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

Anonymous ACL submission

Abstract

Effective personalization of LLMs is critical for a broad range of user-interfacing applications such as virtual assistants and content cu-004 ration. Inspired by the strong in-context learning capabilities of LLMs, we propose Few-Shot Preference Optimization (FSPO), which reframes reward modeling as a meta-learning problem. Under this framework, an LLM learns to quickly adapt to a user via a few labeled preferences from that user, constructing a per-011 sonalized reward function for them. Additionally, since real-world preference data is scarce 013 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 par-017 ticular, to successfully transfer from synthetic 019 data to real users, we find it crucial for the data to exhibit both high diversity and coherent, self-consistent structure. We evaluate FSPO 021 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, along with a controlled human study. Overall, FSPO achieves an 87% Alpaca Eval winrate on average in generating responses that are personalized to synthetic users and a 72% winrate with real human users in open-ended question answering. 031

1 Introduction

032

034

042

As language models 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 may be reflected through human-annotated preference datasets; yet, 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 personal preferences of the users and not the aggregated preferences of all users? 043

045

047

049

051

054

055

057

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

077

078

079

Addressing this challenge requires a shift from modeling a singular aggregate reward function to modeling a distribution of reward functions (Sorensen et al., 2024; Jang et al., 2023) that capture the diversity of human preferences. 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. However, how can this be effectively done for open-ended question answering and transfer to 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 meta-learning for language modeling (Chen et al., 2022; Min et al., 2022; Yu et al., 2024), we fine-tune the model with a meta-learning objective, using preferencelearning objectives such as IPO (Gheshlaghi Azar et al., 2023). We additionally propose user description chain-of-thought (COT), allowing the model to leverage additional inference-compute for better reward modeling and the model's instruction fol-



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.

lowing capabilities for better response generation.

However, to learn a model that effectively personalizes to real people, we need to collect a diverse preference dataset spanning diverse users. One natural approach to do this is to curate data from humans, but this curation is difficult and timeconsuming. In contrast, in this work, we propose instantiating this dataset synthetically, and present careful design decisions to generate a dataset that is diverse and structured, following task construction considerations from the meta-learning literature (Hsu et al., 2019; Yin et al., 2019).

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

Personalized learning of preferences. Prior research has explored personalization through various methods. One approach is distributional alignment, which focuses on 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 et al., 2024), yet they do not explicitly optimize for individual preference adaptation. Another strategy involves 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 multiplechoice questions, such as GPO (Zhao et al., 2024). Although effective in structured survey settings, these methods have not been validated for openended 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.

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

161

163

Algorithms for preference learning. LLMs are typically fine-tuned via supervised next-token prediction on high-quality responses and later refined with human preference data (Casper et al., 2023; Ouyang et al., 2022). This process can use onpolicy reinforcement learning methods like RE-INFORCE (Sutton et al., 1999) or PPO (Schulman et al., 2017), which optimize a reward model

121

122

with a KL constraint. Alternatively, supervised 164 fine-tuning may be applied to a curated subset of 165 preferred responses (Dubois et al., 2024b) or iter-166 atively to preferred completions as in ReST (Gul-167 cehre et al., 2023). Other methods, such as 168 DPO (Rafailov et al., 2023), IPO (Gheshlaghi Azar 169 et al., 2023), and KTO (ContextualAI, 2024), learn 170 directly from human preferences without an ex-171 plicit reward model, with recent work exploring iterative preference modeling applications (Yuan 173 et al., 2024; Chen et al., 2024). 174

176

177

178

179

180

181

182

184

185

186

187

188

189

190

191

192

193

194

195

199

201

206

209

210

211

212

Black-box meta-learning. FSPO is an instance of black-box 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 metalearning 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**

Preference fine-tuning algorithms, such as Reinforcement Learning from Human Feedback (RLHF) and Direct Preference Optimization (DPO), typically involve two main stages (Ouyang et al., 2022; Ouyang et al., 2022): Supervised Fine-Tuning (SFT) and Preference Optimization (DPO/RLHF). First, a pre-trained model is finetuned on high-quality data from the target task using Supervised Fine-Tuning (SFT). This process produces a reference model, denoted as π_{ref} . The purpose of this stage is to bring the responses from a particular domain in distribution with supervised learning. To further refine π_{ref} according to human preferences, a preference dataset $\mathcal{D}_{\text{pref}} = \{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}_w, \mathbf{y}^{(i)}_l)\}$ is collected. In this dataset, $\mathbf{x}^{(i)}$ represents a prompt or input context, $\mathbf{y}_{w}^{(i)}$ is the preferred response, and $\mathbf{y}_{l}^{(i)}$ is the less preferred response. These responses are typically sampled from the output distribution of π_{ref} and are labeled based on human feedback.

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 y_1

over y_2 , given a prompt x, as:

$$p^{*}(\mathbf{y}_{1} \succ \mathbf{y}_{2} \mid \mathbf{x}) = \frac{e^{r^{*}(\mathbf{x}, \mathbf{y}_{1})}}{e^{r^{*}(\mathbf{x}, \mathbf{y}_{1})} + e^{r^{*}(\mathbf{x}, \mathbf{y}_{2})}} \quad (1)$$

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

The objective of preference fine-tuning is to optimize the policy π_{θ} to maximize the expected reward r^* . However, directly optimizing r^* is often impractical due to model limitations or noise in reward estimation. Therefore, a reward model r_{ϕ} is trained to approximate r^* . To prevent the finetuned policy π_{θ} from deviating excessively from the reference model π_{ref} , a Kullback-Leibler (KL) divergence constraint is imposed. This leads to the following fine-tuning objective:

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

In this equation, the regularization term weighted by β controls how much π_{θ} diverges from π_{ref} , based on the reverse KL divergence constraint. This constraint ensures that the 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:

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

where σ 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:

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

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

The Few-Shot Preference Optimization 4 (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} _

3

213

214

215

216

217

218

219

236

237

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235



Figure 2: User Description Chain-of-Thought (COT).

257

259

260

262

264

265

269

270

271

273

274

275

276

279

281

285

 $\{(\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 scorerids $\mathbf{S}^{(i)}$ of each user for differentiating users that have labeled preferences in our dataset: $\mathcal{D}_{\text{pref}} =$ $\{(\mathbf{x}^{(i)}, \mathbf{y}_w^{(i)}, \mathbf{y}_l^{(i)}, \mathbf{S}^{(i)})\}$. Now consider each user as a task instance, where the objective is to learn an effective reward function for that user using the user's set of preferences. This can be naturally instantiated as a black-box meta-learning objective, where meta-learning is done over users (also referred to as a task in meta-learning). Meta-learning should enable rapid personalization, i.e. adaptability to new users with just a few preferences.

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 θ as:

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

293
$$\mathcal{L}_{pref}^{\theta} = ||h_{\pi_{\theta}}^{y_w, y_l} - (2\beta)^{-1}||_2^2$$

$$h_{\pi_{\theta}}^{y_w, y_l} = \log \frac{\pi_{\theta}(y_w|x)}{\pi_{ref}(y_w|x)} - \log \frac{\pi_{\theta}(y_l|x)}{\pi_{ref}(y_l|x)}$$
(6)

Following black-box meta-learning approaches, FSPO receives as input a sequence of preferences $D_i^{fewshot} \sim D_i$ 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 prediction c. To make preferences compatible with a pre-trained language model, a few-shot prompt is constructed, comprising of preferences from a user and the heldout query as seen in Figure 1. This construction has an added benefit of leveraging a pretrained language model's capabilities for few-shot conditioning (Brown et al., 2020), which can enable some amount of steerage/personalization. This prediction c is implicitly learned by a preference optimization algorithm such as DPO (Rafailov et al., 2023), which parameterizes the reward model as $\beta \frac{\log \pi_{\theta}(y|x)}{\log \pi_{ref}(y|x)}$. This parameterization enables us to leverage the advantages of preference optimization algorithms such as eliminating policy learning instabilities and computational burden of on-policy sampling, learning an effective model with a simple classification objective.

Algorithm 1 Overview of FSPO

- **Require:** for each unique user $S^{(i)}$, a dataset of preferences $\mathcal{D} := (x, y_1, y_2, c)_i$, and optionally user description $y_{S^{(i)}}$ for COT, $\forall i$
- 1: while not done do
- 2: Sample training user $S^{(i)}$ (or minibatch)
- 3: Sample a subset of preferences from the user $\mathcal{D}_{i}^{fewshot} \sim \mathcal{D}_{i}$
- 4: Sample held-out preference examples $D_i^{heldout} \sim \mathcal{D}_i \setminus \mathcal{D}_i^{fewshot}$
- 5: if COT then
- 6: Use eq. (5) and eq. (6) to predict the loss on the user description $y_{\mathcal{S}^{(i)}}$.
- 7: Conditioning on $\mathcal{D}_{i}^{fewshot}$ (optionally $y_{\mathcal{S}(i)}$), use eq. (5) and eq. (6) to predict the loss on the held-out preference example $D_{i}^{heldout}$
- 8: Update learner parameters θ , using gradient of loss on $D_i^{heldout}$
- 9: return π_{θ}

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

295

296

297

298

299

300

301

302

303

304

305

306

307

309

310

311

312

313

314

315

316

325

318

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 fewshot 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.

339

341

344

345

347

351

355

357

364

367

371

373

374

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

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 opensource 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 posi-

tive 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).

spired by the subreddit "Explain Like I'm 5" where users answer questions at a very basic level appropriate for a 5 year old. Here we study the ability of the model to personalize a pedagogical explanation to a user's education background. We construct two variants of the task. The first variant is ELIXeasy where users are one of 5 education levels (elementary school, middle school, high school, college, expert) and the goal of the task is to explain a question such as "How are beaches formed?" to a user of that education background. The second, more realistic variant is ELIX-hard, which consists of question answering at a high school to university level. Here, users may have different levels of expertise in different domains. For example, a PhD student in Computer Science may have a very different educational background from an undergraduate studying studying Biology, allowing

376

377

378

379

380

381

382

383

386

387

389

390

391

392

393

394

395

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

for preferences from diverse users (550 users). **Roleplay.** The Roleplay task tackles general question answering across a wide set of users, following PRISM (Kirk et al., 2024) and PERSONA Bench (Castricato et al., 2024) to study personalization representative of the broad human population. We start by identifying three demographic traits (age, geographic location, and gender) that humans differ in that can lead to personalization. For each trait combination, we generate 30 personas, leading to 1,500 total personas. To more accurately model the distribution of questions, we split our questions into two categories: global and specific. Global questions are general where anyone may ask it, but specific questions revolve around a trait, for example an elderly person asking about retirement or a female asking about breast cancer screening.

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 annota-

375 **ELIX.** The Explain Like I'm X (ELIX) task is in-



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.

tion 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-toreal transfer in non-language domains like robotics (Makoviychuk et al., 2021) and selfdriving 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 question dataset, while preserving the stylistic choices and beliefs of human writers. For the reviews dataset, we compiled a list of popular media from sites such as Goodreads, IMDb, and MyAnimeList. For the Roleplay dataset, we prompted GPT-40 to generate questions all users would ask (global) or questions only people with a specific trait would ask (specific). This allows us to have questions that are more consistent with the distribution of questions people may ask.

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

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

460

461

462

463

427

428

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.

501

502

503

506

507

510

511

512

513

514

515

516

517

518

519

524

525

527

529

531

534

535

536

537

539

541

542

548

552

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 finetuned 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 metalearning 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.

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

Additionally, as mentioned above, in some cases, such as with the Roleplay dataset, the user description is underspecified, leading to challenges in labeling consistent preferences. For example, if a user description does not have information about dietary preferences, inconsistency may arise for



Figure 4: **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.



Figure 5: Flowchart of roleplay dataset generation

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

labeling preferences about that topic. For instance, in one preference pair, vegan cake recipes may be preferred but in another, steakhouses are preferred for date night. To fix this, we take an iterative process to constructing user descriptions. Firstly, we start with a seed set of user descriptions generated from the trait attributes. After generating questions and responses based on these seed descriptions, we take a set of question and response pairs. For each pair, we iteratively refine the user description by prompting a model like GPT4-o to either label the preference pair or if the user description is insufficient, to randomly choose a preference and append information to the description so a future scorer would make the same decision. Finally, we utilize the updated user description to relabel preferences for the set of questions and responses allocated to that user with the labeling scheme above. This fix for underspecification also helps the COT prediction as predicting an underspecified user persona, can 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 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, meta-learning 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 4. We outline our full dataset generation process in Figure 5.

Method	Winrate (%)
Base (Llama 3.2 3B instruct)	50.0
IPO	72.4
Few-shot Prompting	63.2
Few-shot Pref-FT	62.8
FSPO (ours)	82.6
FSPO + COT (ours)	90.3
Oracle (prompt w/ g.t. persona)	90.9

Table 1: Automatic Winrates on Roleplay (1500 users)

Baseline Method	Winrate (%)
FSPO vs Base	71.2
FSPO vs SFT	72.3

Table 2: Roleplay: Human Eval Winrates

7 Experimental Evaluation

586

587

588

592

593

594

599

Baselines. We compare FSPO against four baselines: (1) a base model generating user-agnostic responses, (2) few-shot 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) and (4) prompting with an oracle user description following Persona Steering (Cheng et al., 2023). 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 (2) zero-shot predicts the preferred response and (3) is optimized with SFT to increase the likelihood on the preferred response. In (4), the base model is prompted with the oracle user description, representing an upper bound on FSPO's performance.

Synthetic winrates. We first generate automated win rates using the modified AlpacaEval procedure 606 from Section 6. In the ELIX task in Table 3, we study two levels of difficulty (easy, hard), where we find a consistent improvement of FSPO over baselines. Next, in Table 4 for the Review task, on both Trained and Interpolated Users, FSPO allows for 611 better performance on held-out questions. Finally, 612 in Table 1, we study Roleplay, scaling to 1500 real users, seeing a win rate of 82.6% on both held-out 614 users and questions. Additionally, COT closes the 615 gap to the oracle response, showing effective re-616 covery of the user description. In appendix A.1, 617 618 sample generations from FSPO show effective personalization to the oracle user description. Given 619 this result, can we personalize to real people? Preliminary human study. We evaluate our model trained on the Roleplay task by personalizing re-622

Method	ELIX-easy	ELIX-hard
Base	50.0	50.0
Few-shot Prompted	92.4	81.4
Few-shot Pref-FT	91.2	82.9
FSPO (Ours)	97.8	91.8

Table 3: GPT-40 Winrates on ELIX-easy and ELIX-hard

Method	Trained	Interpolated
Base (Llama 3.2 3B instruct)	50.0	50.0
Few-shot Prompted (4-shot)	66.6	61.9
Few-shot Pref-FT (4-shot)	66.5	66.1
FSPO (4-shot, Ours)	78.4	71.3
Few-shot Prompted (8-shot)	69.1	59.1
Few-shot Pref-FT (8-shot)	65.6	70.7
FSPO (8-shot, Ours)	80.4	73.6

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

Table 4: Review Winrates - Trained and Interpolated Users

sponses for real human participants. We build a data collection app (Figure 7), interacting with a user in two stages. First, we ask participants to label preference pairs, used as the few-shot examples in FSPO. Then, for held out questions, we show a user a set of two responses: (1) a response from FSPO personalized based on their preferences and (2) a baseline response. Prolific is used to recruit a diverse set of study participants, evenly split across genders and continents, corresponding to the traits used to construct user descriptions. Question and response order is randomized to remove confounding factors. We evaluate with 25 users and 11 questions. As seen in Figure 2, 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 openended 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, thereby addressing the 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.

9 Limitations and Potential Risks

657

There are several limitations and potential risks. One limitation pertains to the ethical and fairness 659 considerations of personalization. While FSPO improves inclusivity by modeling diverse preferences, the risk of reinforcing user biases (echo chambers) 662 or inadvertently amplifying harmful viewpoints re-663 quires careful scrutiny. Future work should explore mechanisms to balance personalization with ethical 665 safeguards, ensuring that models remain aligned with fairness principles while respecting user individuality. Additionally, our human study was 669 preliminary with control over the questions that a user may ask, format normalization where format-670 ting details such as markdown are removed, and 671 view normalization comparing the same number of 672 viewpoints for both FSPO and the baselines. How-673 ever, to the best of our knowledge, we are the first 674 675 approach to perform such a human study for personalization to open-ended question answering. Future work should do further ablations with human 677 evaluation for personalization. Additionally, due to 678 compute constraints, we work with models in the 679 parameter range of 3B (specifically Llama 3.2 Instruct 3B) with a limited context window of 128K, 681 and without context optimization such as sequence 682 parallelism (Li et al., 2022; Yang et al., 2024), further limiting the effective context window. It is an open question on how fine-tuning base models 686 with better long-context and reasoning capabilities would help with FSPO for personalization, such as 687 688 the 2M context window of Gemini Flash Thinking models, especially in the case of COT.

References

690

693

702

703

704

705

710

711

712

713

714

715

716

717

718

720

721

722

723

724

725

728

729

730

731

732

733

734

735

736

738

739

740

741

742

744

745

746

- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal Ndousse, Kamile Lukosuite, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemi Mercado, Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna Kravec, Sheer El Showk, Stanislav Fort, Tamera Lanham, Timothy Telleen-Lawton, Tom Conerly, Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. 2022. Constitutional ai: Harmlessness from ai feedback. Preprint, arXiv:2212.08073.
 - Ralph Allan Bradley and Milton E. Terry. 1952. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324– 345.
 - Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam Mc-Candlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. *Preprint*, arXiv:2005.14165.
 - Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémy Scheurer, Javier Rando, Rachel Freedman, Tomasz Korbak, David Lindner, Pedro Freire, Tony Tong Wang, Samuel Marks, Charbel-Raphael Segerie, Micah Carroll, Andi Peng, Phillip Christoffersen, Mehul Damani, Stewart Slocum, Usman Anwar, Anand Siththaranjan, Max Nadeau, Eric J Michaud, Jacob Pfau, Dmitrii Krasheninnikov, Xin Chen, Lauro Langosco, Peter Hase, Erdem Biyik, Anca Dragan, David Krueger, Dorsa Sadigh, and Dylan Hadfield-Menell. 2023. Open problems and fundamental limitations of reinforcement learning from human feedback. *Transactions on Machine Learning Research*. Survey Certification.
 - Louis Castricato, Nathan Lile, Rafael Rafailov, Jan-Philipp Fränken, and Chelsea Finn. 2024. Persona: A reproducible testbed for pluralistic alignment. *Preprint*, arXiv:2407.17387.
 - Yanda Chen, Ruiqi Zhong, Sheng Zha, George Karypis, and He He. 2022. Meta-learning via language model in-context tuning. *Preprint*, arXiv:2110.07814.

Zixiang Chen, Yihe Deng, Huizhuo Yuan, Kaixuan Ji, and Quanquan Gu. 2024. Self-Play Fine-Tuning Converts Weak Language Models to Strong Language Models. *arXiv e-prints*, arXiv:2401.01335. 747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

773

774

775

776

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

- Myra Cheng, Esin Durmus, and Dan Jurafsky. 2023. Marked personas: Using natural language prompts to measure stereotypes in language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1504–1532, Toronto, Canada. Association for Computational Linguistics.
- ContextualAI. 2024. Human-centered loss functions (halos).
- Yan Duan, John Schulman, Xi Chen, Peter L Bartlett, Ilya Sutskever, and Pieter Abbeel. 2016. RI@: Fast reinforcement learning via slow reinforcement learning. *arXiv preprint arXiv:1611.02779*.
- Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B. Hashimoto. 2024a. Length-controlled alpacaeval: A simple way to debias automatic evaluators. *Preprint*, arXiv:2404.04475.
- Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2024b. Alpacafarm: A simulation framework for methods that learn from human feedback. *Preprint*, arXiv:2305.14387.
- Mohammad Gheshlaghi Azar, Mark Rowland, Bilal Piot, Daniel Guo, Daniele Calandriello, Michal Valko, and Rémi Munos. 2023. A General Theoretical Paradigm to Understand Learning from Human Preferences. *arXiv e-prints*, arXiv:2310.12036.
- Goodreads. 2025. Goodreads: Book reviews, recommendations, and discussion. Accessed: 2025-02-15.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet,

Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, 808 Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Jun-811 teng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal 814 Lakhotia, Lauren Rantala-Yeary, Laurens van der 815 Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, 816 Louis Martin, Lovish Madaan, Lubo Malo, Lukas 817 Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kam-822 badur, Mike Lewis, Min Si, Mitesh Kumar Singh, 823 Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Va-826 sic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, 832 Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan 833 Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye 837 Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Syd-841 ney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vítor Albiero, Vladan Petro-847 vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xi-849 aofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Gold-851 schlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, 852 Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, 853 Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, 855 856 Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, 857 Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit San-859 gani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew 861 Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi 867 Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic,

Brian Gamido, Britt Montalvo, Carl Parker, Carly 869 Burton, Catalina Mejia, Ce Liu, Changhan Wang, 870 Changkyu Kim, Chao Zhou, Chester Hu, Ching-871 Hsiang Chu, Chris Cai, Chris Tindal, Christoph Fe-872 ichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, 873 Daniel Kreymer, Daniel Li, David Adkins, David 874 Xu, Davide Testuggine, Delia David, Devi Parikh, 875 Diana Liskovich, Didem Foss, Dingkang Wang, Duc 876 Le, Dustin Holland, Edward Dowling, Eissa Jamil, 877 Elaine Montgomery, Eleonora Presani, Emily Hahn, 878 Emily Wood, Eric-Tuan Le, Erik Brinkman, Este-879 ban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, 880 Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat 881 Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, 883 Gada Badeer, Georgia Swee, Gil Halpern, Grant 884 Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanaz-886 eri, Han Zou, Hannah Wang, Hanwen Zha, Haroun 887 Habeeb, Harrison Rudolph, Helen Suk, Henry As-888 pegren, Hunter Goldman, Hongyuan Zhan, Ibrahim 889 Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis, 890 Irina-Elena Veliche, Itai Gat, Jake Weissman, James 891 Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jen-893 nifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy 894 Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe 895 Cummings, Jon Carvill, Jon Shepard, Jonathan Mc-896 Phie, Jonathan Torres, Josh Ginsburg, Junjie Wang, 897 Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khan-898 delwal, Katayoun Zand, Kathy Matosich, Kaushik 899 Veeraraghavan, Kelly Michelena, Kegian Li, Ki-900 ran Jagadeesh, Kun Huang, Kunal Chawla, Kyle 901 Huang, Lailin Chen, Lakshya Garg, Lavender A, 902 Leandro Silva, Lee Bell, Lei Zhang, Liangpeng 903 Guo, Licheng Yu, Liron Moshkovich, Luca Wehrst-904 edt, Madian Khabsa, Manav Avalani, Manish Bhatt, 905 Martynas Mankus, Matan Hasson, Matthew Lennie, 906 Matthias Reso, Maxim Groshev, Maxim Naumov, 907 Maya Lathi, Meghan Keneally, Miao Liu, Michael L. 908 Seltzer, Michal Valko, Michelle Restrepo, Mihir Pa-909 tel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, 910 Mike Macey, Mike Wang, Miquel Jubert Hermoso, 911 Mo Metanat, Mohammad Rastegari, Munish Bansal, 912 Nandhini Santhanam, Natascha Parks, Natasha 913 White, Navyata Bawa, Nayan Singhal, Nick Egebo, 914 Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich 915 Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, 916 Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin 917 Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pe-918 dro Rittner, Philip Bontrager, Pierre Roux, Piotr 919 Dollar, Polina Zvyagina, Prashant Ratanchandani, 920 Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel 921 Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu 922 Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, 923 Raymond Li, Rebekkah Hogan, Robin Battey, Rocky 924 Wang, Russ Howes, Ruty Rinott, Sachin Mehta, 925 Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara 926 Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, 927 Satadru Pan, Saurabh Mahajan, Saurabh Verma, 928 Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lind-929 say, Shaun Lindsay, Sheng Feng, Shenghao Lin, 930 Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, 931 Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, 932

933 Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve 934 Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj 937 Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary 951 DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. 2024. The llama 3 herd of models. Preprint, arXiv:2407.21783.

> Caglar Gulcehre, Tom Le Paine, Srivatsan Srinivasan, Ksenia Konyushkova, Lotte Weerts, Abhishek Sharma, Aditya Siddhant, Alex Ahern, Miaosen Wang, Chenjie Gu, et al. 2023. Reinforced selftraining (rest) for language modeling. *arXiv preprint arXiv:2308.08998*.

957

960

961

962

963

964

965

966

967

968 969

971

972

973

974

975

976

978

979

982 983

985

987

991

- Kyle Hsu, Sergey Levine, and Chelsea Finn. 2019. Unsupervised learning via meta-learning. In *International Conference on Learning Representations*.
- IMDb. 2025. Imdb: Ratings, reviews, and where to watch the best movies & tv shows. Accessed: 2025-02-15.
- Hakan Inan, Kartikeya Upasani, Jianfeng Chi, Rashi Rungta, Krithika Iyer, Yuning Mao, Michael Tontchev, Qing Hu, Brian Fuller, Davide Testuggine, and Madian Khabsa. 2023. Llama guard: Llm-based input-output safeguard for human-ai conversations. *Preprint*, arXiv:2312.06674.
- Joel Jang, Seungone Kim, Bill Yuchen Lin, Yizhong Wang, Jack Hessel, Luke Zettlemoyer, Hannaneh Hajishirzi, Yejin Choi, and Prithviraj Ammanabrolu. 2023. Personalized soups: Personalized large language model alignment via post-hoc parameter merging. *Preprint*, arXiv:2310.11564.
- Hannah Rose Kirk, Alexander Whitefield, Paul Röttger, Andrew Bean, Katerina Margatina, Juan Ciro, Rafael Mosquera, Max Bartolo, Adina Williams, He He, Bertie Vidgen, and Scott A. Hale. 2024. The prism alignment dataset: What participatory, representative and individualised human feedback reveals about the subjective and multicultural alignment of large language models. *Preprint*, arXiv:2404.16019.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. *Preprint*, arXiv:2309.06180.

Yoonho Lee, Jonathan Williams, Henrik Marklund, Archit Sharma, Eric Mitchell, Anikait Singh, and Chelsea Finn. 2024. Test-time alignment via hypothesis reweighting. *Preprint*, arXiv:2412.08812. 992

993

994

995

996

997

998

999

1001

1002

1003

1004

1005

1006

1007

1008

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1029

1030

1031

1032

1033

1034

1035

1036

1037

- Haoran Li, Qingxiu Dong, Zhengyang Tang, Chaojun Wang, Xingxing Zhang, Haoyang Huang, Shaohan Huang, Xiaolong Huang, Zeqiang Huang, Dongdong Zhang, Yuxian Gu, Xin Cheng, Xun Wang, Si-Qing Chen, Li Dong, Wei Lu, Zhifang Sui, Benyou Wang, Wai Lam, and Furu Wei. 2024. Synthetic data (almost) from scratch: Generalized instruction tuning for language models. *Preprint*, arXiv:2402.13064.
- Shenggui Li, Fuzhao Xue, Chaitanya Baranwal, Yongbin Li, and Yang You. 2022. Sequence parallelism: Long sequence training from system perspective. *Preprint*, arXiv:2105.13120.
- Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts.
 2011. Learning word vectors for sentiment analysis. In Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies, pages 142–150, Portland, Oregon, USA. Association for Computational Linguistics.
- Viktor Makoviychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles Macklin, David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, and Gavriel State. 2021. Isaac gym: High performance gpu-based physics simulation for robot learning. *Preprint*, arXiv:2108.10470.
- Nicole Meister, Carlos Guestrin, and Tatsunori Hashimoto. 2024. Benchmarking distributional alignment of large language models. *Preprint*, arXiv:2411.05403.
- Igor Melnyk, Youssef Mroueh, Brian Belgodere, Mattia Rigotti, Apoorva Nitsure, Mikhail Yurochkin, Kristjan Greenewald, Jiri Navratil, and Jerret Ross. 2024. Distributional preference alignment of llms via optimal transport. *Preprint*, arXiv:2406.05882.
- Sewon Min, Mike Lewis, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2022. Metaicl: Learning to learn in context. *Preprint*, arXiv:2110.15943.
- Nikhil Mishra, Mostafa Rohaninejad, Xi Chen, and Pieter Abbeel. 2018. A simple neural attentive metalearner. *Preprint*, arXiv:1707.03141.
- MyAnimeList. 2025. Myanimelist: Track, discover, and discuss anime & manga. Accessed: 2025-02-15.
- OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, 1039 Adam Perelman, Aditya Ramesh, Aidan Clark, 1040 AJ Ostrow, Akila Welihinda, Alan Hayes, Alec 1041 Radford, Aleksander Mądry, Alex Baker-Whitcomb, 1042 Alex Beutel, Alex Borzunov, Alex Carney, Alex 1043 Chow, Alex Kirillov, Alex Nichol, Alex Paino, Alex 1044 Renzin, Alex Tachard Passos, Alexander Kirillov, 1045 Alexi Christakis, Alexis Conneau, Ali Kamali, Allan 1046 Jabri, Allison Moyer, Allison Tam, Amadou Crookes, 1047

Amin Tootoochian, Amin Tootoonchian, Ananya Kumar, Andrea Vallone, Andrej Karpathy, Andrew Braunstein, Andrew Cann, Andrew Codispoti, Andrew Galu, Andrew Kondrich, Andrew Tulloch, Andrey Mishchenko, Angela Baek, Angela Jiang, Antoine Pelisse, Antonia Woodford, Anuj Gosalia, Arka Dhar, Ashley Pantuliano, Avi Nayak, Avital Oliver, Barret Zoph, Behrooz Ghorbani, Ben Leimberger, Ben Rossen, Ben Sokolowsky, Ben Wang, Benjamin Zweig, Beth Hoover, Blake Samic, Bob McGrew, Bobby Spero, Bogo Giertler, Bowen Cheng, Brad Lightcap, Brandon Walkin, Brendan Quinn, Brian Guarraci, Brian Hsu, Bright Kellogg, Brydon Eastman, Camillo Lugaresi, Carroll Wainwright, Cary Bassin, Cary Hudson, Casey Chu, Chad Nelson, 1063 Chak Li, Chan Jun Shern, Channing Conger, Charlotte Barette, Chelsea Voss, Chen Ding, Cheng Lu, Chong Zhang, Chris Beaumont, Chris Hallacy, Chris Koch, Christian Gibson, Christina Kim, Christine Choi, Christine McLeavey, Christopher Hesse, Claudia Fischer, Clemens Winter, Coley Czarnecki, Colin Jarvis, Colin Wei, Constantin Koumouzelis, Dane 1069 Sherburn, Daniel Kappler, Daniel Levin, Daniel Levy, David Carr, David Farhi, David Mely, David Robinson, David Sasaki, Denny Jin, Dev Valladares, Dimitris Tsipras, Doug Li, Duc Phong Nguyen, Duncan Findlay, Edede Oiwoh, Edmund Wong, Ehsan Asdar, Elizabeth Proehl, Elizabeth Yang, Eric Antonow, Eric Kramer, Eric Peterson, Eric Sigler, Eric Wallace, Eugene Brevdo, Evan Mays, Farzad Khorasani, Felipe Petroski Such, Filippo Raso, Francis Zhang, Fred von Lohmann, Freddie Sulit, Gabriel Goh, Gene Oden, Geoff Salmon, Giulio Starace, Greg Brockman, Hadi Salman, Haiming Bao, Haitang Hu, Hannah Wong, Haoyu Wang, Heather Schmidt, Heather Whitney, Heewoo Jun, Hendrik Kirchner, Henrique Ponde de Oliveira Pinto, Hongyu Ren, Huiwen Chang, Hyung Won Chung, Ian Kivlichan, Ian O'Connell, Ian O'Connell, Ian Osband, Ian Silber, Ian Sohl, Ibrahim Okuyucu, Ikai Lan, Ilya Kostrikov, Ilya Sutskever, Ingmar Kanitscheider, Ishaan Gulrajani, Jacob Coxon, Jacob Menick, Jakub 1090 Pachocki, James Aung, James Betker, James Crooks, James Lennon, Jamie Kiros, Jan Leike, Jane Park, Jason Kwon, Jason Phang, Jason Teplitz, Jason Wei, Jason Wolfe, Jay Chen, Jeff Harris, Jenia Varavva, Jessica Gan Lee, Jessica Shieh, Ji Lin, Jiahui Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joanne Jang, Joaquin Quinonero Candela, Joe Beutler, Joe Landers, Joel Parish, Johannes Heidecke, John Schulman, Jonathan Lachman, Jonathan McKay, Jonathan Uesato, Jonathan Ward, Jong Wook Kim, Joost Huizinga, Jordan Sitkin, Jos Kraaijeveld, Josh Gross, Josh Kaplan, Josh Snyder, Joshua Achiam, Joy Jiao, Joyce Lee, Juntang Zhuang, Justyn Harriman, Kai Fricke, Kai Hayashi, Karan Singhal, Katy Shi, Kavin Karthik, Kayla Wood, Kendra Rimbach, Kenny Hsu, 1105 Kenny Nguyen, Keren Gu-Lemberg, Kevin Button, Kevin Liu, Kiel Howe, Krithika Muthukumar, Kyle Luther, Lama Ahmad, Larry Kai, Lauren Itow, Lauren Workman, Leher Pathak, Leo Chen, Li Jing, Lia Guy, Liam Fedus, Liang Zhou, Lien Mamitsuka, Lilian Weng, Lindsay McCallum, Lindsey Held, Long Ouyang, Louis Feuvrier, Lu Zhang, Lukas Kon-

1048

1049

1050

1052

1055

1056

1057

1058

1059

1060

1066

1068

1070

1071

1073

1075

1076

1078

1079

1080

1081

1082

1083

1084

1085

1086

1087

1088

1089

1091

1092

1093

1094

1095

1096

1097

1098

1099

1100

1101

1102

1103

1104

1106

1107

1108

1109

1110

1111

draciuk, Lukasz Kaiser, Luke Hewitt, Luke Metz, 1112 Lyric Doshi, Mada Aflak, Maddie Simens, Madelaine 1113 Boyd, Madeleine Thompson, Marat Dukhan, Mark 1114 Chen, Mark Gray, Mark Hudnall, Marvin Zhang, 1115 Marwan Aljubeh, Mateusz Litwin, Matthew Zeng, 1116 Max Johnson, Maya Shetty, Mayank Gupta, Meghan 1117 Shah, Mehmet Yatbaz, Meng Jia Yang, Mengchao 1118 Zhong, Mia Glaese, Mianna Chen, Michael Jan-1119 ner, Michael Lampe, Michael Petrov, Michael Wu, 1120 Michele Wang, Michelle Fradin, Michelle Pokrass, 1121 Miguel Castro, Miguel Oom Temudo de Castro, 1122 Mikhail Pavlov, Miles Brundage, Miles Wang, Mi-1123 nal Khan, Mira Murati, Mo Bavarian, Molly Lin, 1124 Murat Yesildal, Nacho Soto, Natalia Gimelshein, Na-1125 talie Cone, Natalie Staudacher, Natalie Summers, 1126 Natan LaFontaine, Neil Chowdhury, Nick Ryder, 1127 Nick Stathas, Nick Turley, Nik Tezak, Niko Felix, 1128 Nithanth Kudige, Nitish Keskar, Noah Deutsch, Noel 1129 Bundick, Nora Puckett, Ofir Nachum, Ola Okelola, 1130 Oleg Boiko, Oleg Murk, Oliver Jaffe, Olivia Watkins, 1131 Olivier Godement, Owen Campbell-Moore, Patrick 1132 Chao, Paul McMillan, Pavel Belov, Peng Su, Pe-1133 ter Bak, Peter Bakkum, Peter Deng, Peter Dolan, 1134 Peter Hoeschele, Peter Welinder, Phil Tillet, Philip 1135 Pronin, Philippe Tillet, Prafulla Dhariwal, Qiming 1136 Yuan, Rachel Dias, Rachel Lim, Rahul Arora, Ra-1137 jan Troll, Randall Lin, Rapha Gontijo Lopes, Raul 1138 Puri, Reah Miyara, Reimar Leike, Renaud Gaubert, 1139 Reza Zamani, Ricky Wang, Rob Donnelly, Rob 1140 Honsby, Rocky Smith, Rohan Sahai, Rohit Ramchan-1141 dani, Romain Huet, Rory Carmichael, Rowan Zellers, 1142 Roy Chen, Ruby Chen, Ruslan Nigmatullin, Ryan 1143 Cheu, Saachi Jain, Sam Altman, Sam Schoenholz, 1144 Sam Toizer, Samuel Miserendino, Sandhini Agar-1145 wal, Sara Culver, Scott Ethersmith, Scott Gray, Sean 1146 Grove, Sean Metzger, Shamez Hermani, Shantanu 1147 Jain, Shengjia Zhao, Sherwin Wu, Shino Jomoto, Shi-1148 rong Wu, Shuaiqi, Xia, Sonia Phene, Spencer Papay, 1149 Srinivas Narayanan, Steve Coffey, Steve Lee, Stew-1150 art Hall, Suchir Balaji, Tal Broda, Tal Stramer, Tao 1151 Xu, Tarun Gogineni, Taya Christianson, Ted Sanders, 1152 Tejal Patwardhan, Thomas Cunninghman, Thomas 1153 Degry, Thomas Dimson, Thomas Raoux, Thomas 1154 Shadwell, Tianhao Zheng, Todd Underwood, Todor 1155 Markov, Toki Sherbakov, Tom Rubin, Tom Stasi, 1156 Tomer Kaftan, Tristan Heywood, Troy Peterson, Tyce 1157 Walters, Tyna Eloundou, Valerie Qi, Veit Moeller, 1158 Vinnie Monaco, Vishal Kuo, Vlad Fomenko, Wayne 1159 Chang, Weiyi Zheng, Wenda Zhou, Wesam Manassra, 1160 Will Sheu, Wojciech Zaremba, Yash Patil, Yilei Qian, 1161 Yongjik Kim, Youlong Cheng, Yu Zhang, Yuchen 1162 He, Yuchen Zhang, Yujia Jin, Yunxing Dai, and 1163 Yury Malkov. 2024. Gpt-40 system card. Preprint, 1164 arXiv:2410.21276. 1165 1166

Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. Preprint, arXiv:2203.02155.

1167

1168

1169

1170

1171

1172

Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. arXiv e-prints, arXiv:2203.02155.

1174

1175

1176

1177

1179

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1203

1204

1205

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225

1226

1227

1228

1229

- Sriyash Poddar, Yanming Wan, Hamish Ivison, Abhishek Gupta, and Natasha Jaques. 2024. Personalizing reinforcement learning from human feedback with variational preference learning. *Preprint*, arXiv:2408.10075.
- Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. 2025. Qwen2.5 technical report. *Preprint*, arXiv:2412.15115.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D Manning, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. *arXiv preprint arXiv:2305.18290*.
- Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy Lillicrap. 2016. Oneshot learning with memory-augmented neural networks. *Preprint*, arXiv:1605.06065.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal Policy Optimization Algorithms. *arXiv e-prints*, arXiv:1707.06347.
- Omar Shaikh, Michelle Lam, Joey Hejna, Yijia Shao, Michael Bernstein, and Diyi Yang. 2024. Show, don't tell: Aligning language models with demonstrated feedback. *Preprint*, arXiv:2406.00888.
- Anand Siththaranjan, Cassidy Laidlaw, and Dylan Hadfield-Menell. 2024. Distributional preference learning: Understanding and accounting for hidden context in rlhf. *Preprint*, arXiv:2312.08358.
- Taylor Sorensen, Jared Moore, Jillian Fisher, Mitchell Gordon, Niloofar Mireshghallah, Christopher Michael Rytting, Andre Ye, Liwei Jiang, Ximing Lu, Nouha Dziri, Tim Althoff, and Yejin Choi. 2024. A roadmap to pluralistic alignment. *Preprint*, arXiv:2402.05070.
- Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. 1999. Policy gradient methods for reinforcement learning with function approximation. In *Advances in Neural Information Processing Systems*, volume 12. MIT Press.
- Gemma Team, Morgane Riviere, Shreya Pathak, 1231 Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupati-1232 raju, Léonard Hussenot, Thomas Mesnard, Bobak 1233 Shahriari, Alexandre Ramé, Johan Ferret, Peter 1234 Liu, Pouya Tafti, Abe Friesen, Michelle Casbon, 1235 Sabela Ramos, Ravin Kumar, Charline Le Lan, Sammy Jerome, Anton Tsitsulin, Nino Vieillard, 1237 Piotr Stanczyk, Sertan Girgin, Nikola Momchev, 1238 Matt Hoffman, Shantanu Thakoor, Jean-Bastien Grill, 1239 Behnam Neyshabur, Olivier Bachem, Alanna Wal-1240 ton, Aliaksei Severyn, Alicia Parrish, Aliya Ah-1241 mad, Allen Hutchison, Alvin Abdagic, Amanda 1242 Carl, Amy Shen, Andy Brock, Andy Coenen, An-1243 thony Laforge, Antonia Paterson, Ben Bastian, Bilal 1244 Piot, Bo Wu, Brandon Royal, Charlie Chen, Chintu 1245 Kumar, Chris Perry, Chris Welty, Christopher A. 1246 Choquette-Choo, Danila Sinopalnikov, David Wein-1247 berger, Dimple Vijaykumar, Dominika Rogozińska, Dustin Herbison, Elisa Bandy, Emma Wang, Eric 1249 Noland, Erica Moreira, Evan Senter, Evgenii Elty-1250 shev, Francesco Visin, Gabriel Rasskin, Gary Wei, 1251 Glenn Cameron, Gus Martins, Hadi Hashemi, Hanna 1252 Klimczak-Plucińska, Harleen Batra, Harsh Dhand, 1253 Ivan Nardini, Jacinda Mein, Jack Zhou, James Svens-1254 son, Jeff Stanway, Jetha Chan, Jin Peng Zhou, Joana 1255 Carrasqueira, Joana Iljazi, Jocelyn Becker, Joe Fer-1256 nandez, Joost van Amersfoort, Josh Gordon, Josh 1257 Lipschultz, Josh Newlan, Ju yeong Ji, Kareem Mo-1258 hamed, Kartikeya Badola, Kat Black, Katie Mil-1259 lican, Keelin McDonell, Kelvin Nguyen, Kiranbir 1260 Sodhia, Kish Greene, Lars Lowe Sjoesund, Lau-1261 ren Usui, Laurent Sifre, Lena Heuermann, Leti-1262 cia Lago, Lilly McNealus, Livio Baldini Soares, 1263 Logan Kilpatrick, Lucas Dixon, Luciano Martins, 1264 Machel Reid, Manvinder Singh, Mark Iverson, Mar-1265 tin Görner, Mat Velloso, Mateo Wirth, Matt Davi-1266 dow, Matt Miller, Matthew Rahtz, Matthew Watson, 1267 Meg Risdal, Mehran Kazemi, Michael Moynihan, 1268 Ming Zhang, Minsuk Kahng, Minwoo Park, Mofi 1269 Rahman, Mohit Khatwani, Natalie Dao, Nenshad 1270 Bardoliwalla, Nesh Devanathan, Neta Dumai, Nilay 1271 Chauhan, Oscar Wahltinez, Pankil Botarda, Parker 1272 Barnes, Paul Barham, Paul Michel, Pengchong 1273 Jin, Petko Georgiev, Phil Culliton, Pradeep Kup-1274 pala, Ramona Comanescu, Ramona Merhej, Reena 1275 Jana, Reza Ardeshir Rokni, Rishabh Agarwal, Ryan 1276 Mullins, Samaneh Saadat, Sara Mc Carthy, Sarah 1277 Cogan, Sarah Perrin, Sébastien M. R. Arnold, Se-1278 bastian Krause, Shengyang Dai, Shruti Garg, Shruti 1279 Sheth, Sue Ronstrom, Susan Chan, Timothy Jor-1280 dan, Ting Yu, Tom Eccles, Tom Hennigan, Tomas 1281 Kocisky, Tulsee Doshi, Vihan Jain, Vikas Yadav, 1282 Vilobh Meshram, Vishal Dharmadhikari, Warren 1283 Barkley, Wei Wei, Wenming Ye, Woohyun Han, 1284 Woosuk Kwon, Xiang Xu, Zhe Shen, Zhitao Gong, 1285 Zichuan Wei, Victor Cotruta, Phoebe Kirk, Anand 1286 Rao, Minh Giang, Ludovic Peran, Tris Warkentin, 1287 Eli Collins, Joelle Barral, Zoubin Ghahramani, Raia 1288 Hadsell, D. Sculley, Jeanine Banks, Anca Dragan, 1289 Slav Petrov, Oriol Vinyals, Jeff Dean, Demis Hass-1290 abis, Koray Kavukcuoglu, Clement Farabet, Elena 1291 Buchatskaya, Sebastian Borgeaud, Noah Fiedel, Ar-1292 mand Joulin, Kathleen Kenealy, Robert Dadashi, 1293 and Alek Andreev. 2024. Gemma 2: Improving 1294

1295open language models at a practical size. Preprint,1296arXiv:2408.00118.

1297

1298

1299

1302

1303

1304

1305

1307

1309

1310

1311 1312

1313 1314

1315

1316

1317

1318

1319

1320

1322

1326

1328

1329

1330

1331

1332

1333

1334 1335

1336

1337

1338

1339

1340 1341

1342

- Joshua Tobin, Lukas Biewald, Rocky Duan, Marcin Andrychowicz, Ankur Handa, Vikash Kumar, Bob McGrew, Jonas Schneider, Peter Welinder, Wojciech Zaremba, and Pieter Abbeel. 2018. Domain randomization and generative models for robotic grasping. *Preprint*, arXiv:1710.06425.
- Jane X Wang, Zeb Kurth-Nelson, Dhruva Tirumala, Hubert Soyer, Joel Z Leibo, Remi Munos, Charles Blundell, Dharshan Kumaran, and Matt Botvinick. 2016. Learning to reinforcement learn. *arXiv preprint arXiv:1611.05763*.
 - Amy Yang, Jingyi Yang, Aya Ibrahim, Xinfeng Xie, Bangsheng Tang, Grigory Sizov, Jeremy Reizenstein, Jongsoo Park, and Jianyu Huang. 2024. Context parallelism for scalable million-token inference. *Preprint*, arXiv:2411.01783.
- Ze Yang, Yun Chen, Jingkang Wang, Sivabalan Manivasagam, Wei-Chiu Ma, Anqi Joyce Yang, and Raquel Urtasun. 2023. Unisim: A neural closed-loop sensor simulator. *Preprint*, arXiv:2308.01898.
- Mingzhang Yin, George Tucker, Mingyuan Zhou, Sergey Levine, and Chelsea Finn. 2019. Metalearning without memorization. *arXiv preprint arXiv:1912.03820*.
- Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T. Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. 2024. Metamath: Bootstrap your own mathematical questions for large language models. *Preprint*, arXiv:2309.12284.
- Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Xian Li, Sainbayar Sukhbaatar, Jing Xu, and Jason Weston. 2024. Self-Rewarding Language Models. *arXiv e-prints*, arXiv:2401.10020.
- Lunjun Zhang, Arian Hosseini, Hritik Bansal, Mehran Kazemi, Aviral Kumar, and Rishabh Agarwal. 2024. Generative verifiers: Reward modeling as next-token prediction. *Preprint*, arXiv:2408.15240.
- Siyan Zhao, John Dang, and Aditya Grover. 2024. Group preference optimization: Few-shot alignment of large language models. *Preprint*, arXiv:2310.11523.
- Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie, Chuyue Sun, Jeff Huang, Cody Hao Yu, Shiyi Cao, Christos Kozyrakis, Ion Stoica, Joseph E. Gonzalez, Clark Barrett, and Ying Sheng. 2024. Sglang: Efficient execution of structured language model programs. *Preprint*, arXiv:2312.07104.

1345

1347

1348

1349

1350

1351

1352

1353

1354

1355

1356

1357

1358

1359

1360

1361

1362

1363

1364

1365

1367

1368

1370

1371

1372

1373

1374

1375

1376

1378

1379

1380

1381

1382

1383

1384

1385

1386

1388

1390

1391

1392

1393

1394

1395

1396

1397

1399

1400

1401

1402

1403

1404

1405

1406

1407

A Appendix

1346 A.1 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.1.1 Roleplay CoT

QUESTION:

What are some good volunteer opportunities for $\operatorname{\mathsf{me}}\nolimits?$

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 longterm financial security. She is at a point in her career where opportunities for advancement might be time-sensitive. She enjoys movies that challenge perspectives and spark discussions. She prioritizes well-being over being fully informed. She is interested in the van life movement. She appreciates creativity and variety in meal options. She enjoys engaging in creative and tactile activities and prefers traditional learning methods. She values compassionate and responsible choices. She enjoys exploring new and dynamic experiences. She is concerned about fast food. She values mental well-being and mindfulness. She enjoys attending events and conferences to meet new people. She appreciates practicality and health. She enjoys shows with dark humor and suspenseful storytelling.

COT GENERATED PERSONA:

Based on the examples provided, the user appears to value a more holistic and experiential approach to learning and personal growth. They seem to prioritize human connection, nature, and personal development over technology and material wealth.

This user is likely someone who is interested in personal stories, emotions, and experiences, as evident from their appreciation for Sylvia Plath's poetry and the suggestion to weave storytelling into public speaking. They may also be someone who values authenticity, simplicity, and community, as seen in their suggestion of house swapping and volunteering as alternative vacation options.

In terms of lifestyle, this user likely prefers a balanced and flexible approach to life, allowing for a mix of work, personal growth, and leisure activities. They may also be someone who values environmental sustainability and social responsibility, as hinted at by their consideration of eco-friendly home improvements and sustainable financial planning. Overall, this user seems to be a thoughtful, empathetic, and nature-loving individual who prioritizes human connections, experiences, and personal growth over material possessions and technological advancements.

1408

1409

1410

1411

1412

1413

1414

1415 1416

1417

1418

1419

1420

1421

1422

1423

1424 1425

1426

1427

1428

1429

1430

1431

1432

1433

1434

1435

1436

1437

1438

1439

1440

1441

1442

1443

1444

1445

1446

1447

1448

1449

1450

1451

1452

1453

1454

1455

1456

1457

1458

1459

1460

1461

1462

1463

1464

1465

1466

1467

1468

1469

1470

1471

1472

FINAL RESPONSE:

Based on your affinity for human connections, personal growth, and environmental sustainability, I'd like to suggest the following volunteer opportunities that might resonate with you:

Community Garden or Urban Farm: Many organizations in your area offer community gardens or urban farms where you can help cultivate and maintain a shared space, promoting sustainability and food security. This allows you to connect with nature and contribute to your community's well-being.

Nature Reserves or Conservation Efforts: Consider volunteering with local wildlife conservation organizations, such as bird watching tours, tree planting initiatives, or habitat restoration projects. This will enable you to work outdoors, connect with others who share your passion for nature, and contribute to preserving the environment.

Community Centers or Environmental Education Programs: Many community centers offer various programs, such as gardening workshops, environmental education classes, or sustainable living workshops. By volunteering at these centers, you can share your knowledge, learn from others, and make a positive impact on your community.

Animal Shelters or Wildlife Rehabilitation Centers: If you're an animal lover, consider volunteering at a local animal shelter or wildlife rehabilitation center. You'll have the opportunity to care for animals, interact with like-minded individuals, and contribute to animal welfare and conservation efforts.

Neighborhood Events or Community Festivals: Look for opportunities to volunteer at local events, such as festivals, farmers' markets, or street fairs. This allows you to engage with your community, promote sustainability, and build connections with neighbors through shared passion for environmental causes.

These opportunities will not only nurture your spirit of volunteering but also help you foster meaningful connections with nature, others, and your community while nurturing a sense of personal growth and fulfillment.

A.2 Additional Details on Dataset Construction and Evaluation

Dataset DetailsWe include a representative sub-
set of the roleplay task and a jupyter notebook147314741474

Name	Value
Learning Rate (SFT/Pref-FT)	$1e^{-5}, 1e^{-6}, 1e^{-7}$
Learning Rate (IPO)	$1e^{-5}, \mathbf{1e^{-6}}, 1e^{-7}$
Beta (IPO)	$0.1, 0.05, 0.01, \boldsymbol{0.005}, 0.001$
Number of Shots	4, 8
Model Name	Llama 3.2 3B Instruct (Grattafiori et al., 2024)

Table 5: Sweep over hyperparameters for FSPO, recommended hyperparameters in bold.

Personalized Responses (ELIX-easy)



Elementary School Student



Personalized Responses (Reviews)



Interpolated User: Concise + Negative

Interpolated User: Verbose + Negative

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



Pre-compute responses for all possible preference selections

Figure 7: An overview of the Human Study Interface. First, users label a set of preferences. Then, a set of personalized answers are provided, conditioned on label preferences.

demonstrating shot construction for training and
evaluation in the supplementary material (due to
size restrictions). We will release the full datasets
for each task in the final release of the paper.

1479Evaluation Prompt for Synthetic Preferences1480We used GPT-40 as a Judge using a modified vari-1481ant of the Alpaca Eval (Dubois et al., 2024b,a)1482Prompt to be aware of a user description when1483scoring preference examples.

Here is the system prompt:

You are a highly efficient assistant, who evaluates and selects the best large language model (LLMs) based on the quality of their responses to a given instruction. This process will be used to create a leaderboard reflecting the most accurate and human-preferred answers.

Here is the user prompt:

You are tasked with evaluating the outputs of multiple large language models to determine which model produces the best response from a human perspective.

Instructions

1484

1485

1486

1487

1488

1489

1490

1491

1492

1493

1494

1495

1496

1497

1498

1499

1501

1503

1504

1505

1506

1507

1508

1509

You will receive:

 A **User Instruction**: This is the query or task provided to the models.
 Model Outputs: Unordered responses from different models, each identified by a unique model identifier.
 A **User Description**: This describes the user's preferences or additional context to guide your evaluation.

Your task is to:

1. Evaluate the outputs based on quality and 1510 relevance to the users instruction and 1511 1512 description. 2. Select the best output that meets the user's 1513 1514 needs. ## Input Format 1517 ### User Instruction {QUESTION} 1519 1520 ### Model Outputs 1521 - Model "m": {RESPONSE_A} 1522 - Model "M": {RESPONSE_B} 1523 1524 ### User Description 1525 {USER_DESCRIPTION} 1526 1527 ## Task 1528 1529 From the provided outputs, determine which model 1530 produces the best response. Output only the 1531 model identifier of the best response (either 'm' 1532

model identifier of the best response (either 'm'
or 'M') with no additional text, quotes, spaces,
or new lines.

1533 1534

1535

1536

Best Model Identifier

Additional Human Study Details As shown in Alpaca Eval 2.0 (Dubois et al., 2024a), several bi-1538 ases can affect the evaluation of language models 1539 such as length, format, and more. For this rea-1540 son, we took action to normalize both FSPO and 1541 baselines in 3 different categories. First, length 1542 is an evaluation bias. For this reason, we com-1543 puted the average length of responses from FSPO 1544 and prompted the base model during evaluation 1545 to keep its responses around the average length in 1546

words (≈ 250 words). For the SFT baseline, we 1547 found that this was consistent with FSPO since it 1548 was fine-tuned on the same preference dataset. Ad-1549 ditionally, due to context length restrictions and 1550 the instruction following abilities of smaller opensource LLMs, we decided to have formatting be 1552 consistent as paragraphs rather than markdown for 1553 the Roleplay task. Thus, we similarly prompted the 1554 Base model with this behavior. Finally, a differing 1555 number of views can also skew the evaluation, as 1556 a large proportion of users seem to prefer direct answers. Additionally, if more views are presented, 1558 a user may prefer just one of the many views pro-1559 vided, skewing evaluation. Thus, we ensure that 1560 when two responses are compared, they have the 1561 same number of views. In future, work, it would be interesting to consider how to relax some of 1563 the design decisions needed for the human study. 1564 We additionally provide screenshots of the human 1565 study interface in Figure 7. 1566

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

1569

1570

1571

1574

1576

1577

1578

1581

1582

1583

1584

1585

1586

1588

1589

1590

1591

1592

1593

1594

1595

1597

"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.

A.3 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. Additional hyperparameters can be found in Table 5. 1598

1599

1602

1603

1605

1606

1607

1608

1609

1610

1611

1612

1613

1614

1615

1616

1617

1619

1620

1621

1622

1623

1624

1625

1626

1627

1628

1630

1631

1632

1634

1635

1636

1637

1638

1640

1641

1642

1643

1644

1645

1646

1648

A.4 Details for ARR checklist

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.2.

For models, we used a combination of opensource and closed-source models. The models that we used for sampling data 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 (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 license, the Gemma 2 family of models (Team et al., 2024) (Gemma 2 2b, Gemma 2 9b, and Gemma 2 27b) with the gemma license, and the OpenAI (OpenAI et al., 2024) family of models (GPT4o, GPT4omini) with the OpenAI API License (based off of the MIT License). We used SGLang (Zheng et al., 2024) and VLLM (Kwon et al., 2023) for model inference. For training, we used 1 node of A100 GPUs (8 GPUs) for 8 hours for each experiment with FSDP. Cumulatively, we used approximately 4000 hours of GPU hours for ablations over dataset, architecture design and other details.

With respect to the dataset, for questions for the review dataset, we sourced media names from IMDb (IMDb, 2025), Goodreads (Goodreads, 2025), and MyAnimeList (MyAnimeList, 2025). We define the domains in more detail in section 5. Seed questions for ELIX were human generated, sourced from Prolific. The dataset is entirely in English, with some artifacts of Chinese from the Qwen model family, which will be filtered out for the final release of the dataset. None of this data has identifying information about individual people or offensive content as the dataset was sourced from instruction and safety-tuned models, with each step of the dataset having a manual check of the inputs and outputs. Before final release, we will verify the dataset with Llama 3 Guard (Inan et al., 2023).

In terms of statistics of the dataset, the review dataset has 130K train/dev examples and 32.4K

1649 test examples, the ELIX-easy dataset has 235K train/dev examples and 26.1K test examples, the 1650 ELIX-hard dataset has 267K train/dev examples 1651 and 267K test examples, and the roleplay dataset 1652 has 362K train/dev examples and 58.2K test exam-1653 1654 ples, with a total of 1.378 million examples. The final release will be in this ballpark but may be 1655 adjusted for quality and safety purposes. 1656

1657

1658

1659

1660

1661

1662

1663

1664

1665

1666

1667

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).

AI assistants were used in this work, with tools like copilot used for coding and ChatGPT, Claude and Gemini were used for writing assistance, latex table creation, and figure construction.

Each of the artifacts above was consistent with its intended use and the dataset should be usable outside of research contexts.