# FSPO: Few-Shot Preference Optimization of Synthetic Preference Data Elicits LLM Personalization to Real Users

Anonymous Authors<sup>1</sup>

#### Abstract

Effective personalization of LLMs is critical for a broad range of user-interfacing applications such as virtual assistants and content curation. Inspired by the strong in-context learning capabilities of LLMs, we propose few-shot preference optimization (FSPO), an algorithm for LLM personalization that reframes reward modeling as a metalearning problem. Under FSPO, an LLM learns to quickly infer a personalized reward function for a user via a few labeled preferences. Since real-world preference data is scarce and challenging to collect at scale, we propose careful design choices to construct synthetic preference datasets for personalization, generating over 1M synthetic personalized preferences using publicly available LLMs. In particular, to successfully transfer from synthetic data to real users, we find it crucial for the data to exhibit both high diversity and coherent, self-consistent structure. We evaluate FSPO on personalized open-ended generation for up to 1,500 synthetic users across across three domains: movie reviews, pedagogical adaptation based on educational background, and general question answering. We also run a controlled human study. Overall, FSPO achieves an 87% Alpaca Eval winrate in generating responses that are personalized to synthetic users and a 72% winrate with real human users in open-ended question answering.

#### 1. Introduction

As large language models (LLMs) increasingly interact with a diverse user base, it becomes important for models to generate responses that align with individual user preferences. People exhibit a wide range of preferences and beliefs shaped by their cultural background, personal experience, and individual values. These diverse preferences are humanannotated preference datasets; however, current preferences optimization techniques like reinforcement learning from human feedback (RLHF) largely focus on optimizing a *single* model based on preferences aggregated over the entire population. This approach may neglect minority viewpoints, embed systematic biases into the model, and ultimately lead to worse performance compared to personalized models. Can we create language models that can adaptively align with the personal preferences of each user instead of the aggregated preferences of all users?

Addressing this challenge requires a shift from modeling a singular aggregate reward function to modeling a distribution of reward functions that captures the diversity of human preferences (Sorensen et al., 2024; Jang et al., 2023). By doing so, we can enable personalization in language models, allowing them to generate a wide range of responses tailored to individual subpopulations. This approach not only enhances user satisfaction but also promotes inclusivity by acknowledging and respecting the varied perspectives that exist within any user base. Despite this problem's importance, to our knowledge LLM personalization has yet to be achieved for open-ended question answering with real users.

In this paper, we introduce few-shot preference optimization (FSPO), a novel framework designed to model diverse subpopulations in preference datasets to elicit personalization in language models for open-ended question answering. At a high level, FSPO leverages in-context learning to adapt to new subpopulations. This adaptability is crucial for practical applications, where user preferences can be dynamic and multifaceted. Inspired by past work on black-box metalearning for language modeling (Chen et al., 2022; Min et al., 2022; Yu et al., 2024), we fine-tune the model in a meta-learning setup using preference-learning objectives such as IPO (Gheshlaghi Azar et al., 2023). To further improve personalized generation, we additionally propose user description chain-of-thought (COT), which allows the model to leverage additional inference-time compute for better reward modeling and instruction following.

Learning a model that effectively personalizes to real people requires training on a realistic, user-stratified preference dataset. One natural approach to consider is to curate such

<sup>&</sup>lt;sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

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



Figure 1: Overview of FSPO. N previously collected preferences are fed into the LLM along with the current query, allowing the LLM to personalize its response to the query using the past preferences.

data from humans, but this is difficult and time-consuming.
Instead, we propose instantiating this dataset synthetically,
and present careful design decisions inspired from the metalearning literature (Hsu et al., 2019; Yin et al., 2019) to
generate a dataset that is both diverse and structured.

074 To evaluate the efficacy of our approach, we construct a 075 set of three semi-realistic domains to study personalization: 076 (1) Reviews, studying the generation ability of models for 077 reviews of movies, TV shows, and books that are consis-078 tent with a user's writing style, (2) Explain Like I'm X 079 (ELIX): studying the generation ability of models for re-080 sponses that are consistent with a user's education level, and 081 (3) Roleplay: studying the generation ability of models for 082 responses that are consistent with a user's description, with 083 effective transferability to a real human-study. Here we find 084 that FSPO outperforms an unpersonalized model on average by 87%. We additionally perform a controlled human study 086 showcasing a winrate of 72% of FSPO over unpersonalized 087 models.

By addressing limitations of existing reward modeling techniques, our work paves the way for more inclusive and personalized LLMs. We believe that FSPO represents a significant step toward models that better serve the needs of all users, respecting the rich diversity of human preferences.

#### 2. Related Work

094

095

096

109

066

067

Personalized learning of preferences. Prior research has 097 explored personalization through various methods. One 098 approach is distributional alignment, which focuses on 099 100 matching model outputs to broad target distributions rather than tailoring them to individual user preferences. For example, some prior work have concentrated on aligning model-generated distributions with desired statistical properties (Siththaranjan et al., 2024; Meister et al., 2024; Melnyk 104 105 et al., 2024), yet they do not explicitly optimize for individual preference adaptation. Another strategy involves 106 explicitly modeling a distribution of rewards (Lee et al., 2024; Poddar et al., 2024). However, these methods suffer

from sample inefficiency during both training and inference (Rafailov et al., 2023; Gheshlaghi Azar et al., 2023). Additionally, these approaches have limited evaluations: Lee et al. (2024) focuses solely on reward modeling, while Poddar et al. (2024) tests with a very limited number of artificial users (e.g helpfulness user and honest user). Other works have investigated personalization in multiple-choice questions, such as GPO (Zhao et al., 2024). Although effective in structured survey settings, these methods have not been validated for open-ended personalization tasks. Similarly, Shaikh et al. (2024) explores personalization via explicit human corrections, but relying on such corrections is expensive and often impractical to scale. Finally, several datasets exist for personalization, such as Prism (Kirk et al., 2024) and Persona Bench (Castricato et al., 2024). Neither of these datasets demonstrate that policies trained on these benchmarks lead to effective personalization. Unlike these prior works which study personalization based off of human values and controversial questions, we instead study more general questions that a user may ask.

Algorithms for preference learning. LLMs are typically fine-tuned via supervised next-token prediction on highquality responses and later refined with human preference data (Casper et al., 2023; Ouyang et al., 2022). This process can use on-policy reinforcement learning methods like REINFORCE (Sutton et al., 1999) or PPO (Schulman et al., 2017), which optimize a reward model with a KL constraint. Alternatively, supervised fine-tuning may be applied to a curated subset of preferred responses (Dubois et al., 2024b) or iteratively to preferred completions as in ReST (Gulcehre et al., 2023). Other methods, such as DPO (Rafailov et al., 2023), IPO (Gheshlaghi Azar et al., 2023), and KTO (ContextualAI, 2024), learn directly from human preferences without an explicit reward model, with recent work exploring iterative preference modeling applications (Yuan et al., 2024).

**Black-box meta-learning.** FSPO is an instance of blackbox meta-learning, which has been studied in a wide range of domains spanning image classification (Santoro et al., 2016; Mishra et al., 2018), language modeling (Chen et al.,
2022; Min et al., 2022; Yu et al., 2024), and reinforcement
learning (Duan et al., 2016; Wang et al., 2016). Black-box
meta-learning is characterized by the processing of task
contexts and queries using generic sequence operations like
recurrence or self-attention, instead of specifically designed
adaptation mechanisms.

#### **3. Preliminaries and Notation**

117

144

145

159 160

119 Preference fine-tuning algorithms, such as reinforcement 120 learning from human feedback (RLHF) and direct pref-121 erence optimization (DPO), typically involve two main 122 stages (Ouyang et al., 2022; Ouyang et al., 2022): su-123 pervised fine-tuning (SFT) and preference optimization 124 (DPO/RLHF). First, a pre-trained model is fine-tuned on 125 high-quality data from the target task using SFT. This pro-126 cess produces a reference model, denoted as  $\pi_{ref}$ . The pur-127 pose of this stage is to bring the responses from a partic-128 ular domain in distribution with supervised learning. To 129 further refine  $\pi_{ref}$  according to human preferences, a pref-130 erence dataset  $\mathcal{D}_{\text{pref}} = \{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}_w, \mathbf{y}^{(i)}_l)\}$  is collected. In 131 this dataset,  $\mathbf{x}^{(i)}$  represents a prompt or input context,  $\mathbf{y}^{(i)}_w$ 132 is the preferred response, and  $\mathbf{y}_{l}^{(i)}$  is the less preferred re-133 sponse. These responses are typically sampled from the 134 output distribution of  $\pi_{ref}$  and are labeled based on human 135 feedback. 136

Most fine-tuning pipelines assume the existence of an underlying reward function  $r^*(\mathbf{x}, \cdot)$  that quantifies the quality of responses. A common approach to modeling human preferences is the Bradley-Terry (BT) model (Bradley and Terry, 1952), which expresses the probability of preferring response  $\mathbf{y}_1$  over  $\mathbf{y}_2$ , given a prompt  $\mathbf{x}$ , as:

$$p^*(\mathbf{y}_1 \succ \mathbf{y}_2 \mid \mathbf{x}) = \frac{e^{r^*(\mathbf{x}, \mathbf{y}_1)}}{e^{r^*(\mathbf{x}, \mathbf{y}_1)} + e^{r^*(\mathbf{x}, \mathbf{y}_2)}}$$
(1)

147 Here,  $p^*(\mathbf{y}_1 \succ \mathbf{y}_2 \mid \mathbf{x})$  denotes the probability that  $\mathbf{y}_1$  is 148 preferred over  $\mathbf{y}_2$  given  $\mathbf{x}$ .

149 The objective of preference fine-tuning is to optimize the 150 policy  $\pi_{\theta}$  to maximize the expected reward  $r^*$ . However, 151 directly optimizing  $r^*$  is often impractical due to model 152 limitations or noise in reward estimation. Therefore, a re-153 ward model  $r_{\phi}$  is trained to approximate  $r^*$ . To prevent 154 the fine-tuned policy  $\pi_{\theta}$  from deviating excessively from 155 the reference model  $\pi_{ref}$ , a Kullback-Leibler (KL) diver-156 gence constraint is imposed. This leads to the following 157 fine-tuning objective: 158

$$\max \mathbb{E}[r^*(x,y)] - \beta D_{\mathrm{KL}}(\pi \parallel \pi_{\mathrm{ref}})$$
(2)

<sup>161</sup> In this equation, the regularization term weighted by  $\beta$  controls how much  $\pi_{\theta}$  diverges from  $\pi_{ref}$ , based on the reverse KL divergence constraint. This constraint ensures that the



Figure 2: User description chain-of-thought (COT). Prediction is a two-stage process: first predicting a (synthetic) user description from the few-shot preferences and next predicting the response.

updated policy remains close to the reference model while improving according to the reward function.

**Reward model training.** To fine-tune the large language model (LLM) policy  $\pi_{\theta}(\mathbf{y} \mid \mathbf{x})$ , the Bradley-Terry framework allows for either explicitly learning a reward model  $r_{\phi}(\mathbf{x}, \mathbf{y})$  or directly optimizing preferences. Explicit reward models are trained using the following classification objective:

$$\max_{\phi} \mathbb{E}_{\mathcal{D}_{\text{pref}}}\left[\log \sigma\left(r_{\phi}(\mathbf{x}, \mathbf{y}_w) - r_{\phi}(\mathbf{x}, \mathbf{y}_l)\right)\right]$$
(3)

where  $\sigma$  is the logistic function, used to map the difference in rewards to a probability. Alternatively, contrastive learning objectives such as Direct Preference Optimization (Rafailov et al., 2023) and Implicit Preference Optimization (Gheshlaghi Azar et al., 2023) utilize the policy's log-likelihood  $\log \pi_{\theta}(\mathbf{y} \mid \mathbf{x})$  as an implicit reward:

$$r_{\theta}(\mathbf{x}, \mathbf{y}) = \beta \log \left( \pi_{\theta}(\mathbf{y} \mid \mathbf{x}) / \pi_{\text{ref}}(\mathbf{y} \mid \mathbf{x}) \right)$$
(4)

This approach leverages the policy's log probabilities to represent rewards, thereby simplifying the reward learning process.

# 4. The Few-Shot Preference Optimization (FSPO) Framework

**Personalization as a meta-learning problem.** Generally, for fine-tuning a model with RLHF a preference dataset of the form:  $\mathcal{D}_{pref} = \{(\mathbf{x}^{(i)}, \mathbf{y}_w^{(i)}, \mathbf{y}_l^{(i)})\}$  is collected, where x is a prompt,  $y_w$  is a preferred response, and  $y_l$  is a dispreferred response. Here, preferences from different users are aggregated to learn the preferences over a population. However, through this aggregation, individual user preferences are marginalized, leading to the model losing personalized values or beliefs due to population-based preference learning and RLHF algorithms such as DPO as seen in prior work (Siththaranjan et al., 2024).

How can we incorporate user information when learning from preference datasets? In this work, we have a weak requirement to collect scorer-ids  $S^{(i)}$  of each user for differentiating users that have labeled preferences in our dataset:

165  $\mathcal{D}_{\text{pref}} = \{(\mathbf{x}^{(i)}, \mathbf{y}_w^{(i)}, \mathbf{y}_l^{(i)}, \mathbf{S}^{(i)})\}$ . Now consider each user 166 as a task instance, where the objective is to learn an effective 167 reward function for that user using the user's set of pref-168 erences. This can be naturally instantiated as a black-box 169 meta-learning objective, where meta-learning is done over 170 users (also referred to as a task in meta-learning). Meta-171 learning should enable rapid personalization, i.e. adaptabil-172 ity to new users with just a few preferences. 173

More formally, consider that each unique user  $S^{(i)}$ 's reward function is characterized by a set of preferences with prompt and responses  $(x, y_1, y_2)$ , and preference label c (indicating if  $y_1 \succ y_2$  or  $y_1 \prec y_2$ ). Given a distribution over users  $S = P(S^{(i)})$ , a meta-learning objective can be derived to minimize its expected loss with respect to  $\theta$  as:

174

175

176

177

178

179

180

181

182

183

184

185

186

187

188

189 190

191

$$\min_{\theta} \mathbb{E}_{\mathcal{S}^{(i)} \sim \mathcal{S}} \left[ \mathbb{E}_{(x, y_1, y_2, c) \sim \mathcal{D}_i, \{(x, y_1, y_2, c)\}_1^N \sim D_i} \left[ \mathcal{L}_{pref}^{\theta} \left( x, y_1, y_2, c | \{(x, y_1, y_2, c)\}_1^N \right) \right] \right] (5)$$

where  $D_i$  is a distribution over preference tuples  $(x, y_1, y_2, c)$  for each user  $S^{(i)}$ , and  $\mathcal{L}_{pref}^{\theta}$  is a preference learning objective such as DPO (Rafailov et al., 2023) or IPO (Gheshlaghi Azar et al., 2023):

$$\mathcal{L}_{pref}^{\theta} = ||h_{\pi_{\theta}}^{y_{w},y_{l}} - (2\beta)^{-1}||_{2}^{2},$$

$$h_{\pi_{\theta}}^{y_{w},y_{l}} = \log \frac{\pi_{\theta}(y_{w}|x)}{\pi_{ref}(y_{w}|x)} - \log \frac{\pi_{\theta}(y_{l}|x)}{\pi_{ref}(y_{l}|x)}$$
(6)

where  $y_w$  and  $y_l$  are the preferred and dispreferred responses (respectively) according to the responses  $y_1, y_2$  and class label c in the preference dataset.

196 Following black-box meta-learning approaches, FSPO re-197 ceives as input a sequence of preferences  $D_i^{fewshot} \sim D_i$ 198 from a User  $S^{(i)}$ . This is followed by an unlabeled, held-out preference  $(x, y_1, y_2) \sim \mathcal{D}_i \setminus \mathcal{D}_i^{fewshot}$  for which it outputs its 199 200 prediction c. To make preferences compatible with a pre-201 trained language model, a few-shot prompt is constructed, 202 comprising of preferences from a user and the held-out query as seen in Figure 1. This construction has an added 204 benefit of leveraging a pretrained language model's capabilities for few-shot conditioning (Brown et al., 2020), which 206 can enable some amount of steerage/personalization. This prediction c is implicitly learned by a preference optimiza-208 tion algorithm such as DPO (Rafailov et al., 2023), which 209 parameterizes the reward model as  $\beta \frac{\log \pi_{\theta}(y|x)}{\log \pi_{ref}(y|x)}$ . This pa-210 rameterization enables us to leverage the advantages of pref-211 erence optimization algorithms such as eliminating policy 212 learning instabilities and computational burden of on-policy 213 214 sampling, learning an effective model with a simple classification objective. 215

User description chain-of-thought (COT). If provided
with a description of the user (potentially synthetically generated), FSPO can be converted to a two-step prediction

problem as seen in Figure 2. In the first step, conditioned on user few-shot preferences, the user description is generated, then conditioned on the prompt, few-shot preferences, and generated user description, a response can then be generated. This prediction of the user description is an interpretable summarization of the fewshot preferences and a better representation to condition on for response generation. Similar to the rationale generated in Zhang et al. (2024) for verifiers, the COT prediction can be viewed as using additional inference-compute for better reward modeling. Additionally, this formulation leverages the instruction following ability of LLMs (Ouyang et al., 2022) for response generation.

User representation through preference labels. From an information-theoretic perspective, the few-shot binary preferences can be seen as a N-bit representation of the user, representing up to  $2^N$  different personas or reward functions. There are several ways to represent users: surveys, chat histories, or other forms of interaction that reveal hidden preferences. We restrict our study to such a N-bit user representation, as such a constrained representation can improve the performance when transferring reward models learned on synthetic personalities to real users. We defer the study of less constrained user representations to future work.

We summarize FSPO in Algorithm 1. Next, we will discuss domains to study FSPO.

Algorithm 1 Overview of few-Shot preference optimization (FSPO).

- 0: Input: For each unique user S<sup>(i)</sup>, a dataset of preferences D := (x, y<sub>1</sub>, y<sub>2</sub>, c)<sub>i</sub>, and optionally user description y<sub>S<sup>(i)</sup></sub> for COT, ∀i
- 0: **Output:** Learned policy  $\pi_{\theta}$
- 0: while not done do
- 0: Sample training user  $S^{(i)}$  (or minibatch)
- 0: Sample a subset of preferences from the user  $\mathcal{D}_i^{fewshot} \sim \mathcal{D}_i$
- 0: Sample held-out preference examples  $D_i^{heldout} \sim D_i \setminus D_i^{fewshot}$
- 0: if COT then
- 0: Use Equation (5) and Equation (6) to predict the loss on the user description  $y_{S^{(i)}}$
- 0: end if
- 0: Conditioning on  $\mathcal{D}_i^{fewshot}$  (optionally  $y_{\mathcal{S}^{(i)}}$ ), use Equation (5) and Equation (6) to predict the loss on the held-out preference example  $D_i^{heldout}$
- 0: Update learner parameters  $\theta$ , using gradient of loss on  $D_i^{heldout}$

#### 0: end while

0: **Return**  $\pi_{\theta} = 0$ 

#### 5. Domains to Study Personalization

To study personalization with FSPO we construct a benchmark across 3 domains ranging from generating personalized movie reviews (**Reviews**), generating personalized responses based off a user's education background (**ELIX**), and personalizing for general question answering (**Roleplay**). We open-source preference datasets and evaluation protocols from each of these tasks for future work looking to study personalization (sample in supplementary).

**Reviews.** The Reviews task is inspired by the IMDB dataset (Maas et al., 2011), containing reviews for movies. We curate a list of popular media such as movies, TV shows, anime, and books for a language model to review. We consider two independent axes of variation for users: sentiment (positive and negative) and conciseness (concise and verbose). Here being able to pick up the user is crucial as the users from the same axes (e.g positive and negative) would have opposite preferences, making this *difficult* to learn with any population based RLHF method. We also study the steerability of the model considering the axes of verbosity and sentiment in tandem (e.g positive + verbose).

ELIX. The Explain Like I'm X (ELIX) task is inspired 243 by the subreddit "Explain Like I'm 5" where users answer 244 questions at a very basic level appropriate for a 5 year old. Here we study the ability of the model to personalize a ped-246 agogical explanation to a user's education background. We 247 construct two variants of the task. The first variant is ELIX-248 easy where users are one of 5 education levels (elementary 249 school, middle school, high school, college, expert) and the 250 goal of the task is to explain a question such as "How are 251 beaches formed?" to a user of that education background. 252 The second, more realistic variant is ELIX-hard, which 253 consists of question answering at a high school to university 254 level. Here, users may have different levels of expertise in 255 different domains. For example, a PhD student in computer 256 science may have a very different educational background 257 from an undergraduate studying studying biology, allowing 258 for preferences from diverse users (550 users). 259

**Roleplay.** The Roleplay task tackles general question an-261 swering across a wide set of users, following PRISM (Kirk et al., 2024) and PERSONA Bench (Castricato et al., 2024) 263 to study personalization representative of the broad human 264 population. We start by identifying three demographic traits 265 (age, geographic location, and gender) that humans differ in 266 that can lead to personalization. For each trait combination, 267 we generate 30 personas, leading to 1,500 total personas. 268 To more accurately model the distribution of questions, we 269 split our questions into two categories: global and specific. 270 Global questions are general where anyone may ask it, but 271 specific questions revolve around a trait, for example an 272 elderly person asking about retirement or a female asking 273 about breast cancer screening. 274

One crucial detail for each task is the construction of a preference dataset that spans multiple users. But how should one construct such a dataset that is realistic and effective?

# 6. Sim2Real: Synthetic Preference Data Transfers to Real Users

Collecting personalized data at scale presents significant challenges, primarily due to the high cost and inherent unreliability of human annotation. Curating a diverse set of users to capture the full spectrum of real-world variability further complicates the process, often limiting the scope and representativeness of the data. Synthetically generating data using a language model (Li et al., 2024; Bai et al., 2022) is a promising alternative, since it can both reduce costly human data generation and annotation and streamline the data curation process. Can we generate diverse user preference data using language models in a way that transfers to real people?

We draw inspiration from simulation-to-real transfer in nonlanguage domains like robotics (Makoviychuk et al., 2021) and self-driving cars (Yang et al., 2023), where the idea of domain randomization (Tobin et al., 2018) has been particularly useful in enabling transfer to real environments. Domain randomization enables efficient adaptation to novel test scenarios by training models in numerous simulated environments with varied, randomized properties.

But why is this relevant to personalization? As mentioned previously, each user can be viewed as a different "environment" to simulate as each user has a unique reward function that is represented by their preferences. To ensure models trained on synthetic data generalize to real human users, we employ domain randomization to simulate a diverse set of synthetic preferences. However, diversity alone isn't sufficient to learn a personalized LM. As studied in prior work (Hsu et al., 2019; Yin et al., 2019), it is crucial that the task distribution in meta-learning exhibits sufficient structure to rule out learning shortcuts that do not generalize. But how can we elicit both **diversity** and **structure** in our preference datasets?

**Encouraging diversity.** Diversity of data is crucial to learning a reward function that generalizes across prompts. Each domain has a slightly different generation setup as described in Section 5, but there are some general design decisions that are shared across all tasks to ensure diversity.

One source of diversity is in the questions used in the preferences. We use a variety of strategies to procure questions for the three tasks. For question selection for ELIX, we first sourced questions from human writers and then synthetically augmented the set of questions by prompting GPT-40 (OpenAI et al., 2024) with subsets of these human-generated questions. This allows us to scalably augment the human



Figure 3: **Overview of domain randomization techniques.** View-conditioning (left) decomposes a given question into multiple viewpoints, allowing for diverse response generation. Iterative persona generation (right) allows for better structure by removing underspecification of the persona by iteratively refining a persona if it is insufficient to make a preference prediction.



299

300

301

302

303

304

Figure 4: Flowchart of Roleplay dataset generation: Starting from a set of traits, a seed persona is constructed and a set of specific questions about that trait. Then responses are constructed with View-Conditioning. The seed personas are then iteratively refined to not be underspecified. Finally, the refined persona is used to score consistent preferences.

question dataset, while preserving the stylistic choices and 305 beliefs of human writers. For the reviews dataset, we com-306 piled a list of popular media from sites such as Goodreads, 307 IMDb, and MyAnimeList. For the Roleplay dataset, we 308 prompted GPT-40 to generate questions all users would ask 309 (global) or questions only people with a specific trait would ask (specific). This allows us to have questions that are 311 more consistent with the distribution of questions people 312 313 may ask.

314 Additionally, having a diversity of responses is crucial for 315 not only training the model on many viewpoints but also 316 reward labeling, allowing for greater support over the set 317 of possible responses for a question. To achieve diverse re-318 sponses, we employ two strategies: Persona Steering (Cheng 319 et al., 2023) and view conditioning. For ELIX and Reviews, 320 we use persona steering by prompting the model with a question and asking it to generate an answer for a randomly 322 selected persona. For Roleplay, the user description was 323 often underspecified so responses generated with persona 324 steering were similar. Therefore, we considered a multi-325 turn approach to generating a response. First, we asked the model to generate different viewpoints that may be possible 327 for a question. Then, conditioned on each viewpoint inde-328 pendently, we prompted the model with the question and 329



Disagreement Matrix

Figure 5: **Disagreement matrix across 5 users in Roleplay.** Here we plot the disagreement of preferences for 5 users. There is a mix of users with high and low disagreement.

the viewpoint and asked it to answer the question adhering to the viewpoint presented. For example, if you consider the question, "How can I learn to cook a delicious meal?", one viewpoint here could be "watching a youtube video", better suited for a younger, more tech savvy individual, whereas viewpoints such as "using a recipe book" or "taking a cooking class" may be better for an older population or those who would have the time or money to spend on a cooking class. This allowed for more diversity in the responses and resulting preferences.

Finally, we sampled responses from an ensemble of models with a high temperature, including those larger than the base model we fine-tuned such as Llama 3.3 70b (Grattafiori et al., 2024) and Gemma 2 27b (Team et al., 2024), allowing for better instruction following abilities of the fine-tuned model, than the Llama 3.2 3B we fine-tune.

**Encouraging task structure.** Meta-learning leverages a shared latent structure across tasks to adapt to a new task

quickly. The structure can be considered as similar feature
representations, function families, or transition dynamics
that the meta-learning algorithm can discover and leverage.
For a preference dataset, this structure can be represented
as the distribution of preferences across different users and
is controlled by the scoring function and the distribution of
responses.

337 One thing we controlled to enable better structure is the scor-338 ing function used to generate synthetic preferences. Firstly, 339 we wanted to ensure consistent preference labeling. We 340 use AI Feedback (Bai et al., 2022) to construct this, using 341 relative pairwise feedback for preference labels, akin to Al-342 pacaEval (Dubois et al., 2024b), as an alternative to absolute 343 rubric based scoring, which we found to be noisy and inaccurate. The preference label along with being conditioned 345 on the prompt, response, and general guidance on scoring, 346 is now also conditioned on the scoring user description and 347 additional scoring guidelines for user-aware preference labeling. Additionally, due to context length constraints, many 349 responses for our preference dataset are shorter than the in-350 struct model that we fine-tune from. Therefore, we prompt 351 the model to ignore this bias. Furthermore, we provide each 352 preference example to the model twice, flipping the order 353 of the responses, and keeping filtering out responses that 354 are not robust to order bias for both training and evaluation 355 (win rates). 356

357 Additionally, as mentioned above, in some cases, such as 358 with the Roleplay dataset, the user description is underspeci-359 fied, leading to challenges in labeling consistent preferences. 360 For example, if a user description does not have informa-361 tion about dietary preferences, inconsistency may arise for 362 labeling preferences about that topic. For instance, in one 363 preference pair, vegan cake recipes may be preferred but in another, steakhouses are preferred for date night. To fix this, 365 we take an iterative process to constructing user descriptions. Firstly, we start with a seed set of user descriptions gener-367 ated from the trait attributes. After generating questions and 368 responses based on these seed descriptions, we take a set of 369 question and response pairs. For each pair, we iteratively re-370 fine the user description by prompting a model like GPT-40 371 to either label the preference pair or if the user description 372 is insufficient, to randomly choose a preference and append 373 information to the description so a future scorer would make 374 the same decision. Finally, we utilize the updated user de-375 scription to relabel preferences for the set of questions and 376 responses allocated to that user with the labeling scheme 377 above. This fix for underspecification also helps the COT 378 prediction as predicting an underspecified user persona, can 379 lead to ambiguous generated descriptions.

Finally, we desire structured relationships between users.
To ensure this, we analyzed the disagreement (average difference of preference labels) of user's preferences across

384

prompts to understand where users agreed and disagreed, and regenerated data if this disagreement was too high across users. By having users with some overlap, metalearning algorithms can learn how to transfer knowledge effectively from one user to another. A sample disagreement plot for a subset of users in the Roleplay task can be found in Figure 5. We outline our full dataset generation process in Figure 4 in the Roleplay Task, starting from just a simple set of demographic traits.

Strategy	Mean Similarity $(\downarrow)$	Median Similarity $(\downarrow)$
Temp. $= 0.3$	0.96	0.97
Temp. $= 1.0$	0.94	0.95
Persona Steering (ours)	0.81	0.82
View Steering (ours)	0.78	0.78
Ensemble + View (ours)	0.71	0.73

Table 1: Comparison of diversity-inducing strategies as evaluated under ALOE (Wu et al., 2024).

Evaluating diversity and structure. We evaluate our design decisions with the following vignettes. For diversity, we measure semantic similarity using the dense score from the BGE-M3 model, following ALOE (Wu et al., 2024), on 100 randomly sampled prompts and 10 responses per prompt in the Roleplay task. As seen in Table 1, our proposed steering and ensembling mechanisms result in the base Llama 3.2 3B Instruct model exhibiting significantly reduced mean similarity. For structure, we estimate binary Shannon entropy of the preference label before and after iterative refinement. We condition on the persona and an unlabeled preference tuple (prompt and responses) and sample a preference label with a fixed temperature of 1.0 on 100 randomly sampled prompts from the Roleplay task with 100 pairs of personas and 10 samples per prompt. We use GPT-40 as the scoring model. Iterative persona refinement causes the entropy to drop from 0.64 nats to 0.13 nats, validating the efficacy of this approach in inducing better persona-prompt-response consistency.

#### 7. Experimental Evaluation

**Baselines.** We compare FSPO against five baselines: (1) a base model generating user-agnostic responses, (2) fewshot prompting with a base model, following Meister et al. (2024), (3) few-shot supervised fine-tuning (Pref-FT) based off the maximum likelihood objective from GPO (Zhao et al., 2024), (4) prompting with an oracle user description following Persona Steering (Cheng et al., 2023), and (5) Rewards-in-Context (Yang et al., 2024b). Specifically, for (1) we use a standard instruct model that is prompted solely with the query, resulting in unconditioned responses. For (2) and (3), the base instruct model is provided with the same few-shot personalization examples as in FSPO, but

FSPO:	Few-Shot	<b>Optimization</b>	of Synthetic	Preferences	Effectively	Personalizes to	o Real Users

Method	Traine	d Interpolated
Llama 3.2 3B Instruct	50.0	50.0
4-shot Prompted	66.6	61.9
4-shot Pref-FT	66.5	66.1
4-shot FSPO (Ours)	78.4	71.3
8-shot Prompted	69.1	59.1
8-shot Pref-FT	65.6	70.7
8-shot FSPO (Ours)	80.4	73.6
8-shot FSPO + COT (C	Ours) 92.3	84.6
Table 2: Re	eview Winrates	
Method	ELIX-easy	ELIX-hard
Llama 3.2 3B Instruct	50.0	50.0
Few-shot Prompted	92.4	81.4

#### Table 4: GPT-40 Winrates on ELIX

Few-shot Pref-FT

FSPO (Ours)

403

404

405

406

407 (2) zero-shot predicts the preferred response and (3) is opti408 mized with SFT to increase the likelihood on the preferred
409 response. In (4), the base model is prompted with the oracle,
410 ground truth user description, representing an upper bound
411 on FSPO's performance.

91.2

97.8

82.9

91.8

412 Synthetic winrates. We first generate automated win rates 413 using the modified AlpacaEval procedure from Section 6. 414 In the ELIX task in Table 4, we study two levels of difficulty 415 (easy, hard), where we find a consistent improvement of 416 FSPO over baselines. Next, in Table 2 for the Review task, 417 on both Trained and Interpolated Users, FSPO allows for bet-418 ter performance on held-out questions. Finally, in Table 3, 419 we study Roleplay, scaling to 1500 real users, seeing a win 420 rate of 82.6% on both held-out users and questions. Also, 421 COT closes the gap to the oracle response, effectively recov-422 ering the ground-truth user description. In Appendix A.2, 423 sample generations from FSPO show effective personaliza-424 tion to the oracle user description. Given this result, can we 425 personalize to real people? 426

427 Preliminary human study. We evaluate our model trained 428 on the Roleplay task by personalizing responses for real hu-429 man participants. We build a data collection app (Figure 7), 430 interacting with a user in two stages. First, we ask partici-431 pants to label preference pairs, used as the few-shot exam-432 ples in FSPO. Then, for held out questions, we show a user a set of two responses: (1) a response from FSPO personal-433 ized based on their preferences and (2) a baseline response. 434 435 Prolific is used to recruit a diverse set of study participants, 436 evenly split across genders and continents, corresponding 437 to the traits used to construct user descriptions. Question 438 and response order is randomized to remove confounding 439

Method	Winrate (%)
Llama 3.2 3B Instruct	50.0
IPO	72.4
Few-shot Prompting	63.2
Few-shot Pref-FT (GPO (Zhao et al., 2024))	62.8
RIC (Yang et al., 2024b)	53.3
FSPO (Ours, DPO)	81.3
FSPO (Ours, IPO)	82.6
FSPO + COT (Ours)	90.3
Oracle (prompt w/ g.t. persona)	90.9

Table 3: GPT-4c	Winrates	on Rolepla	y (1500 users)
-----------------	----------	------------	----------------

<b>Baseline Method</b>	Winrate (%)
FSPO vs Base	71.2
FSPO vs SFT	72.3

Table 5: Roleplay: Human Eval Winrates

factors. We evaluate with 25 users and 11 questions. As seen in Figure 5, we find that FSPO has a 71% win rate over the Base model and a 72% win rate over an SFT model trained on diverse viewpoints from the preference dataset.

#### 8. Discussion and Conclusion

We introduce FSPO, a novel framework for eliciting personalization in language models for open-ended question answering that models a distribution of reward functions to capture diverse human preferences. Our approach leverages meta-learning for rapid adaptation to each user, addressing limitations of conventional reward modeling techniques that learn from aggregated preferences. Through rigorous evaluation in 3 domains, we demonstrate that FSPO's generations are consistent with user context and preferred by real human users. Our findings also underscore the importance of diversity and structure in synthetic personalized preference datasets to bridge the Sim2Real gap. Overall, FSPO is a step towards developing more inclusive, user-centric language models.

#### 8.1. Limitations, Potential Risks, and Societal Impact

Key limitations include ethical concerns like reinforcing user biases, requiring future work on safeguards. The preliminary human study, though a first for personalized openended QA, had controlled elements and needs further ablation. Computational constraints restricted us to smaller models with limited context; investigating larger, more capable models is an open question, especially for Chain-of-Thought processes.

#### 440 **References**

- 441 Y. Bai, S. Kadavath, S. Kundu, A. Askell, J. Kernion, 442 A. Jones, A. Chen, A. Goldie, A. Mirhoseini, C. McK-443 innon, C. Chen, C. Olsson, C. Olah, D. Hernan-444 dez, D. Drain, D. Ganguli, D. Li, E. Tran-Johnson, 445 E. Perez, J. Kerr, J. Mueller, J. Ladish, J. Landau, 446 K. Ndousse, K. Lukosuite, L. Lovitt, M. Sellitto, N. El-447 hage, N. Schiefer, N. Mercado, N. DasSarma, R. Lasenby, 448 R. Larson, S. Ringer, S. Johnston, S. Kravec, S. E. 449 Showk, S. Fort, T. Lanham, T. Telleen-Lawton, T. Con-450 erly, T. Henighan, T. Hume, S. R. Bowman, Z. Hatfield-451 Dodds, B. Mann, D. Amodei, N. Joseph, S. McCandlish, 452 T. Brown, and J. Kaplan. Constitutional ai: Harmlessness 453 from ai feedback, 2022. URL https://arxiv.org/ 454 abs/2212.08073. 455
- R. A. Bradley and M. E. Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.
  ISSN 00063444. URL http://www.jstor.org/ stable/2334029.
- T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Ka-462 plan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sas-463 try, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, 464 T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, 465 C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, 466 S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, 467 A. Radford, I. Sutskever, and D. Amodei. Language 468 models are few-shot learners, 2020. 469
- 470 S. Casper, X. Davies, C. Shi, T. K. Gilbert, J. Scheurer, 471 J. Rando, R. Freedman, T. Korbak, D. Lindner, P. Freire, 472 T. T. Wang, S. Marks, C.-R. Segerie, M. Carroll, 473 A. Peng, P. Christoffersen, M. Damani, S. Slocum, 474 U. Anwar, A. Siththaranjan, M. Nadeau, E. J. Michaud, 475 J. Pfau, D. Krasheninnikov, X. Chen, L. Langosco, 476 P. Hase, E. Bivik, A. Dragan, D. Krueger, D. Sadigh, 477 and D. Hadfield-Menell. Open problems and funda-478 mental limitations of reinforcement learning from hu-479 man feedback. Transactions on Machine Learning Re-480 search, 2023. ISSN 2835-8856. URL https:// 481 openreview.net/forum?id=bx24KpJ4Eb. Sur-482 vey Certification. 483
- 484
  485
  486
  486
  487
  487
  487
  488
  489
  480
  480
  480
  481
  481
  482
  483
  484
  485
  486
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
  487
- Y. Chen, R. Zhong, S. Zha, G. Karypis, and H. He. Metalearning via language model in-context tuning, 2022.
  URL https://arxiv.org/abs/2110.07814.

488

M. Cheng, E. Durmus, and D. Jurafsky. Marked personas:
 Using natural language prompts to measure stereotypes

in language models. In A. Rogers, J. Boyd-Graber, and N. Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1504–1532, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.84. URL https:// aclanthology.org/2023.acl-long.84/.

- ContextualAI. Human-centered loss functions (halos), 2024. URL https://github.com/ContextualAI/ HALOs.
- Y. Duan, J. Schulman, X. Chen, P. L. Bartlett, I. Sutskever, and P. Abbeel. Fast reinforcement learning via slow reinforcement learning. *arXiv preprint arXiv:1611.02779*, 2016.
- Y. Dubois, B. Galambosi, P. Liang, and T. B. Hashimoto. Length-controlled alpacaeval: A simple way to debias automatic evaluators, 2024a. URL https://arxiv. org/abs/2404.04475.
- Y. Dubois, X. Li, R. Taori, T. Zhang, I. Gulrajani, J. Ba, C. Guestrin, P. Liang, and T. B. Hashimoto. Alpacafarm: A simulation framework for methods that learn from human feedback, 2024b.
- M. Gheshlaghi Azar, M. Rowland, B. Piot, D. Guo, D. Calandriello, M. Valko, and R. Munos. A General Theoretical Paradigm to Understand Learning from Human Preferences. *arXiv e-prints*, art. arXiv:2310.12036, Oct. 2023. doi: 10.48550/arXiv.2310.12036.
- Goodreads. Goodreads: Book reviews, recommendations, and discussion, 2025. URL https://www. goodreads.com/. Accessed: 2025-02-15.
- A. Grattafiori, A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Vaughan, A. Yang, A. Fan, A. Goyal, A. Hartshorn, A. Yang, A. Mitra, A. Sravankumar, A. Korenev, A. Hinsvark, A. Rao, A. Zhang, A. Rodriguez, A. Gregerson, A. Spataru, B. Roziere, B. Biron, B. Tang, B. Chern, C. Caucheteux, C. Nayak, C. Bi, C. Marra, C. McConnell, C. Keller, C. Touret, C. Wu, C. Wong, C. C. Ferrer, C. Nikolaidis, D. Allonsius, D. Song, D. Pintz, D. Livshits, D. Wyatt, D. Esiobu, D. Choudhary, D. Mahajan, D. Garcia-Olano, D. Perino, D. Hupkes, E. Lakomkin, E. AlBadawy, E. Lobanova, E. Dinan, E. M. Smith, F. Radenovic, F. Guzmán, F. Zhang, G. Synnaeve, G. Lee, G. L. Anderson, G. Thattai, G. Nail, G. Mialon, G. Pang, G. Cucurell, H. Nguyen, H. Korevaar, H. Xu, H. Touvron, I. Zarov, I. A. Ibarra, I. Kloumann, I. Misra, I. Evtimov, J. Zhang, J. Copet, J. Lee, J. Geffert, J. Vranes, J. Park, J. Mahadeokar, J. Shah, J. van der Linde, J. Billock, J. Hong, J. Lee, J. Fu, J. Chi, J. Huang, J. Liu, J. Wang, J. Yu, J. Bitton, J. Spisak, J. Park, J. Rocca,

495 J. Johnstun, J. Saxe, J. Jia, K. V. Alwala, K. Prasad, K. Up-496 asani, K. Plawiak, K. Li, K. Heafield, K. Stone, K. El-497 Arini, K. Iyer, K. Malik, K. Chiu, K. Bhalla, K. Lakho-498 tia, L. Rantala-Yeary, L. van der Maaten, L. Chen, 499 L. Tan, L. Jenkins, L. Martin, L. Madaan, L. Malo, 500 L. Blecher, L. Landzaat, L. de Oliveira, M. Muzzi, M. Pa-501 supuleti, M. Singh, M. Paluri, M. Kardas, M. Tsim-502 poukelli, M. Oldham, M. Rita, M. Pavlova, M. Kambadur, 503 M. Lewis, M. Si, M. K. Singh, M. Hassan, N. Goyal, 504 N. Torabi, N. Bashlykov, N. Bogoychev, N. Chatterji, 505 N. Zhang, O. Duchenne, O. Celebi, P. Alrassy, P. Zhang, 506 P. Li, P. Vasic, P. Weng, P. Bhargava, P. Dubal, P. Kr-507 ishnan, P. S. Koura, P. Xu, Q. He, Q. Dong, R. Srini-508 vasan, R. Ganapathy, R. Calderer, R. S. Cabral, R. Sto-509 jnic, R. Raileanu, R. Maheswari, R. Girdhar, R. Patel, 510 R. Sauvestre, R. Polidoro, R. Sumbaly, R. Taylor, R. Silva, 511 R. Hou, R. Wang, S. Hosseini, S. Chennabasappa, 512 S. Singh, S. Bell, S. S. Kim, S. Edunov, S. Nie, S. Narang, 513 S. Raparthy, S. Shen, S. Wan, S. Bhosale, S. Zhang, 514 S. Vandenhende, S. Batra, S. Whitman, S. Sootla, S. Col-515 lot, S. Gururangan, S. Borodinsky, T. Herman, T. Fowler, 516 T. Sheasha, T. Georgiou, T. Scialom, T. Speckbacher, 517 T. Mihaylov, T. Xiao, U. Karn, V. Goswami, V. Gupta, 518 V. Ramanathan, V. Kerkez, V. Gonguet, V. Do, V. Vo-519 geti, V. Albiero, V. Petrovic, W. Chu, W. Xiong, W. Fu, 520 W. Meers, X. Martinet, X. Wang, X. Wang, X. E. Tan, 521 X. Xia, X. Xie, X. Jia, X. Wang, Y. Goldschlag, Y. Gaur, 522 Y. Babaei, Y. Wen, Y. Song, Y. Zhang, Y. Li, Y. Mao, Z. D. 523 Coudert, Z. Yan, Z. Chen, Z. Papakipos, A. Singh, A. Sri-524 vastava, A. Jain, A. Kelsey, A. Shajnfeld, A. Gangidi, 525 A. Victoria, A. Goldstand, A. Menon, A. Sharma, A. Boe-526 senberg, A. Baevski, A. Feinstein, A. Kallet, A. Sangani, 527 A. Teo, A. Yunus, A. Lupu, A. Alvarado, A. Caples, 528 A. Gu, A. Ho, A. Poulton, A. Ryan, A. Ramchandani, 529 A. Dong, A. Franco, A. Goyal, A. Saraf, A. Chowd-530 hury, A. Gabriel, A. Bharambe, A. Eisenman, A. Yazdan, 531 B. James, B. Maurer, B. Leonhardi, B. Huang, B. Lovd, 532 B. D. Paola, B. Paranjape, B. Liu, B. Wu, B. Ni, B. Han-533 cock, B. Wasti, B. Spence, B. Stojkovic, B. Gamido, 534 B. Montalvo, C. Parker, C. Burton, C. Mejia, C. Liu, 535 C. Wang, C. Kim, C. Zhou, C. Hu, C.-H. Chu, C. Cai, 536 C. Tindal, C. Feichtenhofer, C. Gao, D. Civin, D. Beaty, 537 D. Kreymer, D. Li, D. Adkins, D. Xu, D. Testuggine, 538 D. David, D. Parikh, D. Liskovich, D. Foss, D. Wang, 539 D. Le, D. Holland, E. Dowling, E. Jamil, E. Montgomery, 540 E. Presani, E. Hahn, E. Wood, E.-T. Le, E. Brinkman, 541 E. Arcaute, E. Dunbar, E. Smothers, F. Sun, F. Kreuk, 542 F. Tian, F. Kokkinos, F. Ozgenel, F. Caggioni, F. Kanayet, 543 F. Seide, G. M. Florez, G. Schwarz, G. Badeer, G. Swee, 544 G. Halpern, G. Herman, G. Sizov, Guangyi, Zhang, 545 G. Lakshminarayanan, H. Inan, H. Shojanazeri, H. Zou, 546 H. Wang, H. Zha, H. Habeeb, H. Rudolph, H. Suk, 547 H. Aspegren, H. Goldman, H. Zhan, I. Damlaj, I. Moly-548 bog, I. Tufanov, I. Leontiadis, I.-E. Veliche, I. Gat, 549

J. Weissman, J. Geboski, J. Kohli, J. Lam, J. Asher, J.-B. Gaya, J. Marcus, J. Tang, J. Chan, J. Zhen, J. Reizenstein, J. Teboul, J. Zhong, J. Jin, J. Yang, J. Cummings, J. Carvill, J. Shepard, J. McPhie, J. Torres, J. Ginsburg, J. Wang, K. Wu, K. H. U, K. Saxena, K. Khandelwal, K. Zand, K. Matosich, K. Veeraraghavan, K. Michelena, K. Li, K. Jagadeesh, K. Huang, K. Chawla, K. Huang, L. Chen, L. Garg, L. A, L. Silva, L. Bell, L. Zhang, L. Guo, L. Yu, L. Moshkovich, L. Wehrstedt, M. Khabsa, M. Avalani, M. Bhatt, M. Mankus, M. Hasson, M. Lennie, M. Reso, M. Groshev, M. Naumov, M. Lathi, M. Keneally, M. Liu, M. L. Seltzer, M. Valko, M. Restrepo, M. Patel, M. Vyatskov, M. Samvelyan, M. Clark, M. Macey, M. Wang, M. J. Hermoso, M. Metanat, M. Rastegari, M. Bansal, N. Santhanam, N. Parks, N. White, N. Bawa, N. Singhal, N. Egebo, N. Usunier, N. Mehta, N. P. Laptev, N. Dong, N. Cheng, O. Chernoguz, O. Hart, O. Salpekar, O. Kalinli, P. Kent, P. Parekh, P. Saab, P. Balaji, P. Rittner, P. Bontrager, P. Roux, P. Dollar, P. Zvyagina, P. Ratanchandani, P. Yuvraj, Q. Liang, R. Alao, R. Rodriguez, R. Ayub, R. Murthy, R. Nayani, R. Mitra, R. Parthasarathy, R. Li, R. Hogan, R. Battey, R. Wang, R. Howes, R. Rinott, S. Mehta, S. Siby, S. J. Bondu, S. Datta, S. Chugh, S. Hunt, S. Dhillon, S. Sidorov, S. Pan, S. Mahajan, S. Verma, S. Yamamoto, S. Ramaswamy, S. Lindsay, S. Lindsay, S. Feng, S. Lin, S. C. Zha, S. Patil, S. Shankar, S. Zhang, S. Zhang, S. Wang, S. Agarwal, S. Sajuyigbe, S. Chintala, S. Max, S. Chen, S. Kehoe, S. Satterfield, S. Govindaprasad, S. Gupta, S. Deng, S. Cho, S. Virk, S. Subramanian, S. Choudhury, S. Goldman, T. Remez, T. Glaser, T. Best, T. Koehler, T. Robinson, T. Li, T. Zhang, T. Matthews, T. Chou, T. Shaked, V. Vontimitta, V. Ajayi, V. Montanez, V. Mohan, V. S. Kumar, V. Mangla, V. Ionescu, V. Poenaru, V. T. Mihailescu, V. Ivanov, W. Li, W. Wang, W. Jiang, W. Bouaziz, W. Constable, X. Tang, X. Wu, X. Wang, X. Wu, X. Gao, Y. Kleinman, Y. Chen, Y. Hu, Y. Jia, Y. Qi, Y. Li, Y. Zhang, Y. Zhang, Y. Adi, Y. Nam, Yu, Wang, Y. Zhao, Y. Hao, Y. Qian, Y. Li, Y. He, Z. Rait, Z. DeVito, Z. Rosnbrick, Z. Wen, Z. Yang, Z. Zhao, and Z. Ma. The llama 3 herd of models, 2024.

- C. Gulcehre, T. L. Paine, S. Srinivasan, K. Konyushkova, L. Weerts, A. Sharma, A. Siddhant, A. Ahern, M. Wang, C. Gu, et al. Reinforced self-training (rest) for language modeling. *arXiv preprint arXiv:2308.08998*, 2023.
- K. Hsu, S. Levine, and C. Finn. Unsupervised learning via meta-learning. In *International Conference on Learning Representations*, 2019.
- IMDb. Imdb: Ratings, reviews, and where to watch the best movies & tv shows, 2025. URL https://www.imdb.com/. Accessed: 2025-02-15.

- J. Jang, S. Kim, B. Y. Lin, Y. Wang, J. Hessel, L. Zettlemoyer, H. Hajishirzi, Y. Choi, and P. Ammanabrolu.
  Personalized soups: Personalized large language model alignment via post-hoc parameter merging, 2023. URL https://arxiv.org/abs/2310.11564.
- H. R. Kirk, A. Whitefield, P. Röttger, A. Bean, K. Margatina, J. Ciro, R. Mosquera, M. Bartolo, A. Williams, H. He, B. Vidgen, and S. A. Hale. The prism alignment dataset: What participatory, representative and individualised human feedback reveals about the subjective and multicultural alignment of large language models, 2024. URL https://arxiv.org/abs/2404.16019.
- W. Kwon, Z. Li, S. Zhuang, Y. Sheng, L. Zheng, C. H. Yu,
  J. E. Gonzalez, H. Zhang, and I. Stoica. Efficient memory
  management for large language model serving with pagedattention, 2023. URL https://arxiv.org/abs/
  2309.06180.
- Y. Lee, J. Williams, H. Marklund, A. Sharma, E. Mitchell, A. Singh, and C. Finn. Test-time alignment via hypothesis reweighting, 2024. URL https://arxiv.org/ abs/2412.08812.
- H. Li, Q. Dong, Z. Tang, C. Wang, X. Zhang, H. Huang,
  S. Huang, X. Huang, Z. Huang, D. Zhang, Y. Gu,
  X. Cheng, X. Wang, S.-Q. Chen, L. Dong, W. Lu, Z. Sui,
  B. Wang, W. Lam, and F. Wei. Synthetic data (almost)
  from scratch: Generalized instruction tuning for language models, 2024. URL https://arxiv.org/
  abs/2402.13064.
- 582 S. Li, F. Xue, C. Baranwal, Y. Li, and Y. You. Sequence
  583 parallelism: Long sequence training from system perspec584 tive, 2022. URL https://arxiv.org/abs/2105.
  585 13120.
- 586 A. L. Maas, R. E. Daly, P. T. Pham, D. Huang, A. Y. Ng, and 587 C. Potts. Learning word vectors for sentiment analysis. 588 In D. Lin, Y. Matsumoto, and R. Mihalcea, editors, Pro-589 ceedings of the 49th Annual Meeting of the Association 590 for Computational Linguistics: Human Language Tech-591 nologies, pages 142-150, Portland, Oregon, USA, June 592 2011. Association for Computational Linguistics. URL 593 https://aclanthology.org/P11-1015/. 594
- V. Makoviychuk, L. Wawrzyniak, Y. Guo, M. Lu, K. Storey, M. Macklin, D. Hoeller, N. Rudin, A. Allshire, A. Handa, and G. State. Isaac gym: High performance gpubased physics simulation for robot learning, 2021. URL https://arxiv.org/abs/2108.10470.
- N. Meister, C. Guestrin, and T. Hashimoto. Benchmarking distributional alignment of large language models, 2024. URL https://arxiv.org/abs/2411.05403.

- I. Melnyk, Y. Mroueh, B. Belgodere, M. Rigotti, A. Nitsure, M. Yurochkin, K. Greenewald, J. Navratil, and J. Ross. Distributional preference alignment of llms via optimal transport, 2024. URL https://arxiv.org/abs/ 2406.05882.
- S. Min, M. Lewis, L. Zettlemoyer, and H. Hajishirzi. Metaicl: Learning to learn in context, 2022. URL https://arxiv.org/abs/2110.15943.
- N. Mishra, M. Rohaninejad, X. Chen, and P. Abbeel. A simple neural attentive meta-learner, 2018. URL https://arxiv.org/abs/1707.03141.
- MyAnimeList. Myanimelist: Track, discover, and discuss anime & manga, 2025. URL https:// myanimelist.net/. Accessed: 2025-02-15.
- OpenAI, :, A. Hurst, A. Lerer, A. P. Goucher, A. Perelman, A. Ramesh, A. Clark, A. Ostrow, A. Welihinda, A. Hayes, A. Radford, A. Mądry, A. Baker-Whitcomb, A. Beutel, A. Borzunov, A. Carney, A. Chow, A. Kirillov, A. Nichol, A. Paino, A. Renzin, A. T. Passos, A. Kirillov, A. Christakis, A. Conneau, A. Kamali, A. Jabri, A. Moyer, A. Tam, A. Crookes, A. Tootoochian, A. Tootoonchian, A. Kumar, A. Vallone, A. Karpathy, A. Braunstein, A. Cann, A. Codispoti, A. Galu, A. Kondrich, A. Tulloch, A. Mishchenko, A. Baek, A. Jiang, A. Pelisse, A. Woodford, A. Gosalia, A. Dhar, A. Pantuliano, A. Nayak, A. Oliver, B. Zoph, B. Ghorbani, B. Leimberger, B. Rossen, B. Sokolowsky, B. Wang, B. Zweig, B. Hoover, B. Samic, B. McGrew, B. Spero, B. Giertler, B. Cheng, B. Lightcap, B. Walkin, B. Quinn, B. Guarraci, B. Hsu, B. Kellogg, B. Eastman, C. Lugaresi, C. Wainwright, C. Bassin, C. Hudson, C. Chu, C. Nelson, C. Li, C. J. Shern, C. Conger, C. Barette, C. Voss, C. Ding, C. Lu, C. Zhang, C. Beaumont, C. Hallacy, C. Koch, C. Gibson, C. Kim, C. Choi, C. McLeavey, C. Hesse, C. Fischer, C. Winter, C. Czarnecki, C. Jarvis, C. Wei, C. Koumouzelis, D. Sherburn, D. Kappler, D. Levin, D. Levy, D. Carr, D. Farhi, D. Mely, D. Robinson, D. Sasaki, D. Jin, D. Valladares, D. Tsipras, D. Li, D. P. Nguyen, D. Findlay, E. Oiwoh, E. Wong, E. Asdar, E. Proehl, E. Yang, E. Antonow, E. Kramer, E. Peterson, E. Sigler, E. Wallace, E. Brevdo, E. Mays, F. Khorasani, F. P. Such, F. Raso, F. Zhang, F. von Lohmann, F. Sulit, G. Goh, G. Oden, G. Salmon, G. Starace, G. Brockman, H. Salman, H. Bao, H. Hu, H. Wong, H. Wang, H. Schmidt, H. Whitney, H. Jun, H. Kirchner, H. P. de Oliveira Pinto, H. Ren, H. Chang, H. W. Chung, I. Kivlichan, I. O'Connell, I. O'Connell, I. Osband, I. Silber, I. Sohl, I. Okuyucu, I. Lan, I. Kostrikov, I. Sutskever, I. Kanitscheider, I. Gulrajani, J. Coxon, J. Menick, J. Pachocki, J. Aung, J. Betker, J. Crooks, J. Lennon, J. Kiros, J. Leike, J. Park, J. Kwon, J. Phang, J. Teplitz, J. Wei, J. Wolfe, J. Chen, J. Harris, J. Varavva,

605 J. G. Lee, J. Shieh, J. Lin, J. Yu, J. Weng, J. Tang, J. Yu, 606 J. Jang, J. Q. Candela, J. Beutler, J. Landers, J. Parish, 607 J. Heidecke, J. Schulman, J. Lachman, J. McKay, J. Ue-608 sato, J. Ward, J. W. Kim, J. Huizinga, J. Sitkin, J. Kraai-609 jeveld, J. Gross, J. Kaplan, J. Snyder, J. Achiam, J. Jiao, 610 J. Lee, J. Zhuang, J. Harriman, K. Fricke, K. Hayashi, 611 K. Singhal, K. Shi, K. Karthik, K. Wood, K. Rim-612 bach, K. Hsu, K. Nguyen, K. Gu-Lemberg, K. Button, 613 K. Liu, K. Howe, K. Muthukumar, K. Luther, L. Ah-614 mad, L. Kai, L. Itow, L. Workman, L. Pathak, L. Chen, 615 L. Jing, L. Guy, L. Fedus, L. Zhou, L. Mamitsuka, 616 L. Weng, L. McCallum, L. Held, L. Ouyang, L. Feuvrier, 617 L. Zhang, L. Kondraciuk, L. Kaiser, L. Hewitt, L. Metz, 618 L. Doshi, M. Aflak, M. Simens, M. Boyd, M. Thompson, 619 M. Dukhan, M. Chen, M. Gray, M. Hudnall, M. Zhang, 620 M. Aljubeh, M. Litwin, M. Zeng, M. Johnson, M. Shetty, 621 M. Gupta, M. Shah, M. Yatbaz, M. J. Yang, M. Zhong, 622 M. Glaese, M. Chen, M. Janner, M. Lampe, M. Petrov, 623 M. Wu, M. Wang, M. Fradin, M. Pokrass, M. Cas-624 tro, M. O. T. de Castro, M. Pavlov, M. Brundage, 625 M. Wang, M. Khan, M. Murati, M. Bavarian, M. Lin, 626 M. Yesildal, N. Soto, N. Gimelshein, N. Cone, N. Stau-627 dacher, N. Summers, N. LaFontaine, N. Chowdhury, 628 N. Ryder, N. Stathas, N. Turley, N. Tezak, N. Felix, N. Kudige, N. Keskar, N. Deutsch, N. Bundick, N. Puck-629 630 ett, O. Nachum, O. Okelola, O. Boiko, O. Murk, O. Jaffe, 631 O. Watkins, O. Godement, O. Campbell-Moore, P. Chao, 632 P. McMillan, P. Belov, P. Su, P. Bak, P. Bakkum, P. Deng, 633 P. Dolan, P. Hoeschele, P. Welinder, P. Tillet, P. Pronin, 634 P. Tillet, P. Dhariwal, Q. Yuan, R. Dias, R. Lim, R. Arora, 635 R. Troll, R. Lin, R. G. Lopes, R. Puri, R. Miyara, 636 R. Leike, R. Gaubert, R. Zamani, R. Wang, R. Don-637 nelly, R. Honsby, R. Smith, R. Sahai, R. Ramchandani, 638 R. Huet, R. Carmichael, R. Zellers, R. Chen, R. Chen, 639 R. Nigmatullin, R. Cheu, S. Jain, S. Altman, S. Schoen-640 holz, S. Toizer, S. Miserendino, S. Agarwal, S. Culver, 641 S. Ethersmith, S. Gray, S. Grove, S. Metzger, S. Her-642 mani, S. Jain, S. Zhao, S. Wu, S. Jomoto, S. Wu, Shuaiqi, 643 Xia, S. Phene, S. Papay, S. Narayanan, S. Coffey, S. Lee, 644 S. Hall, S. Balaji, T. Broda, T. Stramer, T. Xu, T. Gogi-645 neni, T. Christianson, T. Sanders, T. Patwardhan, T. Cun-646 ninghman, T. Degry, T. Dimson, T. Raoux, T. Shad-647 well, T. Zheng, T. Underwood, T. Markov, T. Sherbakov, 648 T. Rubin, T. Stasi, T. Kaftan, T. Heywood, T. Peterson, 649 T. Walters, T. Eloundou, V. Qi, V. Moeller, V. Monaco, 650 V. Kuo, V. Fomenko, W. Chang, W. Zheng, W. Zhou, 651 W. Manassra, W. Sheu, W. Zaremba, Y. Patil, Y. Qian, 652 Y. Kim, Y. Cheng, Y. Zhang, Y. He, Y. Zhang, Y. Jin, 653 Y. Dai, and Y. Malkov. Gpt-40 system card, 2024. URL 654 https://arxiv.org/abs/2410.21276. 655

L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. L. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, J. Schulman, J. Hilton, F. Kelton, L. Miller,

M. Simens, A. Askell, P. Welinder, P. Christiano, J. Leike, and R. Lowe. Training language models to follow instructions with human feedback. *arXiv e-prints*, art. arXiv:2203.02155, Mar. 2022. doi: 10.48550/arXiv.2203. 02155.

- L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. L. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, J. Schulman, J. Hilton, F. Kelton, L. Miller, M. Simens, A. Askell, P. Welinder, P. Christiano, J. Leike, and R. Lowe. Training language models to follow instructions with human feedback, 2022.
- S. Poddar, Y. Wan, H. Ivison, A. Gupta, and N. Jaques. Personalizing reinforcement learning from human feedback with variational preference learning, 2024. URL https://arxiv.org/abs/2408.10075.
- Qwen, :, A. Yang, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Li, D. Liu, F. Huang, H. Wei, H. Lin, J. Yang, J. Tu, J. Zhang, J. Yang, J. Yang, J. Zhou, J. Lin, K. Dang, K. Lu, K. Bao, K. Yang, L. Yu, M. Li, M. Xue, P. Zhang, Q. Zhu, R. Men, R. Lin, T. Li, T. Tang, T. Xia, X. Ren, X. Ren, Y. Fan, Y. Su, Y. Zhang, Y. Wan, Y. Liu, Z. Cui, Z. Zhang, and Z. Qiu. Qwen2.5 technical report, 2025. URL https://arxiv.org/abs/2412.15115.
- R. Rafailov, A. Sharma, E. Mitchell, S. Ermon, C. D. Manning, and C. Finn. Direct preference optimization: Your language model is secretly a reward model. *arXiv preprint arXiv:2305.18290*, 2023.
- A. Santoro, S. Bartunov, M. Botvinick, D. Wierstra, and T. Lillicrap. One-shot learning with memory-augmented neural networks, 2016. URL https://arxiv.org/ abs/1605.06065.
- J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal Policy Optimization Algorithms. *arXiv e-prints*, art. arXiv:1707.06347, July 2017. doi: 10.48550/arXiv.1707.06347.
- O. Shaikh, M. Lam, J. Hejna, Y. Shao, M. Bernstein, and D. Yang. Show, don't tell: Aligning language models with demonstrated feedback, 2024. URL https:// arxiv.org/abs/2406.00888.
- A. Siththaranjan, C. Laidlaw, and D. Hadfield-Menell. Distributional preference learning: Understanding and accounting for hidden context in rlhf, 2024. URL https: //arxiv.org/abs/2312.08358.
- T. Sorensen, J. Moore, J. Fisher, M. Gordon, N. Mireshghallah, C. M. Rytting, A. Ye, L. Jiang, X. Lu, N. Dziri, T. Althoff, and Y. Choi. A roadmap to pluralistic alignment, 2024. URL https://arxiv.org/abs/ 2402.05070.

660	R. S. Sutton, D. McAllester, S. Singh, and Y. Mansour.
661	Policy gradient methods for reinforcement learning
662	with function approximation. In S. Solla, T. Leen,
663	and K. Müller, editors, Advances in Neural Infor-
664	mation Processing Systems, volume 12. MIT Press,
665	1999. URL https://proceedings.neurips.
666	cc/paper_files/paper/1999/file/
667	464d828b85b0bed98e80ade0a5c43b0f-Paper.
668	pdf.

670 G. Team, M. Riviere, S. Pathak, P. G. Sessa, C. Hardin, 671 S. Bhupatiraju, L. Hussenot, T. Mesnard, B. Shahriari, 672 A. Ramé, J. Ferret, P. Liu, P. Tafti, A. Friesen, M. Cas-673 bon, S. Ramos, R. Kumar, C. L. Lan, S. Jerome, A. Tsit-674 sulin, N. Vieillard, P. Stanczyk, S. Girgin, N. Mom-675 chev, M. Hoffman, S. Thakoor, J.-B. Grill, B. Neyshabur, 676 O. Bachem, A. Walton, A. Severyn, A. Parrish, A. Ah-677 mad, A. Hutchison, A. Abdagic, A. Carl, A. Shen, 678 A. Brock, A. Coenen, A. Laforge, A. Paterson, B. Bas-679 tian, B. Piot, B. Wu, B. Royal, C. Chen, C. Kumar, 680 C. Perry, C. Welty, C. A. Choquette-Choo, D. Sinopal-681 nikov, D. Weinberger, D. Vijaykumar, D. Rogozińska, 682 D. Herbison, E. Bandy, E. Wang, E. Noland, E. Mor-683 eira, E. Senter, E. Eltyshev, F. Visin, G. Rasskin, G. Wei, 684 G. Cameron, G. Martins, H. Hashemi, H. Klimczak-685 Plucińska, H. Batra, H. Dhand, I. Nardini, J. Mein, 686 J. Zhou, J. Svensson, J. Stanway, J. Chan, J. P. Zhou, 687 J. Carrasqueira, J. Iljazi, J. Becker, J. Fernandez, J. van 688 Amersfoort, J. Gordon, J. Lipschultz, J. Newlan, J. yeong 689 Ji, K. Mohamed, K. Badola, K. Black, K. Millican, K. Mc-690 Donell, K. Nguyen, K. Sodhia, K. Greene, L. L. Sjoe-691 sund, L. Usui, L. Sifre, L. Heuermann, L. Lago, L. Mc-692 Nealus, L. B. Soares, L. Kilpatrick, L. Dixon, L. Martins, 693 M. Reid, M. Singh, M. Iverson, M. Görner, M. Velloso, 694 M. Wirth, M. Davidow, M. Miller, M. Rahtz, M. Wat-695 son, M. Risdal, M. Kazemi, M. Moynihan, M. Zhang, 696 M. Kahng, M. Park, M. Rahman, M. Khatwani, N. Dao, 697 N. Bardoliwalla, N. Devanathan, N. Dumai, N. Chauhan, 698 O. Wahltinez, P. Botarda, P. Barnes, P. Barham, P. Michel, 699 P. Jin, P. Georgiev, P. Culliton, P. Kuppala, R. Comanescu, 700 R. Merhej, R. Jana, R. A. Rokni, R. Agarwal, R. Mullins, S. Saadat, S. M. Carthy, S. Cogan, S. Perrin, S. M. R. Arnold, S. Krause, S. Dai, S. Garg, S. Sheth, S. Ronstrom, 703 S. Chan, T. Jordan, T. Yu, T. Eccles, T. Hennigan, T. Ko-704 cisky, T. Doshi, V. Jain, V. Yadav, V. Meshram, V. Dhar-705 madhikari, W. Barkley, W. Wei, W. Ye, W. Han, W. Kwon, 706 X. Xu, Z. Shen, Z. Gong, Z. Wei, V. Cotruta, P. Kirk, A. Rao, M. Giang, L. Peran, T. Warkentin, E. Collins, J. Barral, Z. Ghahramani, R. Hadsell, D. Sculley, J. Banks, 709 A. Dragan, S. Petrov, O. Vinyals, J. Dean, D. Has-710 sabis, K. Kavukcuoglu, C. Farabet, E. Buchatskaya, S. Borgeaud, N. Fiedel, A. Joulin, K. Kenealy, R. Dadashi, 712 and A. Andreev. Gemma 2: Improving open lan-713 guage models at a practical size, 2024. URL https: 714

//arxiv.org/abs/2408.00118.

- J. Tobin, L. Biewald, R. Duan, M. Andrychowicz, A. Handa, V. Kumar, B. McGrew, J. Schneider, P. Welinder, W. Zaremba, and P. Abbeel. Domain randomization and generative models for robotic grasping, 2018. URL https://arxiv.org/abs/1710.06425.
- J. X. Wang, Z. Kurth-Nelson, D. Tirumala, H. Soyer, J. Z. Leibo, R. Munos, C. Blundell, D. Kumaran, and M. Botvinick. Learning to reinforcement learn. *arXiv* preprint arXiv:1611.05763, 2016.
- S. Wu, M. Fung, C. Qian, J. Kim, D. Hakkani-Tur, and H. Ji. Aligning llms with individual preferences via interaction. arXiv preprint arXiv:2410.03642, 2024.
- A. Yang, J. Yang, A. Ibrahim, X. Xie, B. Tang, G. Sizov, J. Reizenstein, J. Park, and J. Huang. Context parallelism for scalable million-token inference, 2024a. URL https://arxiv.org/abs/2411.01783.
- R. Yang, X. Pan, F. Luo, S. Qiu, H. Zhong, D. Yu, and J. Chen. Rewards-in-context: Multi-objective alignment of foundation models with dynamic preference adjustment. *International Conference on Machine Learning*, 2024b.
- Z. Yang, Y. Chen, J. Wang, S. Manivasagam, W.-C. Ma, A. J. Yang, and R. Urtasun. Unisim: A neural closedloop sensor simulator, 2023. URL https://arxiv. org/abs/2308.01898.
- M. Yin, G. Tucker, M. Zhou, S. Levine, and C. Finn. Meta-learning without memorization. *arXiv preprint arXiv:1912.03820*, 2019.
- L. Yu, W. Jiang, H. Shi, J. Yu, Z. Liu, Y. Zhang, J. T. Kwok, Z. Li, A. Weller, and W. Liu. Metamath: Bootstrap your own mathematical questions for large language models, 2024. URL https://arxiv.org/abs/2309. 12284.
- W. Yuan, R. Yuanzhe Pang, K. Cho, X. Li, S. Sukhbaatar, J. Xu, and J. Weston. Self-Rewarding Language Models. *arXiv e-prints*, art. arXiv:2401.10020, Jan. 2024. doi: 10.48550/arXiv.2401.10020.
- L. Zhang, A. Hosseini, H. Bansal, M. Kazemi, A. Kumar, and R. Agarwal. Generative verifiers: Reward modeling as next-token prediction, 2024. URL https://arxiv. org/abs/2408.15240.
- S. Zhao, J. Dang, and A. Grover. Group preference optimization: Few-shot alignment of large language models, 2024. URL https://arxiv.org/abs/2310.11523.

715	L. Zheng, L. Yin, Z. Xie, C. Sun, J. Huang, C. H. Yu,
716	S. Cao, C. Kozyrakis, I. Stoica, J. E. Gonzalez, C. Bar-
717	rett, and Y. Sheng. Sglang: Efficient execution of struc-
718	tured language model programs, 2024. URL https:
719	//arxiv.org/abs/2312.07104.
720	
721	
722	
723	
724	
725	
726	
727	
728	
720	
720	
730	
731	
732	
733	
734	
735	
736	
737	
738	
739	
740	
741	
742	
743	
744	
745	
746	
747	
748	
749	
750	
751	
752	
753	
754	
755	
756	
757	
758	
759	
760	
761	
762	
763	
764	
765	
766	
767	
768	
769	

# 770 A. Appendix

787

794

807

814

815 816

817

818 819

# A.1. Limitations, Potential Risks, and Societal Impact

772 There are several limitations and potential risks. One limitation pertains to the ethical and fairness considerations of 773 personalization. While FSPO improves inclusivity by modeling diverse preferences, the risk of reinforcing user biases 774 (echo chambers) or inadvertently amplifying harmful viewpoints requires careful scrutiny. Future work should explore 775 mechanisms to balance personalization with ethical safeguards, ensuring that models remain aligned with fairness principles 776 while respecting user individuality. Additionally, our human study was preliminary with control over the questions that 777 a user may ask, format normalization where formatting details such as markdown are removed, and view normalization 778 comparing the same number of viewpoints for both FSPO and the baselines. However, to the best of our knowledge, we 779 are the first approach to perform such a human study for personalization to open-ended question answering. Future work 780 should do further ablations with human evaluation for personalization. Additionally, due to compute constraints, we work 781 with models in the parameter range of 3B (specifically Llama 3.2 Instruct 3B) with a limited context window of 128K, and 782 without context optimization such as sequence parallelism (Li et al., 2022; Yang et al., 2024a), further limiting the effective 783 context window. It is an open question on how fine-tuning base models with better long-context and reasoning capabilities 784 would help with FSPO for personalization, such as the 2M context window of Gemini Flash Thinking models, especially in 785 the case of COT. 786

# A.1.1. ADDITIONAL ABLATIONS

Percentage of Preference Data	Winrate (%)
10%	70.1
25%	69.5
50%	78.3
100%	82.6 (reported)

Table 6: Ablation study showing the effect of varying the percentage of preference data per user for FSPO without COT on
 roleplay task winrates with held-out synthetic users.

Number of Few-Shot Examples	Winrate (%)
1	65.7
2	69.3
4	72.1
8	82.6 (reported)

Table 7: Ablation study showing the impact of varying the number of few-shot examples per user on roleplay task performance.

# <sup>808</sup> A.2. Sample Personalized Responses

We provide sample responses from FSPO in Figure 6 across the 3 tasks that were studied (ELIX, Reviews, and Roleplay). We additionally include the oracle scoring description for each response, demonstrating FSPO's ability to personalize. Additionally, for the Review task, we show responses for interpolated users which were not found in the preference dataset, showing the ability of FSPO to generalize and be controllable.

# A.2.1. ROLEPLAY COT

#### QUESTION :

What are some good volunteer opportunities for me?

# TRUE PERSONA:

A 53-year-old woman living in North America, an avid traveler and writer who
 shares her adventures through a popular travel blog. She loves spontaneous
 adventures and values long-term financial security. She is at a point in her
 career where opportunities for advancement might be time-sensitive. She

825	enjoys movies that challenge perspectives and spark discussions. She
826	prioritizes well-being over being fully informed. She is interested in the
827	van life movement. She appreciates creativity and variety in meal options.
828	She anious angaging in creative and tactile activities and prefers
020	she enjoys engaging in creative and tactific activities and prefers
829	traditional learning methods. She values compassionate and responsible
830	choices. She enjoys exploring new and dynamic experiences. She is concerned
831	about fast food. She values mental well-being and mindfulness. She enjoys
832	attending events and conferences to meet new people. She appreciates
022	prostigative and health. She arises shows with dark humar and successful
033	practicality and nearth. She enjoys shows with dark number and suspenserul
834	storytelling.
835	
836	COT GENERATED PERSONA:
837	Based on the examples provided the user appears to value a more holistic and
020	a second the examples provided, the user appears to varie a more instruction and
020	experiential approach to learning and personal growth. They seem to
839	prioritize human connection, nature, and personal development over technology
840	and material wealth.
841	
842	This user is likely someone who is interested in personal stories emotions and
012	averagion and a suident from their expension for Coluin Dist, emotions, and
043	experiences, as evident from their appreciation for Sylvia Plath s poetry and
844	the suggestion to weave storytelling into public speaking. They may also be
845	someone who values authenticity, simplicity, and community, as seen in their
846	suggestion of house swapping and volunteering as alternative vacation options
8/7	suggestion of nearest starting and totaliseting as attended to prove
0.40	•
848	
849	In terms of lifestyle, this user likely prefers a balanced and flexible approach
850	to life, allowing for a mix of work, personal growth, and leisure activities.
851	They may also be someone who values environmental sustainability and social
852	responsibility as hinted at by their consideration of eco-friendly home
052	responsibility, as infinited at by the in consideration of eco-intendry none
833	improvements and sustainable financial planning. Overall, this user seems to
854	be a thoughtful, empathetic, and nature-loving individual who prioritizes
855	human connections, experiences, and personal growth over material possessions
856	and technological advancements.
857	
050	
0.50	FINAL RESPONSE.
859	Based on your affinity for human connections, personal growth, and environmental
860	sustainability, I'd like to suggest the following volunteer opportunities
861	that might resonate with you:
862	
962	. Community Condon on Unken Formers Many enconingstions in your enco. offer
005	**Community Garden of Orban Farm**. Many organizations in your area offer
864	community gardens or urban farms where you can help cultivate and maintain a
865	shared space, promoting sustainability and food security. This allows you to
866	connect with nature and contribute to your community's well-being.
867	
868	A Natura Pasaryas or Conservation Efforts and Consider valuateoring with local
000	**Nature Reserves of Conservation Efforts **. Consider volunteering with focal
809	wildlife conservation organizations, such as bird watching tours, tree
870	planting initiatives, or habitat restoration projects. This will enable you
871	to work outdoors, connect with others who share your passion for nature, and
872	contribute to preserving the environment
872	controute to preserving the environment.
073	
8/4	**Community Centers or Environmental Education Programs**: Many community centers
875	offer various programs, such as gardening workshops, environmental education
876	classes, or sustainable living workshops. By volunteering at these centers.
877	you can share your knowledge learn from others and make a nositive impact
011	jou can share your knowledge, rearn from others, and make a positive impact

on your community.

# FSPO: Few-Shot Optimization of Synthetic Preferences Effectively Personalizes to Real Users

	Name	Value	
	Learning Rate (SFT/Pref-FT)	$1e^{-5} 1e^{-6} 1e^{-7}$	
	Learning Rate (JPO)	$1e^{-5}$ $1e^{-6}$ $1e^{-7}$	
	Beta (IPO)	0 1 0 05 0 01 <b>0 005</b> 0 001	
	Number of Shots	4 8	
	Model Name	Llama 3.2.3B Instruct (Grattafiori et al. 2024)	
		Elana 5.2 5D Instruct (Grananori et al., 2024)	
Ta	able 8: Sweep over hyperparame	eters for FSPO, recommended hyperparameters in bold	d.
**Animal She	lters or Wildlife Reh	abilitation Centers **: If you're an	animal lover
, conside	r volunteering at a l	ocal animal shelter or wildlife reha	abilitation
center. Y	ou'll have the opport	unity to care for animals, interact	with like –
minded in	dividuals, and contrib	bute to animal wentare and conservat	ion errorts
**Neighborhoo	od Events or Community	y Festivals **: Look for opportunities	s to
volunteer	at local events, such	h as festivals, farmers' markets, or	street
fairs. Th	is allows you to engaged	ge with your community, promote sust	ainability,
	a connections with ne	ignoors through shared passion for	
Chvironine	chuses.		
These opport	unities will not only	nurture your spirit of volunteering	but also
help you	foster meaningful con	nections with nature, others, and yo	our
community	while nurturing a set	nse of personal growth and fulfillme	nt.
A.3. Additional D Dataset Details construction for tr datasets for each ta	<b>Petails on Dataset Construction</b> We include a representative su aining and evaluation in the sup ask in the final release of the pap	n and Evaluation ubset of the roleplay task and a jupyter notebook de pplementary material (due to size restrictions). We w per.	emonstrating shot vill release the full
<b>Evaluation Prom</b> Eval (Dubois et al. Here is the system	<b>pt for Synthetic Preferences</b> , 2024b;a) Prompt to be aware o	We used GPT-40 as a Judge using a modified vari of a user description when scoring preference example	ant of the Alpaca es.
Vou are a t	promption and accieta	nt who avaluates and calasts the h	ast large
language instructio most acc	model (LLMs) based on on. This process will urate and human-prefe	the quality of their responses to a be used to create a leaderboard ref rred answers.	given given lecting the
Here is the user pr	ompt:		
You are taske determine	ed with evaluating the which model produces	e outputs of multiple large language the best response from a human pers	models to pective.
## Instructio	on s		
You will reco 1. A **User 2. **Model O by a uniq	eive: Instruction**: This is Putputs**: Unordered reue model identifier.	s the query or task provided to the esponses from different models, each	models. identified

# Personalized Responses (ELIX-easy)

Q: How does a thermostat control temperature?	Q: What makes diamonds so hard?
Oh boy, are you ready to learn about thermostats? So, you know how sometimes your room gets too hot or too cold? That's because of something called temperature. Temperature is like a measure of how warm or cool things are. And thermostats are like superherces that help control the temperature in your home. Imagine you're playing with a thermostat in your nom. It's like a special button that makes the temperature go up or down. But how doe it octuols work?	Diamonds are renowned for their exceptional hardness, which is attributed to their unique crystal structure and chemical composition. Below is a breakdown of the factors contributing to diamond hardness: <b>Key Factors Behind Diamond Hardness</b> • Crystal Structure: Diamonds possess a face-centered
<ul> <li>Network in the control of the entrol parts:</li> <li>Sensors: These are like special eyes that look at the room's temperature. They tell the thermostat how warm or cool it is.</li> <li>Control Unit: This is like the brain of the thermostat. It looks at the sensor's readings and decides what to do next.</li> <li>Heating/Cooling System: This is like the arm that actually moves. It heats or cools the room based on what the control unit tells it to do.</li> <li>Here's how it works:</li> <li>The sensors send a message to the control unit, like "Hey, it's 75 degrees in here!"</li> <li>If the control unit thinks the room is too cold, it sends a signal to the heating system, saying "Turn on the heater!"</li> </ul>	<ul> <li>cubic (FCC) crystal structure. Each carbon atom is bonded to four neighboring carbon atoms in a tetrahedral arrangement. This structure provides unparalleled strength and rigidity, enabling diamonds to resist deformation and fracture.</li> <li>Chemical Composition: Composed solely of carbon atoms arranged in a repeating pattern, diamonds owe their remarkable hardness to this carbon-rich composition.</li> <li>Tetravalent Bonds. Every carbon atom in a diamond is bonded to four neighbors through strong covalent bonds. These tetravalent bonds create a robust framework that supports the diamond's rigid crystal structure.</li> <li>[]</li> </ul>
Elementary School Student	Expert
Elementary School Student Personalized Res	Expert
Elementary School Student <b>Personalized Res</b> Write a review for the Action movie, The Dark Knight. The Dark Knight is a thrilling ride, but its flaws hold it back from being a	Expert  Sponses (Reviews)  Write a review for the Action movie, The Dark Knight.  "A Flop of Catastrophic Proportions: "The Dark Knight Hises" Falls Wardfully Shart of ite Branque Permages**
Elementary School Student <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>Decision</b> <b>De</b>	Expert  Sponses (Reviews)  Write a review for the Action movie, The Dark Knight  **A Flop of Catastrophic Proportions: "The Dark Knight Rises" Falls Worklully Short of its Pompous Promises"  The Dark Knight Rises," Christopher Nolan's vaunted swan song to the Batman franchise, is a calamitous, soulless, and creatively bankrupt exercise in superfrectio excess, a monstrosity of monumental proportions that has left me questioning the very fabric of cinematic reality.
Elementary School Student <b>December</b> <b>December</b> Write a review for the Action movie, The Dark Knight is a thrilling ride, but its flaws hold it back from being a ruly great movie. The film's pacing is well-balanced, and the action scenes are intense and well-executed. The cast, led by Christian Bala and Heath Ledger, deliver solid performances. However, the plot is somewhat predictable, and the character's molviations are not always clear. Overall, The Dark Knight is a fun, but forgettable, superhero film.	Expert  Constant of the excession of the
Elementary School Student  Decision Strategies  Decision Strategies Decision Strategies  Decision Strategies Decision Strategies Decision Strategies Decision Strategies Decision Strategies Decision	Expert  Constant of the entropy of t
Elementary School Student  Description  Desc	Expert  Expert Expe

Figure 6: Sample Personalized Response for ELIX (top) and Reviews (bottom).

```
990
             Label preferences
                                                                            Personalized answers conditioned on labeled preferences
991
992
                               of accomplishment. By eliminating distractions and dedicating your full menta
sk, you allow your brain to process information more effectively and arrive at
993
                               reduces cognitive load, minimizing mental fatigue and burnout. Reme
re things, but rather doing one thing exceptionally well.
994
995
                                                                                      What are some fun activities for kids?
                         ow can I learn to play a musical inst
996
997
998
999
1000
1001
1002
1003
1004
1005
                                                                                        2 3 4 5
                             mend a good restaurant near me?
1006
                                                                                        plain why you chose this answer (1-2 sent
                               hat kind of food
1007
1008
1009
                                             Pre-compute responses for all possible preference selections
1012
       Figure 7: An overview of the Human Study Interface. First, users label a set of preferences. Then, a set of personalized answers are
1013
       provided, conditioned on label preferences.
1014
1015
       3. A **User Description **: This describes the user's preferences or additional
1016
            context to guide your evaluation.
1017
1018
       Your task is to:
       1. Evaluate the outputs based on quality and relevance to the user's instruction
            and description.
       2. Select the best output that meets the user's needs.
1022
1023
       ## Input Format
1024
       ### User Instruction
1026
       {QUESTION}
1027
1028
       ### Model Outputs
1029
       - Model "m": {RESPONSE A}
       - Model "M": {RESPONSE B}
       ### User Description
       {USER_DESCRIPTION}
1034
1035
       ## Task
1036
1037
       From the provided outputs, determine which model produces the best response.
1038
            Output only the model identifier of the best response (either 'm' or 'M')
1039
            with no additional text, quotes, spaces, or new lines.
1040
1041
       ## Best Model Identifier
1042
1043
1044
```

1045 Additional Human Study Details As shown in Alpaca Eval 2.0 (Dubois et al., 2024a), several biases can affect the evaluation of language models such as length, format, and more. For this reason, we took action to normalize both FSPO 1047 and baselines in 3 different categories. First, length is an evaluation bias. For this reason, we computed the average length of 1048 responses from FSPO and prompted the base model during evaluation to keep its responses around the average length in 1049 words ( $\approx 250$  words). For the SFT baseline, we found that this was consistent with FSPO since it was fine-tuned on the same preference dataset. Additionally, due to context length restrictions and the instruction following abilities of smaller 1051 open-source LLMs, we decided to have formatting be consistent as paragraphs rather than markdown for the Roleplay 1052 task. Thus, we similarly prompted the Base model with this behavior. Finally, a differing number of views can also skew 1053 the evaluation, as a large proportion of users seem to prefer direct answers. Additionally, if more views are presented, a 1054 user may prefer just one of the many views provided, skewing evaluation. Thus, we ensure that when two responses are compared, they have the same number of views. In future, work, it would be interesting to consider how to relax some of the 1055 1056 design decisions needed for the human study. We additionally provide screenshots of the human study interface in Figure 7.

Below is the full text of instructions given to the participants:

"This is a study about personalization. You will be asked to read a set of 20 questions (9 on the first page, 11 on the second page). For each question, there are two responses. Please select the response that you prefer. Make this selection based on your individual preferences and which response you find the most helpful. Read the entire response and think carefully before making your selection."

We utilize the demographic information that Prolific provides for each user such as their age group, continent and gender to chose questions but do not store that information about the user. We collect no identifying information about the user and will not make any of the individual preferences from a user public. We pay each user a fair wage subject to the current region that we reside in. We received consent from the people whose data we are using and curating as the very first question in our survey. The demographic and geographic characteristics of the annotator population is exactly the same as Prolific. We do no filtering of this at all.

# 70 A.4. Training Details and Hyperparameters for FSPO and baselines

Similar to DPO (Rafailov et al., 2023) and IPO (Gheshlaghi Azar et al., 2023), we trained FSPO in a two stage manner.
The first stage is Fewshot Pref-FT, increasing the likelihood of the preferred response. The second stage is Fewshot IPO,
initialized from the checkpoint of Fewshot Pref-FT. One epoch of the dataset was performed for each stage. For the IPO
baseline, we followed a similar procedure. Additional hyperparameters can be found in Table 8.

#### 1076 1077 A.5. Additional Details of Setup for Reproducability

1069

We used both code, models, and data as scientific artifacts. In particular, for code, we built off of the codebase from Rafailov et al. (2023), with an Apache 2.0 license. We additionally adapted our evaluation script from Alpaca EVAL, including the prompt, and other criterion for evaluation and normalization. We have reported the implementation details for synthetic evaluation in Section 6 and human study evaluation in Section A.3.

For models, we used a combination of open-source and closed-source models. The models that we used for sampling data 1083 are the Llama family of models (Grattafiori et al., 2024) (Llama 3.2 3b, Llama 3.1 8b, Llama 3.3 70b) with the llama license 1084 (3.1, 3.2, 3.3), the Qwen family of models (Qwen et al., 2025) (Qwen 2.5 3b, Qwen 2.5 32b, Qwen 2.5 72b) with the qwen 1085 license, the Gemma 2 family of models (Team et al., 2024) (Gemma 2 2b, Gemma 2 9b, and Gemma 2 27b) with the gemma 1086 license, and the OpenAI (OpenAI et al., 2024) family of models (GPT4o, GPT4o-mini) with the OpenAI API License 1087 (based off of the MIT License). We used SGLang (Zheng et al., 2024) and VLLM (Kwon et al., 2023) for model inference. 1088 For training, we used 1 node of A100 GPUs (8 GPUs) for 8 hours for each experiment with FSDP. Cumulatively, we used 1089 approximately 4000 hours of GPU hours for ablations over dataset, architecture design and other details. 1090

1091 With respect to the dataset, for questions for the review dataset, we sourced media names from IMDb (IMDb, 2025), 1092 Goodreads (Goodreads, 2025), and MyAnimeList (MyAnimeList, 2025). We define the domains in more detail in section 5. 1093 Seed questions for ELIX were human generated, sourced from Prolific. The dataset is entirely in English, with some artifacts 1094 of Chinese from the Qwen model family, which will be filtered out for the final release of the dataset. None of this data 1095 has identifying information about individual people or offensive content as the dataset was sourced from instruction and 1096 safety-tuned models, with each step of the dataset having a manual check of the inputs and outputs. In terms of statistics 1097 of the dataset, the review dataset has 130K train/dev examples and 32.4K test examples, the ELIX-easy dataset has 235K 1098 train/dev examples and 26.1K test examples, the ELIX-hard dataset has 267K train/dev examples and 267K test examples, 1099

	FSPO: Few-Shot Optimization of Synthetic Preferences Effectively Personalizes to Real Users
1100 1101 1102	and the roleplay dataset has 362K train/dev examples and 58.2K test examples, with a total of 1.378 million examples. For our statistics, we reported the average winrate % for each method on both synthetic and human evals, following prior work in alignment like AlpacaFarm (Dubois et al., 2024b).
1103 1104 1105	Each of the artifacts above was consistent with its intended use and the code, models, and datasets should be usable outside of research contexts.
1106	
1107	
1108	
1109	
1111	
1112	
1113	
1114	
1115	
1117	
1118	
1119	
1120	
1121	
1122	
1123	
1125	
1126	
1127	
1128	
1129	
1131	
1132	
1133	
1134	
1135	
1137	
1138	
1139	
1140	
1141	
1142	
1144	
1145	