COPL: COLLABORATIVE PREFERENCE LEARNING FOR PERSONALIZING LLMS

Youngbin Choi¹, Seunghyuk Cho¹, Minjong Lee², MoonJeong Park¹, Yesong Ko², Jungseul Ok^{1,2}, Dongwoo Kim^{1,2}

¹Graduate School of Artificial Intelligence, POSTECH,

²Department of Computer Science and Engineering, POSTECH,

{choi.youngbin, shhj1998, minjong.lee, mjeongp, yesong.ko, jungseul, dongwoo.kim}@postech.ac.kr

Abstract

Personalizing large language models (LLMs) is important for aligning outputs with diverse user preferences, yet existing methods struggle with flexibility and generalization. We propose CoPL (Collaborative Preference Learning), a graphbased collaborative filtering framework that models user-response relationships to enhance preference estimation, particularly in sparse annotation settings. By integrating a mixture of LoRA experts, CoPL efficiently fine-tunes LLMs while dynamically balancing shared and user-specific preferences. Additionally, an optimization-free adaptation strategy enables generalization to unseen users without fine-tuning. Experiments on UltraFeedback-P demonstrate that CoPL outperforms existing personalized reward models, effectively capturing both common and controversial preferences, making it a scalable solution for personalized LLM alignment.

1 INTRODUCTION

Large language models (LLMs) have rapidly expanded across diverse applications, from customer service and tutoring to creative content generation Shi et al. (2024); Molina et al. (2024); Venkatraman et al. (2024). As increasing numbers of users with varied backgrounds interact with LLMs, accounting for diverse preferences has become essential. Most reward models rely on the Bradley-Terry-Luce (BTL) framework (Bradley & Terry, 1952), which learns preferences from pairwise comparisons provided by humans. However, earlier studies largely assumed a single, uniform preference and neglected the diversity of user preferences (Siththaranjan et al., 2024; Li et al., 2024). This limitation has led to growing interest in personalized reward models (Sorensen et al., 2024).

There are two different approaches to utilizing the BTL framework for personalized reward models. The first approach has explored combining multiple reward models, each trained for a specific preference and later aggregated (Jang et al., 2023; Oh et al., 2024). However, this approach relies on pre-trained models for different preference types, reducing flexibility. Another work introduces user latent variables into a single BTL framework, learning personalized representations from user annotations (Chen et al., 2024a; Poddar et al., 2024; Li et al., 2024). While this method captures individual preferences, the latent variable model does not explicitly account for relationships between users sharing similar responses. As a result, it struggles to generalize in sparse annotation settings.

To address these limitations, we propose Collaborative Preference Learning (CoPL), which constructs a user-response bipartite graph from pairwise annotations and uses a graph-based collaborative filtering (GCF) framework for personalized reward modeling. Unlike approaches that model each user separately, GCF on the graph structure allows preference signals to propagate across users, enabling to exploit multi-hop relationships among users and responses (Wang et al., 2019; He et al., 2020). As a result, CoPL can capture users' diverse preferences even in sparse annotation settings.

Based on the user embedding, we develop an LLM-based reward model that can predict the preference score of a user given input text. We adopt the mixture of LoRA experts (MoLE) (Chen et al.,

Correspondence to: Dongwoo Kim (dongwoo.kim@postech.ac.kr)

2023; 2024c; Liu et al., 2024) that allows parameter efficient fine-tuning while routing different users to different paths based on the learned embedding. Specifically, we develop a user preference-aware gating function that dynamically selects the experts in the forward pass, making the LLM predict a personalized preference.

While the reward model can predict preferences for users included in the training set, the model cannot handle newly participated *unseen* users whose embeddings are unknown. To estimate the preferences of unseen users, we propose an optimization-free adaptation method. Given a few annotations from an unseen user, we exploit the existing graph to find users with similar preferences and aggregate their embeddings to represent the unseen user.

Experimental results demonstrate that CoPL consistently outperforms existing personalized reward models in both seen and unseen users. Especially, CoPL generalizes to unseen users, maintaining high accuracy with only a few provided annotations. Embedding visualizations show that CoPL clusters users with similar preferences more closely than competing baselines. Further ablation studies confirm that both GCF and MoLE contribute significantly to performance.

2 RELATED WORK

In this section, we summarize relevant lines of research, such as personalized alignment and preference learning with sparse interactions.

Personalized alignment. With the growth of generative models, alignment has emerged as a crucial strategy for mitigating undesirable outcomes, such as biased or harmful outputs, and ensuring that the model works with human preference (Dai et al., 2023; Yang et al., 2024a). Alignment methods often rely on reward models. They typically build on the BTL framework, which relies on pairwise comparisons from various annotators. However, previous research has often focused on the average preference of annotators (Achiam et al., 2023), ignoring the diverse preferences.

To address preference diversity, recent works (Jang et al., 2023; Oh et al., 2024; Yang et al., 2024b) view this problem as a soft clustering problem, where user-specific preferences are treated as mixtures of predefined preference types. Although this approach effectively handles diverse preferences, it relies on specifying several preference types in advance.

Another line of work introduces user latent variable in the BTL framework (Poddar et al., 2024; Li et al., 2024; Chen et al., 2024a). Although extending the BTL framework with latent user variables can address diverse preferences, the main challenge lies in obtaining user representations. One approach is to treat each user embedding as learnable parameters, (Li et al., 2024; Chen et al., 2024a), and the other strategy is to train an encoder that infers embeddings from the small set of annotated pairs provided by each user (Poddar et al., 2024).

Preference learning with sparse interactions. Preference learning with sparse interactions is a well-studied challenge in recommendation systems, where each user typically interacts with only a small fraction of the available items. Despite these limited interactions, the system should infer the preference of each user and recommend additional items accordingly (He & Chua, 2017; Chen et al., 2020; Li et al., 2022; Lin et al., 2022). Collaborative filtering (CF) is a widely adopted solution that assumes users with similar interaction histories will exhibit similar preferences.

Graph-based CF (GCF) (Wang et al., 2019; He et al., 2020) has been considered one of the most advanced algorithms for a recommendation system. GCF leverages graph neural networks (GNNs) to capture preference through the connectivity among users and items. Many GCFs are developed based on an implicit feedback assumption (Rendle et al., 2012), where an edge between a user and an item reveals a preferable relation. Whereas in our setting, users provide explicit feedback given a pair of responses, making direct application of GCF unsuitable.

3 PROBLEM FORMULATION

We aim to develop a reward model that can capture diverse user preferences from a limited set of preference annotations. Instead of directly defining a user's preference, we collect pairwise comparisons indicating which item a user prefers. Let $\mathcal{U} = \{1, \dots, U\}$ be a set of users and \mathcal{X} be

a space of LLM's responses. To estimate the preferences of users, we first curate a survey set $S = \{(q_i, a_i, b_i)\}_{i=1}^R$ consisting of predefined questions q_i and two different responses $a_i, b_i \in \mathcal{X}$ from LLMs. For each user u, we first randomly sample N_u number of survey items and then collect the preferences over the response pairs, resulting in *preference dataset* \mathcal{D}_u . We use $(a \succ b) \in \mathcal{D}_u$ to denote that user u prefers response a over the response b. Given these pairwise preferences, we aim to learn a numerical reward function

$$f(u,r): \mathcal{U} \times \mathcal{X} \to \mathbb{R},\tag{1}$$

where f(u, r) represents a scalar *preference score* of response r for user u. The model is trained to satisfy

f(u,a) > f(u,b)

for all u and preference pairs $a \succ b$ observed in the data.

Following previous works (Li et al., 2024; Poddar et al., 2024), we consider the Bradly-Terry-Luce (BTL) choice model (Bradley & Terry, 1952) with maximum likelihood estimation to train the reward function. The likelihood of user u prefers item a over b can be defined using the BTL model as

$$p(a \succ b \mid u) = \frac{\exp(f(u, a))}{\exp(f(u, a)) + \exp(f(u, b))}.$$

Conversely, if b was chosen over a, i.e., $a \prec b$, the likelihood is

$$p(b \succ a \mid u) = 1 - p(a \succ b \mid u).$$

Through the maximum likelihood estimation with preference data for all users, one can learn the reward function f to make the reward function align with user preference. In the case of the universal preference model, user u is ignored in Eq. (1) (Chen et al., 2024b; Achiam et al., 2023; Dai et al., 2023; Bai et al., 2022). In practice, the user u is replaced by a user embedding (Poddar et al., 2024; Li et al., 2024; Chen et al., 2024a).

4 Method

In this section, we describe our Collaborative Preference Learning (CoPL). We first learn user embeddings based on GCF with the preference data. We then train the reward model based on the learned user embeddings. Finally, we provide an optimization-free adaptation strategy to obtain embeddings of users who are unseen during training.

4.1 USER REPRESENTATION LEARNING

Users who share similar preferences are likely to respond to similar responses. When the number of annotated responses is very small, it is unlikely to annotate the same responses between users. However, if we exploit multi-hop relations between users and responses, we may estimate user preference accurately. In fact, the exploitation of the relationship between users and items is the key idea behind graph-based collaborative filtering (GCF).

The preference dataset for all users can be naturally converted into a bipartite graph, where each user and response is represented as a node, and an edge between a user and a response represents the user's preference over the response. The edge can have two different types: positive or negative, indicating whether a user prefers the response or not.

Given a bipartite graph, we design a message-passing algorithm to update user and response representations. Let $e_u \in \mathbb{R}^d$ be an embedding vector of user u, and $e_r \in \mathbb{R}^d$ be an embedding vector of response r. Since there are two different edge types, we use different parameterizations for each type. Let \mathcal{N}_u^+ be a set of positive edges and \mathcal{N}_u^- be a set of negative edges from user u. Similarly, we can define \mathcal{N}_r^+ and \mathcal{N}_r^- for response r. Given user and response embeddings at layer ℓ , the message passing computes a message from neighborhood responses to the user as

$$\begin{split} \boldsymbol{m}_{u}^{+} &= \sum_{r \in \mathcal{N}_{u}^{+}} \alpha_{u,r} \Big(W_{1}^{(\ell)} \boldsymbol{e}_{r}^{(\ell)} + W_{2}^{(\ell)} (\boldsymbol{e}_{r}^{(\ell)} \odot \boldsymbol{e}_{u}^{(\ell)}) \Big), \\ \boldsymbol{m}_{u}^{-} &= \sum_{r \in \mathcal{N}_{u}^{-}} \beta_{u,r} \Big(W_{3}^{(\ell)} \boldsymbol{e}_{r}^{(\ell)} + W_{4}^{(\ell)} (\boldsymbol{e}_{r}^{(\ell)} \odot \boldsymbol{e}_{u}^{(\ell)}) \Big), \\ \boldsymbol{m}_{u}^{(\ell)} &= W_{\text{self}}^{(\ell)} \boldsymbol{e}_{u}^{(\ell)} + \boldsymbol{m}_{u}^{+} + \boldsymbol{m}_{u}^{-}, \end{split}$$
(2)

where $W_1^{(\ell)}, W_2^{(\ell)}, W_3^{(\ell)}, W_4^{(\ell)}, W_{self}^{(\ell)} \in \mathbb{R}^{d \times d}$ are parameter matrices, \odot is element-wise multiplication, and $\alpha_{u,r}$ and $\beta_{u,r}$ are normalization factors, set to $\frac{1}{\sqrt{|\mathcal{N}_u^+||\mathcal{N}_r^+|}}$ and $\frac{1}{\sqrt{|\mathcal{N}_u^-||\mathcal{N}_r^-|}}$, respectively.

Then, the user embedding is updated with the aggregated message $m_u^{(\ell)}$:

$$\boldsymbol{e}_{u}^{(\ell+1)} = \psi(\boldsymbol{m}_{u}^{(\ell)}), \tag{3}$$

where $\psi(\cdot)$ is a non-linear activation. The response embedding $e_r^{(\ell)}$ is updated with analogous process. We randomly initialize the user and response embeddings at the first layer and then fine-tune the embeddings through training. The update steps for the response embeddings are provided in Appendix A.

After L propagation steps, user and response embeddings accumulate information from their local neighborhood. Given the final user embedding $e_u^{(L)}$ and response embedding $e_r^{(L)}$, we use the inner product between the embeddings as a predicted preference :

$$s_{u,r} = \left(\boldsymbol{e}_{u}^{(L)}\right)^{\top} \left(\boldsymbol{e}_{r}^{(L)}\right). \tag{4}$$

With the score function, the GNN is trained on preference data D_u for all users by minimizing the following loss function:

$$\mathcal{L}_{\text{GCF}}(\theta) := \sum_{u \in \mathcal{U}} \sum_{(a \succ b) \in \mathcal{D}_u} -\log \sigma \left(s_{u,a} - s_{u,b} \right) + \lambda \|\theta\|_2^2, \tag{5}$$

where $\sigma(\cdot)$ denotes a sigmoid function, λ is a regularization hyper-parameter and θ represents all trainable parameters, including weights of the propagation layers and initial embeddings of the users $e_u^{(0)}$ and responses $e_r^{(0)}$.

4.2 PERSONALIZED REWARD MODEL WITH USER REPRESENTATIONS

Based on the learned user embeddings $e_u^{(L)}$, we build a reward model that can accommodate the preferences of diverse users. We use an LLM-based reward function:

$$f_{\phi}(\boldsymbol{e}_{u}, r) : \mathbb{R}^{d} \times \mathcal{X} \to \mathbb{R}$$

$$\tag{6}$$

where f is an LLM parameterized by ϕ taking user embedding e_u and the response r as inputs and predicts preference score. Unlike the response, the user embedding is not used as an input token. Instead, it is used in the gating mechanism described below. To learn the reward model, we can employ the BTL model, resulting in the maximum likelihood objective:

$$\mathcal{L}_{\mathsf{RM}}(\phi) = \sum_{u} \sum_{(a \succ b) \in \mathcal{D}_u} \log p_\phi(a \succ b \mid \boldsymbol{e}_u) \tag{7}$$

However, naively optimizing this objective starting from a pretrained LLM requires fine-tuning billions of parameters. Moreover, different preferences of users result in conflicting descent directions of the model parameters, resembling a multi-task learning scenario.

Mixture of LoRA experts for personalized reward function. For an efficient parameter update while minimizing the negative effect of diverse preferences, we adopt the mixture of LoRA experts (MoLE) (Hu et al., 2021; Liu et al., 2024) into our framework. MoLE is proposed to maximize the benefit of the mixture of experts (MoE) while maintaining efficient parameter updates. With MoLE, the model parameter matrix W is decomposed into pretrained and frozen W_0 and trainable ΔW ,

i.e., $W = W_0 + \Delta W$. ΔW is further decomposed into a shared LoRA expert $A_s \in \mathbb{R}^{d_{\text{out}} \times n}, B_s \in \mathbb{R}^{n \times d_{\text{in}}}$, which is used across all users, and M individual LoRA experts $\{A_i, B_i\}_{i=1}^M$ with the same dimensionality of the shared expert. Formally, this can be written as

$$\Delta W_u = A_s B_s + \sum_{i=1}^M w_i A_i B_i,\tag{8}$$

where $w_i \in [0, 1]$ denotes the importance of expert *i*.

To adopt the different preferences of users, we define a user-dependent gating mechanism to model the importance parameter w_i . For each user u, a gating function $g : \mathbb{R}^d \to \mathbb{R}^M$ maps $e_u^{(L)}$ to expert-selection logits:

$$\mathbf{z} = g(\boldsymbol{e}_u^{(L)}). \tag{9}$$

We convert these logits z into gating weight w_i by selecting the top one expert from the logits:

$$w_{i} = \begin{cases} \frac{\exp(z_{i}/\tau)}{\sum_{j=1}^{M} \exp(z_{j}/\tau)} & \text{if } i = \arg\max_{i} z_{i} \\ 0 & \text{otherwise,} \end{cases}$$
(10)

where τ is a temperature parameter. In practice, one can use top-k experts, but we could not find a significant difference in our experiments. For computational efficiency, we keep the top one expert.

4.3 Optimization-free User Adaptation

While we can predict a preference score of unseen responses for a known user, the reward model trained in Section 4.2 cannot be used to predict the preference of users who have not been observed during training. To estimate the embeddings of unseen users, we propose an optimization-free adaptation approach.

Let u^* be an unseen user who annotates a small set of response pairs. Under the assumption that users who have similar responses have similar preferences, we can estimate the embedding of an unseen user by taking an embedding of users with similar tastes. For example, if both user u^* and u share positive preference over the same response r, then we can use the embedding of uto approximate that of u^* . Based on this intuition, we propose the following optimization-free adaptation strategy for unseen user embedding:

$$\boldsymbol{e}_{u^{*}}^{(L)} = \sum_{u \in \mathcal{N}_{u^{*}}^{+}(k)} w_{u,u^{*}} \boldsymbol{e}_{u}^{(L)}, \tag{11}$$

where $\mathcal{N}_{u^*}^+(k)$ is a set of k-hop neighborhood¹ of user u^* connected by only positive edges, and w_{u,u^*} is a normalized alignment score between u and u^* . The normalized alignment score w_{u,u^*} is defined as

$$w_{u,u^*} = \frac{\exp(\gamma_{u,u^*}/\kappa)}{\sum_{\tilde{u}\in\mathcal{N}_{u^*}^+(k)}\exp(\gamma_{\tilde{u},u^*}/\kappa)},$$
$$\gamma_{u,u^*} = \sum_{(a\succ b)\in\mathcal{D}_{u^*}}\log\sigma(s_{u,a} - s_{u,b}),$$

where

where $s_{u,i}$ is an inner product between user and response embeddings, κ is a temperature parameter, and γ_{u,u^*} is an alignment score between user u and u^* . Intuitively, γ_{u,u^*} measures how well the *predicted preference* of user u aligns with the *annotated preference* provided by user u^* . If the preferences of both users align well, γ_{u,u^*} is large. Consequently, their embeddings become similar to each other. By collecting embeddings of well-aligned neighborhood users, we can obtain embeddings of user u^* without having further optimization.

5 EXPERIMENTS

In this section, we aim to show whether reward models can accurately learn user preferences in sparse annotation scenarios. Specifically, we examine situations where many users contribute only a few annotated pairs.

 $^{^{1}}k$ must be an even number to aggregate only the user embeddings.

5.1 EXPERIMENTAL SETTINGS

Datasets. We employ the UltraFeedback-P (UF-P) dataset (Poddar et al., 2024), which is explicitly designed to capture diverse user preferences from UltraFeedback (Cui et al., 2023). Unlike traditional reward modeling datasets that assume a single dominant preference, UF-P explicitly builds diverse preference groups through fine-grained scores across multiple preference attributes about response from UltraFeedback.

UF-P is created by grouping users based on distinct preference priorities, including helpfulness, honesty, instruction-following, and truthfulness. This dataset consists of two environments, each with a different number of groups. First, UF-P-2 consists of two user groups, each prioritizing either helpfulness or honesty. UF-P-4 expands to four groups, each concentrating on a different attribute. We provide a detailed explanation of the construction of the UF-P dataset from UltraFeedback in Appendix C.1.

While UF-P supports personalized reward modeling, it does not inherently reflect scenarios where a large number of users each provides only a handful of annotations. To reflect our target scenario, we generate a modified version of UF-P with 10,000 users evenly distributed across different preference groups and a survey set of 25,993 pairs.

Specifically, we construct four experimental environments based on UF-P-2 and UF-P-4:

- UF-P-2-ALL: In two preference groups, each user contributes exactly 8 annotations.
- UF-P-2-AVG: In two preference groups, each user contributes 8 annotations on average.
- UF-P-4-ALL: In four preference groups, each user contributes exactly 16 annotations.
- UF-P-4-AVG: In four preference groups, each user contributes 16 annotations on average.

For UF-P-2-AVG and UF-P-4-AVG, we randomly sample the number of annotations from a uniform distribution over $1 \sim 15$ and $1 \sim 31$, respectively.

Since UF-P-4 encompasses a broader range of preferences, users provide more annotations to capture this added complexity. These configurations enable us to rigorously evaluate how reward models perform under sparse user annotations, a critical challenge for large-scale personalized alignment in practical settings.

Notably, our experimental environments remain consistent with previous work (Poddar et al., 2024), but more closely mirror our target environments. Specifically, Poddar et al. (2024) infers user preferences from a small, predefined pool of unannotated pairs, so all users must be evaluated within that limited query set. In contrast, we consider a much broader range of unannotated pairs, allowing the model to capture preferences across diverse contexts and better adapt to real-world personalized alignment scenarios.

Baselines. We evaluate six baselines to benchmark. First, we use a uniform preference model (Uniform) trained on all annotations via BTL. Additionally, we consider four personalized reward models: I2E, $I2E_{proxy}$ (Li et al., 2024), VPL (Poddar et al., 2024), and PAL Chen et al. (2024a). Finally, we include an Oracle, which has access to user group information and all annotations in the survey set and trains a separate reward function in Eq. (1) for each preference group. The details of each model are provided in the Appendix B.

Training and evaluation details. For reward function training, we utilize two LLM backbones: gemma-2b-it and gemma-7b-it (Team et al., 2024). Our model uses one shared LoRA, eight LoRA experts, each with a rank of eight, and a two-layer MLP for the gating function. The other baselines, e.g., Uniform, I2E, VPL, PAL, and Oracle, use a LoRA rank of 64. Other training details, such as hyper-parameters and model architecture, are provided in Appendix C.2.

We report reward model accuracy on unseen test pairs that are not in the survey set. We define a correct prediction as assigning a higher score to the preferred response. We evaluate performance for both seen and unseen users. For seen user experiments, each user is assigned 10 test pairs, and accuracy is calculated over all seen users. We fix the number of unseen users at 100, evenly distributed across preference groups. To adapt the reward model for each unseen user, we provide 8

Table 1: Accuracy of reward models on unseen annotated pairs. The Seen user results report per-
formance for all users encountered during training in the upper block of the table. The Unseen user
results report performance for 100 new users, evenly distributed across preference groups. Unseen
users provide 8 annotations under UF-P-2-ALL/AVG and 16 annotations under UF-P-4-ALL/AVG.
Bold represents the best result, except with Oracle. All experiments run on three seeds.

		Gemma-2b-it			Gemma-7b-it				
		UF	-P-2	UF-P-4		UF-P-2		UF-P-4	
		ALL	AVG	ALL	AVG	ALL	AVG	ALL	AVG
	Oracle	$64.53_{\pm 0.14}$	$64.53_{\pm 0.14}$	$61.52_{\pm 0.13}$	$61.52_{\pm 0.13}$	$66.80_{\pm 0.17}$	$66.80_{\pm 0.17}$	$62.17_{\pm 0.09}$	$62.17_{\pm 0.09}$
Seen	Uniform I2E I2E _{proxy} VPL PAL CoPL	$\begin{array}{c} 61.82 {\pm} 0.16 \\ 61.48 {\pm} 0.18 \\ 61.43 {\pm} 0.56 \\ 61.11 {\pm} 0.16 \\ 59.95 {\pm} 0.04 \\ \textbf{63.81} {\pm} 0.16 \end{array}$	$\begin{array}{c} 61.82 {\scriptstyle \pm 0.16} \\ 61.49 {\scriptstyle \pm 0.70} \\ 61.33 {\scriptstyle \pm 0.61} \\ 61.86 {\scriptstyle \pm 0.84} \\ 61.53 {\scriptstyle \pm 0.22} \\ \textbf{63.45} {\scriptstyle \pm 0.38} \end{array}$	$\begin{array}{c} 56.15 {\scriptstyle \pm 0.22} \\ 57.21 {\scriptstyle \pm 0.37} \\ 56.78 {\scriptstyle \pm 0.14} \\ 56.04 {\scriptstyle \pm 1.71} \\ 56.95 {\scriptstyle \pm 0.13} \\ \textbf{62.57} {\scriptstyle \pm 0.38} \end{array}$	$\begin{array}{c} 56.15_{\pm 0.22} \\ 57.44_{\pm 0.37} \\ 57.14_{\pm 0.31} \\ 56.77_{\pm 0.38} \\ 57.37_{\pm 0.14} \\ \textbf{62.08}_{\pm 0.27} \end{array}$	$\begin{array}{c} 61.96 {\scriptstyle \pm 0.07} \\ 62.10 {\scriptstyle \pm 0.28} \\ 62.03 {\scriptstyle \pm 0.30} \\ 62.39 {\scriptstyle \pm 0.10} \\ 62.59 {\scriptstyle \pm 0.06} \\ \textbf{63.90} {\scriptstyle \pm 0.07} \end{array}$	$\begin{array}{c} 61.96 {\scriptstyle \pm 0.07} \\ 61.43 {\scriptstyle \pm 0.23} \\ 62.27 {\scriptstyle \pm 0.09} \\ 62.59 {\scriptstyle \pm 0.24} \\ 62.47 {\scriptstyle \pm 0.13} \\ \textbf{63.48} {\scriptstyle \pm 0.13} \end{array}$	$\begin{array}{c} 56.80_{\pm 0.12} \\ 57.90_{\pm 0.21} \\ 57.54_{\pm 0.16} \\ 58.87_{\pm 0.25} \\ 57.17_{\pm 0.22} \\ \textbf{62.90}_{\pm 0.05} \end{array}$	$\begin{array}{c} 56.80 {\pm} 0.12 \\ 58.50 {\pm} 0.09 \\ 58.12 {\pm} 0.14 \\ 57.55 {\pm} 1.00 \\ 56.27 {\pm} 0.13 \\ \textbf{61.93} {\pm} 0.02 \end{array}$
Unseen	Oracle Uniform I2E I2E _{proxy} VPL PAL CoPL	$\begin{array}{c} 64.66_{\pm 1.10} \\ 62.82_{\pm 0.59} \\ 61.67_{\pm 0.82} \\ 62.30_{\pm 0.54} \\ 60.83_{\pm 0.40} \\ 59.83_{\pm 0.69} \\ \textbf{63.92}_{\pm 0.54} \end{array}$	$\begin{array}{c} 64.66_{\pm 1.10} \\ 62.82_{\pm 0.59} \\ 59.52_{\pm 0.51} \\ 61.70_{\pm 0.63} \\ 62.62_{\pm 0.49} \\ 61.71_{\pm 0.31} \\ \textbf{63.26}_{\pm 0.51} \end{array}$	$\begin{array}{c} 61.33_{\pm 0.35} \\ 55.65_{\pm 0.61} \\ 56.42_{\pm 0.41} \\ 56.00_{\pm 1.15} \\ 54.03_{\pm 1.54} \\ 57.07_{\pm 0.22} \\ \textbf{61.62}_{\pm 0.10} \end{array}$	$\begin{array}{c} 61.33_{\pm 0.35}\\ 55.65_{\pm 0.61}\\ 56.75_{\pm 0.68}\\ 56.50_{\pm 0.34}\\ 56.13_{\pm 0.57}\\ 57.13_{\pm 0.33}\\ \textbf{61.97}_{\pm 0.35}\end{array}$	$\begin{array}{c} 67.43 \scriptstyle \pm 0.65 \\ 62.23 \scriptstyle \pm 0.06 \\ 62.62 \scriptstyle \pm 0.95 \\ 61.99 \scriptstyle \pm 0.33 \\ 62.69 \scriptstyle \pm 0.99 \\ 63.08 \scriptstyle \pm 0.73 \\ \textbf{64.08} \scriptstyle \pm 0.71 \end{array}$	$\begin{array}{c} 67.43_{\pm 0.65}\\ 62.23_{\pm 0.06}\\ 61.88_{\pm 0.21}\\ 62.84_{\pm 0.40}\\ 63.67_{\pm 0.12}\\ 62.52_{\pm 0.58}\\ \textbf{64.38}_{\pm 1.00}\end{array}$	$\begin{array}{c} 62.01_{\pm 0.04} \\ 57.02_{\pm 0.27} \\ 57.62_{\pm 0.92} \\ 57.69_{\pm 0.70} \\ 58.49_{\pm 1.22} \\ 57.15_{\pm 0.48} \\ \textbf{62.77}_{\pm 1.32} \end{array}$	$\begin{array}{c} 62.01_{\pm 0.04}\\ 57.02_{\pm 0.27}\\ 58.12_{\pm 0.98}\\ 57.73_{\pm 0.32}\\ 56.85_{\pm 0.84}\\ 56.44_{\pm 0.67}\\ \textbf{62.08}_{\pm 0.64}\end{array}$



Figure 1: T-SNE visualization of seen user embeddings in UF-P-4-AVG with gemma-2b-it. Points are colored by their preference group. Our method clusters users in the same group more effectively, whereas other baselines fail to cluster users by their preference groups in user embedding space.

annotations in UF-P-2-ALL/AVG and 16 annotations in UF-P-4-ALL/AVG, followed by evaluation on 50 test pairs per unseen user. CoPL uses 2-hop neighbors for unseen user adaptation.

5.2 Results

Table 1 presents accuracy for both seen and unseen users. CoPL consistently outperforms other baselines, except for Oracle, in both seen user and unseen user experiments. Notably, CoPL is comparable with Oracle in UF-P-4-ALL/AVG. In unseen user experiments, CoPL achieves accuracy comparable to the seen user setting, indicating the effectiveness of our unseen user adaptation.

Fig. 1 visualizes the embedding space of seen users in UF-P-4-AVG, which is the most challenging environment in these experiments, and demonstrates that GNN-based representation learning can capture preference similarity between users even when each user provides few annotations.

5.3 ANALYSIS

Analysis of performance in UF-P-2. In Table 1, all models appear capable of representing diverse preferences, surprisingly including the uniform models in UF-P-2-ALL/AVG. To investigate further, we divide the test pairs of UF-P-2 into *common* and *controversial* categories, where common pairs have identical annotations from both preference groups, and controversial pairs differ. Focusing on the seen user results in UF-P-2-ALL with gemma-2b-it from Table 1, we break down the accuracy in Table 2. The results indicate that baselines, except Oracle, struggle with controversial pairs, suggesting a tendency to capture only the common preference across all users. By contrast,

Table 2: Accuracy of reward models on UF-P-2-ALL with gemma-2b-it, broken down by pair type. *Common* refers to pairs for which the two preference groups provide the same preference label, *Controversial* refers to pairs labeled differently by the two groups, and *Total* encompasses all pairs. These categories reflect how diverse user preferences affect the performance of reward models.

	Oracle	Uniform	I2E	I2E _{proxy}	VPL	PAL	CoPL
Common Controversial	$\begin{array}{c} 71.86_{\pm 0.14} \\ 57.68_{\pm 0.27} \end{array}$	$\begin{array}{c} \textbf{74.52}_{\pm 0.45} \\ 49.86_{\pm 0.30} \end{array}$	$\begin{array}{c} 73.94_{\pm 0.21} \\ 49.61_{\pm 0.05} \end{array}$	$\begin{array}{c} 74.15_{\pm 1.53} \\ 49.86_{\pm 0.06} \end{array}$	$\begin{array}{c} 72.73_{\pm 1.00} \\ 50.26_{\pm 0.44} \end{array}$	$\begin{array}{c} 70.82_{\pm 0.17} \\ 49.79_{\pm 0.12} \end{array}$	$\begin{array}{c} 71.23_{\pm 1.63} \\ \textbf{56.89}_{\pm 1.56} \end{array}$
Total	$64.53_{\pm 0.14}$	$61.82_{\pm 0.16}$	$61.48{\pm 0.18}$	$61.59_{\pm 0.79}$	$61.11_{\pm 0.32}$	$59.95_{\pm 0.04}$	$\textbf{63.81}_{\pm 0.15}$



Figure 2: Accuracy of unseen user adaptation as the number of provided annotation sets increases, evaluated on UF-P-2/4-AVG with gemma-2b-it. 2-hop and 4-hop indicates 2-hop and 4-hop adaptation, respectively.



Figure 3: Expert allocation at layers 2 and 3 in UF-P-4-ALL with gemma-2b-it. Colors indicate preference groups. Users with similar preference groups are mapped to the same expert.

our method achieves comparable performance to Oracle on controversial pairs while preserving high accuracy on common pairs.

Effect of the number of annotations in unseen user adaptation. Fig. 2 shows accuracy as the number of provided annotations increases in UF-P-2-AVG and UF-P-4-AVG. We observe that additional annotations lead to more accurate preference predictions for unseen users in general. However, in practice, even eight annotations are sufficient, enabling accurate inference of each user's preference. We also compare two-hop and four-hop adaptations, but there is no significant difference.

Ablation study of CoPL. Table 3 presents an ablation study of CoPL, focusing on GNNderived user embeddings and the MoLE architecture. When GNN embeddings are removed, user representations become learnable parameters. Without MoLE, user embeddings are projected into the token space and passed as an additional token to the reward model. The results indicate that components of CoPL are effective. Specifically, GNN-based embeddings are a crucial component of CoPL, and the MoLE architecture further enhances accuracy. Notably, CoPL uses fewer activated parameters than w/o MoLE (n = 64). Table 3: Ablation study of CoPL in UF-P-2/4-ALL with gemma-2b-it. *w/o GNN embedding* replaces user embeddings from GNN with learnable user embeddings. *w/o MoLE* removes the MoLE architecture and projects user embeddings into the token space. The symbol n denotes the LoRA rank. All experiments run on three seeds.

	UF-P-2-ALL	UF-P-4-ALL
CoPL	$\textbf{63.81}_{\pm 0.16}$	$\textbf{62.57}_{\pm 0.38}$
w/o GNN embedding	62.09 ± 0.38	56.75 ± 0.30
w/o MoLE $(n = 64)$	62.69 ± 0.86	62.28 ± 0.33
w/o MoLE $(n = 16)$	$62.43_{\pm 0.69}$	$62.13_{\pm 0.12}$

UF	-P-2	UF-P-4		
ALL	AVG	ALL	AVG	
$84.84_{\pm 0.83}$	$84.32_{\pm 0.09}$	$90.01_{\pm 0.35}$	$87.74_{\pm 0.19}$	

Table 5: Test accuracy of the GNN. We evaluate the model using the same users from training but with annotation pairs not reflected in the graph. All experiments run on three seeds.

Fig. 3 depicts expert allocation, the user-conditioned gating mechanism partitions users differently at each layer. We observe users with the same preferences tend to be routed to the same expert.

Ablation study of unseen user adaptation.

We conduct an ablation study to evaluate the effectiveness of the unseen user adaptation strategy, comparing it to two baselines, Naive Avg and User Opt. Naive Avg assigns each unseen user embedding as the unweighted average of 2-hop seen user embeddings. User Opt replaces $e_u^{(L)}$ with a parameterized embedding learned by minimizing Equation (5) on the provided annotations. Table 4 reports results in UF-P-4-ALL/AVG with gemma-2b-it, showing that CoPL outperforms both alternatives while achieving better computational efficiency than the optimization-based User Opt.

Fig. 4 illustrates that naive averaging places unseen users away from identical preference

Table 4: Accuracy of unseen-user adaptation in UF-P-4-ALL/AVG with gemma-2b-it. *Naive* Avg. computes the unseen user's embedding as the unweighted average of 2-hop neighbors. *User* Opt. represents an optimization-based approach that learns a parameterized user embedding by maximizing the likelihood of the given annotations. All experiments run on three seeds.

	UF-P-4-ALL	UF-P-4-AVG
CoPL	$61.62_{\pm0.10}$	$\textbf{61.97}_{\pm 0.35}$
Naive Avg. User Opt.	$\begin{array}{c} 59.91_{\pm 0.59} \\ 59.24_{\pm 0.71} \end{array}$	$\begin{array}{c} 59.39_{\pm 0.50} \\ 59.45_{\pm 0.72} \end{array}$

group users, whereas our method clusters them closely with users who share the same preferences.

Training reward models with GNN. Table 5 reports GNN accuracy on seen users and responses for test pairs excluded from the training dataset. The results demonstrate that GNN can accurately predict labels for unannotated pairs with sparse annotations.

Table 6 examines the impact of training with GNN-based pseudo labels, allowing the model to leverage additional preference data. Although the pseudo-labeled pairs increase the dataset size, performance is slightly worse than using only user-provided annotations, suggesting that noise degrades model accuracy.

Table 6: Accuracy of reward model trained by using a pre-trained GNN in UF-P-2/4-ALL with gemma-2b-it. The *pseudo-label* trains a reward model on all seen user-response pairs, with annotations provided by GNN-predicted labels. The *user-specific* refers to a model trained with pseudo labels for each user. 10 users per group are sampled. All experiments run on three seeds.

	UF-P-2-ALL	UF-P-4-ALL
CoPL Pseudo label	$\begin{array}{c} 63.81_{\pm 0.16} \\ 62.77_{\pm 0.70} \end{array}$	$\begin{array}{c} 62.57_{\pm 0.38} \\ 62.26_{\pm 0.27} \end{array}$
Oracle User-specific	$\begin{array}{c} 64.53_{\pm 0.14} \\ 58.09_{\pm 1.73} \end{array}$	$ \begin{array}{r} 61.52_{\pm 0.13} \\ 55.30_{\pm 3.30} \end{array} $

To investigate the effect of noise further, a userspecific reward model is trained on pseudo la-

bels for a random sample of 10 users per group. The results are worse than the Oracle, indicating that noisy labels introduce training instability. This observation aligns with Wang et al. (2024), which notes that noisy preference labels can lead to training instability and performance degradation.

6 CONCLUSION

In this work, we introduced CoPL, a novel approach for personalizing LLMs through graph-based collaborative filtering and MoLE. Unlike existing methods that treat user preferences independently or require predefined clusters, our approach leverages multi-hop user-response relationships to improve preference estimation, even in sparse annotation settings. By integrating user embeddings into the reward modeling process with MoLE, CoPL effectively predicts an individual preference.

REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- 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.
- Daiwei Chen, Yi Chen, Aniket Rege, and Ramya Korlakai Vinayak. Pal: Pluralistic alignment framework for learning from heterogeneous preferences, 2024a. URL https://arxiv.org/abs/2406.08469.
- Lei Chen, Le Wu, Richang Hong, Kun Zhang, and Meng Wang. Revisiting graph based collaborative filtering: A linear residual graph convolutional network approach. In *Proceedings of the AAAI conference on artificial intelligence*, 2020.
- Lu Chen, Rui Zheng, Binghai Wang, Senjie Jin, Caishuang Huang, Junjie Ye, Zhihao Zhang, Yuhao Zhou, Zhiheng Xi, Tao Gui, et al. Improving discriminative capability of reward models in rlhf using contrastive learning. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 15270–15283, 2024b.
- Shaoxiang Chen, Zequn Jie, and Lin Ma. Llava-mole: Sparse mixture of lora experts for mitigating data conflicts in instruction finetuning mllms. *arXiv preprint arXiv:2401.16160*, 2024c.
- Zeren Chen, Ziqin Wang, Zhen Wang, Huayang Liu, Zhenfei Yin, Si Liu, Lu Sheng, Wanli Ouyang, and Jing Shao. Octavius: Mitigating task interference in mllms via lora-moe. In *ICLR*, 2023.
- Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and Maosong Sun. Ultrafeedback: Boosting language models with high-quality feedback, 2023.
- Josef Dai, Xuehai Pan, Ruiyang Sun, Jiaming Ji, Xinbo Xu, Mickel Liu, Yizhou Wang, and Yaodong Yang. Safe rlhf: Safe reinforcement learning from human feedback. *arXiv preprint arXiv:2310.12773*, 2023.
- Xiangnan He and Tat-Seng Chua. Neural factorization machines for sparse predictive analytics, 2017. URL https://arxiv.org/abs/1708.05027.
- Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yongdong Zhang, and Meng Wang. Lightgcn: Simplifying and powering graph convolution network for recommendation. In *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval*, pp. 639–648, 2020.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- 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. *arXiv preprint arXiv:2310.11564*, 2023.
- Jiacheng Li, Tong Zhao, Jin Li, Jim Chan, Christos Faloutsos, George Karypis, Soo-Min Pantel, and Julian McAuley. Coarse-to-fine sparse sequential recommendation. In *Proceedings of the 45th international ACM SIGIR conference on research and development in information retrieval*, pp. 2082–2086, 2022.
- Xinyu Li, Zachary C Lipton, and Liu Leqi. Personalized language modeling from personalized human feedback. *arXiv preprint arXiv:2402.05133*, 2024.

- Zihan Lin, Changxin Tian, Yupeng Hou, and Wayne Xin Zhao. Improving graph collaborative filtering with neighborhood-enriched contrastive learning. In *Proceedings of the ACM Web Conference 2022*, WWW '22, pp. 2320–2329, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450390965. doi: 10.1145/3485447.3512104. URL https://doi.org/10.1145/3485447.3512104.
- Qidong Liu, Xian Wu, Xiangyu Zhao, Yuanshao Zhu, Derong Xu, Feng Tian, and Yefeng Zheng. When moe meets llms: Parameter efficient fine-tuning for multi-task medical applications. In Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '24. Association for Computing Machinery, 2024. ISBN 9798400704314. URL https://doi.org/10.1145/3626772.3657722.
- I Loshchilov. Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101, 2017.
- Ismael Villegas Molina, Audria Montalvo, Benjamin Ochoa, Paul Denny, and Leo Porter. Leveraging llm tutoring systems for non-native english speakers in introductory cs courses. *arXiv preprint arXiv:2411.02725*, 2024.
- Minhyeon Oh, Seungjoon Lee, and Jungseul Ok. Active preference-based learning for multidimensional personalization, 2024. URL https://arxiv.org/abs/2411.00524.
- Sriyash Poddar, Yanming Wan, Hamish Ivison, Abhishek Gupta, and Natasha Jaques. Personalizing reinforcement learning from human feedback with variational preference learning. *arXiv preprint arXiv:2408.10075*, 2024.
- Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. Bpr: Bayesian personalized ranking from implicit feedback. *arXiv preprint arXiv:1205.2618*, 2012.
- Jingzhe Shi, Jialuo Li, Qinwei Ma, Zaiwen Yang, Huan Ma, and Lei Li. Chops: Chat with customer profile systems for customer service with llms. *arXiv preprint arXiv:2404.01343*, 2024.
- Anand Siththaranjan, Cassidy Laidlaw, and Dylan Hadfield-Menell. Distributional preference learning: Understanding and accounting for hidden context in rlhf. In *ICLR*, 2024. URL https://arxiv.org/abs/2312.08358.
- Taylor Sorensen, Jared Moore, Jillian Fisher, Mitchell Gordon, Niloofar Mireshghallah, Christopher Michael Rytting, Andre Ye, Liwei Jiang, Ximing Lu, Nouha Dziri, et al. A roadmap to pluralistic alignment. *arXiv preprint arXiv:2402.05070*, 2024.
- Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024.
- Saranya Venkatraman, Nafis Irtiza Tripto, and Dongwon Lee. Collabstory: Multi-Ilm collaborative story generation and authorship analysis. *arXiv preprint arXiv:2406.12665*, 2024.
- Binghai Wang, Rui Zheng, Lu Chen, Zhiheng Xi, Wei Shen, Yuhao Zhou, Dong Yan, Tao Gui, Qi Zhang, and Xuan-Jing Huang. Reward modeling requires automatic adjustment based on data quality. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 4041– 4064, 2024.
- Xiang Wang, Xiangnan He, Meng Wang, Fuli Feng, and Tat-Seng Chua. Neural graph collaborative filtering. In *Proceedings of the 42nd international ACM SIGIR conference on Research and development in Information Retrieval*, pp. 165–174, 2019.
- Kai Yang, Jian Tao, Jiafei Lyu, Chunjiang Ge, Jiaxin Chen, Weihan Shen, Xiaolong Zhu, and Xiu Li. Using human feedback to fine-tune diffusion models without any reward model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 8941–8951, June 2024a.
- Rui Yang, Xiaoman Pan, Feng Luo, Shuang Qiu, Han Zhong, Dong Yu, and Jianshu Chen. Rewardsin-context: Multi-objective alignment of foundation models with dynamic preference adjustment. *arXiv preprint arXiv:2402.10207*, 2024b.

APPENDIX

A MESSAGE PASSING FOR RESPONSE EMBEDDINGS

Given user and response embeddings at layer $\ell,$ a message from neighborhood users to the response as

$$\begin{split} \boldsymbol{m}_{r}^{+} &= \sum_{u \in \mathcal{N}_{r}^{+}} \alpha_{u,r} \Big(\hat{W}_{1}^{(\ell)} \boldsymbol{e}_{u}^{(\ell)} + \hat{W}_{2}^{(\ell)} (\boldsymbol{e}_{u}^{(\ell)} \odot \boldsymbol{e}_{r}^{(\ell)}) \Big), \\ \boldsymbol{m}_{r}^{-} &= \sum_{u \in \mathcal{N}_{r}^{-}} \beta_{u,r} \Big(\hat{W}_{3}^{(\ell)} \boldsymbol{e}_{u}^{(\ell)} + \hat{W}_{4}^{(\ell)} (\boldsymbol{e}_{u}^{(\ell)} \odot \boldsymbol{e}_{r}^{(\ell)}) \Big), \\ \boldsymbol{m}_{r}^{(\ell)} &= \hat{W}_{\text{self}}^{(\ell)} \boldsymbol{e}_{r}^{(\ell)} + \boldsymbol{m}_{r}^{+} + \boldsymbol{m}_{r}^{-}, \end{split}$$
(12)

where $\hat{W}_1^{(\ell)}, \hat{W}_2^{(\ell)}, \hat{W}_3^{(\ell)}, \hat{W}_4^{(\ell)}, \hat{W}_{self}^{(\ell)} \in \mathbb{R}^{d \times d}$ are parameter matrices, \odot is element-wise multiplication, and $\alpha_{u,r}$ and $\beta_{u,r}$ are normalization factors, set to $\frac{1}{\sqrt{|\mathcal{N}_u^+| \cdot |\mathcal{N}_r^+|}}$ and $\frac{1}{\sqrt{|\mathcal{N}_u^-| \cdot |\mathcal{N}_r^-|}}$, respectively.

Then, the response embedding is updated with the aggregated message $m_r^{(\ell)}$:

$$\boldsymbol{e}_{r}^{(\ell+1)} = \psi(\boldsymbol{m}_{r}^{(\ell)}), \tag{13}$$

where $\psi(\cdot)$ is a non-linear activation.

B METHOD BASELINES

Uniform. The uniform model is a standard approach for pairwise preference comparisons. We train the uniform model with all annotation pairs, which will capture the common preference.

Oracle. For an oracle model of our setting, we train the model with the true group membership of all users. A separate uniform model is trained for each group by aggregating annotations from the users in that group.

I2E (Li et al., 2024). I2E is a framework that uses DPO to personalize LLM. However, it can be easily extended to reward modeling. I2E trains a model that maps the user index into a learnable embedding. It appends each user embedding as an additional input token to the LLM, providing user-specific signals for reward prediction.

 $I2E_{proxy}$ (Li et al., 2024). A variant of I2E that introduces N proxy embeddings. A weighted combination of these proxies forms the final user embedding, which is passed to the LLM for reward prediction. In our experiments, we use N = 10.

VPL (**Poddar et al., 2024**). Variational Preference Learning (VPL) encodes user-specific annotations into user embeddings. The user embeddings are then combined with sentence representations via an MLP to predict reward scores. To capture the user preferences effectively, VPL uses a variational approach that maps the user annotations into a prior distribution.

PAL (Chen et al., 2024a). Pluralistic Alignment (PAL) applies an ideal-point model, where the distance between the user and the response determines the reward. The ideal point of the user is represented by N proxies, set to N = 10 in this work. Among variants of PAL, we use PAL-A with logistic loss.

C EXPERIMENTAL DETAILS

In this section, we provide a detailed explanation of dataset construction and hyper-parameters.

C.1 ULTRAFEEDBACK-P

Poddar et al. (2024) proposes the Ultrafeedback-P (UF-P) benchmark for personalized reward modeling, based on the Ultrafeedback (UF) dataset Cui et al. (2023), which provides response pairs rated on four attributes: helpfulness, honesty, instruction following, and truthfulness. In UF-P, each attribute corresponds to a distinct preference. For instance, a user belonging to the helpfulness group annotates pairs, solely considering the helpfulness score.

UF-P-2. This version employs only two attributes and removes pairs that both user groups label identically, focusing on controversial cases where preferences differ.

UF-P-4. All four attributes are retained as preference dimensions, which allows for partial agreement among groups and hence increases complexity. Although Poddar et al. (2024) also excludes pairs fully agreed upon by all users, the remaining set is larger and exhibits more variety than UF-P-2.

In Poddar et al. (2024), each user is given a small context sample from a limited set of unannotated pairs to infer the user's preference. In contrast, we leverage every available pair in the dataset to infer each user's preferences. For our dataset construction, we use UF-P-4 dataset.

C.2 HYPER-PARAMETERS

We describe the training details of GNN, a reward model, and unseen user adaptation, such as model architecture and hyper-parameters.

GNN. The model consists of four message-passing layers, each with user and response embeddings of dimension 512. We use Leaky ReLU as non-linear activation function to update user and response embeddings. Training proceeds for 300 epochs using the AdamW optimizer Loshchilov (2017) with a learning rate of 1×10^{-4} and a cosine scheduler with warmup ratio 0.1. The batch size is 1024, and all experiments are conducted on an RTX 4090 GPU.

Reward models. CoPL comprises an LLM backbone and a MoLE adapter. We use gemma-2b-it or gemma-7b-it as the LLM backbone. MoLE includes one shared expert and eight LoRA experts with a rank of eight. A two-layer MLP with a hidden dimension of 256 and ReLU activation serves as the gating mechanism, with a temperature set to 1.

We train the reward models using the AdamW optimizer with a learning rate of 5×10^{-5} and a cosine scheduler with warmup ratio 0.03. Four GPUs, such as RTX6000ADA, L40S, and A100-PCIE-40GB, are employed with a batch size of 32 per GPU for gemma-2b-it and 16 per GPU for gemma-7b-it.

Baseline models use LoRA with rank 64. They also trained with an AdamW optimizer and a cosine scheduler with a warmup ratio 0.03. We search the learning rate from $[1 \times 10^{-4}, 5 \times 10^{-5}, 1 \times 10^{-5}, 5 \times 10^{-6}]$.

User adaptation. We use two-hop seen user and 0.07 as temperature for unseen user adaptation of CoPL. For I2E, each learnable user representation is mapped into each user. For I2E_{proxy} and PAL, user representations are determined by N = 10 proxies. Adapting to an unseen user requires parameter optimization for unseen users, typically through several gradient steps. To optimize the parameters for unseen users, 50 gradient steps are applied during adaptation.



Figure 4: T-SNE visualization of seen and unseen user embeddings in UF-P-4-AVG. *Naive Avg.* computes unseen user embeddings as the unweighted mean of 2-hop neighbor embeddings. *User Opt.* represents an optimization-based approach that learns a parameterized user embedding by maximizing the likelihood of the given annotations. Colors indicate preference groups, and points with black edges represent unseen users. Unseen users adapted by our method align with their respective preference groups.