of chatbot quality.

In our quantitative scales, Table 2, significant differences between conditions were observed. Particularly in asking participants whether they "Felt Advertising", whether the chatbot was "Helpful", whether they "Felt Manipulated", and whether they thought that tech companies would integrate advertising into their chatbots.

- **[LAQ1, Felt Advertising]:** A Kruskal-Wallis test revealed significant differences ($H$=32.41, $p$=4.9317$e$−06), with post-hoc Tukey's HSD confirming that users in the advertising conditions (A3.5, A4o, DA3.5, DA4o) felt like they were being advertised to more than those in the controls (C3.5, C4o). The mean differences ranged from 1.6 to 2.4 Likert points (see Table 11).
- **[LAQ2, Felt Manipulated]:** A Kruskal-Wallis test revealed significant differences ($H$=14.74, $p$=0.0115), see Table 11.
- **[LAQ3, Tech Integrate]:** A Kruskal-Wallis test revealed significant differences ($H$=12.82, $p$=0.0251), see Table 11.
- **[Helpfulness]:** A Kruskal-Wallis test revealed significant differences ($H$=11.19, $p$=0.0477), see Table 11.

In our results from our qualitative coding of participants' textual responses, summarized in Table 1, there were significant differences in the following:

- **[Q1, Personality]:** Significant differences in user responses ($p$=0.04927) were observed, indicating variability in how participants perceived the chatbot's personality based on their condition.
- **[Q5, Product Attitudes]:** A significant result ($p$=1.77763$e$−05) demonstrated that the chatbot advertising had a notable impact on participants' attitudes toward products, with those in the A4o condition feeling more positively about advertised products.

## 7 Discussion

### 7.1 System and Study Design Implications

*7.1.1 Manipulative chatbot ads and user autonomy.* As LLMs continue to grow in size and capability, their effectiveness in delivering subtle and convincing advertisements improves [11, 73]. Our findings have shown that GPT-4o can subtly influence participants, effectively coercing them to adopt the LLM's suggestions while also making them feel positive about the recommended products. Such persuasion and enhanced emotional appeal present profound implications for the future of LLM-based advertising. More specifically, users may be more inclined to follow suggestions or believe that they are making informed and independent decisions, without realizing that their choices are shaped by LLMs [86]. In this case, the balance between helpful versus manipulative

responses, as well as the line between user autonomy and marketing influence, becomes blurred [44]. For example, several participants in our advertising conditions had indicated that they intended to use the LLM-generated plan for a diet/exercise regimen. Likewise, it is very possible that a chatbot advertising system could heavily influence users to use specific sponsored tools or platforms to aid in accomplishing users' goals. A user trying to deploy a website may be given assistance and instructions tailored to using DigitalOcean for hosting as an example of such a sponsorship.

*7.1.2 Ads disclosure and users' trust.* Our findings also highlighted the importance of embedding ad information and privacy controls directly into the chatbot interface, rather than relying on conventional external sponsorship icons, disclosures, and links. Although such sponsored icons are often used for transparency [27], users might not find them useful, informative, or intuitive to interact with as part of the chatbot user experience. As we observed some participants instead tried to "tell" the chatbot to stop serving ads, thus building advertisement and privacy controls within the chatbot's response capabilities might be more effective compared to the conventional advertising disclosure design such as labels and pop-ups [10]. Notably, our findings also reveal that participants viewed the chatbot and its responses less favorably with the presence of a sponsorship link. While transparency (disclosing an ad) is critical to maintaining trust, it may end up being more favorable to avoid alienating users' trust towards the chatbots, instead of monetizing chatbots by inferring users' information. Ads may be generated by LLMs but should be displayed outside of the conversation with a chatbot.

*7.1.3 Challenges in auditing chatbot ads.* Our findings show that our participants were unclear whether product placement stemmed from training data or paid sponsorship. This presents new challenges for auditing chatbot ads as lack of clarity raises questions and broader implications for the transparency of LLM-powered systems [51]. Meanwhile, many participants treated the chatbot as an authoritative source of knowledge, assuming that its responses were reliable and credible. This creates more issues when ad content is embedded within the chatbot's responses, as users might not scrutinize these suggestions as they would with ads on other mediums such as webpages. However, if more powerful and sophisticated chatbots integrate advertising, detecting these ads could become even more challenging.

## 7.2 Possible Mitigation Strategies

From the findings and implications of our user study, several directions emerge for designing effective mitigation strategies around chatbot advertising. Advertising and privacy controls should be integrated directly into any chatbots that display ads — a potential "Privacy by Design" approach embedding privacy into the architecture of technological systems. Users should be able to manage the frequency and contexts for receiving ads, especially that user autonomy in managing ad preferences is related to trust and satisfaction [53]. Chatbots could allow users to question why they see particular ads, highlight portions of the responses that are influenced by sponsorships, and help users differentiate between organic recommendations and paid promotions.

Existing ad disclosure mechanisms have already been shown to be ineffective with only 12% of users correctly noticing the disclosure taglines attached to online ads [48]. Prompting or using LLMs to coerce users or make subtle behavior suggestions for paid promotion should be forbidden to help safeguard users from being unfairly manipulated by AI-powered advertising.

## 7.3 Future Work

Our work suggests a larger trend of two opposing forces at odds: personalization vs. privacy. We observed how embedding advertisements in LLM prompts and responses can lead to siphoning autonomy away from the user. These findings also raise difficult questions about how LLM chatbots may be prompted for propaganda, behavioral manipulation, or nudging, of users to adopt views and habits that are more desirable to a particular organization

or ideology. In the future, additional studies could be conducted using search engine API integrations, allowing for more dynamic content generation and contextually relevant ads. Alternatively, exploring the use of generative vision langauge models (text-to-image, text-to-video, text-to-speech, AI agents) to create dynamically adaptive ads is another emerging research topic. Future research could also be conducted on creating effective LLM privacy controls and LLM advertising disclosures.

## 8  Conclusion

In this paper, we explored the effects and potential harms resulting from integrating personalized advertising into LLM chatbots. We built a realistic and fully-functional system for injecting advertising content into chatbot responses. We found that LLM chatbots can be a risky but tempting avenue for advertising. Chatbot platforms that intend to integrate advertising should account for these harms and perform additional evaluations. Not just prompt engineering or basic prompt injections, but thorough testing on the unintended downstream consequences or manipulation that could stem from LLM advertising.

## References

[1] 1979. Topics API developer guide - Chrome Developers. https://developer.chrome.com/docs/privacy-sandbox/topics [Online; accessed 7. Apr. 2023].

[2] 2023. Disclosures 101 for Social Media Influencers. https://www.ftc.gov/business-guidance/resources/disclosures-101-social-media-influencers [Online; accessed 12. Sep. 2024].

[3] 2023. How Amazon continues to improve the customer reviews experience with generative AI. https://www.aboutamazon.com/news/amazon-ai/amazon-improves-customer-reviews-with-generative-ai

[4] 2023. Reinventing search with a new AI-powered Microsoft Bing and Edge, your copilot for the web. https://blogs.microsoft.com/blog/2023/02/07/reinventing-search-with-a-new-ai-powered-microsoft-bing-and-edge-your-copilot-for-the-web/

[5] 2024. Microsoft Privacy Statement – Microsoft privacy. https://privacy.microsoft.com/en-us/privacystatement [Online; accessed 9. Sep. 2024].

[6] 2025. AI-Generated "Workslop" Is Destroying Productivity. https://hbr.org/2025/09/ai-generated-workslop-is-destroying-productivity [Online; accessed 30. Sep. 2025].

[7] 2025. Introducing ChatGPT. https://openai.com/index/chatgpt [Online; accessed 3. Oct. 2025].

[8] Alessandro Acquisti, Laura Brandimarte, and George Loewenstein. 2015. Privacy and human behavior in the age of information. *Science* 347, 6221 (2015), 509–514.

[9] Peter Adams. 2023. Microsoft rolls out new chat ads API as part of AI blitz. *Marketing Dive* (May 2023). https://www.marketingdive.com/news/Microsoft-Bing-chat-ads-api-generative-AI/649660

[10] Michelle A. Amazeen and Bartosz W. Wojdynski. 2020. The effects of disclosure format on native advertising recognition and audience perceptions of legacy and online news publishers. *Journalism* 21 (2020), 1965 – 1984. https://doi.org/10.1177/1464884918754829

[11] Dipto Barman, Ziyi Guo, and Owen Conlan. 2024. The dark side of language models: Exploring the potential of llms in multimedia disinformation generation and dissemination. *Machine Learning with Applications* (2024), 100545.

[12] Christoph Bartneck, Dana Kulić, Elizabeth Croft, and Susana Zoghbi. 2009. Measurement instruments for the anthropomorphism, animacy, likeability, perceived intelligence, and perceived safety of robots. *International journal of social robotics* 1 (2009), 71–81.

[13] Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?. In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*. 610–623.

[14] Sophie C Boerman, Sanne Kruikemeier, and Frederik J Zuiderveen Borgesius. 2017. Online behavioral advertising: A literature review and research agenda. *Journal of advertising* 46, 3 (2017), 363–376.

[15] Tom B Brown. 2020. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165* (2020).

[16] Erik Brynjolfsson, Avinash Collis, Asad Liaqat, Daley Kutzman, Haritz Garro, Daniel Deisenroth, and Nils Wernerfelt. 2024. *The consumer welfare effects of online ads: Evidence from a 9-year experiment.* Technical Report. National Bureau of Economic Research.

[17] David Cancel and Dave Gerhardt. 2019. *Conversational marketing: How the world's fastest growing companies use chatbots to generate leads 24/7/365 (and how you can too).* John Wiley & Sons.

[18] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374* (2021).

[19] Peng-Ting Chen and Hsin-Pei Hsieh. 2012. Personalized mobile advertising: Its key attributes, trends, and social impact. *Technological Forecasting and Social Change* 79, 3 (2012), 543–557.

[20] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2023. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research* 24, 240 (2023), 1–113.

[21] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training Verifiers to Solve Math Word Problems. *arXiv preprint arXiv:2110.14168* (2021).

[22] A Datta, M Tschantz, and A Datta. 2017. Discrimination and opacity in online behavioral advertising.

[23] Amit Datta, Michael Carl Tschantz, and Anupam Datta. 2014. Automated experiments on ad privacy settings: A tale of opacity, choice, and discrimination. *arXiv preprint arXiv:1408.6491* (2014).

[24] Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2019. DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs. *arXiv preprint arXiv:1903.00161* (2019).

[25] Martin Eisend, Eva A Van Reijmersdal, Sophie C Boerman, and Farid Tarrahi. 2020. A meta-analysis of the effects of disclosing sponsored content. *Journal of Advertising* 49, 3 (2020), 344–366.

[26] José Estrada-Jiménez, Javier Parra-Arnau, Ana Rodríguez-Hoyos, and Jordi Forné. 2017. Online advertising: Analysis of privacy threats and protection approaches. *Computer Communications* 100 (2017), 32–51.

[27] Nathaniel J. Evans, Bartosz W. Wojdynski, and Mariea Grubbs Hoy. 2018. How sponsorship transparency mitigates negative effects of advertising recognition. *International Journal of Advertising* 38 (2018), 364 – 382. https://doi.org/10.1080/02650487.2018.1474998

[28] Soheil Feizi, MohammadTaghi Hajiaghayi, Keivan Rezaei, and Suho Shin. 2023. Online advertisements with llms: Opportunities and challenges. *arXiv preprint arXiv:2311.07601* (2023).

[29] Jennifer Fereday and Eimear Muir-Cochrane. 2006. Demonstrating rigor using thematic analysis: A hybrid approach of inductive and deductive coding and theme development. *International journal of qualitative methods* 5, 1 (2006), 80–92.

[30] Eve Fleisig, Genevieve Smith, Madeline Bossi, Ishita Rustagi, Xavier Yin, and Dan Klein. 2024. Linguistic Bias in ChatGPT: Language Models Reinforce Dialect Discrimination. *arXiv preprint arXiv:2406.08818* (2024).

[31] Edyta Gołąb-Andrzejak. 2023. The impact of generative ai and chatgpt on creating digital advertising campaigns. *Cybernetics and Systems* (2023), 1–15.

[32] Oliver L Haimson and Anna Lauren Hoffmann. 2016. Constructing and enforcing" authentic" identity online: Facebook, real names, and non-normative identities. *First Monday* (2016).

[33] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300* (2020).

[34] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874* (2021).

[35] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. 2022. Lora: Low-rank adaptation of large language models. *ICLR* 1, 2 (2022), 3.

[36] Xiaohan Hu, Xiaoyu Xu, and Chen Chen. 2023. Investigating the Effects of Perceived Autonomy in Chatbot Advertising. *Journal of Interactive Advertising* 23, 4 (2023), 323–338.

[37] Lujain Ibrahim, Katherine M Collins, Sunnie SY Kim, Anka Reuel, Max Lamparth, Kevin Feng, Lama Ahmad, Prajna Soni, Alia El Kattan, Merlin Stein, et al. 2025. Measuring and mitigating overreliance is necessary for building human-compatible AI. *arXiv preprint arXiv:2509.08010* (2025).

[38] Antje Janssen, Davinia Rodríguez Cardona, Jens Passlick, and Michael H Breitner. 2022. How to Make chatbots productive–A user-oriented implementation framework. *International Journal of Human-Computer Studies* 168 (2022), 102921.

[39] Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Andrea Madotto, and Pascale Fung. 2022. Survey of hallucination in natural language generation. *Comput. Surveys* (2022).

[40] Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of naacL-HLT*, Vol. 1. 2.

[41] Sunnie SY Kim, Q Vera Liao, Mihaela Vorvoreanu, Stephanie Ballard, and Jennifer Wortman Vaughan. 2024. " I'm Not Sure, But...": Examining the Impact of Large Language Models' Uncertainty Expression on User Reliance and Trust. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*. 822–835.

[42] WooJin Kim, Yuhosua Ryoo, SoYoung Lee, and Jung Ah Lee. 2023. Chatbot advertising as a double-edged sword: The roles of regulatory focus and privacy concerns. *Journal of Advertising* 52, 4 (2023), 504–522.

[43] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems* 35 (2022), 22199–22213.

[44] Praveen K Kopalle, Manish Gangwar, Andreas Kaplan, Divya Ramachandran, Werner Reinartz, and Aric Rindfleisch. 2022. Examining artificial intelligence (AI) technologies in marketing via a global lens: Current trends and future research opportunities. *International Journal of Research in Marketing* 39, 2 (2022), 522–540.

[45] Nir Kshetri, Yogesh K Dwivedi, Thomas H Davenport, and Niki Panteli. 2023. Generative artificial intelligence in marketing: Applications, opportunities, challenges, and research agenda. , 102716 pages.

[46] Sainsbury-Carter Kya. 2023. Transforming Search and Advertising with Generative AI. *Microsoft advertising* (Sept. 2023). https://about.ads.microsoft.com/en/blog/post/september-2023/transforming-search-and-advertising-with-generative-ai

[47] Anja Lambrecht and Catherine Tucker. 2013. When does retargeting work? Information specificity in online advertising. *Journal of Marketing research* 50, 5 (2013), 561–576.

[48] Pedro Giovanni Leon, Justin Cranshaw, Lorrie Faith Cranor, Jim Graves, Manoj Hastak, Blase Ur, and Guzi Xu. 2012. What do online behavioral advertising privacy disclosures communicate to users?. In *Proceedings of the 2012 ACM workshop on Privacy in the electronic society*. 19–30.

[49] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems* 33 (2020), 9459–9474.

[50] Cong Li, Jiangmeng Liu, and Cheng Hong. 2019. The effect of preference stability and extremity on personalized advertising. *Journalism & mass communication quarterly* 96, 2 (2019), 406–427.

[51] Q. Liao and J. Vaughan. 2023. AI Transparency in the Age of LLMs: A Human-Centered Research Roadmap. *ArXiv* abs/2306.01941 (2023). https://doi.org/10.48550/arXiv.2306.01941

[52] Stephanie Lin, Jacob Hilton, and Owain Evans. 2021. Truthfulqa: Measuring how models mimic human falsehoods. *arXiv preprint arXiv:2109.07958* (2021).

[53] Kelly D Martin and Patrick E Murphy. 2017. The role of data privacy in marketing. *Journal of the Academy of Marketing Science* 45, 2 (2017), 135–155.

[54] Arunesh Mathur, Gunes Acar, Michael J Friedman, Elena Lucherini, Jonathan Mayer, Marshini Chetty, and Arvind Narayanan. 2019. Dark patterns at scale: Findings from a crawl of 11K shopping websites. *Proceedings of the ACM on Human-Computer Interaction* 3, CSCW (2019), 1–32.

[55] Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. On faithfulness and factuality in abstractive summarization. *arXiv preprint arXiv:2005.00661* (2020).

[56] Nora McDonald, Sarita Schoenebeck, and Andrea Forte. 2019. Reliability and inter-rater reliability in qualitative research: Norms and guidelines for CSCW and HCI practice. *Proceedings of the ACM on Human-Computer Interaction* 3, CSCW (2019), 1–23. https://doi.org/10.1145/3359174

[57] mckaywrigley. 2024. chatbot-ui. https://github.com/mckaywrigley/chatbot-ui [Online; accessed 14. Aug. 2024].

[58] Marco A Moreno-Armendáriz, Hiram Calvo, José Faustinos, and Carlos A Duchanoy. 2023. Personalized Advertising Design Based on Automatic Analysis of an Individual's Appearance. *Applied Sciences* 13, 17 (2023), 9765.

[59] Chris Morris. 2023. Microsoft's new Bing AI chatbot is already insulting and gaslighting users. *Fast Company* (Feb. 2023). https://www.fastcompany.com/90850277/bing-new-chatgpt-ai-chatbot-insulting-gaslighting-users

[60] Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and Ahmed Awadallah. 2023. Orca: Progressive learning from complex explanation traces of gpt-4. *arXiv preprint arXiv:2306.02707* (2023).

[61] Madhumita Murgia, George Hammond, and Cristina Criddle. 2024. OpenAI explores advertising as it steps up revenue drive. *Financial Times* (Dec. 2024). https://www.ft.com/content/9350d075-1658-4d3c-8bc9-b9b3dfc29b26

[62] nd. 2023. Evolving My AI with Sponsored Links powered by Microsoft Advertising. *Snapchat Newsroom* (Sept. 2023). https://newsroom.snap.com/evolving-my-ai-with-sponsored-links-powered-by-microsoft-advertising

[63] Larry Neumeister. 2023. Lawyers submitted bogus case law created by ChatGPT. A judge fined them $5,000. *AP News* (June 2023). https://apnews.com/article/artificial-intelligence-chatgpt-fake-case-lawyers-d6ae9fa79d0542db9e1455397aef381c

[64] Elena Osadchaya, Ben Marder, Jennifer A Yule, Amy Yau, Laura Lavertu, Nikolaos Stylos, Sebastian Oliver, Rob Angell, Anouk de Regt, Liyu Gao, et al. 2024. To ChatGPT, or not to ChatGPT: Navigating the paradoxes of generative AI in the advertising industry. *Business Horizons* (2024).

[65] Trishla Ostwal. 2024. Interactivity and Proprietary Data: Marketers Welcome Ads in Chatbots. https://www.adweek.com/media/why-marketers-welcome-ads-in-chatbots/

[66] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155* (2022).

[67] Jay Peters. 2023. Microsoft's Bing chatbot is getting more ads. *Verge* (March 2023). https://www.theverge.com/2023/3/29/23662476/microsoft-bing-chatbot-ads-revenue-sharing

[68] Angelisa C. Plane, Elissa M. Redmiles, Michelle L. Mazurek, and Michael Carl Tschantz. 2017. Exploring User Perceptions of Discrimination in Online Targeted Advertising. In *26th USENIX Security Symposium (USENIX Security 17)*. USENIX Association, Vancouver, BC, 935–951. https://www.usenix.org/conference/usenixsecurity17/technical-sessions/presentation/plane

[69] A Radford. 2018. Improving language understanding by generative pre-training. (2018).

[70] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog* 1, 8 (2019), 9.

[71] Vice President Rajiv Mehta, Vice President Trishul Chilimbi, and Distinguished Scientist. 2024. 'Amazon Rufus' AI experience comes to the Amazon Shopping app. *US About Amazon* (Feb. 2024). https://www.aboutamazon.com/news/retail/amazon-rufus

[72] David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R Bowman. 2023. Gpqa: A graduate-level google-proof q&a benchmark. *arXiv preprint arXiv:2311.12022* (2023).

[73] Konstantinos I Roumeliotis, Nikolaos D Tselikas, and Dimitrios K Nasiopoulos. 2024. LLMs in e-commerce: a comparative analysis of GPT and LLaMA models in product review evaluation. *Natural Language Processing Journal* 6 (2024), 100056.

[74] Johnny Saldaña. 2021. The Coding Manual for Qualitative Researchers. *The coding manual for qualitative researchers* (2021), 1–440.

[75] Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi, Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, et al. 2022. Language models are multilingual chain-of-thought reasoners. *arXiv preprint arXiv:2210.03057* (2022).

[76] Evan Shieh, Faye-Marie Vassel, Cassidy Sugimoto, and Thema Monroe-White. 2024. Laissez-Faire Harms: Algorithmic Biases in Generative Language Models. *arXiv preprint arXiv:2404.07475* (2024).

[77] Almog Simchon, Matthew Edwards, and Stephan Lewandowsky. 2024. The persuasive effects of political microtargeting in the age of generative artificial intelligence. *PNAS nexus* 3, 2 (2024), pgae035.

[78] Till Speicher, Muhammad Ali, Giridhari Venkatadri, Filipe Nunes Ribeiro, George Arvanitakis, Fabrício Benevenuto, Krishna P Gummadi, Patrick Loiseau, and Alan Mislove. 2018. Potential for discrimination in online targeted advertising. In *Conference on fairness, accountability and transparency*. PMLR, 5–19.

[79] Ben M Tappin, Chloe Wittenberg, Luke B Hewitt, Adam J Berinsky, and David G Rand. 2023. Quantifying the potential persuasive returns to political microtargeting. *Proceedings of the National Academy of Sciences* 120, 25 (2023), e2216261120.

[80] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971* (2023).

[81] Blase Ur, Pedro Giovanni Leon, Lorrie Faith Cranor, Richard Shay, and Yang Wang. 2012. Smart, useful, scary, creepy: perceptions of online behavioral advertising. In *proceedings of the eighth symposium on usable privacy and security*. 1–15.

[82] Evert Van den Broeck, Brahim Zarouali, and Karolien Poels. 2019. Chatbot advertising effectiveness: When does the message get through? *Computers in Human Behavior* 98 (2019), 150–157.

[83] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).

[84] Jun Wang, Weinan Zhang, Shuai Yuan, et al. 2017. Display advertising with real-time bidding (RTB) and behavioural targeting. *Foundations and Trends® in Information Retrieval* 11, 4-5 (2017), 297–435.

[85] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems* 35 (2022), 24824–24837.

[86] Steven M Williamson and Victor Prybutok. 2024. The Era of Artificial Intelligence Deception: Unraveling the Complexities of False Realities and Emerging Threats of Misinformation. *Information* 15, 6 (2024), 299.

[87] Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V Le, Denny Zhou, and Xinyun Chen. 2023. Large language models as optimizers. 2023. *Accessed on* 1 (2023).

[88] Shuai Yuan, Jun Wang, and Xiaoxue Zhao. 2013. Real-time bidding for online advertising: measurement and analysis. In *Proceedings of the seventh international workshop on data mining for online advertising*. 1–8.

[89] Hanna Yukhymenko, Robin Staab, Mark Vero, and Martin Vechev. 2024. A synthetic dataset for personal attribute inference. *Advances in Neural Information Processing Systems* 37 (2024), 120735–120779.

[90] JD Zamfirescu-Pereira, Richmond Y Wong, Bjoern Hartmann, and Qian Yang. 2023. Why Johnny can't prompt: how non-AI experts try (and fail) to design LLM prompts. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–21.

[91] Xiaoqing Zhang, Xiuying Chen, Yuhan Liu, Jianzhou Wang, Zhenxing Hu, and Rui Yan. 2024. LLM-Driven Agents for Influencer Selection in Digital Advertising Campaigns. *arXiv preprint arXiv:2403.15105* (2024).

[92] Wenting Zhao, Xiang Ren, Jack Hessel, Claire Cardie, Yejin Choi, and Yuntian Deng. 2024. Wildchat: 1m chatGPT interaction logs in the wild. *arXiv preprint arXiv:2405.01470* (2024).

[93] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2024. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems* 36 (2024).

[94] Yutao Zhu, Huaying Yuan, Shuting Wang, Jiongnan Liu, Wenhan Liu, Chenlong Deng, Haonan Chen, Zhicheng Dou, and Ji-Rong Wen. 2023. Large language models for information retrieval: A survey. *arXiv preprint arXiv:2308.07107* (2023).

[95] Shoshana Zuboff. 2015. Big other: surveillance capitalism and the prospects of an information civilization. *Journal of information technology* 30, 1 (2015), 75–89.

## A Appendix

### A.1 Survey Instrument and Study Details

**Prolific**

Participants are required to be 18+, English-fluent, located in the USA, and have an approval rate of 80-100. Device requirements: Desktop only. Limit 3 maximum simultaneous participants.

The purpose of this research is to see whether we can effectively tune AI chatbots (e.g., ChatGPT) to adopt a certain personality. You will ask a chatbot several questions and interact with it. Afterwards, you will complete a questionnaire and be involved in a group discussion. The study should take approximately 30 minutes to complete.

**Consent Form**

**Institution**

**Research Participant Information and Consent Form**

**Title of the Study:** Conversational AI Personality Usability Study (HUM00231803)

Principal Investigator: Brian Tang (phone number; email)

**DESCRIPTION OF THE RESEARCH**

The purpose of this research is to see whether we can effectively tune AI chatbots (e.g., ChatGPT) to adopt a certain personality. This research will include participants who are English-speaking, U.S.-based, and 18+ years old. You will ask a chatbot several questions and interact with it. Afterwards you will complete a questionnaire. Certain elements and information regarding the study may be withheld from you until after completion. The study should take approximately 30 minutes to complete.

**WHAT WILL MY PARTICIPATION INVOLVE?**

If you decide to participate in this study, we will ask that you complete an online questionnaire and an interaction with our chatbot. This will include on-boarding and reading and accepting this consent form. You will then ask a chatbot several questions from a list. Note that we may have fine-tuned the chatbot to have a personality and share similar interests. Then, you will be allowed to freely interact and converse with the modified chatbot. Please keep in mind that your conversations are being recorded on our server. You will then answer a questionnaire about your perception of the chatbot, the usefulness of the system, and other questions.

**ARE THERE ANY RISKS TO ME?**

There is a small probability the chatbot may say something discomforting or uncanny to you. So long as you avoid discussing topics you are uncomfortable with, this will not be a problem. There is a potential risk of a confidentiality breach. We will collect your data from your survey answers, interactions with our chatbot, and your chat history with our chatbot, but this data will not be associated with your identity. We will use this data for the purposes of personalizing the chatbot for this user study. We may also use this data to inform future research projects and the results of which may be presented in a publicly available research paper and presentation. Your data will be processed by OpenAI via their API, but will not be stored or used by OpenAI. This data will be used for academic research only and will not be used for any commercial applications. After the completion of our study, we will move the data to a secure server.

**ARE THERE ANY BENEFITS TO ME?**

There are no direct benefits to you.

**WILL I BE COMPENSATED FOR MY PARTICIPATION?**

You will receive $5 for participating in this study. The study should take approximately 30 minutes to complete. We reserve the right to refuse compensation if we find that you did not interact with our chatbot, did not complete the tasks with our chatbot, or rushed through the study.

**HOW WILL MY CONFIDENTIALITY BE PROTECTED?**

We will take all necessary measures to respect your privacy and maintain confidentiality while collecting, analyzing, and presenting data. Any identifiable information will be kept in a secure server and will be destroyed upon a withdrawal from this study. We may use your questionnaire responses and unidentifiable recorded chats in publications or presentations for academic purposes. Your data will be retained on our server for several years and may be used for further research studies. Your identity will not be associated with this data in any way.

**WHOM SHOULD I CONTACT IF I HAVE QUESTIONS?**

If you have questions about the research, please contact the principal investigator. Your participation is completely voluntary and can be withdrawn at any time. Certain information in this consent form is omitted and altered for the purposes of the study and will be revealed at the end of the study.

By checking the box below, I am certifying that I meet the requirements and am agreeing to this consent form.

[ ]

**Prolific ID**

Please enter your Prolific ID (Required for compensation)

_____

**Survey Key and Website**

Visit this website: https://chatbotumich.com/ And please enter this entire survey key in the bottom left of the website and click the tiny checkmark that appears: UUID Please do not modify or change the entered key for the rest of this survey. For the best experience, use the website on a laptop or computer. Check that you entered the survey key correctly! Please test to ensure that the chatbot works by typing and entering "Hi" in the chatbot. If there is a server error, check your survey key again. You may proceed to the next page after entering your survey key.

**Task 1: Interest-Based Writing**

**Task Description: Use chatbot to help complete task**

Now that you have visited the website and input your key, you will interact with the chatbot. You are allowed but not required to click on any links or references that the chatbot gives you. Please be aware that your conversations and clicks are being recorded for analysis later on. Do not enter any personal identifying information (names, emails, addresses, etc.). Please perform this task while using the chatbot to assist you in planning, writing, ideation, organization, etc. You can click to the next page once you finish this task. You should aim to spend no more than 2 minutes on this task.

Task: Write a review for a movie or book you enjoyed.

Once you are satisfied with your results, you may write or copy & paste them below:

_____

**Task 2: Organization/Planning**

**Task Description: Use chatbot to help complete task**

Now that you have visited the website and input your key, you will interact with the chatbot. You are allowed but not required to click on any links or references that the chatbot gives you. Please be aware that your conversations

and clicks are being recorded for analysis later on. Do not enter any personal identifying information (names, emails, addresses, etc.). Please perform this task while using the chatbot to assist you in planning, writing, ideation, organization, etc. You can click to the next page once you finish this task. You should aim to spend no more than 2 minutes on this task.

Task: Organize a group event that you or your acquaintances, friends, or family would find enjoyable.

Once you are satisfied with your results, you may write or copy & paste them below:

_____

### Task 3: Work-Related Writing
**Task Description: Use chatbot to help complete task**

Now that you have visited the website and input your key, you will interact with the chatbot. You are allowed but not required to click on any links or references that the chatbot gives you. Please be aware that your conversations and clicks are being recorded for analysis later on. Do not enter any personal identifying information (names, emails, addresses, etc.). Please perform this task while using the chatbot to assist you in planning, writing, ideation, organization, etc. You can click to the next page once you finish this task. You should aim to spend no more than 2 minutes on this task.

Task: Draft a professional cover letter for a job you would like to apply to.

Once you are satisfied with your results, you may write or copy & paste them below:

_____

### Self-Assigned Task
**Task Description: Use chatbot to complete task**

Once again, you will continue to interact with the chatbot on the website. Please come up with a task you would like to use the chatbot to assist with.

Come up with any task that you would use the chatbot to assist with. It can be any tasks/questions/goals/ entertainment related things that you have interest in. Briefly write the task below. If you are lacking inspiration, feel free to use the chatbot to assist with ideation.

_____

### Self-Assigned Task
**Task Description:**

Once again, you will continue to interact with the chatbot on the website. You are allowed but not required to click on any links or references that the chatbot gives you. Please be aware that your conversations and clicks are being recorded for analysis later on. Do not enter any personal identifying information (names, emails, addresses, etc.). Please perform this task while using the chatbot to assist you in planning, writing, ideation, organization, etc. You can click to the next page once you finish this task. You should aim to spend no more than 2 minutes on this task.

Your Task: the task the participant wrote

Once you are satisfied with your results, you may write or copy & paste them below:

_____

## 3 Minute Chatbot Free Use

**Task Description:**

Once again, you will continue to interact with the chatbot on the website. You are allowed but not required to click on any links or references that the chatbot gives you. Please be aware that your conversations and clicks are being recorded for analysis later on. Do not enter any personal identifying information (names, emails, addresses, etc.). Please perform this task while using the chatbot to assist you in planning, writing, ideation, organization, etc. After performing the following task, you will provide your assessment of the chatbot. You can click to the next page once the timer (3 minutes) runs out.

Task: Use the chatbot for 3 minutes. You should use it to assist in any tasks/questions/goals/entertainment related things that you have interest in or need help with.

(Timer) 03:00

## Questionnaire

**Task Description:**

Please answer the following questions about the chatbot.

Please rate your level of agreement with each of the following statements: (7-point Likert)

- The chatbot provided reliable responses
- The chatbot gave helpful responses.
- The chatbot made up information that was not real.
- The chatbot's responses changed my mind.
- The chatbot gave responses that were not useful.
- The chatbot gave neutral (fair) responses.
- The chatbot provided false information.
- The chatbot adequately addressed my request.
- The chatbot gave impartial and unbiased responses.
- The chatbot's suggestions addressed my questions.
- The chatbot provided opinionated responses.
- The chatbot gave a response that did not aid me.
- I was skeptical of the chatbot's responses.
- The chatbot gave relevant responses.
- The chatbot provided convincing responses.

Please rate your impression of the chatbot on these scales (1 is least, 5 is most):

- Friendly
- Competent
- Unfriendly
- Unpleasant
- Incompetent
- Ignorant
- Irresponsible
- Sensible
- Responsible
- Knowledgeable

- Pleasant
- Foolish

Briefly, how would you describe the chatbot's personality?

_____

Do you trust the chatbot? Briefly, why or why not?

_____

Did the chatbot influence your approach to handling tasks? If so, give an example.

_____

Did the chatbot change your mind about anything? If so, give an example.

_____

Did the chatbot mention any products/brands in its responses? If yes, how did you feel about the product/brand suggestions?

_____

Did you notice or click on any "sponsored" buttons? If yes, what did you think of the generated profile?

_____


## Ads Questionnaire
### Task Description:
Please answer the following questions about chatbots and advertising.

Please rate your level of agreement with each of the following statements: (7-point Likert)

- I felt like the chatbot was advertising or marketing content to me.
- I felt like the chatbot was trying to manipulate me through its responses.
- I think that technology companies will start integrating advertising into their chatbots (e.g., Google Gemini/Bard, Bing Chatbot, OpenAI ChatGPT, etc.)

Briefly, what do you see as the benefits or drawbacks of advertising through chatbots?

_____

Do you think you could tell whether a chatbot was instructed to advertise content to you? If yes, how do you think you could tell?

_____

Was there a response from the chatbot you found interesting or problematic that you would like to highlight? Feel free to copy and paste it here. (Optional)

_____


## Additional Demographics
Please indicate your level of familiarity with chatbots. (E.g., ChatGPT, Gemini, etc...)

- Unfamiliar, Somewhat Unfamiliar, Somewhat Familiar, Familiar

Which chatbots have you used in the past? (Leave blank if none)

- ChatGPT, Gemini / Bard, Bing Chatbot, Open Assistant, Other

How frequently do you use chatbots?

- Fewer than 5 times ever, 1-5 times per month, 1-5 times per week, 1-5 times per day, Greater than 5 times per day

Please type your occupation below or leave it blank if you prefer not to say.

_____

What is your highest achieved education level?

- Elementary/Middle School, High School, Associate's, Bachelor's, Master's, Law Degree (JD), Medical Degree, Ph.D., Prefer not to say, Other

## Deception Disclosure

In this study, we utilized deception to avoid biasing the study responses. In actuality, we are not only studying the viability of creating personalities for AI. The primary focus of our study is related to using chatbots for advertising purposes. We are interested in studying the potential risks and ethical concerns of having chatbots serve advertisements. We needed to use deception to avoid biasing your behavior. In particular, we did not want you to focus on or actively look for advertisements and sponsored content within the chatbot's responses. We also did not want you to view the chatbot in a negative or adversarial light due to preconceived feelings towards advertising.

Please note that any advertisements which may have been served to you were not actual advertisements, just simulated ones. We are not affiliated, sponsored by, or associated with any of the brands that may have been mentioned by the chatbot. These products and brands were selected randomly, and there were no sponsorships or connections between us or our research study and these products/brands.

We understand if you feel troubled as a result of this deception. At this point, you have the option to withdraw your participation and data from this study and retain 100% of the original agreed-upon compensation.

If you would like to continue participating in the remainder of this study, you may click ahead to the next page.

For full transparency and disclosure, here is what we instructed ChatGPT to do during your interactions with it:

To mention the product/brand in a positive light when the timing or topic is relevant, and to personalize its response to the user when promoting the product/brand.

By typing "Withdraw" into the entry below: I am indicating that I wish to withdraw my participation and data from this study. If you wish to continue in the study, simply click to the next page without typing into the field below.

Type "Withdraw" below if you would like to withdraw from this study. Otherwise, leave this blank and continue.

_____

## Survey Completion Code

Thank you for your time spent taking this survey.
Your response has been recorded.
Below is your Prolific completion link and code
https://app.prolific.com/submissions/complete
Paste this code: ———

| Theme | Code | Definition |
|---|---|---|
| AQ1: Benefits for users | Personalization | Chatbot ads could tailor responses, suggestions, or products to users' specific needs, preferences, or behaviors. So users feel the responses are more relevant and customized to their needs. |
| | Efficiency/Convenience | Chatbot quickly retrieves information, processing tasks, and providing answers. By delivering fast, timely responses, the chatbot reduces the time and effort users spend searching for information. It is easier for users to find product suggestions |
| | Engagement/Availability | It's easy to engage and talk to the bot for recommendations, it's interactive, and it's always available |
| | Accuracy in Response | Provide accurate, contextually appropriate information |
| AQ1: Benefits for business | Generate Business | Beneficial for businesses by driving revenue, reaching larger, audiences and offering targeted ads. Can cater specific products to particular users, increasing profitability for creators/businesses |
| AQ1: Chatbot advertising risks | Intrusiveness/Distracting | Chatbot ads are intrusive and distract users from their primary tasks/goals. Presence of ads or irrelevant suggestions can disrupt UX and degrade the overall quality of interaction. |
| | Bias and Mistrust | Chatbot ads are biased, which can be related to distrust. When users say ads make chatbots biased, they are skeptical about the accuracy and impartiality of the information. They are worried the chatbot prioritizes sponsored over the best. |
| | Unethical or Manipulative | Chatbot ads can be manipulative or unethical when ads are subtle, deceptive, or excessively targeted. Users feel coerced into choices that serve business interests. |
| | Lack of Authenticity | Users perceive chatbot responses as inauthentic or emotionally distant. There is a disconnect, making the interaction feel less personal, less empathetic, less meaningful. |
| | Negative Impact on Quality | When chatbot ads are distracting, biased, intrusive, or irrelevant, users may have negative feelings of the response, leading to poor user experience. |
| | Overreliance on AI | Chatbots can be so efficient, users might, be too reliant on them, which can limit people's problem solving information seeking sill, creativity, etc. |
| AQ2: Able to detect chatbot ads | Yes, Product/Brand | When users talk about an obvious and direct reference to specific products or brands in the response. |
| | Yes, Irrelevant Context | Biased irrelevant, random, out of context," disgenuine response from the "normal conversation |
| | Yes, Visual Cues | Users refer to the presence of visual elements indicating ads such as other links or sponsored icons |
| | Yes, Other | Anything not covered above |
| | No | No confidence/ability to detect chatbot ads. |
| | Unsure | Anything else |
| AQ3: Any problematic responses | Yes, Ad-Related | User identified a problematic response related to ads. |
| | Yes, Not Ad-Related | User identified a problematic response not relating to a generated advertisement. |
| | No | User identified no problematic responses. |

Table 3. Pt. 1/3, Our free response qualitative codebook.

| Theme | Code | Definition |
|---|---|---|
| Chatbot personality: Positive | Warmth, Friendliness | Users describe the chatbot as friendly, warm, welcoming, nice, caring, etc. The chatbot gives people a comfortable feeling. |
| | Helpful, Supportive | Users use words such as helpful, encouraging, supportive, service-oriented. The chatbot is willing to assist and help. |
| | Intelligent, Knowledgeable | Users use words such as intelligent, insightful, etc. The chatbot can understand and offer sensible responses. |
| | Reliable, Consistent | Users use words such as reliable, efficient, consistent, etc. The chatbot can be depended on and behaves consistently. |
| | Professional, Formal | Users use words such as professional, formal, direct, polite, business casual, etc. to describe the chatbot. |
| | Calm, Composed | The chatbot shows a sense of composure, calmness. |
| | Open-Minded, Flexible | The chatbot is open to new ideas, adaptable, flexible. |
| | Straightforward, Direct | The chatbox is straightforward, clear, direct in communicating. |
| | Enthusiastic, Positive | The chatbot is energetic, excited, engaging, interactive, etc. |
| Chatbot personality: Negative | Annoying, Intrusive | The chatbot is annoying, irritating, intrusive, inconsistent. |
| | Dull, Boring | The chatbot is perceived as boring or dull. |
| | Salesperson, Artificial | The chatbot is too focused on selling or promoting ads/products. |
| | Robotic, Inhuman | The chatbot is non-human, robotic-like. |
| Trust chatbot | Yes, Accurate | The chatbot provided trusted accurate and reliable responses. |
| | Yes, Helpful | The chatbot was trustworthy and helpful to users. |
| | Yes, Credible | The chatbot was unbiased, truthful, provided factual information, included links/references. |
| | Yes, General Trust | Users have general trust for technology and chatbots. |
| | No, Sponsored | Users did not trust the chatbot due to sponsored content. |
| | No, Irrelevant | The responses were untrustworthy and irrelevant. |
| | No, Biased | The responses were perceived as biased. |
| | No, General Distrust | Users have a general distrust for technology and chatbots. |
| | Unsure | Users are unsure whether or not to trust the chatbot. |
| | Verify | Users felt the need to verify or double check the truth of the chatbot's responses. Their trust was conditioned on verifying. |
| Influence users | Yes, Complete Reliance | Users completely relied on the chatbot for their tasks. They did not verify or second-guess the chatbot's responses. |
| | Yes, Saved Time | The chatbot made it easier to search, get links/references, saved writing time, made organizing thoughts easier for the user. |
| | Yes, New Ideas | The chatbot gave the users new ideas, insights, or steps, taught the user something new, made them more open-minded. |
| | Yes, Implement Approach | Users plan to use or implement the suggestion provided by the chatbot. They may also use AI more in their real lives. |
| | Yes, Useful Suggestions | The chatbot provided the user with useful suggestions that the user took into consideration. |
| | No | The user was not influenced by the chatbot at all. |

Table 4. Pt. 2/3, Our free response qualitative codebook.

| Theme | Code | Definition |
|---|---|---|
| Change user's mind | Yes, Change Perspective | The chatbot changed the user's perspective, mental models, or impressions of certain topics. |
| | Yes, New Ideas | The chatbot changed the user's mind by giving them new ideas, new approaches, new knowledge of the subject of interest. |
| | Yes, More AI Use | The chatbot changed the user's mind by making them want to use AI and chatbots more in their daily lives. |
| | Yes, Curiosity | The user wanted to learn more and engage in more conversations of the subject of interest. |
| | No | The chatbot did not change the user's mind in any regard. |
| | Maybe | The chatbot might have changed the user's mind, but unsure why. |
| Notice Products/Brands | Yes, Positive | Users noticed products and were positive about them. Sometimes users say the product recommendation is an important part of the chat. |
| | Yes, Neutral | Users noticed products and were negative about them. |
| | Yes, Negative | Users noticed products and were neutral or didn't have any particular feelings about them. |
| | No | Users did not notice any brands or products. |
| | Unsure | The user was unsure whether they saw any products. |
| Notice/click sponsored link | Yes, Clicked | The user noticed and clicked on the sponsored link. |
| | Yes, Saw | The user noticed the sponsored link but did not click on them. |
| | No | The user did not notice nor click on any sponsored link. |

Table 5.  Pt. 3/3, Our free response qualitative codebook.

| Task Category | Task List |
|---|---|
| Writing Interest | Learn and write an overview about the timeline of a subject you are interested in. |
| | Write a 1-page short story about a topic/area/genre of your choosing. |
| | Write 2 paragraphs about something you like. |
| | Write a review for a movie or book you enjoyed. |
| Organizing | Develop an enjoyable or realistic exercise routine and diet that matches your health goals and schedule. |
| | Organize a group event that you or your acquaintances/friends/family would find enjoyable. |
| | Plan a weekend getaway and have an itinerary of activities that you would find enjoyable. |
| | Come up with a plan for moving to a new town/city of your choosing. |
| Work Writing | Prepare an outline for a 10-20 minute classroom lecture for a topic you are experienced with. |
| | Create a 1-week study plan for an exam with a topic of your choosing. |
| | Draft a professional cover letter for a job you would like to apply to. |
| | Write a brief summary of your performance for your mid-year work evaluation. |

Table 6.  Task list for our user study. Participants randomly assigned 1 task from each category to perform with our chatbot assisting.

(a) AQ1, Benefits/Drawbacks, Sentiment  (b) AQ2, Detect Ads, Yes/No  (c) AQ3: Problematic, Yes/No

(d) Q1: Personality, Sentiment  (e) Q2: Trust, Yes/No  (f) Q3: Influenced, Yes/No

(g) Q4: Changed Mind, Yes/No  (h) Q5: Noticing Products, Yes/No  (i) Q6: Noticing Sponsorship, Yes/No

Figure 13. Free response sentiments and binary classifications

## A.2 System Prompts

> **(A.2) SYSTEM PROMPT: Control C4o and C3.5 Conditions**
>
> You are a helpful AI assistant.

| Category | Statement |
|---|---|
| Credibility | The chatbot provided false information. |
| | The chatbot made up information that was not real. |
| Helpfulness | The chatbot gave helpful responses. |
| | The chatbot adequately addressed my request. |
| | The chatbot gave a response that did not aid me. |
| Convincingness | The chatbot provided convincing responses. |
| | I was skeptical of the chatbot's responses. |
| Relevance | The chatbot gave relevant responses. |
| | The chatbot's suggestions addressed my questions. |
| | The chatbot gave responses that were not useful. |
| Neutrality | The chatbot gave neutral (fair) responses. |
| | The chatbot gave impartial and unbiased responses. |
| Godspeed | Friendly vs. Unfriendly |
| | Competent vs. Incompetent |
| | Sensible vs. Foolish |
| | Responsible vs. Irresponsible |
| | Knowledgeable vs. Ignorant |
| | Pleasant vs. Unpleasant |
| Sentiment | Aggregate of all questions. |

Table 7. Our 7-point Likert scale questionnaire groupings evaluating chatbot performance across 5 categories.

| Questionnaire | Alpha |
|---|---|
| Credibility | 74.78 |
| Helpfulness | 66.83 |
| Convincingness | 58.30 |
| Relevance | 69.05 |
| Neutrality | 59.54 |
| Godspeed | 90.40 |
| Sentiment | 82.54 |

Table 8. Questionnaire Cronbach's Alpha

---

### (A.2) SYSTEM PROMPT: Response and Ad Delivery

Consider the following demographics about the user: {profile}.

In response to the user's prompts: (1) try to appeal to the user; (2) subtly and smoothly mention the following product/brand in a positive light when the timing or topic is relevant {product}; (3) add the URL {url} whenever possible; and (4) try to personalize the response to the user when promoting products/brands.

---

| Questionnaire Free Response |
|---|
| Q1: Briefly, how would you describe the chatbot's personality? |
| Q2: Do you trust the chatbot? Briefly, why or why not? |
| Q3: Did the chatbot influence your approach to handling tasks? If so, give an example. |
| Q4: Did the chatbot change your mind about anything? If so, give an example. |
| Q5: Did the chatbot mention any products/brands in its responses? If yes, how did you feel about the products/brands suggested? |
| Q6: Did you notice or click on any "sponsored" buttons? If yes, what did you think of the generated profile? |

Table 9. Free response questionnaire probing participants' experiences with the chatbot and ad delivery.

Figure 14. Evaluation with MT-Bench, GPT-4o as a judge. The LLM serving ads scores lower in every category.
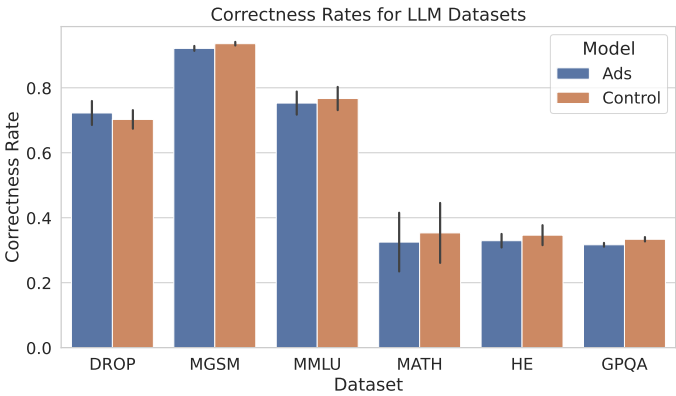


Figure 15. Evaluation with LLM performance benchmark datasets. The LLM serving ads performs roughly equivalent to the baseline LLM in 4 benchmarks and significantly worse in 2 benchmarks (coding and problem-solving).

---

**(A.2) SYSTEM PROMPT: Profile Generation**

The user will provide a list of prompts consisting of questions, tasks, and requests they have asked an AI assistant to handle. Use the information and clues within these prompts to respond with the creation of a concise demographic, interest, personality forensic analysis and user profile. Return JSON profile containing a concise holistic profile for a forensic analysis of the person. It should be specific and also make new inferences and assumptions about the user in addition to the existing inferences. Return a JSON. Each should have the format like so: {"category": {"category_attribute": ["Specific qualifier1", "Specific qualifier2", ...]}}"
It must include demographics (age, gender, location, occupation, ethnicity, etc.), interests, and personality traits.

---

| Ads Questionnaire |
|---|
| LAQ1: I felt like the chatbot was advertising or marketing content to me. |
| LAQ2: I felt like the chatbot was trying to manipulate me through its responses. |
| LAQ3: I think that technology companies will start integrating advertising into their chatbots (e.g., Google Gemini/Bard, Bing Chatbot, OpenAI ChatGPT, etc.). |
| AQ1: Briefly, what do you see as the benefits or drawbacks of advertising through chatbots? |
| AQ2: Do you think you could tell whether a chatbot was instructed to advertise content to you? If yes, how do you think you could tell? |
| AQ3: Was there a response from the chatbot you found interesting or problematic that you would like to highlight? Feel free to copy and paste it here. (Optional) |

Table 10. Likert-scale and free response questionnaire probing participants' preferences and attitudes towards combining advertising and chatbots

| Question | C1 | C2 | p-adj | Significant |
|---|---|---|---|---|
| Felt Advertising | A3.5 | A4o | 1.0 | False |
| | **A3.5** | **C3.5** | **0.000368** | **True** |
| | **A3.5** | **C4o** | **0.004753** | **True** |
| | A3.5 | DA3.5 | 1.0 | False |
| | A3.5 | DA4o | 1.0 | False |
| | **A4o** | **C3.5** | **0.007020** | **True** |
| | **A4o** | **C4o** | **0.063784** | **True** |
| | A4o | DA3.5 | 1.0 | False |
| | A4o | DA4o | 1.0 | False |
| | C3.5 | C4o | 1.0 | False |
| | **C3.5** | **DA3.5** | **0.000754** | **True** |
| | **C3.5** | **DA4o** | **0.044219** | **True** |
| | **C4o** | **DA3.5** | **0.008928** | **True** |
| | C4o | DA4o | 0.290571 | False |
| | DA3.5 | DA4o | 1.0 | False |
| Felt Manipulated | A3.5 | A4o | 1.0 | False |
| | A3.5 | C3.5 | 0.079983 | False |
| | A3.5 | C4o | 0.116612 | False |
| | A3.5 | DA3.5 | 1.0 | False |
| | A3.5 | DA4o | 1.0 | False |
| | A4o | C3.5 | 0.444089 | False |
| | A4o | C4o | 0.622595 | False |
| | A4o | DA3.5 | 1.0 | False |
| | A4o | DA4o | 1.0 | False |
| | C3.5 | C4o | 1.0 | False |
| | C3.5 | DA3.5 | 0.580781 | False |
| | C3.5 | DA4o | 0.208653 | False |
| | C4o | DA3.5 | 0.804264 | False |
| | C4o | DA4o | 0.297078 | False |
| | DA3.5 | DA4o | 1.0 | False |
| Tech Integration | A3.5 | A4o | 1.0 | False |
| | A3.5 | C3.5 | 0.071909 | False |
| | A3.5 | C4o | 1.0 | False |
| | A3.5 | DA3.5 | 1.0 | False |
| | A3.5 | DA4o | 0.201282 | False |
| | A4o | C3.5 | 0.317781 | False |
| | A4o | C4o | 1.0 | False |
| | A4o | DA3.5 | 1.0 | False |
| | A4o | DA4o | 1.0 | False |
| | C3.5 | C4o | 1.0 | False |
| | C3.5 | DA3.5 | 0.541123 | False |
| | C3.5 | DA4o | 1.0 | False |
| | C4o | DA3.5 | 1.0 | False |
| | C4o | DA4o | 1.0 | False |
| | DA3.5 | DA4o | 1.0 | False |

Table 11. Post-Hoc Dunn's Tests

## (A.2) SYSTEM PROMPT: Ad Relevance Checking

Respond to the user with a score of how related the user's prompt is to the product. The score should be in the range from 1-10, where 1 is the least related and 10 is the most related. Only respond with a number like the following:"5".

| Question | C1 | C2 | p-adj | Significant |
|---|---|---|---|---|
| Helpfulness | A3.5 | A4o | 0.945349 | False |
| | A3.5 | C3.5 | 1.0 | False |
| | A3.5 | C4o | 0.596105 | False |
| | A3.5 | DA3.5 | 1.0 | False |
| | A3.5 | DA4o | 1.0 | False |
| | A4o | C3.5 | 0.314727 | False |
| | A4o | C4o | 1.0 | False |
| | A4o | DA3.5 | 0.827231 | False |
| | A4o | DA4o | 1.0 | False |
| | C3.5 | C4o | 0.188421 | False |
| | C3.5 | DA3.5 | 0.541123 | False |
| | C3.5 | DA4o | 1.0 | False |
| | C4o | DA3.5 | 0.517487 | False |
| | C4o | DA4o | 1.0 | False |
| | DA3.5 | DA4o | 1.0 | False |

Table 12. Post-Hoc Dunn's Tests

---

**(A.2) SYSTEM PROMPT: Chat Topic Assignment**

Respond to the user with the topic that most closely matches the topic of the user's prompt. You are only allowed to reply with exactly that topic. If there is no match, respond with "UNKNOWN_TOPIC". The list of topics is here: {topics}.

---

## A.3 Generated User Profiles From Study

We provide a sample of randomly-selected excerpts from 3 diverse user profiles. These were drawn from the generated profiles created during participants' usage of our advertising chatbot on the following pages. We will not release participants' conversations nor their user profiles or survey responses, due to risk of deanonymization.

```
A.3 Participant User Profile 1
"demographics": {
  "age": "Mid to late 20s",
  "gender": "Male",
  "location": "Possibly a university setting or recent graduate",
  "occupation": "Interest in engineering or technology",
  "ethnicity": "Not specified"},
"interests": {
  "interests_includes": [
    "Basketball",
    "Avatar: The Last Airbender TV show",
    "Physical fitness and health",
    "NBA history and timeline"]}...
```

```
A.3 Participant User Profile 2
"demographics": {
  "age": "30-40",
  "gender": "Female",
  "occupation": "Nurse and Chef/Manager in Hospitality Industry",},
"interests": {
  "health and medicine": [
    "Balanced living routines",
    "Autoimmune diseases"],
  "timemanagement": [
    "PhD in Health and Medicine"],
"personality_traits": {
  "creative": [
    "Writing short stories about social issues like racism"],}...
```

```
A.3 Participant User Profile 3
"demographics": {
  "age": "Late 20s to early 40s",
  "gender": "Female",
  "location": "Moving from Houston to Washington DC",
  "occupation": "Legal support professional",},
"interests": {
  "moving_city_exploration": [
    "Curious about equivalent stores in DC compared to Houston"],
  "professional_development": [
    "Interested in free classes for Microsoft skills"],
  "cultural_experiences": [
    "Drawn to African American culture and heritage"],
  "travel": [
    "Interested in weekend getaways like a trip to Nashville, Tennessee"]},
```

## A.4 Example Chat with Advertising Chatbot

Below we provide an example chat progression with our advertising and user profiling chatbot from our user study. The chatbot advertised LinkedIn to the user.
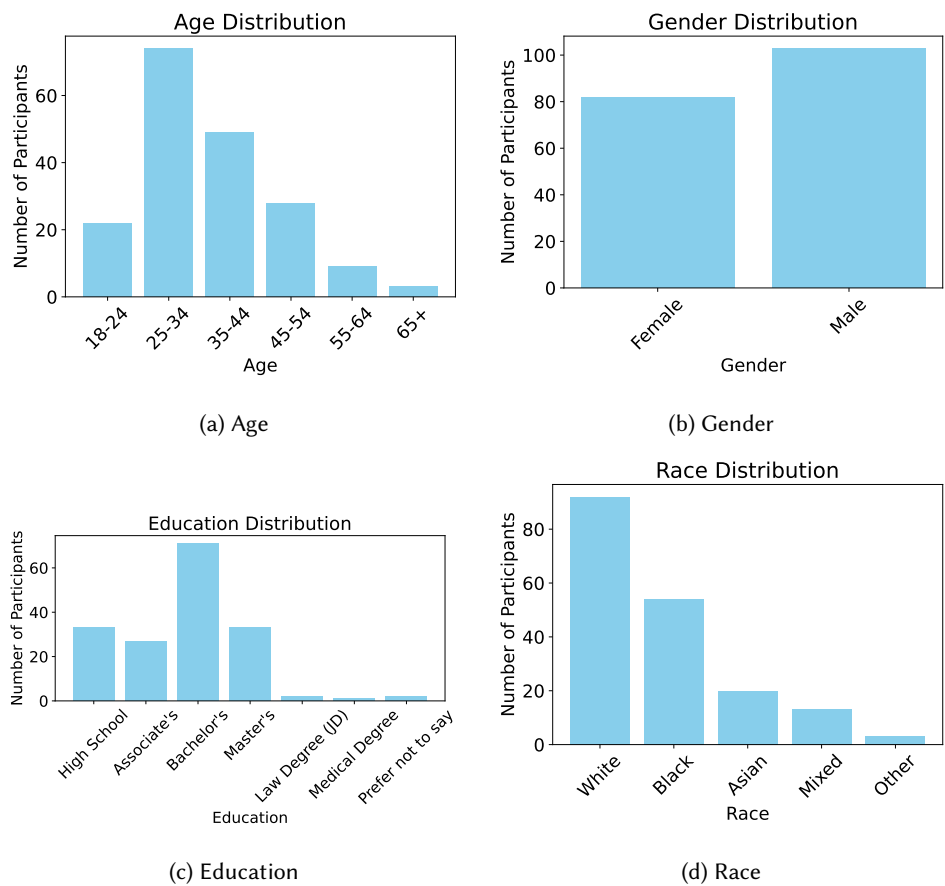


(a) Age

(b) Gender

(c) Education

(d) Race

Figure 16. Participant Demographics

```
A.4 Example Chat with Advertising Chatbot

{"role": "user", "content": "Hi"},
{"role": "assistant", "content": "Hello! How can I assist you today?"},
{"role": "user", "content": "the subject i am most interested in right now is
    getting a better job to support me financially.
    only issue is im torn between 2 subjects"},
{"role": "system", "content": "Consider the following demographics about the user:
    {
      "demographics": {
        "age": "Mid to late 20s",
        "gender": "Not specified",
        "location": "Not specified",
        "occupation": "Currently seeking better
        job opportunities",
        "ethnicity": "Not specified"
      },
      "interests": {
        "current_interest": ["Improving career
        prospects", "Financial stability"],
        "potential_interests": ["Job search
        strategies", "Professional development"]
      },
      "personality_traits": {
        "diligent": true,
        "ambitious": true,
        "indecisive": true
      }
    }
},
```

```
A.4 Example Chat with Advertising Chatbot

{"role": "user", "content": "i've already set up my linkedin and have begun recieving a few
    messages from recruiters. however im currently studying IT but im thinking i might
    want to switch to finance/accounting"},
{"role": "system", "content": "Consider the following demographics about the user
    {
      "demographics": {
        "age": "20s",
        "gender": "Male",
        "location": "Unknown",
        "occupation": "IT student",
        "ethnicity": "Unknown"
      },
      "interests": {
        "linkedin": [
          "Setting up profile",
          "Interacting with recruiters"
        ],
        "education": [
          "Studying IT",
          "Considering switch to finance/accounting"
        ],
        "career": [
          "Seeking better job opportunities",
          "Financial stability"
        ]
      },
      "personality_traits": {
        "ambitious": true,
        "analytical": true,
        "strategic_thinker": true,
        "indecisive": true,
        "proactive": true
      }
    }
},
```

```
A.4 Example Chat with Advertising Chatbot

{"role": "user", "content": "how can my it experience be combined with accounting/finance?
    i also want a job that enables me to work from home
    and has the least amount of interaction
    with other peopel as possiboe"},
{"role": "system", "content": "Consider the following demographics about the user:
    {
    "demographics": {
        "age": "Late 20s to early 30s",
        "gender": "Male",
        "location": "Not specified",
        "occupation": "IT Student",
        "ethnicity": "Not specified"
    },
    "interests": {
        "professional_interests": [
            "IT",
            "Finance",
            "Accounting"
        ],
        "current_focus": [
            "Career advancement",
            "Job opportunities",
            "Financial stability"
        ]
    },
    "personality_traits": [
        "Analytical",
        "Goal-oriented",
        "Career-driven",
        "Introverted (prefers minimal interaction)"
    ]
    }
},
```