
FSPO: Few-Shot Preference Optimization of Synthetic Data in LLMs Elicits Effective Personalization to Real Users

Anikait Singh^{1,◇}, Sheryl Hsu^{1,◇}, Kyle Hsu¹, Eric Mitchell^{1,4}, Stefano Ermon¹, Tatsunori Hashimoto^{1,◆}, Archit Sharma^{1,2,◆} and Chelsea Finn^{1,3,◆}

¹Stanford University, ²Google DeepMind, ³Physical Intelligence, ⁴OpenAI, [◇]Equal Contribution, [◆]Equal Advising

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), 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 personalized reward function for them. Additionally, 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, 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.

1. Introduction

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?

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

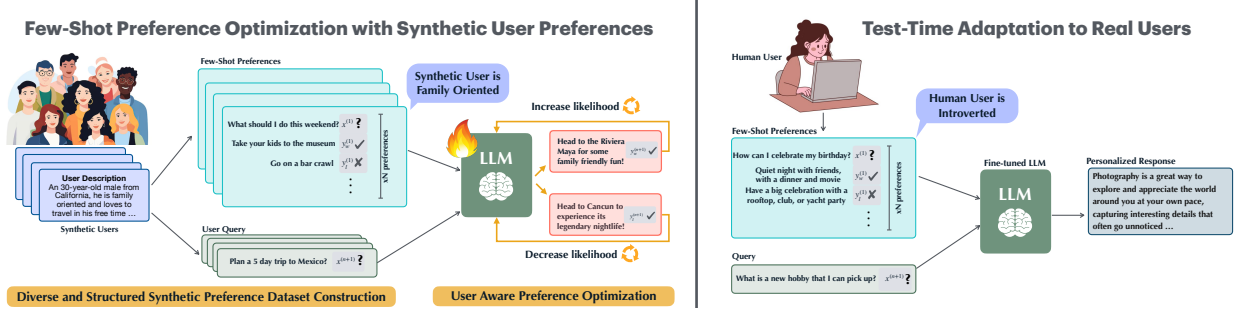


Figure 1: Overview of FSPO. N previously collected preferences are fed into the LLM along with the current query, allowing the LLM to personalize its response to the query using the past preferences.

modeling (Chen et al., 2022; Min et al., 2022; Yu et al., 2024), we fine-tune the model with a meta-learning objective, using preference-learning 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 following 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 time-consuming. 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 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 high-quality 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 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 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

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 fine-tuned 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 $\mathcal{D}_{\text{pref}} = \{(\mathbf{x}^{(i)}, \mathbf{y}_w^{(i)}, \mathbf{y}_l^{(i)})\}$ 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 \mathbf{y}_1 over \mathbf{y}_2 , given a prompt \mathbf{x} , as:

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

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 fine-tuned

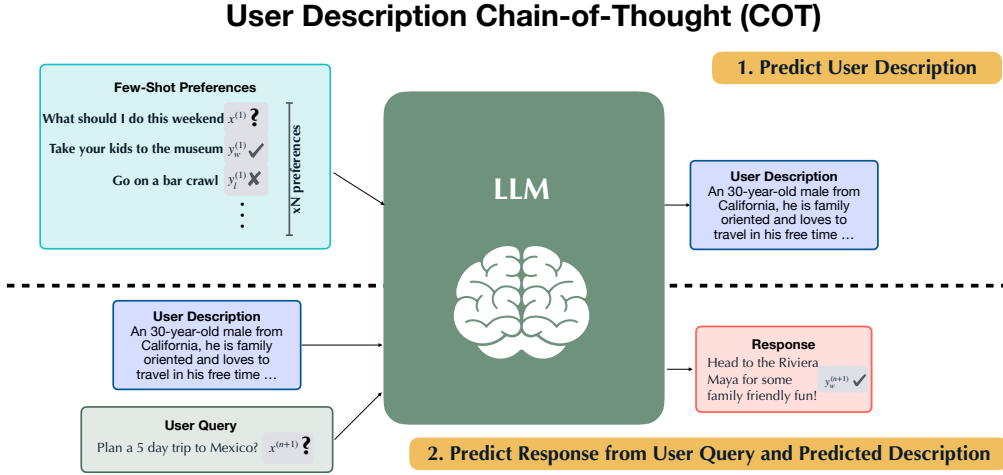


Figure 2: User Description Chain-of-Thought (COT). Prediction is a two-stage process: first predicting a (synthetic) user description from the few-shot preferences and next predicting the response.

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:

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

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 $\pi_\theta(y \mid x)$, the Bradley-Terry framework allows for either explicitly learning a reward model $r_\phi(x, y)$ or directly optimizing preferences. Explicit reward models are trained using the following classification objective:

$$\max_{\phi} \mathbb{E}_{\mathcal{D}_{\text{pref}}} [\log \sigma(r_\phi(x, y_w) - r_\phi(x, y_l))] \quad (3)$$

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 $\log \pi_\theta(y \mid x)$ as an implicit reward:

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

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: $\mathcal{D}_{\text{pref}} = \{(x^{(i)}, y_w^{(i)}, y_l^{(i)})\}$ 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 $\mathbf{S}^{(i)}$ 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 $\mathcal{S}^{(i)}$ ’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 $\mathcal{S} = P(\mathcal{S}^{(i)})$, a meta-learning objective can be derived to minimize its expected loss with respect to θ as:

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

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

$$\mathcal{L}_{\text{pref}}^{\theta} = \|h_{\pi_{\theta}}^{y_w, y_l} - (2\beta)^{-1}\|_2^2, \quad h_{\pi_{\theta}}^{y_w, y_l} = \log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \log \frac{\pi_{\theta}(y_l|x)}{\pi_{\text{ref}}(y_l|x)} \quad (6)$$

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.

Following black-box meta-learning approaches, FSPO receives as input a sequence of preferences $D_i^{\text{fewshot}} \sim D_i$ from a User $\mathcal{S}^{(i)}$. This is followed by an unlabeled, held-out preference $(x, y_1, y_2) \sim D_i \setminus D_i^{\text{fewshot}}$ 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 held-out 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 $\beta \frac{\log \pi_{\theta}(y|x)}{\log \pi_{\text{ref}}(y|x)}$. 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.

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)

- 1: **Input:** For each unique user $\mathcal{S}^{(i)}$, a dataset of preferences $\mathcal{D} := (x, y_1, y_2, c)_i$, and optionally user description $y_{\mathcal{S}^{(i)}}$ for COT, $\forall i$
 - 2: **Output:** Learned policy π_θ
 - 3: **while** not done **do**
 - 4: Sample training user $\mathcal{S}^{(i)}$ (or minibatch)
 - 5: Sample a subset of preferences from the user $\mathcal{D}_i^{\text{fewshot}} \sim \mathcal{D}_i$
 - 6: Sample held-out preference examples $\mathcal{D}_i^{\text{heldout}} \sim \mathcal{D}_i \setminus \mathcal{D}_i^{\text{fewshot}}$
 - 7: **if** COT **then**
 - 8: Use Eq. (5) and Eq. (6) to predict the loss on the user description $y_{\mathcal{S}^{(i)}}$
 - 9: **end if**
 - 10: Conditioning on $\mathcal{D}_i^{\text{fewshot}}$ (optionally $y_{\mathcal{S}^{(i)}}$), use Eq. (5) and Eq. (6) to predict the loss on the held-out preference example $\mathcal{D}_i^{\text{heldout}}$
 - 11: Update learner parameters θ , using gradient of loss on $\mathcal{D}_i^{\text{heldout}}$
 - 12: **end while**
 - 13: **Return** π_θ
-

Takeaways from FSPO framework

FSPO offers an effective approach to personalizing open-ended question answering, by framing personalization as a meta-learning problem, conditioned on few-shot preferences from a user as seen in Figure 1. Additionally, FSPO can be converted to a two-step prediction problem, predicting a user description conditioned on preferences and then a response, leveraging additionally inference-compute and the model’s instruction-tuned prior for better performance as seen Figure 2. We summarize the algorithm framework in Algorithm 1.

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 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 **ELIX-easy** 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 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?

Takeaways from Personalization Domains

We propose a benchmark consisting of 3 domains, where personalization can be studied: (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.

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 non-language 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

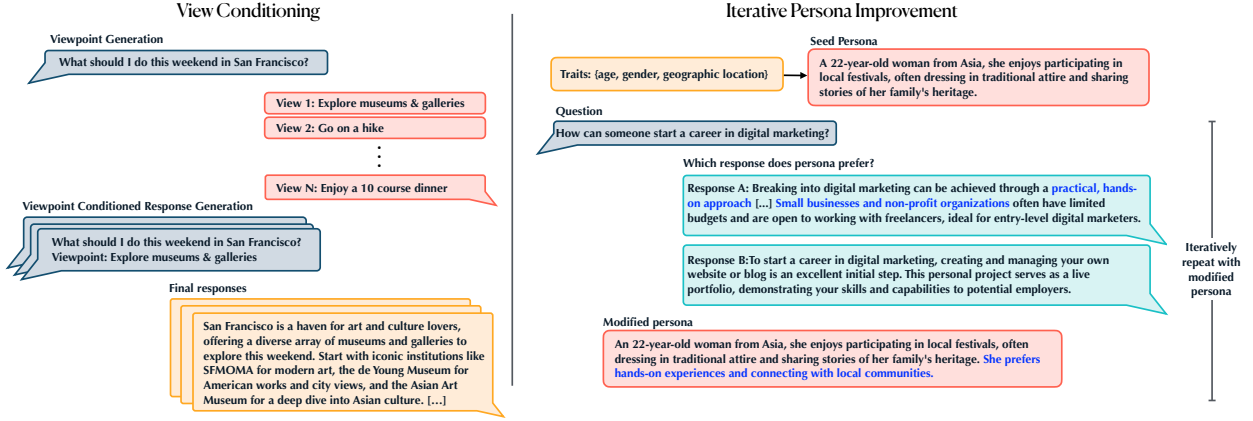


Figure 3: Overview of Domain Randomization Techniques. View-Conditioning (left) decomposes a given question into multiple viewpoints, allowing for diverse response generation. Iterative Persona Generation (right) allows for better structure by removing underspecification of the persona by iteratively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

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-4o (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-4o 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.

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 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

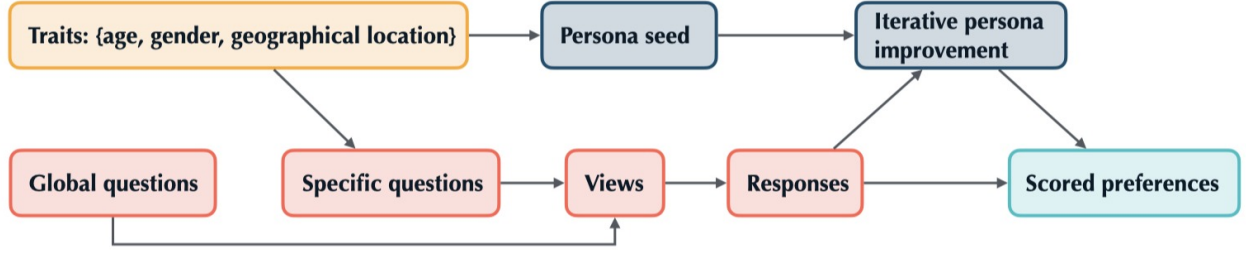


Figure 4: Flowchart of Roleplay dataset generation: Starting from a set of traits, a seed persona is constructed and a set of specific questions about that trait. Then responses are constructed with View-Conditioning. The seed personas are then iteratively refined to not be underspecified. Finally, the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

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.

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 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

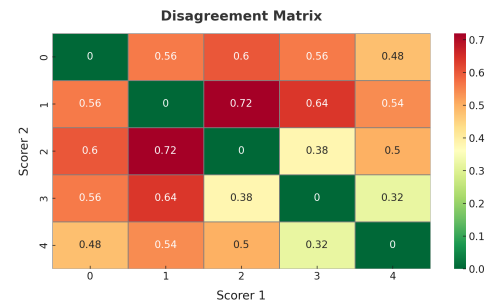


Figure 5: Disagreement Matrix across 5 users in Roleplay. Here we plot the disagreement of preferences for 5 users. There is a mix of users with high and low disagreement.

| 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 3: Roleplay: Human Eval Winrates

| 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 2: GPT-4o 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 |

Table 4: Review Winrates - Trained and Interpolated Users

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 5. We outline our full dataset generation process in Figure 4 in the Roleplay Task, starting from just a simple set of demographic traits.

Takeaways from Sim2Real and Domain Randomization

Since collecting personalized data at scale is challenging, we propose instead to generate diverse synthetic preference datasets that can be transferred to real humans. We study two design decisions to effectively encourage this transfer: (1) **Encouraging Diversity** and (2) **Structured Task Construction** and discuss approaches to instantiate these design choices.

7. Experimental Evaluation

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, ground truth user description, representing an upper bound on

FSPO’s performance.

Synthetic winrates. We first generate automated win rates using the modified AlpacaEval procedure from Section 6. In the ELIX task in Table 2, 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 better performance on held-out questions. Finally, in Table 1, we study Roleplay, scaling to 1500 real users, seeing a win rate of 82.6% on both held-out users and questions. Additionally, COT closes the gap to the oracle response, showing effective recovery of the ground-truth user description. In Appendix A.1, sample generations from FSPO show effective personalization to the oracle user description. Given this result, can we personalize to real people?

Preliminary human study. We evaluate our model trained on the Roleplay task by personalizing responses 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 3, 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.

Takeaways from Experiments

We evaluate FSPO on the 3 tasks discussed and find an **87% Alpaca Eval winrate** on average in generating responses that are personalized to synthetic users. COT also enables us to close the gap to the oracle method, where we prompt with the ground truth persona. Additionally, we run a preliminary, controlled human study, where we find a **72% winrate** with real human users for open-ended question answering.

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, 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

There are several limitations and potential risks. One limitation pertains to the ethical and fairness considerations of personalization. While FSPO improves inclusivity by modeling diverse preferences, the risk of reinforcing user biases (echo chambers) or inadvertently amplifying harmful viewpoints requires careful scrutiny. Future work should explore mechanisms to balance personalization with ethical safeguards, ensuring that models remain aligned with fairness principles while respecting user individuality. Additionally, our human study was preliminary with control over the questions that a user may ask, format normalization where formatting details such as markdown are removed, and view normalization comparing the same number of viewpoints for both FSPO and the baselines.

However, to the best of our knowledge, we are the first approach to perform such a human study for personalization to open-ended question answering. Future work should do further ablations with human evaluation for personalization. Additionally, due to compute constraints, we work with models in the parameter range of 3B (specifically Llama 3.2 Instruct 3B) with a limited context window of 128K, and without context optimization such as sequence 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 with better long-context and reasoning capabilities would help with FSPO for personalization, such as the 2M context window of Gemini Flash Thinking models, especially in the case of COT.

10. Acknowledgements

We would like to thank Yuejiang Liu, Aviral Kumar, Rafael Rafeliov, Violet Xiang, Yoonho Lee, Meihua Dang, and others in Stanford IRIS Lab and Stanford Hashimoto Group for discussions and feedback. AS gratefully acknowledges the support of the NSF Graduate Research Fellowship Program and Toyota Research Institute (TRI) for compute and API credits. TH was supported by a HAI seed grant, gifts from Open Philanthropy, Amazon, Schmidt Sciences, the Tianqiao and Chrissy Chen Foundation and a grant under the NSF CAREER IIS-2338866.

References

- 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. Constitutional ai: Harmlessness from ai feedback, 2022. URL <https://arxiv.org/abs/2212.08073>.
- Ralph Allan Bradley and Milton 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>.
- 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 McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020.
- Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, J  r  my 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. Open problems and fundamental limitations of reinforcement learning from human feedback. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL <https://openreview.net/forum?id=bx24KpJ4Eb>. Survey Certification.
- Louis Castricato, Nathan Lile, Rafael Rafailov, Jan-Philipp Fr  nken, and Chelsea Finn. Persona: A reproducible testbed for pluralistic alignment, 2024. URL <https://arxiv.org/abs/2407.17387>.
- Yanda Chen, Ruiqi Zhong, Sheng Zha, George Karypis, and He He. Meta-learning via language model in-context tuning, 2022. URL <https://arxiv.org/abs/2110.07814>.
- Myra Cheng, Esin Durmus, and Dan Jurafsky. Marked personas: Using natural language prompts to measure stereotypes in language models. In Anna Rogers, Jordan Boyd-Graber, and Naoaki 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>.
- Yan Duan, John Schulman, Xi Chen, Peter L Bartlett, Ilya Sutskever, and Pieter Abbeel. RL  : Fast reinforcement learning via slow reinforcement learning. *arXiv preprint arXiv:1611.02779*, 2016.

- Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B. Hashimoto. Length-controlled alpaca-eval: A simple way to debias automatic evaluators, 2024a. URL <https://arxiv.org/abs/2404.04475>.
- Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. AlpacaFarm: A simulation framework for methods that learn from human feedback, 2024b.
- Mohammad Gheshlaghi Azar, Mark Rowland, Bilal Piot, Daniel Guo, Daniele Calandriello, Michal Valko, and Rémi Munos. A General Theoretical Paradigm to Understand Learning from Human Preferences. *arXiv e-prints*, art. arXiv:2310.12036, October 2023. doi: 10.48550/arXiv.2310.12036.
- Goodreads. Goodreads: Book reviews, recommendations, and discussion, 2025. URL <https://www.goodreads.com/>. 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, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng 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 Lakhotia, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas 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 Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, 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 Vasic, 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, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Rapparthi, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney 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, Vitor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyan Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang,

Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew 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 Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkan Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Rutu Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj 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 DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. The llama 3 herd of models, 2024.
- Caglar Gulcehre, Tom Le Paine, Srivatsan Srinivasan, Ksenia Konyushkova, Lotte Weerts, Abhishek Sharma, Aditya Siddhant, Alex Ahern, Miaosen Wang, Chenjie Gu, et al. Reinforced self-training (rest) for language modeling. *arXiv preprint arXiv:2308.08998*, 2023.
- Kyle Hsu, Sergey Levine, and Chelsea 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.
- Joel Jang, Seungone Kim, Bill Yuchen Lin, Yizhong Wang, Jack Hessel, Luke Zettlemoyer, Hannaneh Hajishirzi, Yejin Choi, and Prithviraj Ammanabrolu. Personalized soups: Personalized large language model alignment via post-hoc parameter merging, 2023. URL <https://arxiv.org/abs/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. 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>.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention, 2023. URL <https://arxiv.org/abs/2309.06180>.
- Yoonho Lee, Jonathan Williams, Henrik Marklund, Archit Sharma, Eric Mitchell, Anikait Singh, and Chelsea Finn. Test-time alignment via hypothesis reweighting, 2024. URL <https://arxiv.org/abs/2412.08812>.
- 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. Synthetic data (almost) from scratch: Generalized instruction tuning for language models, 2024. URL <https://arxiv.org/abs/2402.13064>.
- Shenggui Li, Fuzhao Xue, Chaitanya Baranwal, Yongbin Li, and Yang You. Sequence parallelism: Long sequence training from system perspective, 2022. URL <https://arxiv.org/abs/2105.13120>.
- Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. Learning word vectors for sentiment analysis. In Dekang Lin, Yuji Matsumoto, and Rada Mihalcea, editors, *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 142–150, Portland, Oregon, USA, June 2011. Association for Computational Linguistics. URL <https://aclanthology.org/P11-1015/>.
- Viktor Makoviychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles Macklin, David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, and Gavriel State. Isaac gym: High performance gpu-based physics simulation for robot learning, 2021. URL <https://arxiv.org/abs/2108.10470>.

- Nicole Meister, Carlos Guestrin, and Tatsunori Hashimoto. Benchmarking distributional alignment of large language models, 2024. URL <https://arxiv.org/abs/2411.05403>.
- Igor Melnyk, Youssef Mroueh, Brian Belgodere, Mattia Rigotti, Apoorva Nitsure, Mikhail Yurochkin, Kristjan Greenewald, Jiri Navratil, and Jerret Ross. Distributional preference alignment of llms via optimal transport, 2024. URL <https://arxiv.org/abs/2406.05882>.
- Sewon Min, Mike Lewis, Luke Zettlemoyer, and Hannaneh Hajishirzi. Metaicl: Learning to learn in context, 2022. URL <https://arxiv.org/abs/2110.15943>.
- Nikhil Mishra, Mostafa Rohaninejad, Xi Chen, and Pieter 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, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Mądry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, Alex Nichol, Alex Paino, Alex Renzin, Alex Tachard Passos, Alexander Kirillov, Alexi Christakis, Alexis Conneau, Ali Kamali, Allan Jabri, Allison Moyer, Allison Tam, Amadou Crookes, 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 Gierler, 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, 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 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 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 Quinero 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, 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 Kondraciuk, Lukasz Kaiser, Luke Hewitt, Luke Metz, Lyric Doshi, Mada Aflak, Maddie Simens, Madelaine Boyd, Madeleine Thompson, Marat Dukhan, Mark Chen, Mark Gray, Mark Hudnall, Marvin Zhang, Marwan Aljubeh, Mateusz Litwin, Matthew Zeng, Max Johnson, Maya Shetty, Mayank Gupta, Meghan Shah, Mehmet Yatbaz, Meng Jia Yang, Mengchao Zhong, Mia Glaese, Mianna Chen, Michael Janner, Michael Lampe, Michael Petrov, Michael Wu, Michele Wang, Michelle Fradin, Michelle Pokrass, Miguel Castro, Miguel Oom Temudo de Castro, Mikhail Pavlov, Miles Brundage, Miles Wang, Minal Khan, Mira Murati, Mo Bavarian, Molly Lin, Murat Yesildal, Nacho Soto, Natalia Gimelshein, Natalie Cone, Natalie Staudacher, Natalie Summers, Natan LaFontaine, Neil Chowdhury, Nick Ryder, Nick Stathas, Nick Turley, Nik Tezak, Niko Felix, Nithanth Kudige, Nitish Keskar, Noah Deutsch, Noel Bundick, Nora Puckett, Ofir Nachum, Ola Okelola, Oleg Boiko, Oleg Murk, Oliver Jaffe, Olivia Watkins, Olivier Godement, Owen Campbell-Moore, Patrick Chao, Paul McMillan, Pavel Belov, Peng Su, Peter Bak, Peter Bakkum, Peter Deng, Peter Dolan, Peter Hoeschele, Peter Welinder, Phil Tillet, Philip Pronin, Philippe Tillet, Prafulla Dhariwal, Qiming Yuan, Rachel Dias, Rachel Lim, Rahul Arora, Rajan Troll, Randall Lin, Rapha Gontijo Lopes, Raul Puri, Reah Miyara, Reimar Leike, Renaud Gaubert, Reza Zamani, Ricky Wang, Rob Donnelly, Rob Honsby, Rocky Smith, Rohan Sahai, Rohit Ramchandani, Romain Huet, Rory Carmichael, Rowan Zellers, Roy Chen, Ruby Chen, Ruslan Nigmatullin, Ryan Cheu, Saachi Jain, Sam Altman, Sam Schoenholz, Sam Toizer, Samuel Miserendino, Sandhini Agarwal, Sara Culver, Scott Ethersmith, Scott Gray, Sean Grove, Sean Metzger, Shamez Hermani, Shantanu Jain, Shengjia Zhao, Sherwin Wu, Shino Jomoto, Shirong Wu, Shuaiqi, Xia, Sonia Phene, Spencer Papay, Srinivas Narayanan, Steve Coffey, Steve Lee, Stewart Hall, Suchir Balaji, Tal Broda, Tal Stramer, Tao Xu, Tarun Gogineni, Taya Christianson, Ted Sanders, Tejal Patwardhan, Thomas Cunningham, Thomas Degry, Thomas Dimson, Thomas Raoux, Thomas Shadwell, Tianhao Zheng, Todd Underwood, Todor Markov, Toki Sherbakov, Tom Rubin, Tom Stasi, Tomer Kaftan, Tristan Heywood, Troy Peterson, Tyce Walters, Tyna Eloundou, Valerie Qi, Veit Moeller, Vinnie Monaco, Vishal Kuo, Vlad Fomenko, Wayne Chang, Weiyi Zheng, Wenda Zhou, Wesam Manassra, Will Sheu, Wojciech Zaremba, Yash Patil, Yilei Qian, Yongjik Kim, Youlong Cheng, Yu Zhang, Yuchen He, Yuchen Zhang, Yujia Jin, Yunxing Dai, and Yury Malkov. Gpt-4o system card, 2024. URL <https://arxiv.org/abs/2410.21276>.

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. Training language models to follow instructions with human feedback. *arXiv e-prints*, art. arXiv:2203.02155, March 2022. doi: 10.48550/arXiv.2203.02155.

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. Training language models to follow instructions with human feedback, 2022.

Sriyash Poddar, Yanming Wan, Hamish Ivison, Abhishek Gupta, and Natasha Jaques. Personalizing reinforcement learning from human feedback with variational preference learning, 2024. URL <https://arxiv.org/abs/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. Qwen2.5 technical report, 2025. URL <https://arxiv.org/abs/2412.15115>.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *arXiv preprint arXiv:2305.18290*, 2023.
- Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy Lillicrap. One-shot learning with memory-augmented neural networks, 2016. URL <https://arxiv.org/abs/1605.06065>.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal Policy Optimization Algorithms. *arXiv e-prints*, art. arXiv:1707.06347, July 2017. doi: 10.48550/arXiv.1707.06347.
- Omar Shaikh, Michelle Lam, Joey Hejna, Yijia Shao, Michael Bernstein, and Diyi Yang. Show, don’t tell: Aligning language models with demonstrated feedback, 2024. URL <https://arxiv.org/abs/2406.00888>.
- Anand Siththaranjan, Cassidy Laidlaw, and Dylan Hadfield-Menell. Distributional preference learning: Understanding and accounting for hidden context in rlhf, 2024. URL <https://arxiv.org/abs/2312.08358>.
- Taylor Sorensen, Jared Moore, Jillian Fisher, Mitchell Gordon, Niloofar Miresghallah, Christopher Michael Rytting, Andre Ye, Liwei Jiang, Ximing Lu, Nouha Dziri, Tim Althoff, and Yejin Choi. A roadmap to pluralistic alignment, 2024. URL <https://arxiv.org/abs/2402.05070>.
- Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In S. Solla, T. Leen, and K. Müller, editors, *Advances in Neural Information Processing Systems*, volume 12. MIT Press, 1999. URL https://proceedings.neurips.cc/paper_files/paper/1999/file/464d828b85b0bed98e80ade0a5c43b0f-Paper.pdf.
- Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhatipatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Ferret, Peter Liu, Pouya Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Charline Le Lan, Sammy Jerome, Anton Tsitsulin, Nino Vieillard, Piotr Stanczyk, Sertan Girgin, Nikola Momchev, Matt Hoffman, Shantanu Thakoor, Jean-Bastien Grill, Behnam Neyshabur, Olivier Bachem, Alanna Walton, Aliaksei Severyn, Alicia Parrish, Aliya Ahmad, Allen Hutchison, Alvin Abdagic, Amanda Carl, Amy Shen, Andy Brock, Andy Coenen, Anthony Laforge, Antonia Paterson, Ben Bastian, Bilal Piot, Bo Wu, Brandon Royal, Charlie Chen, Chintu Kumar, Chris Perry, Chris Welty, Christopher A. Choquette-Choo, Danila Sinopalnikov, David Weinberger, Dimple Vijaykumar, Dominika Rogozińska, Dustin Herbison, Elisa Bandy, Emma Wang, Eric Noland, Erica Moreira, Evan Senter, Evgenii Eltyshev, Francesco Visin, Gabriel Rasskin, Gary Wei, Glenn Cameron, Gus Martins, Hadi Hashemi, Hanna Klimczak-Plucińska, Harleen Batra, Harsh Dhand, Ivan Nardini, Jacinda Mein, Jack Zhou, James Svensson, Jeff Stanway, Jetha Chan, Jin Peng Zhou, Joana Carrasqueira, Joana Iljazi, Jocelyn Becker, Joe Fernandez, Joost van Amersfoort, Josh Gordon, Josh Lipschultz, Josh Newlan, Ju yeong Ji, Kareem Mohamed, Kartikeya Badola, Kat Black, Katie Millican, Keelin McDonell, Kelvin Nguyen, Kiranbir Sodhia, Kish Greene, Lars Lowe Sjoesund, Lauren Usui, Laurent Sifre, Lena Heuermann, Leticia Lago, Lilly McNealus, Livio Baldini Soares, Logan Kilpatrick, Lucas Dixon, Luciano Martins, Machel Reid, Manvinder Singh, Mark Iverson, Martin Görner, Mat Velloso, Mateo Wirth, Matt Davidow, Matt Miller, Matthew Rahtz, Matthew Watson, Meg Risdal, Mehran Kazemi, Michael

- Moynihan, Ming Zhang, Minsuk Kahng, Minwoo Park, Mofi Rahman, Mohit Khatwani, Natalie Dao, Nenshad Bardoliwalla, Nesh Devanathan, Neta Dumai, Nilay Chauhan, Oscar Wahltinez, Pankil Botarda, Parker Barnes, Paul Barham, Paul Michel, Pengchong Jin, Petko Georgiev, Phil Culliton, Pradeep Kuppala, Ramona Comanescu, Ramona Merhej, Reena Jana, Reza Ardeshir Rokni, Rishabh Agarwal, Ryan Mullins, Samaneh Saadat, Sara Mc Carthy, Sarah Cogan, Sarah Perrin, Sébastien M. R. Arnold, Sebastian Krause, Shengyang Dai, Shruti Garg, Shruti Sheth, Sue Ronstrom, Susan Chan, Timothy Jordan, Ting Yu, Tom Eccles, Tom Hennigan, Tomas Kocisky, Tulsee Doshi, Vihan Jain, Vikas Yadav, Vilobh Meshram, Vishal Dharmadhikari, Warren Barkley, Wei Wei, Wenming Ye, Woohyun Han, Woosuk Kwon, Xiang Xu, Zhe Shen, Zhitao Gong, Zichuan Wei, Victor Cotruta, Phoebe Kirk, Anand Rao, Minh Giang, Ludovic Peran, Tris Warkentin, Eli Collins, Joelle Barral, Zoubin Ghahramani, Raia Hadsell, D. Sculley, Jeanine Banks, Anca Dragan, Slav Petrov, Oriol Vinyals, Jeff Dean, Demis Hassabis, Koray Kavukcuoglu, Clement Farabet, Elena Buchatskaya, Sebastian Borgeaud, Noah Fiedel, Armand Joulin, Kathleen Kenealy, Robert Dadashi, and Alek Andreev. Gemma 2: Improving open language models at a practical size, 2024. URL <https://arxiv.org/abs/2408.00118>.
- Joshua Tobin, Lukas Biewald, Rocky Duan, Marcin Andrychowicz, Ankur Handa, Vikash Kumar, Bob McGrew, Jonas Schneider, Peter Welinder, Wojciech Zaremba, and Pieter Abbeel. Domain randomization and generative models for robotic grasping, 2018. URL <https://arxiv.org/abs/1710.06425>.
- Jane X Wang, Zeb Kurth-Nelson, Dhruva Tirumala, Hubert Soyer, Joel Z Leibo, Remi Munos, Charles Blundell, Dharshan Kumaran, and Matt Botvinick. Learning to reinforcement learn. *arXiv preprint arXiv:1611.05763*, 2016.
- Amy Yang, Jingyi Yang, Aya Ibrahim, Xinfeng Xie, Bangsheng Tang, Grigory Sizov, Jeremy Reizenstein, Jongsoo Park, and Jianyu Huang. Context parallelism for scalable million-token inference, 2024. URL <https://arxiv.org/abs/2411.01783>.
- Ze Yang, Yun Chen, Jingkang Wang, Sivabalan Manivasagam, Wei-Chiu Ma, Anqi Joyce Yang, and Raquel Urtasun. Unisim: A neural closed-loop sensor simulator, 2023. URL <https://arxiv.org/abs/2308.01898>.
- Mingzhang Yin, George Tucker, Mingyuan Zhou, Sergey Levine, and Chelsea Finn. Meta-learning without memorization. *arXiv preprint arXiv:1912.03820*, 2019.
- Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T. Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions for large language models, 2024. URL <https://arxiv.org/abs/2309.12284>.
- Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Xian Li, Sainbayar Sukhbaatar, Jing Xu, and Jason Weston. Self-Rewarding Language Models. *arXiv e-prints*, art. arXiv:2401.10020, January 2024. doi: 10.48550/arXiv.2401.10020.
- Lunjun Zhang, Arian Hosseini, Hritik Bansal, Mehran Kazemi, Aviral Kumar, and Rishabh Agarwal. Generative verifiers: Reward modeling as next-token prediction, 2024. URL <https://arxiv.org/abs/2408.15240>.
- Siyan Zhao, John Dang, and Aditya Grover. Group preference optimization: Few-shot alignment of large language models, 2024. URL <https://arxiv.org/abs/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. Sglang: Efficient

execution of structured language model programs, 2024. URL <https://arxiv.org/abs/2312.07104>.

A. Appendix

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 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 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.

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.

| Name | Value |
|-----------------------------|--|
| Learning Rate (SFT/Pref-FT) | $1e^{-5}$, $1e^{-6}$, $1e^{-7}$ |
| Learning Rate (IPO) | $1e^{-5}$, $1e^{-6}$, $1e^{-7}$ |
| Beta (IPO) | 0.1, 0.05, 0.01, 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)

| | |
|--|--|
| <p>Q: How does a thermostat control temperature?</p> <p>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 superheroes that help control the temperature in your home. Imagine you're playing with a thermostat in your room. It's like a special button that makes the temperature go up or down. But how does it actually work?</p> <p>Thermostats are made up of three main parts:</p> <ol style="list-style-type: none"> 1. Sensors: These are like special eyes that look at the room's temperature. They tell the thermostat how warm or cool it is. 2. Control Unit: This is like the brain of the thermostat. It looks at the sensor's readings and decides what to do next. 3. 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. <p>Here's how it works:</p> <ul style="list-style-type: none"> • The sensors send a message to the control unit, like "Hey, it's 75 degrees in here!" • If the control unit thinks the room is too cold, it sends a signal to the heating system, saying "Turn on the heater!" <p>[...]</p> | <p>Q: What makes diamonds so hard?</p> <p>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:</p> <p>Key Factors Behind Diamond Hardness</p> <ul style="list-style-type: none"> • Crystal Structure: Diamonds possess a face-centered 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. • Chemical Composition: Composed solely of carbon atoms arranged in a repeating pattern, diamonds owe their remarkable hardness to this carbon-rich composition. • 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. <p>[...]</p> |
|--|--|

Elementary School Student

Expert

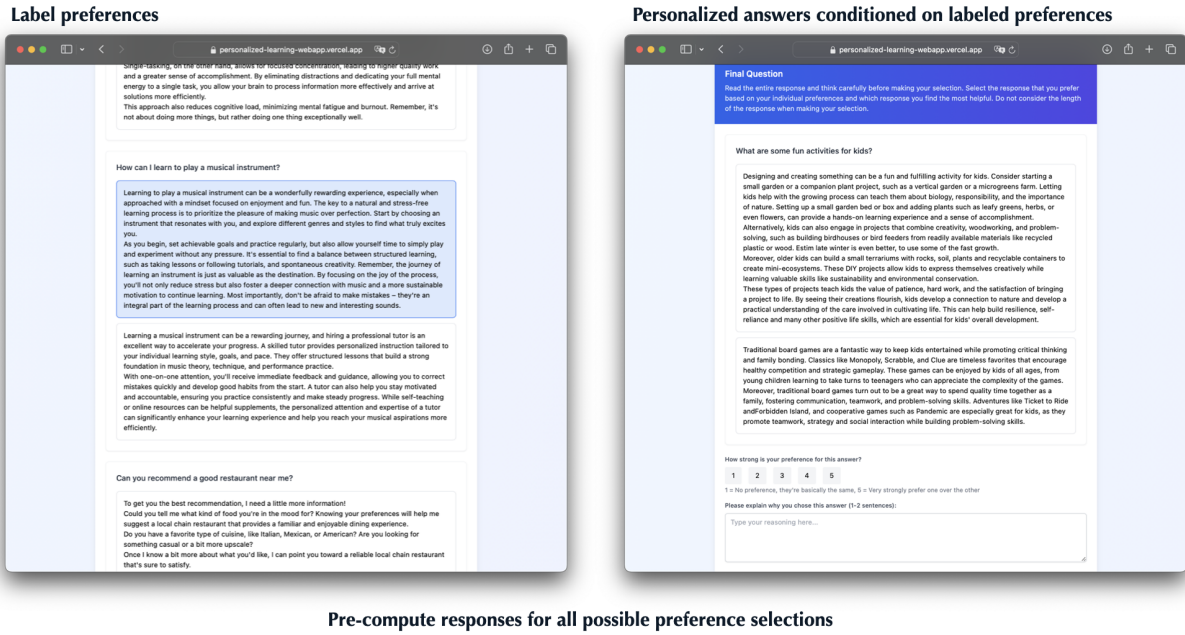
Personalized Responses (Reviews)

| | |
|--|---|
| <p>Write a review for the Action movie, The Dark Knight.</p> <p>The Dark Knight is a thrilling ride, but its flaws hold it back from being a truly great movie. The film's pacing is well-balanced, and the action scenes are intense and well-executed. The cast, led by Christian Bale and Heath Ledger, deliver solid performances. However, the plot is somewhat predictable, and the characters' motivations are not always clear. Overall, The Dark Knight is a fun, but forgettable, superhero film.</p> | <p>Write a review for the Action movie, The Dark Knight.</p> <p><i>"A Flop of Catastrophic Proportions: 'The Dark Knight Rises' Falls Woefully Short of its Pompous Promises"</i></p> <p>"The Dark Knight Rises," Christopher Nolan's vaunted swan song to the Batman franchise, is a calamitous, soulless, and creatively bankrupt exercise in superheroic excess, a monstrosity of monumental proportions that has left me questioning the very fabric of cinematic reality.</p> <p>From its opening frames, it is evident that Nolan has traded his artistic integrity and attention to detail for a saccharine, crowd-pleasing formula that prioritizes style over substance, spectacle over sensibility. The narrative, a jumbled mess of convoluted plot twists, cheap thrills, and cringe-worthy romantic subplots, careens wildly from one paradox to the next, sabotaging any semblance of coherence or logical consistency.</p> <p>Tom Hardy's Bane, the film's ostensibly formidable villain, is a paragon of cartoonish one-dimensionality, his motivations and demeanor reduced to a simplistic, mustache-twirling caricature of menace. The usually reliable Christian Bale, as Bruce Wayne/Batman, fares little better, his character reduced to a puzzle piece of angst-ridden introspection, his entire persona rendered inexplicably thin and unengaging.</p> <p>[...] Continues</p> |
|--|---|

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.

A.2. Additional Details on Dataset Construction and Evaluation

Dataset Details We include a representative subset of the roleplay task and a jupyter notebook 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.

Evaluation Prompt for Synthetic Preferences We used GPT-4o as a Judge using a modified variant of the Alpaca Eval (Dubois et al., 2024b,a) Prompt to be aware of a user description when scoring 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

You will receive:

1. A **User Instruction**: This is the query or task provided to the models.

2. ****Model Outputs****: Unordered responses from different models, each identified by a unique model identifier.
3. 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 relevance to the `users` instruction and description.
2. Select the best output that meets the user's needs.

Input Format

User Instruction
{QUESTION}

Model Outputs
– Model "m": {RESPONSE_A}
– Model "M": {RESPONSE_B}

User Description
{USER_DESCRIPTION}

Task

From the provided outputs, determine which model produces the best response. Output only the model identifier of the best response (either 'm' or 'M') with no additional text, quotes, spaces, or new lines.

Best Model Identifier

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 and baselines in 3 different categories. First, length is an evaluation bias. For this reason, we computed the average length of responses from FSPO and prompted the base model during evaluation to keep its responses around the average length in words (≈ 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 open-source LLMs, we decided to have formatting be consistent as paragraphs rather than markdown for the Roleplay task. Thus, we similarly prompted the Base model with this behavior. Finally, a differing number of views can also skew the evaluation, as a large proportion of users seem to prefer direct answers. Additionally, if more views are presented, a 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 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.

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. For the IPO baseline, we followed a similar procedure. Additional hyperparameters can be found in Table 5.

A.4. Additional Details of Setup for Reproducibility

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 open-source 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, GPT4o-mini) 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. In terms of statistics of the dataset, the review dataset has 130K train/dev examples and 32.4K test examples, the ELIX-easy dataset has 235K train/dev examples and 26.1K test examples, the ELIX-hard dataset has 267K train/dev examples and 267K test examples, 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](#)).

Each of the artifacts above was consistent with its intended use and the code, models, and datasets should be usable outside of research contexts.