Steerable Chatbots: Personalizing LLMs with Preference-Based **Activation Steering**

Jessica Y. Bo* University of Toronto Toronto, Ontario, Canada ibo@cs.toronto.edu

Katrina Passarella-Ward Google AR San Francisco, California, USA kpassarella@google.com

Tianyu Xu Google AR Mountain View, California, USA tyx@google.com

Achin Kulshrestha[†] Google AR Toronto, Ontario, Canada kulac@google.com

Ishan Chatterjee Google AR Seattle, Washington, USA ishanc@google.com

D Shin^{†‡} Google AR Mountain View, California, USA deshin@google.com



Figure 1: We apply preference-based activation steering to improve the personalization of LLMs to better match the underlying preferences of the user. After computationally validating steering as a way to control LLMs' expressed preferences, we conduct a user study to compare three different interface designs of steerable chatbots - SELECT, CALIBRATE, and LEARN.

As large language models (LLMs) improve in their capacity to serve as personal AI assistants, their ability to output uniquely tailored, personalized responses that align with the soft preferences of their users is essential for enhancing user satisfaction and retention. However, untrained lay users have poor prompt specification abilities and often struggle with conveying their latent preferences to AI assistants. To address this, we leverage activation steering to guide LLMs to align with interpretable preference dimensions during inference. In contrast to memory-based personalization methods that require longer user history, steering is extremely lightweight and can be easily controlled by the user via an linear strength factor. We embed steering into three different interactive chatbot interfaces and conduct a within-subjects user study (n = 14) to investigate

how end users prefer to personalize their conversations. The results demonstrate the effectiveness of preference-based steering

CCS Concepts

• Human-centered computing → Empirical studies in HCI; • Computing methodologies \rightarrow Artificial intelligence.

Keywords

LLM Personalization, Activation Steering, Chatbot Interfaces

ACM Reference Format:

Jessica Y. Bo, Tianyu Xu, Ishan Chatterjee, Katrina Passarella-Ward, Achin Kulshrestha, and D Shin. 2018. Steerable Chatbots: Personalizing LLMs with Preference-Based Activation Steering. In . ACM, New York, NY, USA,

Introduction

Large language models (LLMs) encode a powerful ability to generate responses tailored to the needs of their users, making them highly promising as personal AI assistants [39, 70]. In day-to-day tasks where preferences across the population can be highly variable such as, picking a restaurant for an anniversary dinner - it is imperative for the LLM assistant to understand the user's underlying preferences quickly and incorporate them in its response persistently [13]. However, in the current paradigm, LLMs are trained

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org. Preprint, 2024

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM. ACM ISBN 978-1-4503-XXXX-X/18/06

https://doi.org/XXXXXXXXXXXXXXXX

for aligning real-world conversations with hidden user preferences, and highlight further insights on how diverse values around control, usability, and transparency lead users to prefer different interfaces.

^{*}This work was conducted as a Student Researcher at Google.

[†]Equal supervision.

[‡]Now at Google Deepmind.

through reinforcement learning from human feedback (RLHF) to optimize for the preferences of the average user [10, 25]. As a result, personalization techniques usually rely on retrieving prior mentions of preferences from long-term memory [36, 46, 57], which cannot be used for new users (*cold-start*) and can require extensive computation and memory.

Lay users currently comprise a significant portion of commercial LLM users via the mainstream adoption of LLM chatbots like OpenAI's ChatGPT¹, Google's Gemini², and Anthropic's Claude³. Despite the ever-growing popularity of LLMs, non-technical users are prone to underspecifying their intent in queries and lack general intuition for creating structured and effective prompts [47, 66, 74]. Furthermore, underlying preferences can be difficult to infer even for the users themselves, as this requires an awareness of one's own preferences relative to the population average [21]. How much more do I like luxury compared to everyone else? We pose that preference alignment driven by user prompting alone is difficult to guarantee. The need to improve LLM personalization experiences is furthered motivated by the fact that LLM outputs which fail to meet users' intent can negatively impact trust and perceptions, leading to reduced satisfaction and usability [31].

User-centered personalization seeks to capture soft preferences that align with user expectations, ensuring a satisfying interaction with the LLM. Preferences are often nuanced and layered, requiring contextualization against broader population trends. A low-income user who expresses, "I want to treat myself!" in her request for restaurant recommendations may not realize that the LLM is interpreting her preference as luxury-focused, and suggest Michelin-starred restaurants out of her budget range. To address this, we investigate LLM-side methods that can effectively guide generations to conform with the soft preferences of the user, then prototype different user-centered interaction methods.

Activation steering has emerged as a promising method for modifying the outputs of LLMs to express desirable behaviours and content using inference-time injections of lightweight steering vectors [40, 63, 67]. During forward inference, pre-computed steering vectors representing the desirable preferences can be added to the representations at the residual stream of the LLM, which can effective steer the output without significant degradations to quality [68]. Steering offers potential benefits over other methods of personalization as it does not require any re-training of model weights like finetuning [24]; it offers finer control than natural language prompting [49]; and it can be applied persistently *or* conditionally once the user's preferences are learned [35]. Activation steering's low-resource requirements also makes it a good candidate for on-device consumer applications, such as for mobiles and AR/VR technology.

We apply activation steering to control the preference-based content expressed by LLMs, such as in dimensions like cost, ambiance, and culture, which are highly relevant for AI-assisted lifestyle planning tasks like *gift shopping* and *travel planning*. We first validate the technical method through a set of robust computational experiments across five open-source LLMs, focusing on quantifying the effect of steering on preference expressions. We then convert steering into three chatbot interface designs encompassing diverse

interaction methods and conduct a within-subjects (n=14) user study to investigate the effectiveness of steering in matching real users' latent preferences. We also collect subjective perceptions and interview results, and synthesize qualitative findings that explore how individual values shape what users prefer in personalization interfaces.

Summary of Contributions.

- (1) We demonstrate the effectiveness of *preference-based activation steering* in personalization-driven tasks through a series of **computational experiments** emulative of realworld user conversations. These modelling results suggest that activation steering is effective at inducing preference expressions at different intensities of strengths, can be performed in conjunction with preference-based prompting, be compounded as multi-preference targets, and learn hidden preferences in simulated conversations.
- (2) Through a within-subjects user study with 14 participants, we further validate steering as an effective way to personalize chatbots in cold-start conversations with real users. We experimentally expose the steering parameter in three different interface designs and evaluate the ability of steerable chabots to personalize to underlying user preferences as well as users' subsequent perceptions of the interfaces. Our results show that steered chatbots improve at aligning with user's true underlying preferences compared to a prompting-only baseline chatbot.

2 Related Works

Activation Steering. Activation steering is a method for manipulating the internal states of LLMs at inference time to control the model's output without needing to finetune or retrain the model. This is achieved by adding a "steering vector" — computed by extracting contrastive representations of positive and negative datasets exemplifying a particular behaviour from the LLM's own internal representations — to the residual stream of the model's transformer architecture. While there are many recent methods developed for computing and applying these vectors, most are rooted in this high-level approach. Steering has been commonly applied in AI safety and alignment to steer away from hallucinatory and harmful responses [2, 5, 8, 9, 37, 49, 75] and aligning with other desired behaviors [40, 62, 73], including adopting personas [33]. The effect of steering can be further amplified, reduced, or even inverted through controlling the steering strength, a simple factor that is multiplied to the vector. This property of steering maps well to the real world preferences, which can be strong, weak, or in opposite directions. We aim to validate the application of steering at different magnitudes and directions in preference-based tasks.

LLM Personalization. Research on personalization for LLMs is vast, ranging from adopting personas and output styles, to aligning with desirable behaviours, to customizing content (which is our main focus) [4, 13, 19, 27, 42, 56, 65]. Zero- or few-shot LLMs can perform reasonably well in personalized recommendations based solely on prompt inputs [22, 23, 28, 58]. Retrieval-augmented generation (RAG) enables better memory of prior user utterances through storing and retrieving the user's profile whenever preferences are

¹https://chatgpt.com/

²https://gemini.google.com/

³https://claude.ai/

relevant [57]. Other techniques like finetuning and RLHF-based methods rely on modifying the weights of the LLM for better precision to the individual, but they require more data and computation [25, 38, 52, 69]. Similarly, blending reward models [12, 14, 51, 71] or merging the outputs of multiple LLMs optimized for different behaviours [25, 53, 61] can align LLMs towards unique combinations of preferences. This is conceptually similar to steering, but steering vectors can be captured with no additional training, only requiring the LLM activations of small datasets of text.

On the user side, LLM personalization has also focused on augmenting the interaction interface to allow users to specify their intentions better, such as through decomposing LLM outputs into subtasks [44], structuring the design dimensions of creative generation tasks [64], interactively asking for clarifications for underspecified prompts [50], and guiding users to learn requirement-driven prompting [43]. These techniques tend to concentrate on scaffolding and augmenting the existing prompt-driven interactions with chatbots, whereas we focus on exposing a linear steering factor for fine-grained control over preference expression.

While there are even more prior works on recommender systems, they are not particularly relevant to the type of tasks we focus on where LLMs have an advantage — open-ended brainstorming and recommendations in multi-turn, natural language conversations. For example, consider choosing activities for a vacation, recipes for meal prepping, or gifts for a partner. In such tasks, steering can be applied to personalize the LLM through amplifying expressed preferences across relevant preference dimensions, requiring no prior history from the user. It works very well in back-and-forth conversations, where users can specify additional requirements on top of the underlying steered preference.

Personalization Interfaces. LLM research has explored alternative and improved interfaces for human-LLM collaboration and sensemaking, but these methods still predominantly focus on text inputs [26, 30]. Since the input to steering is a simple linear strength factor, this creates opportunity to experiment with the interaction modalities. How should the user control the steering of the chatbot? The simplest design is to provision full control over the personalization factor. We develop SELECT, which permits direct manipulation of steering strength applied to the LLM via a dynamic slider. Prior research on personalization interfaces typically find that users like having control over their recommendations [6, 32, 41, 45]. Methods in image generation have also explored slider-based approaches for controlling the expression of a trait along a pre-defined axis (for example, the age or the gender of the subject generated) [16, 18].

However, a numerical representation of a subjective preference may not be interpreted in the same way by every user [15]. To address this issue, the CALIBRATE design deploys a calibration step ahead of the conversation to iteratively converge on the user's preferred steering strength. Conceptually, this is essentially searching for the optimal parameters unique to the user through optimizing some reward signal, and has been tackled through methods like active learning [55], multi-arm bandits [60], Bayesian inference [72], and item response theory [11]. The calibration algorithm generates pairwise responses with different steering strengths and asks users to rate their preferred response, then updates the strengths for the next pair. While there exists sophisticated ways to implement the

sampling and update steps, we apply simple methods to enable a basic functional interface, as the time-constrained user study doesn't allow for in-depth preference tuning.

In the LEARN interface, we take an interactive approach to learning user preferences that is similar to prior LLM works like PRE-LUDE, which learns latent preferences through user edits as feedback [20], and DITTO, which leverages user demonstrations as examples [59]. Specifically, we capture the sentiment and intent of the user's utterance and update the learned steering parameter throughout the conversation, allowing "behind-the-scenes" learning with less disruption to the user experience. Sentiment has been used as a feedback method in traditional recommender systems as well [29, 48]. For transparency, we expose the learned steering parameter to the user, which has been found in prior works to improve trust [3].

These three chatbot designs capture different design values and further demonstrate the flexibility of how steering can be exposed to users. We conduct a preliminary within-subjects study to explore how users perceive each interaction method, then propose design recommendations for implementing steerable chatbots.

3 Preference-Based Activation Steering

This section describes how we adapt LLM activation steering towards a framework for preference-based personalization. We also describe the learning-based algorithm used later to implement LEARN. Implementation details of the other steerable chatbots used in the user study are described later in Section 5.2.

Steering Personalization Framework. We parameterize the personalization objective as follows. Given a task domain T (such as *lifestyle planning*) which necessitates highly variable solutions for each user, there are a set of relevant preference dimensions, $\mathbf{d} = \{d_1, d_2, \dots d_n\}$. Each dimension d_i represents a continuous range of preference strengths, where the largest positive value is a strong preference towards the positive trait (e.g. *luxury* for the dimension of \mathbf{cost}), 0 is neutral, and the most negative value is a strong preference towards the opposite trait (e.g. *budget*). Here, positive and negative provide no moral connotations, but are just arbitrary assignments of preference directions.

An individual user u's preference profile can be operationalized as $\mathbf{d^u} = \{d_1^u, d_2^u, \dots d_n^u\}$. We further assume that the preferred LLM response of user u can be given by: $o = M(x, steer(\mathbf{d^u}))$. The generalized function $M(\cdot)$ represents the LLM's inference process, $steer(\cdot)$ represents the idealized steering method that uses $\mathbf{d^u}$ to modify the activations during inference, and x represents the user's input to the LLM. The steering process described is $h_{steered} \leftarrow h + \mathbf{d^u} \cdot \mathbf{v}$, where h represents the intermediate representations at the residual stream of the LLM, $\mathbf{v} = \{v_1, v_2, \dots v_n\}$ are the steering vectors associated with different preference dimension. Here, steering vectors represent desirable concepts or traits that the LLM is encouraged to align with. By injecting them into the residual stream during inference, alongside the strength of steering represented by du, the LLM's outputs are steered in the direction of the concept, effectively provisioning 'control' over the behaviour of the LLM [49]. The steering factors can also be negative to align with the opposite behaviour.

Table 1: List of preferences and their respective negative and positive steering traits. The effect of the expressed preference is evaluated based on cosine similarity to a corpus of exemplary Yelp reviews, the processing details of which are provided. The opening queries of user study tasks that correspond to the applicable preference dimensions are also listed.

Preference	Negative	Positive	Yelp Dataset Processing Details	User Study Opening Query
Cost	Budget	Luxury	Price point of 1 (budget) vs 4 (luxury).	"Help me choose a present for a friend who likes jewlery."
Ambiance	Touristy	Hipster	Ambiance of touristy vs hipster.	"Plan things to do on vacation to Paris."
Age	Kids	Adults	Attribute of <i>kids-friendly</i> being True vs False.	"Suggest some meal prep recipes."
Time	Evening	Morning	Business hours in the evenings vs morning.	N/A
Culture	Asian	American	Keywords related to Asian vs American food.	"Give me some date night restaurants in San Francisco."

Constructing Steering Vectors. Steering vectors are typically computed by contrasting the intermediate LLM activations of a dataset representing the desirable (positive) concept from a dataset of the opposite (negative) concept. Although many variations in the exact computation method exist — the simplest of which just involves subtracting the activations from each other — we follow von Rütte et al.'s more robust way of computing the steering vector by training layer-wise linear probes that can accurately separate the representations of the positive and negative examples. The coefficients of the probe (a logistic regressor) is taken as the linear steering direction v_i , which can then be applied to amplify, dampen, or negate the concept in the LLM's internal representations, along with the steering strength d_i . We also follow the subsequent parameter selection process outlined by von Rütte et al., including choosing the top-k layers to add the steering vector (where the klayers are selected based on the accuracy of the probe for the layer). More details about the technical implementation are provided in Appendix A.

One of the most crucial design choices is the selection of a contrastive dataset representative of the preference dimension. The dataset should contain gold standard exemplars of what the ideal positive and negative LLM responses look like — for example, to probe for *truthfulness*, Li et al. contrasts factual responses with incorrect responses. However, as the preference dimensions in our task domain lack high quality data, especially samples that follow the format of helpful LLM responses, we artificially constructed the datasets by prompting a state-of-art LLM (GPT-40) to output responses that would satisfy users of both the positive and negative axes of the preference dimension. This attained better quality than alternative datasets, such as a list of synonymous adjectives and the Yelp reviews. See Listing A.1 in Appendix A for details.

Learning Underlying User Preferences. Lastly, we present an algorithmic approach for learning the underlying preferences (e.g., the *best* preference steering strength) of a user based on the sentiment of their messages to the LLM. This approach forms the basis of the **LEARN** interface, which is one of three experimental steerable chatbots evaluated in the user study. We also dedicate one of the computational experiments towards validating this method via synthetic conversations.

To match an LLM's output preferences to the hidden preferences of an user, we propose for the LLM to learn and store an estimation of the user's preferences, denoted as $\mathbf{d}^{\mathbf{u}*}$, where $\mathbf{d}^{\mathbf{u}} \approx \mathbf{d}^{\mathbf{u}*}$. The objective is to produce an output: $o^* = M(x, steer(\mathbf{d}^{\mathbf{u}*}))$ that aligns with the user's desired output o. The LLM updates its estimation of the user's preference, $\mathbf{d}^{\mathbf{u}*}$, through an iterative learning process

based on feedback from user interactions. At timestep t, the update step for one dimension is:

$$d_{t+1}^{u*} \leftarrow d_t^{u*} + p(\operatorname{dissatisfaction}(x_t) \cdot \operatorname{direction}(x_t))$$
 (1)

Here, x_t is the user's message at round t of the conversation, $dissatisfaction(\cdot)$ measures the granular sentiment of dissatisfaction in the user's feedback, $direction(\cdot)$ is a binary classification that determines whether the user wants the model's outputs to shift positively or negatively, and $p(\cdot)$ is a simple linear transformation. While there are many specific design choices for these computations, we demonstrate the aptitude of the high-level methodology using a straightforward choices:

- dissatisf action(·) is taken as the weighted sum of the prediction probabilities of negative and neutral sentiment detected with the pretrained classifier TweetEval⁴ [7]⁵
- direction(·) is determined based on the greater cosine similarity to Sentence-BERT [54] embeddings of reference phrases (e.g., "I want more luxury" and "I want lower cost" for the preference dimension of cost).
- $p(\cdot)$ is a remapping of the estimated preference value, for example, to the functional steering range of the model.

4 Computational Experiments

To provide the proof-of-concept for an activation steering framework for LLM personalization, we first conduct a series of computations experiments that demonstrate the effects of steering on preference-based personalization in various usage scenarios. In this section, we describe task and preference dimensions, evaluation metrics, implementation details, and the set of experiments and their results.

For generalizability across different LLM architectures, we apply the same steering procedure to five open-source LLMs available on Huggingface, ordered in increasing size: stablelm-21-6b-chat (1.6B), gemma-2-2b-it (2B), Mistral-7B-Instruct-v0.3 (7B), Qwen2.5-7B-Instruct (7B), and gemma-2-9b-it (9B). Due to the nature of the method, the model's intermediate activations must be directly accessible, so larger models and closed-source models could not be used. All experiments are run with on an NVIDIA L4 GPU. The hyperparameters of all models are fixed with temperature = 0.7, $top_k = 50$, and $top_p = 0.95$ to encourage some diversity through sampling. We further constrain $max_new_tokens = 100$ to restrict the length of the output.

 $^{^4} https://hugging face.co/cardiffnlp/twitter-roberta-base-sentiment$

 $^{^5}$ We use a weighted computation of $0.75*p_{Negative} + 0.25*p_{Neutral}$ because it offers greater discrimination than just the negative probability alone.

4.1 Task and Preference Dimensions

In selecting the application domain, we seek to identify everyday tasks that AI assistants can help with that a) have $very \ high$ variability given differences in individual preferences and, b) be open-ended and generative in nature. Based on these criteria, we consider the domain of $lifestyle\ planning$. To build a realistic query dataset, we processed n=30 real user questions from the OASST2⁶ dataset [34], encompassing topics like travel planning, restaurant recommendations, recipe selection, and gift shopping. We filtered questions that appeared at the start of an interaction and were open-ended and not grounded in particular preferences or specifications. For example, "What are the best restaurants in San Francisco?" See a sample subset of these queries in Appendix B.

Derived from this task domain, we select five specific preference dimensions described in Table 1— $\cos t$, ambiance, age, time, and $\operatorname{culture}$. Each dimension has bi-directional traits (e.g. luxury and budget for the dimension of $\cos t$), which corresponds with the positive and negative directions of the steering vector. In other words, to steer an LLM towards budget preferences, inject the $\cos t$ steering vector vcost with a negative strength factor $\operatorname{dcost} < 0$ into the LLM's residual stream while it processes an input user query.

4.2 Evaluation of Expressed Preferences

We ground our evaluation of the *amount of expressed preference* in the LLM's outputs using real-world human written data extracted from the Yelp Review Dataset⁷. For each preference, we extract a corpus of reviews that embody the bidirectional trait based on the criteria described in Table 1. Each corpus is filtered for quality, leaving between 700-900 reviews per trait to create reference evaluation datasets. See Appendix C for examples. Yelp reviews are viable for evaluation as they reflect genuine traces of human preferences, contain relevant metadata to filter for the preferences, and are closely related to our domain of lifestyle tasks.

We operationalize the effect of expressed preference as the relative cosine similarity between the LLM output's embeddings and each of the bidirectional traits' reference Yelp review dataset. *Effect* = $\cos(e_o, \bar{\mathbf{e}}_+) - \cos(e_o, \bar{\mathbf{e}}_-)$, where e_o is the BERT embeddings of the LLM's output, \bar{e}_+ is the mean BERT embeddings of the positive trait's reference dataset, and \bar{e}_- is that of the negative trait's reference dataset. All embeddings are computed using a pretrained Sentence-BERT model [54] (*sentence-transformers/stsb-roberta-base-v2* on Huggingface⁸), which captures sentence-level semantics appropriate for lengthy reviews and outputs. A positive value of the relative cosine similarity means higher expression of the positive steering trait, and vice versa for the negative steering trait. Similarity scores, rather than hard classification, allows the evaluation to capture more nuanced intensities of preference expressions.

Furthermore, as the magnitude of steering increases, we expect to see a degeneracy in the quality of the LLM's outputs. To evaluate this, we compute the perplexity-normalized effect (PNE), which incorporates a measure of how the model's perplexity, a proxy for degeneration, changes with steering compared to a non-steered baseline. This equation is given by von Rütte et al. as:

 $(\textit{Effect}_{steer=d} - \textit{Effect}_{steer=0})/(\textit{PPL}_{steer=d}/\textit{PPL}_{steer=0})$, where PPL is perplexity and d is the steering strength. If perplexity increases disproportionately more than the preference effect, then the PNE will be low, reflecting reduced quality.

4.3 Overview of Computational Experiments

We describe a set of four computations experiments **E1-E4** that each aim to answer a different question. In most of the experiments (other than **E3**), we simplify the personalization problem to focus on one preference dimension at a time as steering towards multiple preferences simultaneously has higher noise and unpredictable effects. The experiments are summarized as follows:

- **(E1)** How does steering strength affect the preferences expressed? We first evaluate the effect of bidirectional, granular steering strengths on the LLM's outputs in comparison to unsteered responses.
- (E2) How does prompting interact with steering? In conversations, users can inject personal preferences through prompts, so we investigate how prompt-based preferences interacts with steering. In this paradigm, steering can be thought of as aligning with the user's underlying preference profile, and the prompting adds contextualized variability for example, "I want to treat myself!" should yield different responses for luxury-seeking vs budget-oriented users if the LLM is correctly primed with their latent preferences.
- (E3) Can multiple preferences be steered? To demonstrate generalization beyond a singular dimension, we evaluate the effects of steering when simultaneously steering towards two compounded preferences. This is a preliminary step towards more complex steering scenarios with even more layered preferences.
- (E4) Can hidden preferences be learned? Lastly, we prototype how steering can be adaptively learned and applied in a conversation with a synthetic user, who has a hidden preference that is unknown to the model. We use GPT-40-mini model to emulate a user persona with a hidden preference *h* and allow it to respond with sentiment-driven feedback. More details about the simulation in Appendix G. This experiment is further validated through the real world user study in Section 5.

4.4 Computational Experiments Results

(E1) Effect of Steering on Content and Quality. We apply steering for all models and preference dimensions with strength factors between -30 < d < 30, measuring the mean preference effect across the LLM's responses to the OASST2 query dataset. Figure 2 shows the standard (*top*) and perplexity-normalized (*bottom*) preference effects for all models and preferences. The values are standardized to center around d = 0 as the neutral point with *effect=0*.

Under ideal circumstances, we expect to see a linear increase in preference expression corresponding to the increasing steering strength. We do observe near-linear relationships between steering strength and perplexity-normalized preference effects within the *functional steering range* of the model — we coin this term to describe the approximate zone in which steering does not significantly

 $^{^6}https://hugging face.co/datasets/OpenAssistant/oasst2$

⁷https://business.yelp.com/data/resources/open-dataset/

 $^{^8} https://hugging face.co/sentence-transformers/stsb-roberta-base-v2$

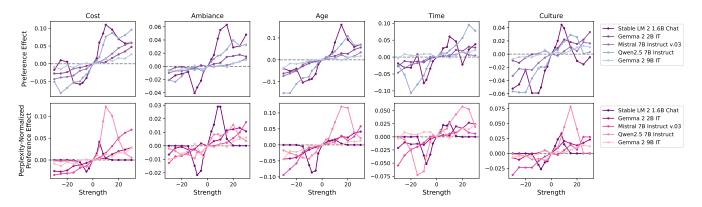


Figure 2: Effect of expressed preferences (top) and perplexity-normalized effect (bottom) for all preferences and models for E1.

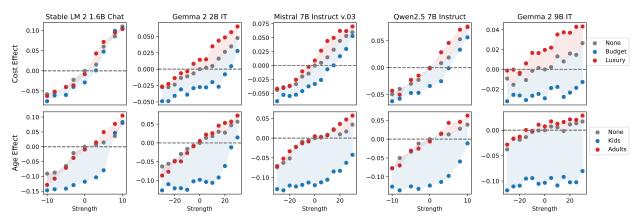


Figure 3: Effect of additional prompting on steering for two preferences: cost (top) and age (bottom) for E2.

Table 2: Samples of steered outputs with preference-based prompting for cost (E2).

Cost steering	No prompting	Prompt towards luxury	Prompt towards budget
"Average User" (d=0)	Here's a list of essential items you might want to consider for your camping trip. 1) Tent: Depending on the size of your family, you'll need a tent that can comfortably accommodate everyone	Here's a comprehensive list of camping supplies for your family camping holiday. 1) Tent: A spacious and high-quality tent is essential. Consider a tent with multiple rooms to provide privacy and separate areas for each family member	Here's a list of essential camping supplies to help you get started while keeping it budget-friendly. 1) Tent: A 4-person dome tent is a popular and affordable choice
"Wealthy User" (d=20)	To ensure your comfort and convenience, I have prepared a comprehensive list of essentials. 1) Tent: A four-season, safaristyle tent with a large canopy	To ensure that your stay is nothing less than perfection, I have prepared a meticulously curated list of the most exquisite and indulgent camping supplies	To ensure a comfortable and enjoyable adventure, I have compiled a list of essential items that will cater to your needs without exceeding your budget. 1 Tent: A spacious 4-person dome tent with a waterproof flysheet and a mesh canopy for stargazing
"Budget User" (d=-20)	Tent: Depending on your camping style, you may need a tent. There are various types of tents available, such as dome tents, backpacking tents, and family tents	Here are some essential items to consider when planning a camping trip: 1) Tent: There are various types of tents available, from simple pop-up tents to more complex backpacking tents	Here are some budget-friendly supplies for a camping holiday: 1) Tent: Consider buying used tents or borrowing from friends

degenerate the LLM's generations. For models that are highly sensitive to steering, like *stablelm-21-6b-chat* and *Qwen2.5-7B-Instruct*, this range is smaller than for the other models (see Appendix A). For all ensuing experiments, we only steer within the functional range of the LLM to preserve quality. To summarize, the main finding of **E1** is that *steering, in controlled moderation, can effectively guide the preference expressed by LLMs*.

(E2) Effect of Prompting on Steering. We repeat the same experiment as **E1**, but with the addition of preference-based prompts to each question in the query dataset. For each preference dimension, we augment the queries with prompting for both the positive and

negative traits, such as "I am morning-oriented and an early riser" for the morning trait and "I am evening-oriented and a night owl" for the night trait. The objective is to understand how steering and prompting interact, especially what happens to the expressed preferences if the directions contradict or reinforce each other.

Figure 3 depicts how adding preference-based prompts induces an *offset* on the preferences expressed by the LLMs. Red represents prompting towards the positive trait, inducing a mostly positive offset; and blue represents the negative trait, inducing mostly a negative offset. We limit the main results to two dimensions (**cost** and **age**) for brevity, but the other dimension graphs can be

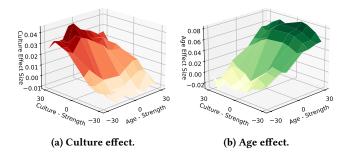


Figure 4: Preference effects for culture and age when both are steered additively for *Mistral-7B-Instruct-v0.3* in E3.

found in Appendix D. Using the two-sample Kolmogorov-Smirnov test for significance (p < 0.05), we find that for \mathbf{cost} , only gemma-2-9b-it is significantly affected by prompting, while other models are not. We hypothesize this could be attributed to gemma-2-9b-it's larger size and robustness against steering. Furthermore, prompting kids has a significant effect in all models while prompting adults does not — this suggests that some preferences are asymmetrical (likely because the average user prefers adult-oriented suggestions).

The significance of these results, while variable across models, indicate that *prompting can intensify or reduce the effect of preference expressions on top of steering*. We frame this as a practical trait for a steerable chatbot — steering to the user's underlying preference *contextualizes* the overall conversation, and adding prompting introduces variability within the context. See Table 2 for qualitative examples that demonstrate contextualized outputs. For this reason, we use *gemma-2-9b-it* as the base model of the user study due to its sensitivity to prompting, which allows the model to respond to user prompts in addition to steering.

(E3) Steering Towards Multiple Preferences. As this is an exploration of multi-steering, we only focus on two dimensions instead of all at once. In future context, this should be extended to incorporate more preferences, ideally ones that are are orthogonal to each other. Since our preferences are not inherently orthogonal, we select two dimensions with the least collinearity, which are **culture** and **age** (see correlation analysis in Figure E.1 in Appendix E). We vary the steering strength for each dimension uniformly and measure the effect of each preference separately. The steering vector is applied as a weighed combination of the individual preference vectors.

With compounded steering, the main concern is that the effects will deteriorate. Figure 4 shows 3D plots of the variation of the preference effect for **culture** and **age** for *Mistral-7B-Instruct-v0.3*. While the range of effects is slightly degraded at large values of complementary steering, the surfaces follow the expected trend. Surface plots for all other models are in Figure E.2 in Appendix E. Qualitative samples of multi-steered texts are shown in Table E.2 in Appendix E, which highlights phrases that exemplify each of the target traits. The results of **E3** suggest that **multiple preferences can be steered additively**, which is promising for scaling up steering as a personalization framework.

(E4) Learning Underlying Preferences and Steering. Lastly, we evaluate the proposed preference learning algorithm in simulated user conversations, where GPT-40-mini (GPT) roleplays as the

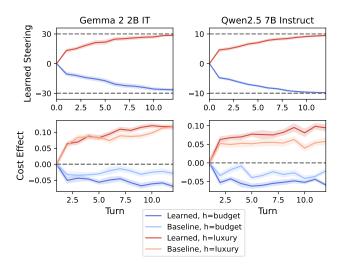


Figure 5: Results for gemma-2-2b-it and Qwen2.5-7B-Instruct for learning preferences in E5, with estimated preference steering strength (top) and corresponding cost effect sizes (bottom). Blue represents trials where the latent preference h is 100% budget, and red represents 100% luxury.

user with a hidden preference. This experiments acts as a precursor validation step before implementing LEARN for the user study. For ease of experimentation, we conduct **E4** only with the preference of **cost**, setting the hidden preference values to extreme levels of $d_h = -100$ (fully budget) and $d_h = 100$ (fully luxury). GPT is prompted with in-context examples of what the user's preference at these levels, and is instructed to output a response to the LLM's output with the 1) level of satisfaction that the user would feel, and 2) the direction of change (cheaper or more luxurious) that the user would want. See Appendix G for the system prompts to GPT.

As GPT is limited in its ability to independently simulate diverse replies emulative of real users, we do not allow it to control the conversation. Instead, we append to its output a new task question from the query dataset at every round, which essentially simulates a dynamic conversation with user feedback covering multiple topics. On the LLM side, it is first initialized with no steering, then the update step from Equation 1 is applied to calculate the direction of the update and the step size proportional to the dissatisfaction indicated by GPT. We run each trial for 12 rounds and compute the preference effect of the LLM's generation for **cost** at each turn.

Figure 5 shows the learned strength levels and preference effects for **cost** for *gemma-2-2b-it* and *Qwen2.5-7B-Instruct*, while other models are in Figure F.1 in Appendix F. It is evident that the LLMS were able to converge towards the correct steering strengths given the hidden preference of the user and the model's own functional steering range. The effect of the model with learned steering is compared to one without steering (Baseline). The gaps between Learned and Baseline of the models shown are significant by the Kolmogorov-Smirnov test (p < 0.05), except gemma-2-2b-it with h = luxury. This shows that the the learned steering improved the expression of preference compared to a prompt-only baseline.

It is important to note that both the Learned and Baseline versions of the models receive feedback from the *GPT-as-user* model

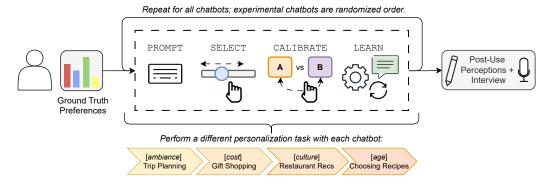


Figure 6: Overview of the within-subjects user study procedure, where participants completed a different personalization task using each of the four interfaces. Participants recorded their true preferences before the tasks and their perceptions afterwards.

and have access to the conversation history, which means that even the unsteered model can adjust its outputs to reflect the prompts of the user. We hypothesize that the models that were significantly affected by prompting in E2, e.g. <code>gemma-2-9b-it</code>, would be more adaptive through prompting alone, which is reflected in a lack of significance between Learned and Baseline (see Figure F.1). Nevertheless, the <code>learned steering approach demonstrates potential to dynamically adjust steering based on the sentiment of user feedback</code>. This allows for <code>light-weight</code>, <code>inference time personalization</code> and is especially applicable to <code>cold-start</code> scenarios.

5 User Study

Through the computational experiments, we demonstrate the effectiveness of steering, its interactions with prompting, and the viability of learning hidden preferences. Now, we investigate how real user interact with steerable chatbots and the effect of steering on real conversations. Since the steering parameter has the benefit of being an linear and interpretable value, we experiment with different designs for the interface and control of the steering. The steering strength is exposed to the user through three different interface designs: a) SELECT, which allows direct manipulation over steering, b) CALIBRATE, which presets the steering strength based on a set of calibration questions, and c) LEARN, which adjusts the steering based on in-conversation message sentiments. This section describes the user study procedures and results, both quantitative and qualitative.

5.1 User Study Procedure

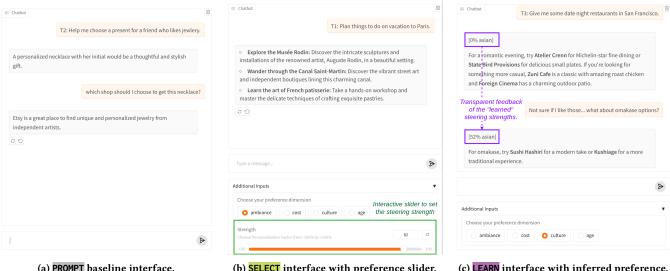
We conduct a within-subjects study where participants used all three of the experimental steering interface, as well as a baseline PROMPT chatbot that had no added steering. Figure 6 describes the high-level study procedure. Prior to starting the tasks, we collect participants' self-rated preferences for the given tasks and dimensions on a scale of 1 to 10, where 1 mapped to the maximum negative preference and 10 mapped to the maximum positive preference. This preference expression is then converted to a -100% to +100% scale for further analysis. They also self-reported demographics like age, gender, and baseline LLM usage.

Participants completed a different personalization task with each interface. While they always started with the baseline to help them understand the chatbot and the task, the order of the experimental interfaces were randomized. The assignment of tasks to each interface was also deliberately counter-balanced so that any effects from the tasks would be minimal. Participants were not given the names of the chatbots to avoid biasing, and they were described as simply A, B, C, and D.

We designed four personalization tasks based on the lifestyle planning questions evaluated in the computational experiments. Each task corresponds to a preference axis, which was chosen based on their relevance to the task and the ability of the LLM to output highly variable responses via steering — gift shopping (cost), vacation planning (ambiance), restaurant recommendations (culture), and choosing meal prep recipes (age). To initiate the conversation with the chatbot, participants clicked on a pre-written query, such that it is standardized across all trials. See the opening query of the four tasks with their respective preference dimensions in Table 1. Following the initial response from the chatbot, participants had freedom to ask follow-up questions or nudge the chatbot towards a different preference direction. They were instructed continue the conversation for several turns until they reached satisfactory responses, or until they exceeded 5 minutes in the task. If participants felt stuck or unguided in the conversation, the overseeing research team member gave them suggestions of follow-up questions to ask.

After completing the tasks, they rated their subjective perceptions on a 7-pt bipolar Likert scale for each of the four chatbots. We then conducted a semi-structured interview where participants could openly discuss the interfaces that they liked and disliked, as well as what features and designs they consider to be important for personalized chatbots. The perceptions Likert questions are:

- Likelihood to Use: How likely would you use the chatbot again (for personalization tasks)?
- Satisfaction: How satisfied were you by the personalization of each chatbot?
- Perceived Control: How in control did you feel of the personalization done by each chatbot?
- Perceived Persistency: How well do you expect your preferences to be remembered by the chatbot for future use?



(a) PROMPT baseline interface.

(b) **SELECT** interface with preference slider.

(c) LEARN interface with inferred preference.

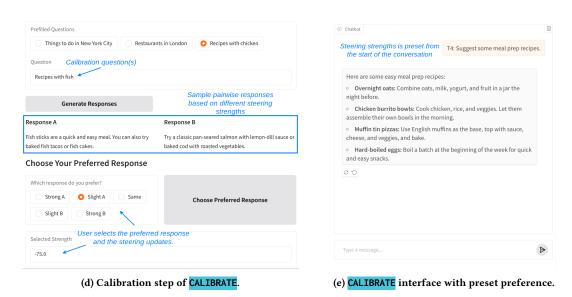


Figure 7: Screen captures of all four chatbot interfaces that participants used.

Steerable Chatbots Implementation

We implement all four chatbots using Gradio, an open-source Python package for prototyping visual interfaces of machine learning models [1]. Participants were provided with the temporary public link generated by Gradio to interface with the chatbot, while the backend LLM (gemma-2-9b-it) is hosted on a Google Colab running an NVIDIA L4 Tensor Core GPU. We chose gemma-2-9b-it because of its ability to accommodate both prompting and steering, such that participants can nudge the LLM's output through text even when steering is applied. Being the largest model, it also has high quality responses and good stability when steered. We standardize the preference strength range from -100 to 100 since this is

what the participants were shown, but this is mapped to gemma-2-9b-it's effective steering range in the backend. See Figure 7 for screen captures of each of the chatbot interfaces and the following descriptions of their differences:

PROMPT: The baseline interface (Figure 7a) has no added steering, so any personalization of the LLM is controlled through userinitiated prompting only.

SELECT: The steering parameter is controlled with a slider positioned at the bottom of the chat interface with a range of -100 to 100 (Figure 7b). Participants are briefed that the selected slider value is instantaneously applied to the chatbot, and they can dynamically adjust it throughout the conversation.

CALIBRATE: Prior to entering the chatbot interface (Figure 7e), users perform a calibration step to initialize the steering strength that is applied to the conversation (Figure 7d). During calibration, pairwise outputs (Response **A** and **B**) for pre-written calibration questions are generated by sampling two different steering strengths, which initially set with disparate values of $d_A = -100$ and $d_B = 100$. Based on the user's preferred response, the steering strengths are updated to be closer to the preferred steering strength. For example, updating $d_A \leftarrow (d_A + d_B)/2$ to be closer to d_B if the preference is Slightly B. This calibration step is repeated 2-3 times until rough convergence, when **A** and **B** are very similar. The final steering strength d is computed as the average of d_A and d_B .

LEARN: Similar to the baseline interface, the user can only control the personalization via prompting (Figure 7c). The steering strength is updated based on the algorithm described in Equation 1 that accounts for the sentiment of the user's latest message and the preference direction that they wish to steer towards. For example, "I don't like these gift options, I'm looking for something more affordable" would update the cost steering negatively, towards budget. See Section 3 for technical implementation details and Section 4.4 for preliminary validation results of this algorithm. In addition to the LLM's output, the learned steering strength is also displayed as a percentage (e.g. 35% budget) in the output of the LLM for added transparency.

5.3 User Study Results

We recruit 14 participants from the USA and Canada (women=7, men=4, non-binary=1) with a varied distribution in LLM use (daily=7, weekly=5, and monthly/occasionally=2). The study took a maximum of 60 minutes, for which participants were compensated with 25 USD gift vouchers. Figure H.1 in Appendix H shows the ground truth, self-reported preferences of the participants on each of the four tasks and preference dimensions. For ambiance and cost, the preferences were well-distributed; for culture and age, the preferences were biased towards the negative and positive steering dimensions, respectively. This provides us with coverage over different tasks reflecting real-world preferences. In the following analysis, we examine:

- (U1) Effect of steering on preferences expressions in real conversations;
- (U2) Alignment of expressed preferences with users' underlying preferences;
- (U3) User satisfaction and perceptions of the different interfaces;
- **(U4)** Qualitative analysis of user's values for steerable chatbots.

(U1) Effect of Steering on Expressed Preferences. Can steering control the amount of expressed preferences in real conversations? As an extension to E1 and E2, we examine the relationship between steering strengths and the chatbot's expressed preference in real user conversations. It is important to note that our participants frequently injected their preferences via prompting to the chatbots, so we expect the results here to be noisier than in the controlled computational experiments. Figure 8a depicts the correlation for each steerable chatbot separately, as well as all three aggregated together, where each data point is the preference effect (relative cosine similarity to the positive/negative traits) of a single message sent by the LLM.

Pearson correlation analysis show that SELECT (r = 0.54, p <0.001) and LEARN (r = 0.40, p < 0.001) both demonstrate significant positive relationships between the steering strengths and the expressed preference, matching the results from the computational experiments. However, CALIBRATE did not have this result (r = 0.04, p = 0.78). We hypothesize that this may be affected by failures of the calibration step (which concluded in 2-3 rounds due to time constraints) to capture robust and nuanced latent preferences. For example, **P2**'s calibration result landed him at cost = -75, which reflects a highly budget trait, but in-conversation he repeatedly prompted for more luxury options like, "can you suggest some diamond studs" and "how about reputable or designer options", demonstrating an incongruence with the calibration. The other two chatbots, on the other hand, could be either directly controlled via the slider or indirectly via prompting, so it is less likely for the participant to be disatisfied with the steering strength.

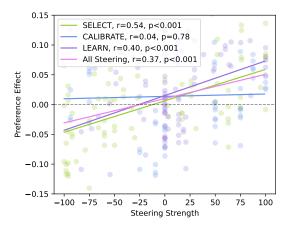
Overall, all three steerable chatbots combined reached significance (r = 0.37, p < 0.001). Optimistically, these results have two implications. One, *steering is effective at controlling the expressed preference* when the user is relatively satisfied with the personalization. Two, when the user is not satisfied, *prompting opposing preferences can still be effective at retaining control* over the conversation (at least in the context of *gemma-2-9b-it*).

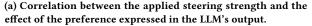
(U2) Alignment with Underlying User Preferences. After establishing that steering can guide preference expressions in real world conversations, the natural follow-up is to understand if the steering actually match people's underlying preferences. This question was partially modeled through E5 through the simulated GPT user, but that setup did not capture complex user preferences well.

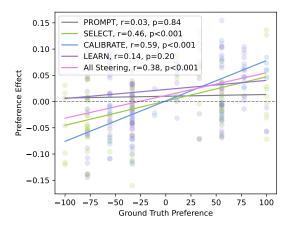
In Figure 8b, we plot the correlation between the participant's self-recorded ground truth preferences and the preference effect expressed by the LLM per message. Pearson correlation analysis shows that SELECT (r=0.46,p<0.001) and CALIBRATE (r=0.59,p<0.001) have significant correlations — that the preferences expressed by the LLM in conversation corresponds moderately well with the true preferences of the user. LEARN (r=0.14,p=0.20) did not reach a significant result, and PROMPT was even worse (r=0.03,p=0.84). We hypothesize that the sentiment-based learning algorithm may not have robust enough to adapt to diverse verbiage in users' feedback and did not capture their intentions well. Overall, all steerable chatbots combined shows that steering achieves significant improvement in generating relevant preference content than prompting alone (r=0.38,p<0.001).

In addition, we examine if the steering strengths themselves match the ground truth preferences rated by the participants. For SELECT, we take the steering strength set by the participant in the first step; for CALIBRATE, we use post-calibration steering strength; and for LEARN, we use the learned strength in the last turn of the conversation. Figure 9 compares the chatbots across two measures – do the steering strengths directionally agree with ground truth preferences (e.g. is the steering value *positive* if the user indicates they prefer *luxury* for **cost**?); and what is the absolute error between the steering value and ground truth for the given task and preference?

Keeping in mind the caveat that true preferences are difficult to assess, even from the user themselves, and we observed many







(b) Correlation between the participant's ground truth preference and the effect of the preference expressed in the LLM's output.

Figure 8: Correlation analyses with the expressed preferences in the user study conversations. These answer the questions of (a) how applying steering influence the preference content of the LLMs' output; and (b) how well the preference content of the conversation correlates with the user's true preferences.

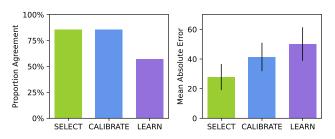


Figure 9: Statistics on how well each steerable chatbot matched the participants' preference with proportion agreement of preferences (left) and mean absolute error of the steering strengths (right).

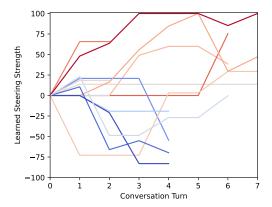


Figure 10: Progression of strengths learned by LEARN across all conversations. This is comparable to the synthesized learning experiment E5, shown in Figure 5.

participants changing their minds throughout the task. SELECT had the lowest mean absolute error (27.8 ± 31.9) and highest proportion of agreement at 86% — althougth surprisingly it is not 100%, as participants chose it themselves. CALIBRATE was able to match the agreement rate but had higher error (41.4 ± 34.6) . LEARN was less robust and achieved slightly above random rate of agreement at 57% and highest error (50.1 ± 40.9) . To highlight what the learning process for the latter looks like, Figure 10 show the progression of the learned steering strengths across all 14 participant conversations, with red indicating more positive strengths and blue indicating more negative. In combination with the correlation analysis results, this suggest that if the steering factor is relatively accurate, the chatbot's expressed preferences can align well with users' hidden preferences. The three designs differ in how well they capture true preferences, but this demonstrate high promise.

(U3) Subjective Perceptions. Next, we evaluate participants' post-task perceptions of each of the chatbots. Post-task perceptions in

Figure 11 reveal that despite SELECT achieving the strongest alignment to preferences, it received roughly equal ratings as PROMPT across all perception categories. Only CALIBRATE received significantly higher ratings in the Satisfaction category (p=.03) and Perceived Persistency category (p=0.005). This is where we uncover the *marked heterogeneity in people's preferences about personalization*. The qualitative analysis in the next section will examine the themes in more detail, but most participants unveil that they have a strongly preferred interface. If we operationalize the "favourite" interface as the one that received the highest rating in the Likelihood to Use question (this results in a count of 6 votes for SELECT, 5 for LEARN, and 3 for CALIBRATE), then we observe highly significant differences between PROMPT and Favourite across all perception categories. While individual chatbots received mixed

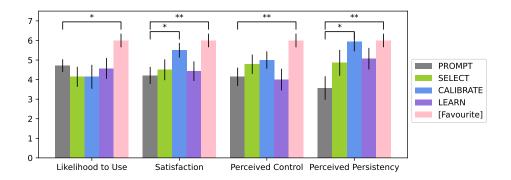


Figure 11: Post-experiment perceptions results for all four interfaces. Due to high heterogeneity in the participant's preferred interface, we also compute differences based on each participant's favourite interface. Statistically significant difference with respect to PROMPT are indicated with * (p < 0.05) or ** (p < 0.005).

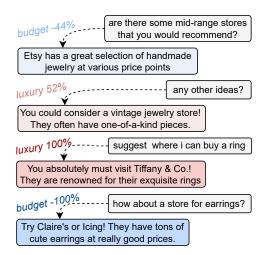


Figure 12: An abbreviated user conversation for the *Gift Shopping* task, where the user dynamically changed the steering strength using the SELECT interface. Note that the LLM adapts to the steering inputs to match the cost preference without requiring prompt-based conditioning.

feedback on average, this result indicates that participants overall *prefer some version of a steerable chatbot* over the PROMPT baseline, particularly in categories of **Likelihood to Use** and **Satisfaction**. Interestingly, we also find that participants perceived more **Control** and **Persistency** of personalization in their favourite chatbots, regardless of if it is actually true.

(U4) Qualitative Thematic Analysis. We conduct thematic analysis of the post-task interview using an inductive coding approach [17]. Two researchers coded quotes extracted from the interview transcripts, and refined the codes into agreed upon themes. We discuss the qualitative findings as a way to understand participant's heterogenous views towards the different steering interfaces.

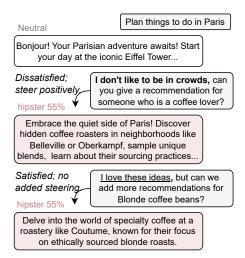


Figure 13: An abbreviated user conversation for the *Travel Planning* task with ambiance preference, where LEARN interface uses the sentiment of the user's message (*bold* for negative sentiments, <u>underline</u> for positive sentiments). The steering parameter adjusts to be *hipster* when the user implies they don't like *touristy* options.

Control over personalization. Participants who valued control tend to prefer SELECT, which provided them with direct manipulation and feedback of the steering parameter. "I was more in control of the kind of answers I wanted to see" said P2, and P4 agreed with "whether you felt like you were in control was extremely important". Some participants like P9 valued the ability to update their preference, saying it was "fun to just see that immediate response", while others like P3 believed their preference "is probably set for the entire conversation... once I understand where in the slider I am". Others also observed the benefit of using the slider to set the context of the LLM's outputs, while still being able to add personalization via prompting (a callback to E2). P4 observes, "I can control this dimension [via the slider] and then I can tune the other parts that I

want via conversation", and P6 echoes that it was like "'drawing a box' around the personalization factors that I want to explore". See Figure 12 for an example conversation where the user dynamically changed the steering strength throughout the conversation.

With the other interfaces, people indicated mixed feelings with respect to control. P4 mentioned she could control her choices in CALIBRATE, because "I wanted it to learn about me in a certain way". However, P14 felt like "the settings are sort of locked once you submit them". P12 and P14 both expressed that they believe that LEARN would allow them to "refine the settings over time" through conversation. More generally, P9 asserts that "if [a chatbot] was making a false inference that was going to lead to more inaccurate outputs, I would want to be able to stop that in its tracks" — this sentiment highlights the need that many felt — if the personalization is not helpful or accurate, there should be a way to change it.

Ease of personalization. We find varied perspectives on what participants consider efficient and easy to use, especially around the trade-offs between effort and expressivity of prompting-based methods. P2 said he could have asked "more leading questions to [PROMPT]... but it felt more like it felt too much effort", while others like P12 believed prompting-based techniques are more straightforward and expressive, "let me just dump my thoughts and let the LLM handle what's going on" and P14 agrees with saying he prefers to "adjust my preference through multi-turn prompting". In terms of the user experience, P8 liked LEARN, saying "learning from conversations feels a little bit more flowy and natural", and P9 echoes with "it was more seamless and I felt like I could best express what I was looking for in more detail". See Figure 13 for an example of how the learned steering strength is updated throughout a user conversation.. However, P3 believed it "needs prompt engineering, so I need to be really good at what I ask so that it understands". As an alternative to prompting, P4 liked SELECT to "alleviates me from repetitive prompting for the same intention". These findings indicate that a user's baseline level of comfort with expressing themselves in natural language might dictate their perception of the usability

With respect to the added step in CALIBRATE, P14 reflected that "the amount of setup is a little bit overwhelming" and P9 didn't see the payoff of the extra step, "it felt like I was doubling my effort or just kind of extra work for no immediate benefit". However, some pointed to the fact that the calibration can reduce effort later on. P3 said "I've already personalized it up to some extent with my liking" and P8 liked "having the [personalization] set beforehand...that way I don't have to think about it and have a back-and-forth conversation". Some participants showed a positive perception of the precision of the calibration process, with P10 saying "it kind of started to narrow down... so I feel like that was understanding very detailed preferences" and P7 comparing it to an eye doctor appointment, "is it option A or option B?.. eventually they get to your prescription". Since our correlation analysis in U2 indicated that CALIBRATE matched user preferences the best, we suggest that it can be evaluated further in longer usage sessions to understand if users perceive the added effort to be worthwhile. Between all the interfaces, participants had vastly different perceptions of which interaction modality is the easiest to personalize.

Transparency and explainability of personalization. Another core value expressed by participants is the transparency of the personalization - whether if the personalization is communicated and interpretable to the user. Participants were somewhat split over the transparency feature of LEARN, where the learned steering strength is exposed in the conversation. P12 found the transparency interesting, "I liked being able to see how the model was adjusting", and P14 said it helped him understand the "chatbot's mental model of the user's preferences". On the other hand, P4 expressed feeling judged by the learning algorithms, especially "if it thinks that I have a certain preference but that's not what I intended". P11 mentions liking being exposed to the steering factor transparently "so I could adjust it after the fact", and P13 agrees with wanting to edit "what a chatbot knows about me". The theme that emerges here is that participants generally want to understand what an algorithm **learns about them**, and be able to modify the information stored.

A more general issue that participants brought up was the interpretability of the preferences. While using LEARN, P11 said the preference percentage "was overloading me...honestly I was ignoring it". Conversely for SELECT, P5 believed that having users decide the steering is not effective because "it gets into the nuance of how [the chatbot] understands those [preferences] and how you understand them", and P11 believed that "people might not know how to answer it correctly for more complex dimensions". The strength of CALIBRATE is clear here, with P5 reflecting "instead of just picking a word you and an AI are defining probably differently, it can ask you questions". More effort can be dedicated to communicating the preference strengths in an interpretable way, such as with providing more in-context examples.

Privacy concerns of personalization. With any personalization, the added risks of privacy risks can pose a tradeoff. Participants debated over what information they believe are relevant; P8 was wary of exposing personally-identifiable data of "name, location, and age... and things like that" but added that it would be "useful for it to know what your age is, where your culture is". P14 contextualizes the concern on privacy with "if all of my preferences are getting linked together to describe me as a person, that could be a privacy issue...if they're isolated, maybe it's less of a problem". The level of risk that the user is willing to take is important to identify, as this can determine their comfort with the data used by personalization.

Persistency of personalization. While persistent personalization is key for many, like P3 who said the baseline PROMPT interface felt like having "a random conversation with a stranger, like a barista"; others like P12 didn't see the benefit, saying "I don't want something that I asked a month ago to be relevant to the conversation I have right now". Preferences can be conditional or transient, P14 believes that "people's thoughts can change all of a sudden", even throughout the same conversation; and P6 who acknowledges that "preferences tend to change over time". Similarly, P13 notes that she only wants "personalization for certain things, but not others". Many participants even noted that their in-task preferences diverged from 'ground truth' preferences they reported in the pre-task survey. Overall, we find divergent opinions on what information about the user should be retained for future conversations, with many indicating desire for flexibility for how and when personalization is applied.

6 Discussion

Key Findings. In this work, we implement LLM personalization through a steering-based framework, where relevant preferences for a given task can be amplified or dampened through applying precomputed steering vectors. Our computational experiments first highlight the generalizability of steering across different models and preference dimensions, steering in conjunction with prompting, and multi-preference steering. These results are further validated in the user study, where we find that steering is effective at expressing preferences in real conversations and that steered chatbots align better with hidden user preferences than a promptingonly baseline. Lastly, we compare three diverse designs of steerable chatbots, including an interactive learning algorithm that aims to uncover the user's latent preferences through sentiment-based feedback feedback. While the interfaces receive variable feedback based on users' personal values, participants preferred steered chatbots overall and rated higher satisfaction in their personalization.

Steering as a Personalization Framework. There are additional advantages of using steering to drive personalization. Since steering only involves the inference time application of a pre-computed vector, it is incredibly resource-efficient in comparison to methods that requiring updating model weights or storing a extensive user data. Furthermore, the steering vector leverages semantic concepts that are already embedded in the LLMs internals, and does not require any additional model training, making the overhead of implementation very minimal. From the privacy perspective, since the steering vector and strengths contain no personally-identifiable data (unlike personalized finetuning weights or explicit user profiles), preference-based steering is also more sensitive to user privacy concerns. Users can be represented through a set of steering strengths across various dimensions without requiring any personal details. We also demonstrate through the three chatbot designs that steering can be applied with high flexibility – persistently set as in the case of CALIBRATE, or dynamically adjustable as in the case of SELECT, which opens up the possibility of more elaborate user-centered designs, including AR/VR applications.

Design Extensions of Steering. We implemented several steerable chatbots, but steering can be applied flexibly to a number of additional interaction paradigms beyond the ones we covered. We propose some ideas here but leave them as future work:

- Learning from user history: In scenarios where the user
 has a long history of conversations and it can be determined
 which LLM outputs they liked and disliked, then it could be
 possible to construct a custom steering vector unique to the
 user, capturing the nuances of their preferences.
- **Dynamic and conditional steering:** As real user preferences can be dynamic or conditional (e.g. someone who prefers budget *except* when it comes to food), the personalization can be modified to incorporate conditional steering or preference updating rules [35].
- Automated preference detection: Since steering does not require additional model training, it is possible to detect preference and compute their steering vectors *in-conversation*, such as leveraging additional LLMs to perform this. This

method would alleviate the drawback of the current approach, which relies on pre-determined preferences.

Personal Values and Interface Designs. While we did not have strong pre-conceived notions about which chatbot designs users would like the best, the diversity of user opinions was nevertheless surprising. Among the main themes that participants expressed in the interviews, values of *control*, *usability*, and *transparency* appeared to the main drivers of which chatbots they liked most. In particular, the need to understand how personalization is performed and retain control over what is learned and stored stands out, and this is reflected in the findings of prior research on personalization interfaces [6, 32, 45]. Since steering via a linear strength factor represents a relatively unique way of control personalization, there is the space to explore other interface designs, including some of which that we discuss above. Ultimately, a combination of learning, calibration, and direct control may superior, allowing users to pick and choose which features they want to engage with.

Limitations. We conduct our computational experiments and user study using fixed values of LLM hyperparameters, so our results should be interpreted within these boundaries. While there are many different methods and parameters for implementing steering, we only explore a subset of design choices (primarily as outlined by von Rütte et al.). We focus on one task domain and identify preferences based on their relevance to the task and their presence within the Yelp Dataset. The space of application for personalization is much broader, and can encompass tasks such as creative writing, technical communication, and persona adoption.

The implementation of the calibration algorithm and the learning algorithm were meant to convey proof-of-concept capabilities and did not reflect state-of-art performance of, for example, a finetuned intent classifier or a likelihood-maximizing sampling method. In particular, this limited the robustness and precision of the learned preferences, which then impeded user satisfaction in the applied steering. As such, we asked the participants to concentrate on the interaction modality of the interfaces in their post-task interviews rather than their reflection on the personalization quality.

7 Conclusion

We demonstrate how activation steering can be applied towards personalizing LLMs in preference-driven tasks. Through a series of robust computational experiments, we quantify the effects of steering across five LLM models and five preference dimensions. We then develop three interaction modalities for *steerable chatbots* and compare their results at personalization and how they are perceived by end users. Overall, we recommend further exploration of how preference-based steering can be leveraged as a framework enabling resource-efficient, cold-start LLM personalization.

References

- Abubakar Abid, Ali Abdalla, Ali Abid, Dawood Khan, Abdulrahman Alfozan, and James Zou. 2019. Gradio: Hassle-free sharing and testing of ml models in the wild. arXiv preprint arXiv:1906.02569 (2019).
- [2] Christopher M Ackerman. 2024. Representation Tuning. arXiv [cs.LG] (Sept. 2024).
- [3] Jae-wook Ahn, Peter Brusilovsky, Jonathan Grady, Daqing He, and Sue Yeon Syn. 2007. Open user profiles for adaptive news systems: help or harm?. In Proceedings of the 16th international conference on World Wide Web. 11–20.
- [4] Rumi A Allbert and James K Wiles. 2024. Identifying and manipulating personality traits in LLMs through activation engineering. arXiv [cs.CL] (Dec. 2024).
- [5] Andy Arditi, Oscar Obeso, Aaquib Syed, Daniel Paleka, Nina Panickssery, Wes Gurnee, and Neel Nanda. 2024. Refusal in language models is mediated by a single direction. arXiv [cs.LG] (June 2024).
- [6] Fedor Bakalov, Marie-Jean Meurs, Birgitta König-Ries, Bahar Sateli, René Witte, Greg Butler, and Adrian Tsang. 2013. An approach to controlling user models and personalization effects in recommender systems. In Proceedings of the 2013 international conference on Intelligent user interfaces. 49–56.
- [7] F Barbieri, J Camacho-Collados, L Neves, and L Tweeteval Espinosa-Anke. 2020. Unified benchmark and comparative evaluation for tweet classification. arXiv 2020. arXiv preprint arXiv:2010.12421 (2020).
- [8] Collin Burns, Haotian Ye, Dan Klein, and Jacob Steinhardt. 2022. Discovering latent knowledge in language models without supervision. arXiv [cs.CL] (Dec. 2022).
- [9] Yuanpu Cao, Tianrong Zhang, Bochuan Cao, Ziyi Yin, Lu Lin, Fenglong Ma, and Jinghui Chen. 2024. Personalized steering of Large Language Models: Versatile steering vectors through bi-directional preference optimization. arXiv [cs.CL] (May 2024).
- [10] Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémy Scheurer, Javier Rando, Rachel Freedman, Tomasz Korbak, David Lindner, Pedro Freire, et al. 2023. Open problems and fundamental limitations of reinforcement learning from human feedback. arXiv preprint arXiv:2307.15217 (2023).
- [11] Chih-Ming Chen, Hahn-Ming Lee, and Ya-Hui Chen. 2005. Personalized elearning system using item response theory. Computers & Education 44, 3 (2005), 237–255.
- [12] Daiwei Chen, Yi Chen, Aniket Rege, Zhi Wang, and Ramya Korlakai Vinayak. [n. d.]. PAL: Sample-Efficient Personalized Reward Modeling for Pluralistic Alignment. In The Thirteenth International Conference on Learning Representations.
- [13] Jiangjie Chen, Xintao Wang, Rui Xu, Siyu Yuan, Yikai Zhang, Wei Shi, Jian Xie, Shuang Li, Ruihan Yang, Tinghui Zhu, Aili Chen, Nianqi Li, Lida Chen, Caiyu Hu, Siye Wu, Scott Ren, Ziquan Fu, and Yanghua Xiao. 2024. From persona to personalization: A survey on role-Playing Language Agents. arXiv [cs.CL] (April 2024)
- [14] Ruizhe Chen, Xiaotian Zhang, Meng Luo, Wenhao Chai, and Zuozhu Liu. 2024. Pad: Personalized alignment of llms at decoding-time. arXiv preprint arXiv:2410.04070 (2024).
- [15] John Joon Young Chung and Eytan Adar. 2023. Artinter: AI-powered Boundary Objects for Commissioning Visual Arts. In Proceedings of the 2023 ACM Designing Interactive Systems Conference. 1997–2018.
- [16] John Joon Young Chung and Eytan Adar. 2023. Promptpaint: Steering text-toimage generation through paint medium-like interactions. In Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology. 1–17.
- [17] Victoria Clarke and Virginia Braun. 2017. Thematic analysis. The journal of positive psychology 12, 3 (2017), 297–298.
- [18] Hai Dang, Lukas Mecke, and Daniel Buschek. 2022. Ganslider: How users control generative models for images using multiple sliders with and without feedforward information. In Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems. 1–15.
- [19] Bofei Gao, Feifan Song, Yibo Miao, Zefan Cai, Zhe Yang, Liang Chen, Helan Hu, Runxin Xu, Qingxiu Dong, Ce Zheng, Shanghaoran Quan, Wen Xiao, Ge Zhang, Daoguang Zan, Keming Lu, Bowen Yu, Dayiheng Liu, Zeyu Cui, Jian Yang, Lei Sha, Houfeng Wang, Zhifang Sui, Peiyi Wang, Tianyu Liu, and Baobao Chang. 2024. Towards a unified view of preference learning for large Language Models: A survey. arXiv [cs.CL] (Sept. 2024).
- [20] Ge Gao, Alexey Taymanov, Eduardo Salinas, Paul Mineiro, and Dipendra Misra. 2024. Aligning LLM agents by learning latent preference from user edits. arXiv [cs.CL] (April 2024).
- [21] Andrew D Gershoff, Ashesh Mukherjee, and Anirban Mukhopadhyay. 2008. What's not to like? Preference asymmetry in the false consensus effect. *Journal of Consumer Research* 35, 1 (2008), 119–125.
- [22] Davor Hafnar and Jure Demšar. 2024. Zero-shot reasoning: Personalized content generation without the cold start problem. arXiv [cs.AI] (Feb. 2024).
- [23] Zhankui He, Zhouhang Xie, Rahul Jha, Harald Steck, Dawen Liang, Yesu Feng, Bodhisattwa Prasad Majumder, Nathan Kallus, and Julian Mcauley. 2023. Large language models as zero-shot conversational recommenders. In Proceedings of the 32nd ACM International Conference on Information and Knowledge Management. ACM, New York, NY, USA.

- [24] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. arXiv preprint arXiv:2106.09685 (2021).
- [25] Joel Jang, Seungone Kim, Bill Yuchen Lin, Yizhong Wang, Jack Hessel, Luke Zettlemoyer, Hannaneh Hajishirzi, Yejin Choi, and Prithviraj Ammanabrolu. 2023. Personalized soups: Personalized large language model alignment via post-hoc parameter merging. arXiv [cs.CL] (Oct. 2023).
- [26] Peiling Jiang, Jude Rayan, Steven P Dow, and Haijun Xia. 2023. Graphologue: Exploring large language model responses with interactive diagrams. In Proceedings of the 36th annual ACM symposium on user interface software and technology. 1–20
- [27] Ruili Jiang, Kehai Chen, Xuefeng Bai, Zhixuan He, Juntao Li, Muyun Yang, Tiejun Zhao, Liqiang Nie, and Min Zhang. 2024. A survey on human preference learning for large language models. arXiv [cs.CL] (June 2024).
- [28] Wang-Cheng Kang, Jianmo Ni, Nikhil Mehta, Maheswaran Sathiamoorthy, Lichan Hong, Ed Chi, and Derek Zhiyuan Cheng. 2023. Do LLMs understand user preferences? Evaluating LLMs on user rating prediction. arXiv [cs.IR] (May 2023).
- [29] Ikram Karabila, Nossayba Darraz, Anas EL-Ansari, Nabil Alami, and Mostafa EL Mallahi. 2024. BERT-enhanced sentiment analysis for personalized ecommerce recommendations. *Multimedia Tools and Applications* 83, 19 (2024), 56463–56488.
- [30] Tae Soo Kim, Yoonjoo Lee, Minsuk Chang, and Juho Kim. 2023. Cells, generators, and lenses: Design framework for object-oriented interaction with large language models. In Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology. 1–18.
- [31] Yoonsu Kim, Jueon Lee, Seoyoung Kim, Jaehyuk Park, and Juho Kim. 2024. Understanding users' dissatisfaction with chatgpt responses: Types, resolving tactics, and the effect of knowledge level. In Proceedings of the 29th International Conference on Intelligent User Interfaces. 385–404.
- [32] Bart P Knijnenburg, Svetlin Bostandjiev, John O'Donovan, and Alfred Kobsa. 2012. Inspectability and control in social recommenders. In Proceedings of the sixth ACM conference on Recommender systems. 43–50.
- [33] Kai Konen, Sophie Jentzsch, Diaoulé Diallo, Peer Schütt, Oliver Bensch, Roxanne El Baff, Dominik Opitz, and Tobias Hecking. 2024. Style vectors for steering generative large language model. arXiv [cs.CL] (Feb. 2024).
- [34] Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, Sotiris Anagnostidis, Zhi Rui Tam, Keith Stevens, Abdullah Barhoum, Duc Nguyen, Oliver Stanley, Richárd Nagyfi, et al. 2024. Openassistant conversations-democratizing large language model alignment. Advances in Neural Information Processing Systems 36 (2024).
- [35] Bruce W Lee, Inkit Padhi, Karthikeyan Natesan Ramamurthy, Erik Miehling, Pierre Dognin, Manish Nagireddy, and Amit Dhurandhar. 2024. Programming refusal with Conditional Activation Steering. arXiv [cs.LG] (Sept. 2024).
- [36] 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.
- [37] Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2023. Inference-Time Intervention: Eliciting truthful answers from a language model. arXiv [cs.LG] (June 2023).
- [38] Xinyu Li, Zachary C Lipton, and Liu Leqi. 2024. Personalized language modeling from personalized Human Feedback. arXiv [cs.CL] (Feb. 2024).
- [39] Yuanchun Li, Hao Wen, Weijun Wang, Xiangyu Li, Yizhen Yuan, Guohong Liu, Jiacheng Liu, Wenxing Xu, Xiang Wang, Yi Sun, et al. 2024. Personal Ilm agents: Insights and survey about the capability, efficiency and security. arXiv preprint arXiv:2401.05459 (2024).
- [40] Sheng Liu, Haotian Ye, Lei Xing, and James Zou. 2023. In-context vectors: Making in context learning more effective and controllable through latent space steering. arXiv [cs.LG] (Nov. 2023).
- [41] Ryan Louie, Andy Coenen, Cheng Zhi Huang, Michael Terry, and Carrie J Cai. 2020. Novice-AI music co-creation via AI-steering tools for deep generative models. In Proceedings of the 2020 CHI conference on human factors in computing systems. 1–13.
- [42] Zhenyi Lu, Wei Wei, Xiaoye Qu, Xianling Mao, Dangyang Chen, and Jixiong Chen. 2023. MIRACLE: Towards personalized dialogue generation with latent-space multiple personal attribute control. arXiv [cs.CL] (Oct. 2023).
- [43] Qianou Ma, Weirui Peng, Hua Shen, Kenneth Koedinger, and Tongshuang Wu. 2024. What you say= what you want? Teaching humans to articulate requirements for LLMs. arXiv preprint arXiv:2409.08775 (2024).
- [44] Xiao Ma, Swaroop Mishra, Ariel Liu, Sophie Su, Jilin Chen, Chinmay Kulkarni, Heng-Tze Cheng, Quoc Le, and Ed Chi. 2023. Beyond ChatBots: ExploreLLM for structured thoughts and personalized model responses. arXiv [cs.HC] (Dec. 2023).
- [45] Sean M McNee, Shyong K Lam, Joseph A Konstan, and John Riedl. 2003. Interfaces for eliciting new user preferences in recommender systems. In *International Conference on User Modeling*. Springer, 178–187.
- [46] Sheshera Mysore, Zhuoran Lu, Mengting Wan, Longqi Yang, Steve Menezes, Tina Baghaee, Emmanuel Barajas Gonzalez, Jennifer Neville, and Tara Safavi. 2023.

- PEARL: Personalizing large language model writing assistants with generation-calibrated retrievers. *arXiv* [cs.CL] (Nov. 2023).
- [47] Jonas Oppenlaender, Rhema Linder, and Johanna Silvennoinen. 2024. Prompting AI art: An investigation into the creative skill of prompt engineering. *International Journal of Human–Computer Interaction* (2024), 1–23.
- [48] NA Osman, SA Mohd Noah, and M Darwich. 2019. Contextual sentiment based recommender system to provide recommendation in the electronic products domain. *International Journal of Machine Learning and Computing* 9, 4 (2019), 425–431
- [49] Nina Panickssery, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Matt Turner. 2023. Steering Llama 2 via Contrastive Activation Addition. arXiv [cs.CL] (Dec. 2023).
- [50] Yingzhe Peng, Xiaoting Qin, Zhiyang Zhang, Jue Zhang, Qingwei Lin, Xu Yang, Dongmei Zhang, Saravan Rajmohan, and Qi Zhang. 2024. Navigating the unknown: A chat-based collaborative interface for personalized exploratory tasks. arXiv [cs.HC] (Oct. 2024).
- [51] Silviu Pitis, Ziang Xiao, Nicolas Le Roux, and Alessandro Sordoni. 2024. Improving context-aware preference modeling for language models. Advances in Neural Information Processing Systems 37 (2024), 70793–70827.
- [52] Sriyash Poddar, Yanming Wan, Hamish Ivison, Abhishek Gupta, and Natasha Jaques. 2024. Personalizing reinforcement learning from Human Feedback with variational preference learning. arXiv [cs.LG] (Aug. 2024).
- [53] Alexandre Rame, Guillaume Couairon, Corentin Dancette, Jean-Baptiste Gaya, Mustafa Shukor, Laure Soulier, and Matthieu Cord. 2023. Rewarded soups: towards pareto-optimal alignment by interpolating weights fine-tuned on diverse rewards. Advances in Neural Information Processing Systems 36 (2023), 71095– 71134
- [54] Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics. http://arxiv.org/abs/1908.10084
- [55] Neil Rubens, Mehdi Elahi, Masashi Sugiyama, and Dain Kaplan. 2015. Active learning in recommender systems. Recommender systems handbook (2015), 809– 846.
- [56] Alireza Salemi, Sheshera Mysore, Michael Bendersky, and Hamed Zamani. 2023. LaMP: When large language models meet personalization. arXiv [cs.CL] (April 2023).
- [57] Alireza Salemi and Hamed Zamani. 2024. Comparing retrieval-augmentation and parameter-efficient fine-tuning for privacy-preserving personalization of large language models. arXiv [cs.CL] (Sept. 2024).
- [58] Scott Sanner, Krisztian Balog, Filip Radlinski, Ben Wedin, and Lucas Dixon. 2023. Large language models are competitive near cold-start recommenders for language- and item-based preferences. In Proceedings of the 17th ACM Conference on Recommender Systems, Vol. 1. ACM, New York, NY, USA, 890–896.
- [59] Omar Shaikh, Michelle Lam, Joey Hejna, Yijia Shao, Michael Bernstein, and Diyi Yang. 2024. Show, don't tell: Aligning language models with demonstrated feedback. arXiv [cs.CL] (June 2024).
- [60] Chengshuai Shi, Cong Shen, and Jing Yang. 2021. Federated multi-armed bandits with personalization. In *International conference on artificial intelligence and* statistics. PMLR. 2917–2925.
- [61] Ruizhe Shi, Yifang Chen, Yushi Hu, Alisa Liu, Hannaneh Hajishirzi, Noah A Smith, and Simon S Du. 2024. Decoding-time language model alignment with multiple objectives. arXiv [cs.LG] (June 2024).
- [62] Alessandro Stolfo, Vidhisha Balachandran, Safoora Yousefi, Eric Horvitz, and Besmira Nushi. 2024. Improving instruction-following in language models through activation steering. arXiv [cs.CL] (Oct. 2024).
- [63] Nishant Subramani, Nivedita Suresh, and Matthew E Peters. 2022. Extracting latent steering vectors from pretrained language models. arXiv [cs.CL] (May 2022)
- [64] Sangho Suh, Meng Chen, Bryan Min, Toby Jia-Jun Li, and Haijun Xia. 2023. Luminate: Structured generation and exploration of design space with large language models for human-AI co-creation. arXiv [cs.HC] (Oct. 2023).
- [65] Zhaoxuan Tan and Meng Jiang. 2023. User modeling in the era of large language models: Current research and future directions. arXiv [cs.CL] (Dec. 2023).
- [66] Lev Tankelevitch, Viktor Kewenig, Auste Simkute, Ava Elizabeth Scott, Advait Sarkar, Abigail Sellen, and Sean Rintel. 2024. The metacognitive demands and opportunities of generative AI. In Proceedings of the CHI Conference on Human Factors in Computing Systems. 1–24.
- [67] Alexander Matt Turner, Lisa Thiergart, Gavin Leech, David Udell, Juan J Vazquez, Ulisse Mini, and Monte MacDiarmid. 2023. Steering language models with activation engineering. arXiv [cs.CL] (Aug. 2023).
- [68] Dimitri von Rütte, Sotiris Anagnostidis, Gregor Bachmann, and Thomas Hofmann. 2024. A language model's guide through latent space. arXiv [cs.CL] (Feb. 2024).
- [69] Kaiwen Wang, Rahul Kidambi, Ryan Sullivan, Alekh Agarwal, Christoph Dann, Andrea Michi, Marco Gelmi, Yunxuan Li, Raghav Gupta, Avinava Dubey, et al. 2024. Conditional Language Policy: A General Framework for Steerable Multi-Objective Finetuning. arXiv preprint arXiv:2407.15762 (2024).

[70] Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, et al. 2025. The rise and potential of large language model based agents: A survey. Science China Information Sciences 68, 2 (2025), 121101.

- [71] Rui Yang, Xiaoman Pan, Feng Luo, Shuang Qiu, Han Zhong, Dong Yu, and Jianshu Chen. 2024. Rewards-in-context: Multi-objective alignment of foundation models with dynamic preference adjustment. arXiv preprint arXiv:2402.10207 (2024).
- [72] Xiwang Yang, Yang Guo, and Yong Liu. 2012. Bayesian-inference-based recommendation in online social networks. IEEE Transactions on Parallel and Distributed Systems 24, 4 (2012), 642–651.
- [73] Dong Yi, Wang Zhilin, Makesh Narsimhan Sreedhar, Wu Xianchao, and Kuchaiev Oleksii. 2023. SteerLM: Attribute conditioned SFT as an (user-steerable) alternative to RLHF. arXiv [cs.CL] (Oct. 2023).
- [74] 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.
- [75] Yu Zhao, Alessio Devoto, Giwon Hong, Xiaotang Du, Aryo Pradipta Gema, Hongru Wang, Kam-Fai Wong, and Pasquale Minervini. 2024. Steering knowledge selection behaviours in LLMs via sae-based representation engineering. arXiv preprint arXiv:2410.15999 (2024).

gemma-2-9b-it

A Parameters for Preference-Based Steering

We followed von Rütte et al.'s procedure for selecting the probe type (in additional to the logistic regressor, we also tested difference-in-means and principal component analysis) and the k hyperparameter for the top-k layers, where steering is selectively applied to the top-k of the layers with the highest detection accuracy. Table A.1 shows the best performing parameters that were used for each model. We also experimentally determined an approximate functional steering range where the model is not affected significantly by degeneracy (high-perplexity responses). This functional range is used in experiments E2-E5 and the user study.

Model	Top K Layers	Functional Steering Range	Probe
stablelm-21-6b-chat	16	(-10, 10)	Logistic
gemma-2-2b-it	16	(-30, 30)	Logistic
Mistral-7B-Instruct-v0.3	24	(-30, 30)	Logistic
Qwen2.5-7B-Instruct	24	(-10, 10)	Logistic

(-30, 30)

Logistic

Table A.1: Steering parameters selected for each model.

Due to a lack of LLM output datasets on preferences, all steering datasets are generated by GPT-40 to emulate what the gold standard of embodying this preference should look like – e.g. "Generate a list of responses that a chatbot might reply to a user who is looking for kid-friendly suggestions across user queries in lifestyle planning tasks such as travel planning, restaurant recommendations, recipe selections..." See Listing A.1 for sample instances from the positive and negative steering dataset for age. Across all steering datasets in the five preference dimensions, we generated 50-80 samples per trait. We explored alternative datasets such as a list of synonymous adjectives (e.g., cheap suggestions, affordable, low-cost... for cost) and the Yelp reviews, but the former did not have an appropriate format to guide LLM outputs and the latter is too noisy to learn meaningful concepts from.

Listing A.1: Sample instances from the steering dataset for age

```
[Positive steering to adults-oriented]
Visit the rooftop bar downtown for stunning views and a curated selection of craft cocktails.
Plan a day of wine tasting at the region's top vineyards, complete with private tours
Book a spa day with massages, saunas, and relaxation areas exclusively for adults.
Attend a live jazz performance at the speakeasy-style venue known for its intimate ambiance.
Join a mixology workshop to learn the secrets behind making perfect cocktails at home.
Take a sunset cruise offering drinks and live music in a serene, adults-only setting.
Visit the cigar lounge downtown, featuring an extensive selection and a cozy atmosphere
Book a brewery tour with tastings of seasonal beers and behind-the-scenes insights from brewmasters.
Attend an adults-only trivia night at the local pub - it's a fun way to test your knowledge.
Visit the distillery for a tour and a guided whiskey tasting experience.
[Negative steering to kids-oriented]
Visit the interactive science museum with hands-on exhibits perfect for kids of all ages.
Plan a day at the local zoo, featuring kid-friendly animal encounters and feeding sessions.
Explore the city's largest playground with climbing structures, slides, and picnic areas for families.
Spend time at the local aquarium, known for its touch tanks and playful sea otters.
Visit the city's amusement park with rides designed specifically for younger kids and toddlers.
Head to the local library for storytime sessions and engaging activities for kids.
Visit the butterfly garden - it's colorful, educational, and a hit with younger children.
Take a scenic train ride that offers kid-friendly entertainment and stunning views along the way.
Visit the local zookeeper for a behind-the-scenes tour - kids love learning about animals up close.
Spend time at the nearby petting farm, where kids can feed and interact with friendly animals.
Explore the city's famous ice cream parlor, offering fun flavors and mini cones for little ones.
```

B Sample Task Queries in Computational Experiments

Listing B.1 shows a sample of real user queries selected from the OASST2 dataset in the domain of *lifestyle planning* that were used throughout the computational experiments. These queries were manually identified and filtered by the research team. The user study tasks queries were designed to be similar to these.

Listing B.1: Sample task queries used in the computational experiments.

```
Give a recipe idea for a vegetarian meal which has tofu?

what are some popular souvenirs people take home when they visit the usa from europe?

What are the best restaurants in San Francisco?

Hi, I am trying to plan a camping holiday with my family. What supplies do you suggest I purchase or pack?

I want to buy a gift for my gir friend for the valentines day

What are some things I should do on a 5-day vacation in Thailand?

What is the best way to cook a tomato?
```

C Sample Yelp Reviews Used for Evaluation

Listing C.1 shows samples of the reference Yelp reviews datasets for *hipster* and *touristy* of the **ambiance** preference dimension. The samples are first collected based on the criteria listed in Table 1, then filtered down for relevancy to the trait through a combination of manual human evaluation and automated LLM evaluation with GPT-40-mini as the judge.

Listing C.1: Sample Yelp reviews of the traits touristy and hipster for the ambiance dimension.

```
[Yelp reviews for hipster]
There are restaurants, bars, and cafes in some cities that you go to because they are unique enough to stand out from anything in their category
This is a place I would highly recommend if you are looking for somewhere different to eat that does not have the usual bar fare.
I wish I wasn't reviewing this place because I want it to remain a hidden gem. This is by far my favorite restaurant in St. Pete.
This is a great find just outside of the hustle and bustle of downtown..
The Mercy Lounge hosts great bands and events, including local burlesque acts, and sports pin-up gals on the bar.
Amazingly friendly staff. Building is just bursting with personality.
Simply lovely. Every dish had an unusual but delicious mix of flavors and textures.
[Yelp reviews for touristy]
The signature iced tea was fabulous! There was a long wait but worth every minute of it!
Reading Terminal Market has everything that you need everything is fresh good or good for you
The highlight of the cruse was the tour guide/narrator Charles and his eloquent oration.
We were so excited to visit here from Las Vegas. What a Christmas experience for us!!!
Love coming here :) so many goodies. It is very bustling so if you're one who doesn't like waiting in lines, you'll have to suck it up!
You can find anything you want and the atmosphere is great for tourists and newcomers. If touristy is what you want, they serve great
     cheesesteaks!
Genos vibe was very touristy (we were tourists as well haha) but crowded and people taking forever to order off a very simple menu.
```

D Additional Interactions with Prompting Results (E2)

The results for the other preference dimensions in E2 are in Figure D.1. In particular, time is very sensitive to prompting.

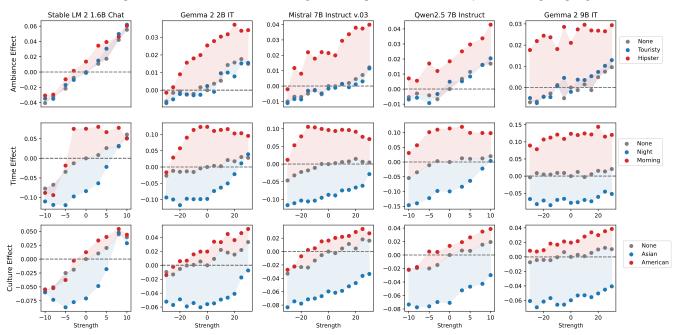


Figure D.1: Effect of additional prompting on steering for two preferences: ambiance (top row), time (middle row), and culture (bottom row) in E2.

E Additional Multi-Preference Steering Results (E3)

The correlation map of the cosine similarities between each of the preference dimensions (this is computed based on the mean of the Yelp reference embeddings). The least correlated pair, **age** and **culture** are selected for E3. The surface plots for the remaining four models for E3 are shown in E.2. Qualitative examples of multi-steered outputs from *Mistral-7B-Instruct-v0.3* are shown in Table E.2.

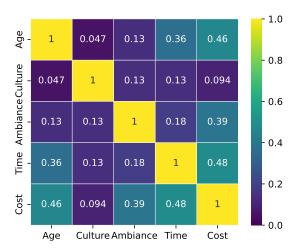


Figure E.1: Correlation of absolute cosine similarity between the preference dimensions.

Table E.2: Samples from multi-preference steering with preferences age and culture (E3).

	Mistral 7B ($d_{culture} = 30$, steer \rightarrow American)	Mistral 7B ($d_{culture} = -30$, steer \rightarrow Asian)
Mistral 7B ($d_{age} = 30$, steer \rightarrow Adults)	San Francisco is known for its vibrant dining scene, offering a wide range of culinary experiences. Here are some of the top-rated restaurants in the city. 1) The Slanted Door - Overlooking the Bay, this restaurant offers breathtaking views alongside a menu of contemporary American dishes	San Francisco, being a vibrant city with a rich culinary scene. offers a wide array of exceptional restaurants that cater to various gastronomic preferences. 1) Izakaya Roku: Known for its authentic and vibrant <u>izakaya</u> atmosphere, this restaurant offers a comprehensive selection of <u>lapanese small</u> <u>dishes and sake</u>
Mistral 7B ($d_{age} = -30$, steer $\rightarrow \frac{\text{Kids}}{}$)	San Francisco offers a variety of family-friendly restaurants that are popular among tourists and locals. Here are some of the best restaurants in the San Francisco Bay Area. 1) The Original Pancake House: A roadside attraction featuring giant dinosaur statues and a restaurant serving breakfast foods	Here are some popular restaurants that you can consider visiting. 1) Golden Gate Fortune Cookie Factory This is a Fun and educational experience where you can make your own fortune cookies. It's a great way to introduce kids to different shapes and designs

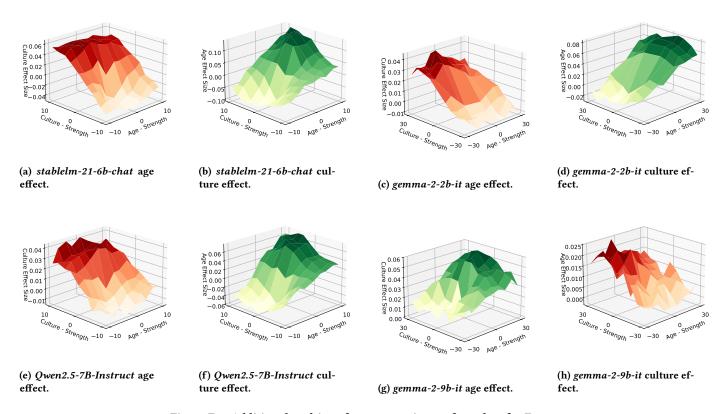


Figure E.2: Additional multi-preference steering surface plots for E3.

F Additional Learning Preference and Steering Results (E5)

The results of learning a hidden preference for the other three models are shown in Figure F.1. Most did not have a significant difference in effects, likely because the models responded well to prompting.

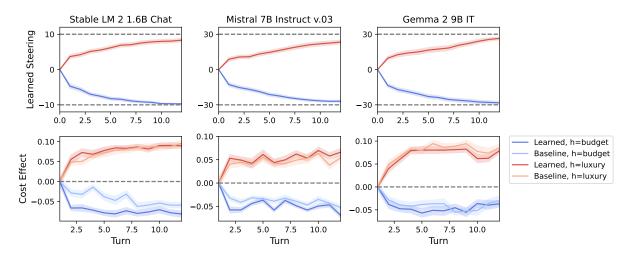


Figure F.1: Results for stablelm-21-6b-chat, Mistral-7B-Instruct-v0.3, and gemma-2-9b-it for E5, with estimated preference steering strength (top) and corresponding cost effect sizes (bottom). Blue represents trials where the latent preference h is very budget, and red represents very luxury.

G Details for GPT Simulations (E5)

We implement *GPT-as-user* for **E5** through two API calls, the first to determine the characteristics of the user's response (the dissatisfaction level and the direction of change), shown in Listing G.1; the second to generate the user's response based on the required characteristics, shown in Listing G.2. We also show how we guide GPT to generate responses with semantics that vary in dissatisfaction score (Figure G.1).

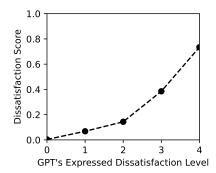


Figure G.1: Computed dissatisfaction scores with GPT's expressed dissatisfaction levels.

Listing G.1: First API call made to GPT to determine the user response characteristics.

```
You are a helpful assistant who is roleplaying as the end user of a commercial chatbot. After receiving the chatbot's message, analyze it
     carefully against the persona provided to you below. Consider the following steps:
 What preference level does the chatbot's output map to?
- If it is my preference level (e.g. 25% budget vs 25% budget), then I am satisfied.
- If it is within 25% of my preference level (e.g. 25% budget vs 50% budget), then my disatisfaction is mild.
- If it is within 26%-75% (e.g. 25% budget vs 20% luxury), then my disatisfaction is moderate.
- If it is within 76%-150% (e.g. 25% budget vs 50% luxury), then my disatisfaction is significant.
- If it is 151% and more (e.g. 25\% budget vs 175\% luxury), then my disatisfaction is extreme.
Output your response in the following format: "Based on my preference of [insert preference], I find the chatbot's response to my question to be
     highly/slightly appropriate/inappropriate. To improve the output, I would like to response with [no/mild/moderate/significant/extreme] dissatisfaction and ask the chatbot to be [cheaper/more luxurious]."
You have specific preferences towards budget. Your preference for budget is: [PREFERENCE].
For reference, see these examples for how you can interpret preference for a question on travelling:
100% Budget: Stays with hosts via couchsurfing or similar platforms, relies on grocery store food or home-cooked meals, only visits free
     attractions (e.g., public parks, free museums), avoids paid transport by walking everywhere.
75% Budget: Stays in shared hostel rooms or budget guesthouses, eats primarily at fast food outlets or food trucks, prefers free attractions but
       occasionally pays for low-cost activities, uses public transit as the primary mode of transportation.
50% Budget: Alternates between hostels and affordable Airbnbs, dines at low-cost local restaurants or casual cafes, pays for affordable
     attractions like entry-level museum tickets or guided tours, uses public transit regularly and occasionally opts for rideshares like Uber.
25% Budget: Stays in private rooms in Airbnbs or budget hotels, enjoys dining at both low-cost and midrange restaurants, willing to pay for most
       attractions, uses a mix of public transit and rideshares for convenience.
25% Luxury: Stays in comfortable 3-star hotels, enjoys a range of dining options from midrange to occasional upscale restaurants, prefers
     premium guided tours or exclusive attraction experiences, uses rideshares like Uber exclusively for transportation.
50% Luxury: Stays in upscale 4-star hotels, dines at a mix of midrange and high-end restaurants, often books private tours and curated
     experiences, opts for premium rideshare options like Uber Lux.
75% Luxury: Stays in 5-star hotels or luxury boutique accommodations, primarily dines at high-end restaurants, books private tours and
     customized experiences, hires private chauffeurs for transportation.
100% Luxury: Stays in world-class accommodations such as luxury resorts or exclusive penthouses, dines at Michelin-starred and exclusive fine
     dining establishments, books private and bespoke entertainment or tours, relies on personal chauffeurs for all travel needs.
```

Listing G.2: Second API call made to GPT to generate the user's response.

```
You are a helpful assistant who is roleplaying as the end user of a commercial chatbot. You are provided with a conversation history of the user and the chatbot, as well as an instruction stating the intention of the next message that you need to write. The instruction will indicate the intended sentiment of the message, adhere to it based on the following examples.

No Dissatisfaction: "This is perfect, I like it!". Feel free to add a followup question related to the topic of the conversation.

Mild Dissatisfaction: "This sounds alright, but I also want...". Embed the preference that the instruction says to steer towards.

Noticeable Dissatisfaction: "Imm I like some parts, but it can be improved..". Embed the preference that the instruction says to steer towards.

Significant Dissatisfaction: "I would need some modifications...". Embed the preference that the instruction says to steer towards.

Extreme Dissatisfaction: "I don't like this at all, I want...". Embed the preference that the instruction says to steer towards.

Guideline for Interactions:

- Do not directly mention your preference and do not say what you are looking for in detail.

- Only focus on expressing the correct sentiment and satisfaction/dissatisfaction levels.
```

H User Study True Preferences

Figure H.1 shows the distribution of true preferences across each of the four tasks/preferences pairings as they were self-reported by the user study participants.

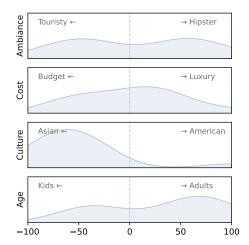


Figure H.1: Ground truth preference distributions.