SWAP-GUIDED PREFERENCE LEARNING FOR PERSONALIZED REINFORCEMENT LEARNING FROM HUMAN FEEDBACK

Anonymous authorsPaper under double-blind review

ABSTRACT

Reinforcement Learning from Human Feedback (RLHF) is a widely used approach to align large-scale AI systems with human values. However, RLHF typically assumes a single, universal reward, which overlooks diverse preferences and limits personalization. Variational Preference Learning (VPL) seeks to address this by introducing user-specific latent variables. Despite its promise, we found that VPL suffers from posterior collapse. While this phenomenon is well known in VAEs, it has not previously been identified in preference learning frameworks. Under sparse preference data and with overly expressive decoders, VPL may cause latent variables to be ignored, reverting to a single-reward model. To overcome this limitation, we propose Swap-guided Preference Learning (SPL). The key idea is to construct fictitious swap annotators and use the mirroring property of their preferences to guide the encoder. SPL introduces three components: (1) swap-guided base regularization, (2) Preferential Inverse Autoregressive Flow (P-IAF), and (3) adaptive latent conditioning. Experiments show that SPL mitigates collapse, enriches user-specific latents, and improves preference prediction. Our code and data are available at https://anonymous.4open.science/r/SPL-0111

1 Introduction

Reinforcement learning from human feedback (RLHF) has emerged as a prominent method for aligning large-scale AI systems with human values in various fields, particularly natural language processing (Ouyang et al., 2022). In RLHF, a reward model is first trained on human comparison data, and then a policy is optimized with reinforcement learning. This approach aligns model behavior more closely with human evaluations, improving performance, accuracy, and fairness across diverse domains (Leike et al., 2018; Ji et al., 2023).

However, most existing RLHF approaches (Christiano et al., 2017; Ouyang et al., 2022) are based on the single-reward assumption that all human preferences can be represented by a universal reward function. This assumption is originated from the Bradley–Terry–Luce (BTL) model (Bradley & Terry, 1952), which is commonly used to model pairwise comparisons and treats preferences as if they were generated from a shared scoring function. While mathematically convenient, this single-reward assumption is problematic in practice. Human preferences are not homogeneous but plural and often diverge across individuals or groups. Recent studies have shown that collapsing diverse perspectives into a single reward function introduces systematic bias in favor of majority preferences, overlooking groups and reducing fairness (Prabhakaran et al., 2021; Feffer et al., 2023; Casper et al., 2023). Consequently, models trained under this assumption may disadvantage underrepresented populations, even when their preferences are valid and important.

To address this issue, researchers have begun exploring what we refer to as personalized alignment (pluralistic alignment) (Sorensen et al., 2024). Instead of forcing all preferences into a single universal reward function, personalized alignment seeks to align different reward functions with different individuals according to their preferences, thereby capturing the heterogeneity of human values. One leading approach is Variational Preference Learning (VPL) (Poddar et al., 2024), which encodes user-specific latent variables from preference data and decodes them into corresponding rewards.

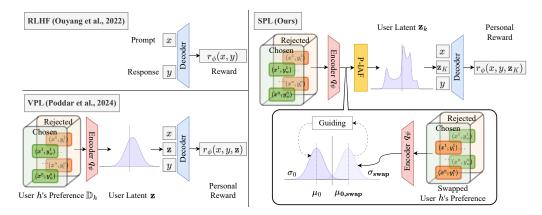


Figure 1: **Overview of SPL.** We propose Swap-guided $Preference\ Learning\ (SPL)$, a new framework for personalized alignment. RLHF (Ouyang et al., 2022) cannot adequately reflect user diversity. To overcome this limitation, VPL (Poddar et al., 2024) encodes text data consisting of a prompt x and response y into a single latent embedding. However, this encoding process is highly prone to collapse. In contrast, SPL leverages the structural properties of preference data through guiding mechanisms and a Preferential Inverse Autoregressive flow, allowing the latent space to capture user-specific characteristics.

This framework allows AI systems to flexibly adapt to diverse users without relying on predefined groupings or rigid categorization.

Despite its promise, we found in our experiments that VPL suffers from practical failure mode: posterior collapse. This phenomenon is sometimes observed in VAEs (Bowman et al., 2016; Chen et al., 2016; He et al., 2019; Lucas et al., 2019; Wang et al., 2021) but has not previously been identified in preference learning frameworks. When combined with a strong reward decoder, this posterior collapse can cause the encoder's latent variable to become uninformative and effectively ignored. In such cases, the latent variable fails to capture user-specific information, and the decoder explains preferences without relying on it. Training then reduces to an implicit single reward model, ignoring minority preferences and undermining the goal of personalized alignment.

To overcome this, we introduce *Swap-guided Preference Learning* (SPL), an expressive variational framework for personalized alignment that explicitly leverages the structural properties of preference pair data. To the best our knowledge, we are the first to report and address posterior collapse in preference learning. Our approach improves user-latent encoding and reward decoding through three key innovations: (i) **Swap-guided Base Regularization**, which encouraging latent space shows *mirrored* characteristics under preference swapping; (ii) **Preferential-Inverse Autoregressive Flow**, which disentangles swap-reversal and swap-invariant signals, conditioning a inverse autoregressive flow on them to yield improved latent representations without collapse; and (iii) **Adaptive Latent Conditioning**, which dynamically adjusts the contribution of the latent variable to reward prediction. Together, these mechanisms consistently reduce posterior collapse and enable more faithful and pluralistic preference modeling.

2 Preliminary fundamentals

Reinforcement Learning from Human Feedback For post-training of Large Language Models (LLM), RLHF relies on a dataset of N human preference pairs, $\mathbb{D} = \{(x^i, y_w^i, y_l^i)\}_{i=1}^N$, where x is a prompt and (y_w, y_l) denote the chosen(winning) and rejected(losing) responses, respectively. RLHF assumes an single universal reward function $r_{\phi}(x, y)$, optimized by maximizing the log-likelihood of observed preferences:

$$\mathbb{E}_{(x,y_w,y_l)\sim\mathbb{D}}\Big[\log p_\phi(y_w \succ y_l \mid x)\Big]. \tag{1}$$

The preference probability $p_{\phi}(y_w \succ y_l \mid x)$ is typically modeled via the Bradley–Terry–Luce (BTL) model (Bradley & Terry, 1952):

$$p_{\phi}(y_w \succ y_l \mid x) = \frac{\exp(r_{\phi}(x, y_w))}{\exp(r_{\phi}(x, y_w)) + \exp(r_{\phi}(x, y_l))} = \sigma(r_{\phi}(x, y_w) - r_{\phi}(x, y_l)), \tag{2}$$

where σ denotes the logistic function. Thus, the reward function r_{ϕ} is trained to explain human-preferred outcomes, and the learned reward model is subsequently used to optimize a policy aligned with human judgments.

Variational Approach for Personalized Alignment A central direction in personalized alignment is to condition reward models and policies on user-specific information (Oh et al., 2024; Poddar et al., 2024; Bose et al., 2025; Shenfeld et al., 2025; Gong et al., 2025). Among these approaches, Variational Preference Learning (VPL) (Poddar et al., 2024) is particularly influential. Inspired by variational autoencoders (VAEs) (Kingma et al., 2013), VPL introduces a user-specific latent variable $z \in \mathbb{R}^d$ inferred from each user h's preference dataset $\mathbb{D}_h = \{(x^i, y_w^i, y_l^i)\}_{i=1}^n \subset \mathbb{D}$. The encoder produces an approximate posterior $q_{\psi}(z \mid \mathbb{D}_h)$, while the decoder predicts rewards for prompt-response pairs (x, y) conditioned on z, denoted as $r_{\phi}(x, y, z)$.

Formally, VPL extends the objective in Eq.(1) by adding a variational regularization term. This yields an evidence lower bound (ELBO):

$$\mathbb{E}_{h \sim \mathbb{H}} \left[\mathbb{E}_{\substack{\boldsymbol{z} \sim q_{\psi}(\boldsymbol{z} \mid \mathbb{D}_{h}) \\ (\boldsymbol{x}, y_{w}, y_{l}) \sim \mathbb{D}_{h}}} [\log p_{\phi}(y_{w} \succ y_{l} \mid \boldsymbol{x}, \boldsymbol{z})] - \beta D_{\mathrm{KL}} [q_{\psi}(\boldsymbol{z} \mid \mathbb{D}_{h}) || p(\boldsymbol{z})] \right], \tag{3}$$

where β is KL divergence weight and the $\log p(z)$ represents the prior distribution's log-density, selected as $\mathcal{N}(\mathbf{0},\mathbf{I})$. This objective maximizes the conditional log-likelihood of preferences while regularizing the user-specific posterior toward the prior, thereby preventing overfitting and encouraging generalizable latent structure. By leveraging z, VPL provides flexibility in modeling personalized traits and has shown strong empirical performance in capturing diverse preferences. However, recent work (Nam et al., 2025) indicates that compressing rich textual preference data into a single latent embedding z remains highly challenging.

Inverse Autoregressive Flow Normalizing flows (Rezende & Mohamed, 2015) is a framework for constructing flexible posterior distributions by applying a sequence of invertible transformations. Among them, Inverse Autoregressive Flow (IAF) (Kingma et al., 2016) is specifically designed to enrich the expressivity of variational posteriors while preserving computational tractability. The procedure begins with a base latent variable $z_0 \in \mathbb{R}^d$ and context vector $c \in \mathbb{R}^{d_c}$ drawn from encoder (i.e., $q_{\psi}(z_0 \mid x) = \mathcal{N}(\mu, \sigma^2)$) with additional output c), followed by a series of parameterized, invertible transformations f_k . After K step transformations, the final variable z_K acquires a more complex distribution:

$$z_0 \sim q_{\psi}(z_0 \mid x), \quad z_k = f_k(z_{k-1}, c), \quad k = 1, \dots, K$$

When each f_k admits a tractable Jacobian determinant, the density of z_K can be computed efficiently via the change-of-variables formula:

$$\log q_{\psi}(\boldsymbol{z}_{K} \mid \boldsymbol{x}) = \log q_{\psi}(\boldsymbol{z}_{0} \mid \boldsymbol{x}) - \sum_{k=1}^{K} \log \det \left| \frac{\partial \boldsymbol{z}_{k}}{\partial \boldsymbol{z}_{k-1}} \right|. \tag{4}$$

In practice, IAF employs autoregressive neural networks to parameterize shift and scale functions:

$$\boldsymbol{z}_k = \mu_k(\boldsymbol{z}_{k-1}, \boldsymbol{c}) + \sigma_k(\boldsymbol{z}_{k-1}, \boldsymbol{c}) \odot \boldsymbol{z}_{k-1}, \tag{5}$$

where μ_k and σ_k are autoregressively conditioned on the preceding dimensions of z_{k-1} . This autoregressive structure ensures a lower-triangular Jacobian, making the determinant easy to compute:

$$\log \det \left| \frac{\partial z_k}{\partial z_{k-1}} \right| = \sum_{j=1}^d \log \left| \sigma_k^j \right|, \tag{6}$$

with σ_k^j denoting the j-th element of the scale function.

As a result, IAF enables parallelizable sampling and yields a substantially richer posterior $q_{\psi}(z_K \mid x)$ that captures inter-dimensional dependencies and non-Gaussian structures (e.g., skewness, heavy tails) beyond the capacity of the base posterior (Kingma et al., 2016; Papamakarios et al., 2021).

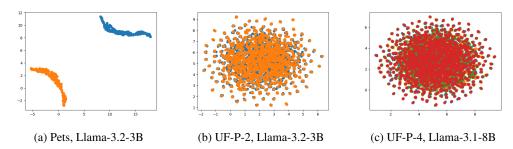


Figure 2: **Posterior collapse in Variational Preference Learning.** We visualize latent embeddings z from the VPL encoder using 2D UMAP (McInnes et al., 2018). Each point denotes a user, colored by their preference type. (a) User preference types are distinctly separated, indicating non-collapse. (b), (c) Latent collapse occurs, making preference types indistinguishable.

3 MOTIVATION

In this section, we explain the posterior collapse that we observed in preference learning and identify some guidance by comparing collapse and non-collapse cases. Fig. 2 illustrates this phenomenon that we observed. In Fig. 2a, two user types (in different colors) are clearly separated in the latent space for a simple dataset *Pets* with the same prompt. In larger, complex datasets *UF-P*, users merge into a single cluster, losing separation as shown in Fig. 2b and 2c. This collapse appears to stem from two factors: (1) noisy and ambiguous human feedback, together with the difficulty of compressing diverse, complex textual preferences in the encoder, often leads to unstable latent learning, which in turn causes the reward decoder to ignore the *z* pathway; and (2) the reward decoder already receives sufficient information from the complete prompt—response pair, allowing it to maximize the likelihood in Eq.(3) without relying on *z*. This leads to the latent variable failing to capture user-specific information and becoming uninformative. Further evidence of posterior collapse is presented in Appendix A.

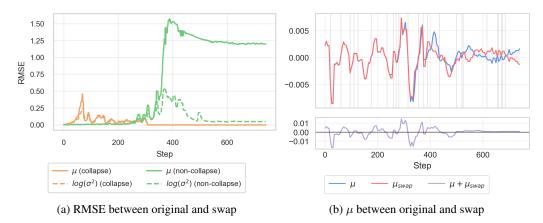


Figure 3: **Differences in posterior distribution between original and swapped inputs.** We test how the encoder's posterior responds when each preference pair is inverted to simulate a user with opposite choices, using the simple dataset *Pets*. (a) Average RMSE between original and swapped inputs across posterior mean μ and log-variance ℓ . Collapse appears in Llama-3.1-8B (orange), where both parameters remain unchanged, whereas Llama-3.2-3B (green) shows distinct behavior. (b) Plot μ vs. μ_{swap} for Llama-3.2-3B; $\mu + \mu_{\text{swap}}$ is in the lower panel. Initially, the curves are similar, but their difference grows and stabilizes as learning continues, resulting in a sign-reversal.

To address the posterior collapse in VPL, we examine the information captured in the posterior distribution when user preferences are successfully encoded, and use this insight to guide the design of an effective user-latent space. To this end, we conduct a simple swap experiment. For a user h with dataset \mathbb{D}_h , suppose the encoder outputs $q_{\psi}(z \mid \mathbb{D}_h) = \mathcal{N}(\mu, \sigma^2)$, where $\mu, \sigma^2 \in \mathbb{R}^d$. We

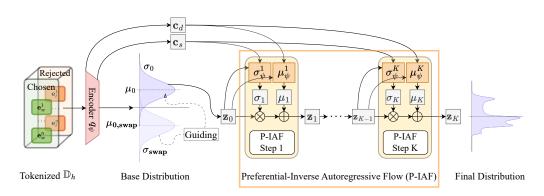


Figure 4: Preference encoding process of SPL

then construct a fictitious user $h_{\rm swap}$ with the opposite preference of h by swapping the chosen and rejected responses in every pair, as shown in the right part of Fig. 1. Feeding these swapped pairs into the encoder yields $q_{\psi}(z \mid \mathbb{D}_{h_{\rm swap}}) = \mathcal{N}(\mu_{\rm swap}, \sigma_{\rm swap}^2)$. Fig. 3a visualizes the RMSE between μ and $\mu_{\rm swap}$, and between $\ell = \log \sigma^2$ and $\ell_{\rm swap} = \log \sigma_{\rm swap}^2$, over the course of training for both collapse and non-collapse cases. In the collapse case, the RMSE converges to zero for both μ and ℓ , i.e., $\mu \approx \mu_{\rm swap}$ and $\ell \approx \ell_{\rm swap}$, indicating that the latent variable carries no user-specific signal and is effectively ignored by the decoder. In the non-collapse case, however, the RMSE of μ converges to a non-zero value, implying a clear separation between the original user h and the fictitious user $h_{\rm swap}$. In particular, μ and $\mu_{\rm swap}$ exhibit a sign-reversal, $\mu \approx -\mu_{\rm swap}$, as shown in Fig. 3b, while the log-variance remains invariant to swaps, $\ell \approx \ell_{\rm swap}$, i.e., the posterior distribution exhibits a "mirrored" distribution when swapped. This structural division implies that μ captures swap-reversal information, whereas ℓ captures swap-invariant information. Such disentanglement makes the latent variable essential for the decoder. In the next section, we use this insight to develop our new preference learning framework.

4 METHOD

We propose *Swap-guided Preference Learning* (SPL), a new framework for preference learning that regularizes the encoder with guidance from preference swapping. This approach consistently reduces posterior collapse while ensuring that user-specific information is faithfully encoded in the latent variable *z*. To achieve this, we introduce three components: (i) Swap-guided Base Regularization, (ii) Preferential Inverse Autoregressive Flow (P-IAF), and (iii) Adaptive Latent Conditioning.

4.1 ENCODING USER PREFERENCES INTO A LATENT

To encourage our SPL to encode user preferences into a latent z, we introduce two strategies in this section. The first is to enforce the output from the encoder to satisfy the mirroring of preference swaps, thereby mitigating posterior collapse. We call the encoder's Gaussian output is termed the base distribution and denoted as z_0 . The second strategy is to transform z_0 using an Inverse Autoregressive Flow (IAF), warping the Gaussian z_0 into a richer distribution z_K . In this second strategy, we also control the flow from the base z_0 to the transformed distribution z_K with guidance from the mirroring of preference swaps. The two strategies are illustrated in Fig. 4. We now explain them one by one.

Swap-guided Base Regularization Based on the mirroring of preference swaps in section 3, the encoder is trained to learn user preferences by generating mirrored distributions for annotators h and h_{swap} . Specifically, given an annotator h and its fictitious opposite annotator h_{swap} , with encoder outputs $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\ell})$ and $\mathcal{N}(\boldsymbol{\mu}_{\text{swap}}, \boldsymbol{\ell}_{\text{swap}})$, respectively, we train the encoder so that the two means $\boldsymbol{\mu}$ and $\boldsymbol{\mu}_{\text{swap}}$ exhibit a sign-reversal, while the two log-variances $\boldsymbol{\ell}$ and $\boldsymbol{\ell}_{\text{swap}}$ remain invariant. This is

achieved by applying the guidance loss \mathcal{L}_{guide} defined by

$$\cos(\boldsymbol{\mu}, \boldsymbol{\mu}_{\text{swap}}) = \frac{\boldsymbol{\mu}^{\top} \boldsymbol{\mu}_{\text{swap}}}{(\|\boldsymbol{\mu}\| + \varepsilon)(\|\boldsymbol{\mu}_{\text{swap}}\| + \varepsilon)}, \quad \cos(\boldsymbol{\ell}, \boldsymbol{\ell}_{\text{swap}}) = \frac{\boldsymbol{\ell}^{\top} \boldsymbol{\ell}_{\text{swap}}}{(\|\boldsymbol{\ell}\| + \varepsilon)(\|\boldsymbol{\ell}_{\text{swap}}\| + \varepsilon)},$$

and define the encoder q_{ψ} training guidance loss as:

$$\mathcal{L}_{\text{guide}} = \mathbb{E}_{h \sim \mathbb{H}} \left[\frac{1}{2} \left(1 + \cos(\boldsymbol{\mu}^h, \boldsymbol{\mu}_{\text{swap}}^h) \right) + \eta \, \frac{1}{2} \left(1 - \cos(\boldsymbol{\ell}^h, \boldsymbol{\ell}_{\text{swap}}^h) \right) \right]. \tag{7}$$

 η balances mean and variance; $\varepsilon > 0$ ensures stability.

Preferential Inverse Autoregressive Flow The next step is to apply IAF to warp the Gaussian z_0 into a multi-modal distribution z_K . Unlike the base regularization, we cannot enforce the mirroring property of preference swaps in this transformation, because z_K is no longer Gaussian and cannot be characterized in terms of mean and variance. In other words, the flow from z_0 to z_K under a standard IAF cannot be directly controlled to satisfy the mirroring property of preference swaps. To address this limitation, we propose Preferential Inverse Autoregressive Flow (P-IAF), which decomposes the context vector c into swap-reversal and swap-invariant components. Intuitively, the swap-reversal context c_d captures the directional preference signals that reflect the mirroring of swaps, while the swap-invariant context c_s captures the background information. Our P-IAF is defined by

$$z_k = f_{\psi}^k(z_{k-1}, c_d, c_s) = \mu_k(z_{k-1}, c_d) + \sigma_k(z_{k-1}, c_s) \odot z_{k-1},$$
 (8)

where k = 1, ..., K. We form c_d and c_s by a swap-reversal and swap-invariant decomposition of the encoder's additional output c (from \mathbb{D}_h) and c_{swap} (from swapped counterpart $\mathbb{D}_{h_{\text{swap}}}$) as follows:

$$oldsymbol{c}_d riangleq rac{1}{2}(oldsymbol{c} - oldsymbol{c}_{ ext{swap}}), \quad oldsymbol{c}_s riangleq rac{1}{2}(oldsymbol{c} + oldsymbol{c}_{ ext{swap}}),$$

which guarantees $c = c_d + c_s$, $c_{\text{swap}} = -c_d + c_s$. By feeding c_d only to the shift function μ_k and c_s only to the scale function σ_k , P-IAF reduces cross-context coupling between swap-reversal and swap-invariant signals, thereby preserving pair-derived user preference more effectively while retaining IAF's expressivity from the autoregressive composition. See Appendix B for details and proof.

Substituting Eq.(8) into Eq.(4) yields the overall log posterior after K flow steps

$$\log q_{\psi}(\boldsymbol{z}_K \mid \mathbb{D}_h) = \log q_{\psi}(\boldsymbol{z}_0 \mid \mathbb{D}_h) - \sum_{k=1}^K \sum_{j=1}^d \log |\sigma_k^j|, \tag{9}$$

and the KL divergence of Eq.(3) is given by¹:

$$D_{\mathrm{KL}} = \mathbb{E}_{h, \mathbb{W}} \left[\log q_{\psi}(\boldsymbol{z}_K \mid \mathbb{D}_h) - \log p(\boldsymbol{z}_K) \right], \tag{10}$$

where $\log q_{\psi}(\boldsymbol{z}_K \mid \mathbb{D}_h)$ is given in Eq.(9).

4.2 DECODING PERSONALIZED REWARDS FROM LATENTS

The decoder scores a prompt–response (x,y) conditioned on the user-latent z_K , yielding $r_{\phi}(x,y,z_K)$, and is trained to satisfy $r_{\phi}(x,y_w,z_K) > r_{\phi}(x,y_l,z_K)$. Extending Eq.(2) about z_K , the decoder training objective over users $h \sim \mathbb{H}$:

$$\mathbb{E}_{h \sim \mathbb{H}} \left[\mathbb{E}_{\substack{z_K \sim q_{\psi}(z_K \mid \mathbb{D}_h) \\ (x, y_w, y_l) \sim \mathbb{D}_h}} \left[\log p_{\phi}(y_w \succ y_l \mid x, z_K) \right] \right]$$
(11)

where $p_{\phi}(y_w \succ y_l \mid x, \mathbf{z}_K) = \sigma(r_{\phi}(x, y_w, \mathbf{z}_K) - r_{\phi}(x, y_l, \mathbf{z}_K))$ which means preference probability conditioned on z_K .

¹For notational simplicity, we denote all learnable parameters by ψ . In practice, ψ includes both (i) encoder parameters ψ_{enc} and (ii) flow parameters $\psi_{\text{flow}} = \{\psi_{\mu_k}, \psi_{\sigma_k}\}_{k=1}^K$ corresponding to the shift and scale function in each flow transformation step k.

Adaptive Latent Conditioning Inspired by feature modulation (Perez et al., 2018), we design a per-user modulation decoder that adapts prompt-response embeddings based on the user-latent embedding z_K , allowing dynamic influence adjustment when predicting input rewards. For example, when the latent embedding provides strong signals of user preference, its contribution to reward prediction is amplified, whereas when the preference signal is uncertain, the contribution is attenuated. Detailed modeling of this adaptive conditioning mechanism is provided in Appendix C.

.

4.3 Objective function of SPL

Maximize the log-likelihood term from Eq.(11) while minimizing the KL divergence term in Eq.(10), the ELBO of SPL is defined across the entire user \mathbb{H} as:

$$ELBO = \underset{h \sim \mathbb{H}}{\mathbb{E}} \left[\underset{\substack{z_K \sim q_{\psi}(z_K \mid \mathbb{D}_h) \\ (x, y_w, y_l) \sim \mathbb{D}_h}}{\mathbb{E}} [\log p_{\phi}(y_w \succ y_l \mid x, z_K)] - \beta(\log q_{\psi}(z_K \mid \mathbb{D}_h) - \log p(z_K)) \right]$$
(12)

We regularize the base posterior $q_{\psi}(z_0 \mid \mathbb{D}_h)$ using the guidance loss in Eq.(7). The final objective minimizes:

$$\mathcal{L}(\phi, \psi) = -\text{ELBO} + \lambda \mathcal{L}_{\text{guide}} \tag{13}$$

where λ controls the strength of the guidance loss term. Consequently, the reward model explicitly conditions on the user-latent z_K , yielding a personalized reward $r_{\phi}(\cdot, \cdot, z_K)$; optimizing the policy under this reward personalizes behavior and thus achieves personalized alignment.

5 EXPERIMENTS

In this section, we evaluate the performance of SPL. First, we examine whether SPL can construct a meaningful latent space without posterior collapse. Second, we evaluate whether SPL effectively improves preference-prediction accuracy. SPL remains stable across different KL divergence weights β , unlike the earlier approach (Poddar et al., 2024) in our experiments. Moreover, SPL consistently outperforms baselines in preference-prediction accuracy. Before presenting these results, we describe our experimental setup.

Baselines We compare our method against the following baselines:

- BTL (Ouyang et al., 2022): The standard RLHF based on Bradley-Terry-Luce model.
- **DPL** (Siththaranjan et al., 2023): Distributional Preference Learning, which captures implicit context across the entire preference dataset and models the reward as a distribution but doesn't consider individual user preferences.
- **VPL** (Poddar et al., 2024): Variational Preference Learning, which employs the user-latent embedding with a simple Gaussian posterior distribution, without swap-guided encoding and latent conditioning.
- VPL-IAF: An extension of VPL with a basic IAF posterior, used to examine the effect of a simple normalizing flow within a variational framework.
- **SPL-IAF**: Identical to SPL but replacing P-IAF with a basic IAF, serving as an ablation to evaluate the contribution of P-IAF.
- **SPL** (Ours): Our proposed method.

For all methods, we use supervised fine-tuned LLMs based on *Llama-3* (Dubey et al., 2024), specifically two variants: *Llama-3.2-3B* and *Llama-3.1-8B*.

Datasets We conduct experiments on two datasets: a simple preference dataset *Pets* and a complex preference dataset *UltraFeedback-P (UF-P)* (Poddar et al., 2024) derived from *Ultrafeedback* (Cui et al., 2023), featuring user types pursuing values like helpfulness, honesty, instruction-following, and truthfulness.

The *Pets* dataset simulates multi-modal user preferences over animals (cats, dogs, birds, and rabbits), capturing consensus in some comparisons (e.g., universally most- and least-preferred pets) and divergence in others (e.g., middle-ranking pets).

The UF-P dataset assumes that each user h belongs to one of several preference types \mathbb{P} (e.g., $p \in \mathbb{P} = \{\text{helpfulness, honesty}\}$). It is constructed from the Ultrafeedback prompt—response data, where responses are labeled by GPT-4 (Achiam et al., 2023) with scores for each type p. For each prompt, the winning and losing responses are selected according to the score associated with the target preference type p. Specifically, UF-P-2 contains two preference types focusing on helpfulness and honesty, while UF-P-4 contains four preference types focusing on helpfulness, honesty, instruction-following, and truthfulness. Due to its diverse preference modes and a wide variety of prompt—response pairs, the UF-P dataset is highly ambiguous and challenging.

In all datasets, one sample \mathbb{D}_h corresponds to a user h with a user type p. This type information is used only when constructing \mathbb{D}_h (to determine winning and losing responses from *Ultrafeedback*) and for qualitative evaluation (to verify whether user types are well-separated). Importantly, the latent embedding relies on user preference data \mathbb{D}_h , not on type p. Additional details about experiments are provided in Appendix E.

5.1 RESULTS

We first demonstrate that our method effectively reduces posterior collapse and encodes a stable latent space. To diagnose collapse quantitatively, we evaluate using the *Active Units* (AU) metric from prior work (Burda et al., 2015). AU counts latent dimensions with variability exceeds a small threshold δ ; a dimension u is considered active if its posterior mean responses show sufficient variability across the evaluation set \mathbb{D}_{eval} .

$$AU = |\{u : \operatorname{Var}_{\mathbb{D}_{\text{eval}}} (\mu_{\psi, u}(\mathbb{D}_{\text{eval}})) > \delta\}|$$

Thus, AU=0 means all latent dimensions are unresponsive across evaluation data. In these runs, the encoder outputs fixed posterior means and variances with AU=0, indicated as *posterior collapse* and shaded gray in tables. Accuracy is the ratio of evaluation samples where predicted rewards match user preferences (i.e., winning responses have higher rewards).

		Active Units [%]			Accuracy [%]				
Model	eta	UF-P-2		UF-P-4		UF-P-2		UF-P-4	
		VPL	SPL	VPL	SPL	VPL	SPL	VPL	SPL
	3.0×10^{-7}	0.00	88.09	0.00	85.35	61.90	62.84	56.83	61.91
Llama-3.2-3B	3.0×10^{-6}	92.09	92.77	0.00	51.76	62.42	63.42	56.75	61.55
	3.0×10^{-5}	19.53	87.60	0.00	21.00	62.49	62.79	56.47	61.48
	3.0×10^{-7}	0.00	94.34	0.00	99.12	62.47	63.54	57.14	62.01
Llama-3.1-8B	3.0×10^{-6}	96.09	96.88	0.00	95.70	62.98	63.65	57.09	62.26
	3.0×10^{-5}	93.95	90.72	0.00	85.35	62.37	63.46	57.15	62.47

Table 1: Active units and accuracy across β

Table 1 shows results for VPL and SPL across a range of KL-divergence weights β under the same random seed. Prior methods require careful tuning of β to avoid collapse. In contrast, SPL exhibits no posterior collapse in any of the tested settings. The advantage is most evident on highly multi-modal preference datasets $\mathit{UF-P-4}$, where VPL collapses under all tested β values, but SPL consistently maintains high AU. Notably, SPL is much less sensitive to β .

Next, Table 2 compares preference-prediction accuracy against baselines. For fairness, we fix $\beta=3\times10^{-6}$ —a setting under which VPL is comparatively more stable—and report the mean \pm standard deviation over three distinct random seeds for all methods. Across all datasets and models, SPL achieves higher preference-prediction accuracy than competing baselines. Simply augmenting VPL with a standard IAF does not yield robust encoding and often fails to prevent collapse. Similarly, simply substituting the flow component of SPL with a standard IAF noticeably reduces

Table 2: Preference-prediction accuracy (%) compared with baselines

Model	Method	Pets	UF-P-2	UF-P-4
Llama-3.2-3B	BTL	57.48 ± 2.37	62.24 ± 0.04	57.05 ± 0.02
Liailia-3.2-3D	DPL	62.02 ± 1.92	62.21 ± 0.03	57.04 ± 0.06
	VPL	99.67 ± 0.38	62.42 ± 0.25	56.99 ± 0.14
	VPL-IAF	$\textbf{100.0} \pm \textbf{0.00}$	62.21 ± 0.17	57.73 ± 1.32
	SPL-IAF	$\textbf{100.0} \pm \textbf{0.00}$	63.10 ± 0.26	59.35 ± 0.22
	SPL (Ours)	$\textbf{100.0} \pm \textbf{0.00}$	$\textbf{63.24} \pm \textbf{0.15}$	$\textbf{61.52} \pm \textbf{0.05}$
Llama-3.1-8B	BTL	60.74 ± 0.49	62.59 ± 0.04	57.42 ± 0.34
Liailia-3.1-6D	DPL	61.03 ± 0.25	62.75 ± 0.02	57.58 ± 0.45
	VPL	75.33 ± 0.63	62.98 ± 0.73	57.14 ± 0.05
	VPL-IAF	$\textbf{100.0} \pm \textbf{0.00}$	63.10 ± 0.26	58.73 ± 0.23
	SPL-IAF	$\textbf{100.0} \pm \textbf{0.00}$	63.27 ± 0.11	60.74 ± 0.40
	SPL (Ours)	$\textbf{100.0} \pm \textbf{0.00}$	$\textbf{63.74} \pm \textbf{0.23}$	$\textbf{62.21} \pm \textbf{0.06}$

accuracy. These results mean swap-guided base regularization and P-IAF are effectively encoding user preference to identifiable user-latent.

We further examine the learned latent spaces qualitatively. Fig. 5 visualizes the encoded user-latent on the *UF-P* dataset (non-collapse cases). SPL yields more compact and distinctly separated embeddings compared to baselines. SPL with a standard IAF (SPL-IAF) prevents collapse and achieves high accuracy but results in a scattered and complex posterior. The standard IAF allocates its modeling capacity toward complex transformations rather than swap-derived properties. In contrast, our P-IAF maintains IAF's expressivity but reduces unnecessary complexity through swap-guided encoding. Further analysis of additional experiments is provided in Appendix D.

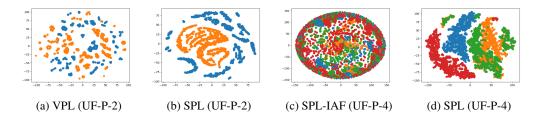


Figure 5: Non-collapsed latent embeddings learned on the UF-P dataset. We visualize latent embeddings z from baselines and SPL (Ours) encoder using 2D t-SNE (Maaten & Hinton, 2008). Each point denotes a user, colored by their preference type.

6 CONCLUSION

We proposed Swap-guided Latent Preference Learning (SPL), a framework that overcomes the failure mode in preference learning on complex textual preference data. Across all experiments, SPL consistently improves prediction accuracy over baselines and prevents collapse. These results suggest that the combination of our base regularization, P-IAF, and adaptive latent conditioning effectively encodes user-specific latent from complex textual preferences—even under sparse preference signals. Consequently, by explicitly conditioning the reward on the user latent and optimizing the policy under this reward, our framework enables user-specific behaviors, achieving personalized alignment.

Limitation Our study focuses on encoding user preferences from independent, single-turn comparison data. This data requirement can be burdensome and may feel unnecessary from a user perspective. We believe our framework can be extended to preferences expressed over natural, multiturn dialogue; we consider this for future work.

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.
- Avinandan Bose, Zhihan Xiong, Yuejie Chi, Simon Shaolei Du, Lin Xiao, and Maryam Fazel. Lore: Personalizing Ilms via low-rank reward modeling. *arXiv preprint arXiv:2504.14439*, 2025.
- Samuel Bowman, Luke Vilnis, Oriol Vinyals, Andrew Dai, Rafal Jozefowicz, and Samy Bengio. Generating sentences from a continuous space. In *Proceedings of the 20th SIGNLL conference on computational natural language learning*, pp. 10–21, 2016.
- 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.
- Yuri Burda, Roger Grosse, and Ruslan Salakhutdinov. Importance weighted autoencoders. *arXiv* preprint arXiv:1509.00519, 2015.
- Stephen Casper, Xander Davies, Claudia Shi, T. Gilbert, J'er'emy Scheurer, Javier Rando, Rachel Freedman, Tomasz Korbak, David Lindner, Pedro J Freire, Tony Wang, Samuel Marks, Charbel-Raphau00ebl Su00e9gerie, Micah Carroll, Andi Peng, Phillip J. K. Christoffersen, Mehul Damani, Stewart Slocum, Usman Anwar, Anand Siththaranjan, Max Nadeau, Eric J. Michaud, J. Pfau, Dmitrii Krasheninnikov, Xin Chen, L. Langosco, Peter Hase, Erdem Biyik, A. Dragan, David Krueger, Dorsa Sadigh, and Dylan Hadfield-Menell. Open problems and fundamental limitations of reinforcement learning from human feedback. In *Trans. Mach. Learn. Res.*, 2023.
- Xi Chen, Diederik P Kingma, Tim Salimans, Yan Duan, Prafulla Dhariwal, John Schulman, Ilya Sutskever, and Pieter Abbeel. Variational lossy autoencoder. *arXiv preprint arXiv:1611.02731*, 2016.
- P. Christiano, Jan Leike, Tom B. Brown, Miljan Martic, S. Legg, and Dario Amodei. Deep reinforcement learning from human preferences. In *Neural Information Processing Systems*, 2017.
- 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. In *International Conference on Machine Learning*, 2023. doi: 10.48550/arXiv.2310.01377.
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Tianyu Liu, et al. A survey on in-context learning. *arXiv preprint arXiv:2301.00234*, 2022.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv e-prints*, pp. arXiv–2407, 2024.
- Michael Feffer, Hoda Heidari, and Zachary Chase Lipton. Moral machine or tyranny of the majority? In *AAAI Conference on Artificial Intelligence*, 2023. doi: 10.48550/arXiv.2305.17319.
- Zhuocheng Gong, Jian Guan, Wei Wu, Huishuai Zhang, and Dongyan Zhao. Latent preference coding: Aligning large language models via discrete latent codes. *arXiv preprint arXiv:2505.04993*, 2025.
- Junxian He, Daniel Spokoyny, Graham Neubig, and Taylor Berg-Kirkpatrick. Lagging inference networks and posterior collapse in variational autoencoders. *arXiv preprint arXiv:1901.05534*, 2019.
- Jiaming Ji, Tianyi Qiu, Boyuan Chen, Borong Zhang, Hantao Lou, Kaile Wang, Yawen Duan, Zhonghao He, Jiayi Zhou, Zhaowei Zhang, Fanzhi Zeng, Kwan Yee Ng, Juntao Dai, Xuehai Pan, Aidan O'Gara, Yingshan Lei, Hua Xu, Brian Tse, Jie Fu, S. McAleer, Yaodong Yang, Yizhou Wang, Song-Chun Zhu, Yike Guo, and Wen Gao. Ai alignment: A comprehensive survey. In *arXiv.org*, 2023.
- Diederik P Kingma, Max Welling, et al. Auto-encoding variational bayes, 2013.

- Durk P Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max Welling. Improved variational inference with inverse autoregressive flow. *Advances in neural information processing systems*, 29, 2016.
 - Jan Leike, David Krueger, Tom Everitt, Miljan Martic, Vishal Maini, and S. Legg. Scalable agent alignment via reward modeling: a research direction. In *arXiv.org*, 2018.
 - James Lucas, George Tucker, Roger B Grosse, and Mohammad Norouzi. Don't blame the elbo! a linear vae perspective on posterior collapse. *Advances in Neural Information Processing Systems*, 32, 2019.
 - Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605, 2008.
 - Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.
 - Hyunji Nam, Yanming Wan, Mickel Liu, Jianxun Lian, and Natasha Jaques. Learning pluralistic user preferences through reinforcement learning fine-tuned summaries. *arXiv* preprint arXiv:2507.13579, 2025.
 - Wen Zheng Terence Ng, Jianda Chen, Yuan Xu, and Tianwei Zhang. Latent embedding adaptation for human preference alignment in diffusion planners. *arXiv preprint arXiv:2503.18347*, 2025.
 - Minhyeon Oh, Seungjoon Lee, and Jungseul Ok. Active preference-based learning for multi-dimensional personalization. In *arXiv.org*, 2024. doi: 10.48550/arXiv.2411.00524.
 - Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. Advances in neural information processing systems, 35: 27730–27744, 2022.
 - George Papamakarios, Eric Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji Lakshminarayanan. Normalizing flows for probabilistic modeling and inference. *Journal of Machine Learning Research*, 22(57):1–64, 2021.
 - Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual reasoning with a general conditioning layer. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
 - 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.
 - Vinodkumar Prabhakaran, Aida Mostafazadeh Davani, and Mark Díaz. On releasing annotator-level labels and information in datasets. *arXiv* (*Cornell University*), 2021.
 - Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in neural information processing systems*, 36:53728–53741, 2023.
 - Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *International conference on machine learning*, pp. 1530–1538. PMLR, 2015.
 - Idan Shenfeld, Felix Faltings, Pulkit Agrawal, and Aldo Pacchiano. Language model personalization via reward factorization. *arXiv preprint arXiv:2503.06358*, 2025.
 - Anand Siththaranjan, Cassidy Laidlaw, and Dylan Hadfield-Menell. Distributional preference learning: Understanding and accounting for hidden context in rlhf. In *International Conference on Learning Representations*, 2023. doi: 10.48550/arXiv.2312.08358.
 - Taylor Sorensen, Jared Moore, Jillian R. Fisher, Mitchell Gordon, Niloofar Mireshghallah, Christopher Rytting, Andre Ye, Liwei Jiang, Ximing Lu, Nouha Dziri, Tim Althoff, and Yejin Choi. A roadmap to pluralistic alignment. In *International Conference on Machine Learning*, 2024. doi: 10.48550/arXiv.2402.05070.

Ruiqi Wang, Dezhong Zhao, Dayoon Suh, Ziqin Yuan, Guohua Chen, and Byung-Cheol Min. Personalization in human-robot interaction through preference-based action representation learning. In 2025 IEEE International Conference on Robotics and Automation (ICRA), pp. 7377–7384. IEEE, 2025.

Yixin Wang, David Blei, and John P Cunningham. Posterior collapse and latent variable non-identifiability. *Advances in neural information processing systems*, 34:5443–5455, 2021.

APPENDIX

A EVIDENCE FOR COLLAPSE VIA POSTERIOR-PRIOR RESPONSE

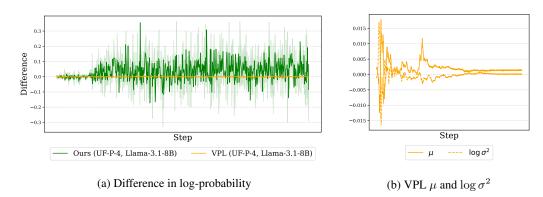


Figure 6: Evidence of posterior collapse in preference learning

We provide evidence that VPL causes the latent to collapse, preventing meaningful encoding during training. Fig. 6a contrasts decoder outputs from a approximate posterior latent and a noise vector ϵ from the prior $\mathcal{N}(\mathbf{0}, I)$. Specifically, we compute

$$[\log p_{\phi}(y_w \succ y_l \mid x, z) - \log p_{\phi}(y_w \succ y_l \mid x, \epsilon)].$$

In non-collapsing runs (e.g., with our SPL), the difference persists, indicating z's informative signal. Conversely, under VPL, the difference remains negligible during training. In this regime, Fig. 6b shows that the encoder's μ and $\ell = \log \sigma^2$ initially carry signal but soon drift toward $\mu \approx 0$ and $\log \sigma^2 \approx 0$, making the posterior almost the same as the prior. The encoded z lacks helpful information for the decoder's reward decision, resulting in a trivial solution that reduces the KL penalty to zero.

B JUSTIFICATION FOR P-IAF

For a prompt-response pair (x, y_w, y_l) and a latent $z \in \mathbb{R}^d$, let us define

$$\Delta r_{\phi}(z) \triangleq r_{\phi}(x, y_w, z) - r_{\phi}(x, y_l, z),$$

which is the part inside the sigmoid σ in Eq.(2). In the swap-guided base regularization, we regularize our encoder q_{ψ} so that $\mu = -\mu_{\text{swap}}$ and $\sigma = \sigma_{\text{swap}}$. Assuming opposite coupling for the fictitious annotator h_{swap} , i.e., $\epsilon_{\text{swap}} = -\epsilon$, the latent samples become

$$oldsymbol{z} = oldsymbol{\mu} + oldsymbol{\sigma} \odot oldsymbol{\epsilon}, \quad oldsymbol{z}_{ ext{swap}} = oldsymbol{\mu}_{ ext{swap}} + oldsymbol{\sigma}_{ ext{swap}} \odot oldsymbol{\epsilon}_{ ext{swap}},$$

respectively. Consequently, we obtain

$$z_0 = -z_{0,\text{swap}}.\tag{14}$$

from the encoder q_{ψ} in Fig. 4 under swapping. Finally, we regularize the base posterior for probability consistency

$$p_{\phi}(y_w \succ y_l \mid x, \mathbf{z}_0) = p_{\phi}(y_l \succ y_w \mid x, \mathbf{z}_{0,\text{swap}}). \tag{15}$$

However, unlike the base posterior z_0 , we cannot directly regularize after transforming $z = z_0$ into z_K :

$$p_{\phi}(y_w \succ y_l \mid x, \mathbf{z}_K) = p_{\phi}(y_l \succ y_w \mid x, \mathbf{z}_{K \text{ swap}}), \tag{16}$$

because IAF entangles dimensions and contexts, so the posterior's mirrored structure need not be preserved after the flow. However, by supplying c_s and c_d to the μ_k and σ_k functions as input

arguments separately, we can obtain similar results as base regularization indirectly. We will show it in this appendix. First, let us define swap probability error of transformed distribution z_K by

$$\delta_p \triangleq \sigma(\Delta r_{\phi}(\mathbf{z}_K)) - \sigma(-\Delta r_{\phi}(\mathbf{z}_{K,\text{swap}})). \tag{17}$$

Further, we assume that (A1) $z \mapsto \Delta r_{\phi}(z)$ is Lipschitz; (A2) opposite coupling is used for the base noise, i.e., $\epsilon = -\epsilon_{\text{swap}}$; (A3) each step k-th scale function σ_k is bounded by $\|\sigma_k(\cdot)\|_{\infty} \leq \rho_k$;

The key idea behind our justification of P-IAF is to demonstrate that the mirroring of preference swaps is realized in the transformed posterior z_K by showing that the swap probability error δ_p of our P-IAF is smaller than that of IAF.

Lemma 1. Let us suppose that z_0 and $z_{0,swap}$ are warped to z_K and $z_{K,swap}$ respectively, by P-IAF. Then, the swap probability error δ_p given in Eq.(17) is bounded by

$$|\delta_p| \le \frac{1}{4} \, \delta_{r,K} + \frac{1}{4} \, L_r \, \|\delta_{z,K}\|.$$
 (18)

where the reward violation $\delta_{r,K} \triangleq |\Delta r_{\phi}(\mathbf{z}_K) + \Delta r_{\phi}(-\mathbf{z}_K)|$ and the latent mismatch $\delta_{z,K} \triangleq \mathbf{z}_{K,swap} + \mathbf{z}_K$.

Proof. For all $a,b\in\mathbb{R}$, the logistic satisfies $|\sigma(a)-\sigma(b)|\leq \frac{1}{4}|a-b|$, the swap probability error δ_p is bounded by

$$\begin{aligned} |\delta_p| &= \left| \sigma(\Delta r_{\phi}(\boldsymbol{z}_K)) - \sigma(-\Delta r_{\phi}(\boldsymbol{z}_{K,\text{swap}})) \right| \\ &\leq \frac{1}{4} \left| \Delta r_{\phi}(\boldsymbol{z}_K) + \Delta r_{\phi}(\boldsymbol{z}_{K,\text{swap}}) \right|. \end{aligned}$$

Using the triangle inequality, we obtain

$$|\delta_{p}| \leq \frac{1}{4} \left| \Delta r_{\phi}(\boldsymbol{z}_{K}) + \Delta r_{\phi}(\boldsymbol{z}_{K,\text{swap}}) \right|$$

$$= \frac{1}{4} \left| \Delta r_{\phi}(\boldsymbol{z}_{K}) + \Delta r_{\phi}(-\boldsymbol{z}_{K}) + \Delta r_{\phi}(\boldsymbol{z}_{K,\text{swap}}) - \Delta r_{\phi}(-\boldsymbol{z}_{K}) \right|$$

$$= \frac{1}{4} \left| \Delta r_{\phi}(\boldsymbol{z}_{K}) + \Delta r_{\phi}(-\boldsymbol{z}_{K}) \right| + \left| \Delta r_{\phi}(\boldsymbol{z}_{K,\text{swap}}) - \Delta r_{\phi}(-\boldsymbol{z}_{K}) \right|$$

$$= \delta_{r,K}$$
(19)

Further, using the Lipschitz assumption (A1),

$$\left|\Delta r_{\phi}(\boldsymbol{z}_{K,\text{swap}}) - \Delta r_{\phi}(-\boldsymbol{z}_{K})\right| \leq L_{r} \|\boldsymbol{z}_{K,\text{swap}} - (-\boldsymbol{z}_{K})\| = L_{r} \|\delta_{z,K}\|.$$

Then we obtain Eq.(18) by combining reward violation and latent mismatch.

Lemma 2. Let us suppose that the base posterior is given by $q_{\psi}(z_0 \mid \mathbb{D}_h) = \mathcal{N}(\mu, \sigma^2)$ and $q_{\psi}(z_0 \mid \mathbb{D}_{h_{swap}}) = \mathcal{N}(\mu_{swap}, \sigma^2_{swap})$. When we sample the latent z_0 and $z_{0,swap}$ based on assumption (A2), the base mismatch defined by

$$\delta_{z,0} \triangleq oldsymbol{z}_{0,\mathit{swap}} + oldsymbol{z}_0 = (oldsymbol{\mu} + oldsymbol{\mu}_\mathit{swap}) + (oldsymbol{\sigma} - oldsymbol{\sigma}_\mathit{swap}) \odot oldsymbol{\epsilon}.$$

And also bounded by

$$\mathbb{E}\|\delta_{z,0}\| \leq \|\mu + \mu_{\text{swap}}\| + \frac{1}{2} \exp(\ell_{\text{max}}^{(\infty)}/2) \|\ell - \ell_{\text{swap}}\|, \tag{20}$$

where $\sigma = exp(\ell/2)$.

Proof.

$$\mathbb{E}\|\delta_{z,0}\| \leq \|\boldsymbol{\mu} + \boldsymbol{\mu}_{\text{swap}}\| + \|\boldsymbol{\sigma} - \boldsymbol{\sigma}_{\text{swap}}\|, \tag{21}$$

since
$$\mathbb{E}\|\boldsymbol{A}\epsilon\| \leq \sqrt{\mathbb{E}\|\boldsymbol{A}\epsilon\|^2} = \|\boldsymbol{A}\|$$
 where $\boldsymbol{A} \triangleq \operatorname{diag}(\boldsymbol{\sigma} - \boldsymbol{\sigma}_{\operatorname{swap}}) \in \mathbb{R}^{d \times d}$.

Moreover, $\sigma = \exp(\ell/2)$, by the mean value theorem for $g(t) = \exp(t/2)$, then, $|g(a) - g(b)| = \frac{1}{2}\exp(\xi/2)|a-b| \le \frac{1}{2}\exp(\max\{a,b\}/2)|a-b|$ for some ξ between a and b.

$$\|\boldsymbol{\sigma} - \boldsymbol{\sigma}_{\text{swap}}\| \le \frac{1}{2} \exp(\boldsymbol{\ell}_{\text{max}}^{(\infty)}/2) \|\boldsymbol{\ell} - \boldsymbol{\ell}_{\text{swap}}\|, \qquad \boldsymbol{\ell}_{\text{max}}^{(\infty)} \triangleq \max\{\|\boldsymbol{\ell}\|_{\infty}, \|\boldsymbol{\ell}_{\text{swap}}\|_{\infty}\}.$$
 (22)

Hence, base regularization Eq.(7) directly decreases the base mismatch $\|\delta_{z,0}\|$.

From now on, we will compute the swap probability errors of our P-IAF and IAF methods one by one. Before deriving the swap probability errors of the two normalizing flows P-IAF and IAF, let us consider context vector c, which is an additional output of encoder q_{ψ} . For further development, we decompose the context vector c into a swap-reversal context c_d and swap-invariant context c_s as:

$$oldsymbol{c}_d = rac{1}{2}(oldsymbol{c} - oldsymbol{c}_{ ext{swap}}), \qquad oldsymbol{c}_s = rac{1}{2}(oldsymbol{c} + oldsymbol{c}_{ ext{swap}}), \qquad oldsymbol{c} = oldsymbol{c}_d + oldsymbol{c}_s,$$

which ensures $c_{d,\text{swap}} = -c_d$ and $c_{s,\text{swap}} = c_s$.

Assumption (A4). Let μ_k, σ_k denote the k-th step shift and scale function. There exist non-negative constants ²

$$L_{\mu,k}^{z}, L_{\mu,k}^{c_d}, L_{\mu,k}^{c_s}, L_{\sigma,k}^{z}, L_{\sigma,k}^{c_d}, L_{\sigma,k}^{c_s}$$

such that, for all (z, c_d, c_s) and (z', c'_d, c'_s) in the valid input space,

$$\|\mu_k(\boldsymbol{z}, \boldsymbol{c}_d, \boldsymbol{c}_s) - \mu_k(\boldsymbol{z}', \boldsymbol{c}_d', \boldsymbol{c}_s')\| \leq L_{\mu,k}^z \|\boldsymbol{z} - \boldsymbol{z}'\| + L_{\mu,k}^{c_d} \|\boldsymbol{c}_d - \boldsymbol{c}_d'\| + L_{\mu,k}^{c_s} \|\boldsymbol{c}_s - \boldsymbol{c}_s'\|, \\ \|\sigma_k(\boldsymbol{z}, \boldsymbol{c}_d, \boldsymbol{c}_s) - \sigma_k(\boldsymbol{z}', \boldsymbol{c}_d', \boldsymbol{c}_s')\| \leq L_{\sigma,k}^z \|\boldsymbol{z} - \boldsymbol{z}'\| + L_{\sigma,k}^{c_d} \|\boldsymbol{c}_d - \boldsymbol{c}_d'\| + L_{\sigma,k}^{c_s} \|\boldsymbol{c}_s - \boldsymbol{c}_s'\|.$$

Lemma 3 (Transformed mismatch of P-IAF). Let us consider a normalizing flow P-IAF given by

$$\boldsymbol{z}_k = \mu_k(\boldsymbol{z}_{k-1}, \boldsymbol{c}_d) + \sigma_k(\boldsymbol{z}_{k-1}, \boldsymbol{c}_s) \odot \boldsymbol{z}_{k-1}$$

If we define the transformed mismatch $\delta_{z,k}$ at the k-th step as:

$$\delta_{z,k} \triangleq z_{k,swap} + z_k$$

Then, the mismatch $\delta_{z,k}$ of P-IAF is bounded by

$$\|\delta_{z,k}\| \leq \left(\rho_{k} + L_{\mu,k}^{z} + L_{\sigma,k}^{z} \|\mathbf{z}_{k-1}\|\right) \|\delta_{z,k-1}\| + \underbrace{\|\delta_{\mu,k}(\mathbf{c}_{d})\|}_{\text{swap-reversal violation }(\mu)} + \underbrace{\|\delta_{\sigma,k}(\mathbf{c}_{s})\|_{\infty} \|\mathbf{z}_{k-1}\|}_{\text{swap-invariant violation }(\sigma)}$$

$$(23)$$

where we define the μ swap-reversal violation and σ swap-invariant violations at step k as:

$$\delta_{\mu,k}(\boldsymbol{c}) \triangleq \mu_k(\boldsymbol{z}_{k-1}, \boldsymbol{c}) + \mu_k(-\boldsymbol{z}_{k-1}, \boldsymbol{c}_{swap}), \qquad \delta_{\sigma,k}(\boldsymbol{c}) \triangleq \sigma_k(\boldsymbol{z}_{k-1}, \boldsymbol{c}) - \sigma_k(-\boldsymbol{z}_{k-1}, \boldsymbol{c}_{swap}). \tag{24}$$

Proof. In the P-IAF, the outputs from the k-th step is given by

$$egin{aligned} oldsymbol{z}_k &= \mu_k(oldsymbol{z}_{k-1}, oldsymbol{c}_d) + \sigma_k(oldsymbol{z}_{k-1}, oldsymbol{c}_s) \odot oldsymbol{z}_{k-1}, \ oldsymbol{z}_{k, ext{swap}} &= \mu_k(oldsymbol{z}_{k-1, ext{swap}}, oldsymbol{c}_{d, ext{swap}}) + \sigma_k(oldsymbol{z}_{k-1, ext{swap}}, oldsymbol{c}_{s, ext{swap}}) \odot oldsymbol{z}_{k-1, ext{swap}}. \end{aligned}$$

Then, the transformed mismatch $\delta_{z,k}$ at the k-th step is given by

$$\begin{split} \delta_{z,k} &= \left[\mu_k(\boldsymbol{z}_{k-1}, \boldsymbol{c}_d) + \mu_k(-\boldsymbol{z}_{k-1}, \boldsymbol{c}_{d, \text{swap}}) \right] \\ &+ \left[\sigma_k(\boldsymbol{z}_{k-1}, \boldsymbol{c}_s) \odot \boldsymbol{z}_{k-1} + \sigma_k(-\boldsymbol{z}_{k-1}, \boldsymbol{c}_{s, \text{swap}}) \odot (-\boldsymbol{z}_{k-1}) \right] \\ &+ \left[\mu_k(\boldsymbol{z}_{k-1, \text{swap}}, \boldsymbol{c}_{d, \text{swap}}) - \mu_k(-\boldsymbol{z}_{k-1}, \boldsymbol{c}_{d, \text{swap}}) \right] \\ &+ \left[\sigma_k(\boldsymbol{z}_{k-1, \text{swap}}, \boldsymbol{c}_{s, \text{swap}}) \odot \boldsymbol{z}_{k-1, \text{swap}} - \sigma_k(-\boldsymbol{z}_{k-1}, \boldsymbol{c}_{s, \text{swap}}) \odot (-\boldsymbol{z}_{k-1}) \right]. \end{split}$$

The first bracket equals $\delta_{\mu,k}(c_d)$ by Eq.(24), hence contributes $\|\delta_{\mu,k}(c_d)\|$.

For the second bracket, by Eq.(24) and $||a \odot b|| \le ||a||_{\infty} ||b||$:

$$\|\delta_{\sigma,k}(\boldsymbol{c}_s)\|_{\infty}\|\boldsymbol{z}_{k-1}\|.$$

For the third bracket, by (A4) in z:

$$\|\mu_k(\boldsymbol{z}_{k-1,\text{swap}}, \boldsymbol{c}_{d,\text{swap}}) - \mu_k(-\boldsymbol{z}_{k-1}, \boldsymbol{c}_{d,\text{swap}})\| \le L_{u,k}^z \|\delta_{z,k-1}\|.$$

²The shift and scale networks are compositions of affine maps and smooth activations; hence they are locally Lipschitz on the working domain considered here.

For the fourth bracket, insert-delete $\sigma_k(z_{k-1,\text{swap}}, c_{s,\text{swap}}) \odot z_{k-1}$:

$$egin{aligned} \sigma_k(oldsymbol{z}_{k-1, ext{swap}}, oldsymbol{c}_{s, ext{swap}}) \odot (oldsymbol{z}_{k-1, ext{swap}} + oldsymbol{z}_{k-1}) \ &+ igl(\sigma_k(oldsymbol{z}_{k-1, ext{swap}}, oldsymbol{c}_{s, ext{swap}}) - \sigma_k(-oldsymbol{z}_{k-1}, oldsymbol{c}_{s, ext{swap}})igr) \odot (-oldsymbol{z}_{k-1}). \end{aligned}$$

Bound the first term using (A3):

$$\|\sigma_k(\boldsymbol{z}_{k-1,\text{swap}}, \boldsymbol{c}_{s,\text{swap}}) \odot (\boldsymbol{z}_{k-1,\text{swap}} + \boldsymbol{z}_{k-1})\| \le \rho_k \|\boldsymbol{z}_{k-1,\text{swap}} + \boldsymbol{z}_{k-1}\| = \rho_k \|\delta_{z,k-1}\|.$$

Bound the second term using (A4) in z:

$$\left(\sigma_{k}(\boldsymbol{z}_{k-1,\text{swap}},\boldsymbol{c}_{s,\text{swap}}) - \sigma_{k}(-\boldsymbol{z}_{k-1},\boldsymbol{c}_{s,\text{swap}})\right) \odot \left(-\boldsymbol{z}_{k-1}\right) \leq L_{\sigma,k}^{z} \left\|\delta_{z,k-1}\right\| \left\|\boldsymbol{z}_{k-1}\right\|$$

Collecting the bounds yields Eq.(23).

Lemma 4 (Transformed mismatch of IAF). Let us consider a normalizing flow IAF given by

$$z_k = \mu_k(z_{k-1}, c) + \sigma_k(z_{k-1}, c) \odot z_{k-1},$$

where $\mathbf{c} = \mathbf{c}_d + \mathbf{c}_s$. Then, the mismatch $\delta_{z,k}$ of IAF is bounded by

$$\|\delta_{z,k}\| \leq (\rho_k + L_{\mu,k}^z + L_{\sigma,k}^z \| \boldsymbol{z}_{k-1} \|) \|\delta_{z,k-1}\| \\ + \underbrace{\|\delta_{\mu,k}(\boldsymbol{c}_d)\|}_{\text{swap-reversal violation }(\mu)} + \underbrace{2L_{\mu,k}^{c_s} \|\boldsymbol{c}_s\|}_{\text{leak }(\boldsymbol{c}_s \to \mu)} + \underbrace{\|\delta_{\sigma,k}(\boldsymbol{c}_s)\|_{\infty} \|\boldsymbol{z}_{k-1}\|}_{\text{swap-invariant violation }(\sigma)} + \underbrace{2L_{\sigma,k}^{c_d} \|\boldsymbol{c}_d\| \|\boldsymbol{z}_{k-1}\|}_{\text{leak }(\boldsymbol{c}_d \to \sigma)}.$$

$$(25)$$

Proof. The derivation mirrors the proof in Eq.(23) except for two aspects. First, we replace $\delta_{\mu,k}(\mathbf{c}_d)$ and $\delta_{\sigma,k}(\mathbf{c}_s)$ with $\delta_{\mu,k}(\mathbf{c})$ and $\delta_{\sigma,k}(\mathbf{c})$ respectively using definition Eq.(24). Decompose these terms via insert–delete step:

$$\begin{split} \delta_{\mu,k}(\boldsymbol{c}) &= \underbrace{\mu_k(\boldsymbol{z}_{k-1},\boldsymbol{c}_d,0) + \mu_k(-\boldsymbol{z}_{k-1},\boldsymbol{c}_{d,\text{swap}},0)}_{\delta_{\mu,k}(\boldsymbol{c}_d)} \\ &+ \underbrace{\left[\mu_k(\boldsymbol{z}_{k-1},\boldsymbol{c}_d,\boldsymbol{c}_s) - \mu_k(\boldsymbol{z}_{k-1},\boldsymbol{c}_d,0)\right]}_{\Delta_1^{(s)}} + \underbrace{\left[\mu_k(-\boldsymbol{z}_{k-1},\boldsymbol{c}_{d,\text{swap}},\boldsymbol{c}_{s,\text{swap}}) - \mu_k(-\boldsymbol{z}_{k-1},\boldsymbol{c}_{d,\text{swap}},0)\right]}_{\Delta_2^{(s)}}, \\ \delta_{\sigma,k}(\boldsymbol{c}) &= \underbrace{\sigma_k(\boldsymbol{z}_{k-1},0,\boldsymbol{c}_s) - \sigma_k(-\boldsymbol{z}_{k-1},0,\boldsymbol{c}_{s,\text{swap}})}_{\delta_{\sigma,k}(\boldsymbol{c}_s)} \\ &+ \underbrace{\left[\sigma_k(\boldsymbol{z}_{k-1},\boldsymbol{c}_d,\boldsymbol{c}_s) - \sigma_k(\boldsymbol{z}_{k-1},0,\boldsymbol{c}_s)\right]}_{\Delta_1^{(d)}} + \underbrace{\left[\sigma_k(-\boldsymbol{z}_{k-1},0,\boldsymbol{c}_{s,\text{swap}}) - \sigma_k(-\boldsymbol{z}_{k-1},\boldsymbol{c}_{d,\text{swap}},\boldsymbol{c}_{s,\text{swap}})\right]}_{\Delta_2^{(d)}}. \end{split}$$

For the Δ terms, by the (A4),

$$\begin{aligned} \|\Delta_1^{(s)}\| &\leq L_{\mu,k}^{c_s} \|\boldsymbol{c}_s\|, \quad \|\Delta_2^{(s)}\| \leq L_{\mu,k}^{c_s} \|\boldsymbol{c}_s\|, \\ \|\Delta_1^{(d)}\| &\leq L_{\sigma,k}^{c_d} \|\boldsymbol{c}_d\|, \quad \|\Delta_2^{(d)}\| \leq L_{\sigma,k}^{c_d} \|\boldsymbol{c}_d\|. \end{aligned}$$

Therefore, we obtain

$$\|\delta_{\mu,k}(\boldsymbol{c})\| \leq \|\delta_{\mu,k}(\boldsymbol{c}_d)\| + 2 L_{\mu,k}^{c_s} \|\boldsymbol{c}_s\|, \\ \|\delta_{\sigma,k}(\boldsymbol{c})\| \leq \|\delta_{\sigma,k}(\boldsymbol{c}_s)\| + 2 L_{\sigma,k}^{c_d} \|\boldsymbol{c}_d\|.$$

Collecting the bounds yields Eq.(25).

Bound-Level Comparison between P-IAF and IAF Assume (A1)–(A4) hold and that P-IAF and IAF share the same architecture and training hyperparameters so that they admit the same upper bounds on the local Lipschitz constants $\{L_{\mu,k}^z, L_{\mu,k}^{c_d}, L_{\sigma,k}^{c_s}, L_{\sigma,k}^{c_d}, L_{\sigma,k}^{c_s}\}$, the same scale bounds ρ_k , the same reward Lipschitz constant L_r , and the same initial mismatch $\|\delta_{z,0}\|$. By Lemma 3 and Lemma 4, the IAF per-step bound contains two additional non-negative leak terms, $2L_{u,k}^{c_s}\|c_s\|$

and $2L_{\sigma,k}^{c_d} \|\mathbf{c}_d\| \|\mathbf{z}_{k-1}\|$, that are absent in P-IAF. Consequently, by induction over K starting from Lemma 2,

$$\mathrm{UB}\big(\|\delta_{z,K}\|\big)^{(\text{P-IAF})} \, \leq \, \mathrm{UB}\big(\|\delta_{z,K}\|\big)^{(\mathrm{IAF})},$$

where $UB(\cdot)$ denotes the upper bound under the shared constants above.

Suppose in addition that the reward violation $\delta_{r,K}$ admits a common bound across the two flows, i.e., $\delta_{r,K}^{\text{(P-IAF)}} \leq C_r$ and $\delta_{r,K}^{\text{(IAF)}} \leq C_r$ for some $C_r \geq 0$. Combining this with Lemma 1, yields the following bound-level comparison.

$$UB(|\delta_p|)^{(P-IAF)} \le UB(|\delta_p|)^{(IAF)}.$$

Thus, P-IAF attains a tighter bound on $|\delta_p|$ than the IAF bound.

Remarks (i) Swap-guided base regularization to encoder output reduces $\|\mu + \mu_{\text{swap}}\|$ and $\|\ell - \ell_{\text{swap}}\|$, thereby directly decreasing the expected base mismatch in Eq.(20). (ii) P-IAF's swap-guided split $(c_d \to \mu_k, c_s \to \sigma_k)$ eliminates leak terms in a shared context, tightening the swap-probability error bound.

C DETAILS OF ADAPTIVE LATENT CONDITIONING

We detail the adaptive latent conditioning applied in the decoder. As illustrated in Fig. 7a, the user-latent z is mapped to a scale γ and shift β that modulate the incoming tokenized prompt-response embedding e in a FiLM-style manner. Concretely, the decoder computes a latent conditioned representation by applying dimension-wise scaling and shifting to e using γ and β , respectively.

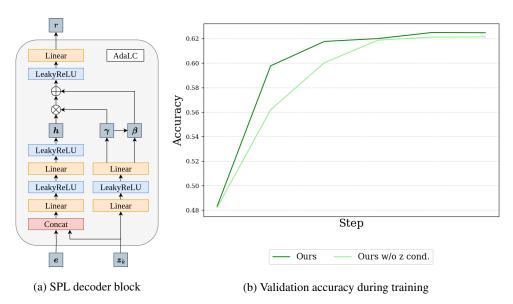


Figure 7: Details and training dynamics of the SPL decoder

Empirically, we also observe a training acceleration effect from adaptive latent conditioning. Fig. 7b reports preference-prediction accuracy evaluated periodically during training on *UF-P-4* with *Llama-3.1-8B*. The curve indicates that adaptive latent conditioning improves early-stage accuracy, suggesting that the reward model captures preferences more quickly with fewer samples. This is beneficial in data-scarce settings with minority preferences.

Adaptive latent conditioning benefits inference as well. When the user-latent encodes preference information clearly (i.e., is low-uncertainty), the decoder leverages it via the modulation to personalize the reward. When the latent is uncertain or uninformative, the decoder naturally reduces the effective contribution of z, behaving closer to the base model. This adaptability ensures robustness across users with different feedback levels and consistency, while allowing strong personalization when reliable signals are available.

D ADDITIONAL EXPERIMENTS

Effect of Components We analyze the impact of three components—(i) the base regularization \mathcal{L}_{guide} , (ii) P-IAF, and (iii) adaptive latent conditioning—via an ablation on *UF-P-4* with *Llama-3.1-8B*. Table 3 summarizes results.

Table 3: Effect of each component

$\mathcal{L}_{ ext{guide}}$	P-IAF	z cond.	Acc. [%]	Active Units [%]
\checkmark			57.18	0.00
	\checkmark		59.10	11.03
		\checkmark	56.95	0.00
\checkmark		\checkmark	56.87	0.00
\checkmark	\checkmark		62.14	94.24
	\checkmark	\checkmark	62.08	93.07
\checkmark	\checkmark	\checkmark	62.26	95.70

Our ablations show that neither the base regularization nor P-IAF alone is sufficient for preference encoding; their combination is the principal source of improvement. The base regularization helps align the base posterior with the intended swap-aware constraints, while P-IAF supplies the expressivity to maintain these constraints through the flow. Together they mitigate collapse and yield informative user-latents. Adaptive latent conditioning contributes only when the encoding retains signal; paired with P-IAF, it amplifies subtle preference cues. The full model—base regularization + P-IAF + adaptive latent conditioning—achieves the best overall performance.

Effect of P-IAF Depth We explore the effect of P-IAF depth K on SPL. Each step updates the entire latent vector, allowing P-IAF to model high-dimensional structures using fewer steps. Using this property, we limit the range to shallow stacks and evaluate $K \in \{1, 2, 4\}$. Table 4 indicates that K = 2 yields the best performance. Thus, K = 2 is our default for all experiments. A single step prevents collapse and improves accuracy. With K = 4, performance drops, suggesting unnecessary expressivity reduces preference-prediction accuracy, similar to standard IAF.

Table 4: Effect of P-IAF depth K

\overline{K}	Accuracy [%]	Active Units [%]
1	60.58	91.60
2	62.26	95.70
4	61.92	93.16

Preference Learning with fewer preference pairs We evaluate preference learning on the *UF-P-4* dataset when only a few preference pairs are provided to the model. Specifically, compared to the default setting (n=8), we randomly supply fewer pairs $(n\in\{2,3,4\})$ to *Llama-3.2-3B* and measure preference-prediction accuracy. As summarized in Table 5, SPL effectively encodes user preferences even under such limited preference signal. By contrast, VPL mitigates collapse with fewer pairs but captures user preferences poorly, resulting in accuracy similar to standard RLHF.

Table 5: Accuracy and active units with fewer preference pairs

Model	Method	Accuracy [%]	Active Units [%]
Llama-3.2-3B	BTL VPL SPL (Ours)	56.94 56.92 58.12	31.35 61.13

E IMPLEMENTATION DETAILS

972

973 974

975976

977

978979

989 990 991

1007

E.1 Hyperparameter settings

We detail the hyperparameters used in our experiments. Table 6 specifies the settings for generating the *Pets* and *UF-P* datasets. Table 7 specifies the training and evaluation hyperparameters. All experiments were run on a single NVIDIA RTX 4090 GPU and completed within two days.

Table 6: Hyperparameters for data generation

Hyperparameter	Value
Token embedding dimension	3072 (Llama-3.2-3B-instruct), 4096 (Llama-3.1-8B-instruct)
Max length	1024
Preference pairs per sample n	8
survey size for <i>UF-P</i>	16
Token data type	bfloat16
Training samples in dataset	4,000 (for <i>Pets</i>), 55,636 (for <i>UF-P-2</i>), 111,272 (for <i>UF-P-4</i>)
Evaluation samples in dataset	400 (for <i>Pets</i>), 6,042 (for <i>UF-P-2</i>), 12,084 (for <i>UF-P-4</i>)

Table 7: Hyperparameters for experiments

Hyperparameter	Value
Encoder input dimension	3072 (Llama-3.2-3B-instruct), 4096 (Llama-3.1-8B-instruct)
Latent dimension d	1024
Learning rate	1.0×10^{-4}
Learning rate scheduler	cosine with 3% warm-up steps
Epoch	2
P-IAF flow step K	2
Batch size	32 (for <i>Pets</i>), 64 (for <i>UF-P</i>)
Optimizer	AdamW(with weight decay = 0.001)
KL Divergence weight β	1.0×10^{-4} (for <i>Pets</i>), 3.0×10^{-6} (for <i>UF-P</i>)
KL annealing scheduler	cosine cyclical from 0 to 1 (period = $10,000$ steps)
Guidance balancing weight η	0.1
Active units threshold δ	0.005

E.2 ALGORITHMS

1008 Algorithm 1 Swap-guided Preference Learning (SPL) 1009 1010 **Require:** Preference Data $\mathbb{D} = \{\mathbb{D}_{h_1}, \cdots, \mathbb{D}_{h_m}\}$ 1011 **Require:** Encoder q_{ψ} , K-step P-IAF $F_{K_{\psi}}$, Reward Model r_{ϕ} , prior $p(\boldsymbol{z}_k)$ 1012 1: **while** not done **do** Sample $\mathbb{D}_h = \{x^i, y_w^i, y_l^i\}_{i=1}^n \sim \mathbb{D}$ 1013 2: Tokenize $e_{(\cdot)}^i = \text{LLM}^{\text{SFT}}(x^i, y_{(\cdot)}^i)$ 1014 3: 1015 4: Compute $\mu, \ell, c = q_{\psi}(\{e_w^i, e_l^i\}_{i=1}^n), \mu_{\text{swap}}, \ell_{\text{swap}}, c_{\text{swap}} = q_{\psi}(\{e_l^i, e_w^i\}_{i=1}^n)$ 1016 5: Sample $z_0 \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\ell})$ 1017 Compute $c_d = \frac{1}{2}(c - c_{\text{swap}}), c_s = \frac{1}{2}(c + c_{\text{swap}})$ 6: 1018 7: Compute $\mathbf{z}_K = F_{K_{ab}}(\mathbf{z}_0, \mathbf{c}_d, \mathbf{c}_s)$ 1019 8: Compute rewards: $r_w = r_\phi(\boldsymbol{e}_w, \boldsymbol{z}_K)$ and $r_l = r_\phi(\boldsymbol{e}_l, \boldsymbol{z}_K)$ 9: Compute reconstruction loss: $\mathcal{L}_{\text{recon}} = -\log(\sigma(r_w - r_l))$ 1020 1021 10: Compute KL-loss: $\mathcal{L}_{KL} = \beta \cdot D_{KL} (\log q_{\psi}(\boldsymbol{z}_k \mid \mathbb{D}_h) || p(\boldsymbol{z}_K))$ Compute guidance loss: $\mathcal{L}_{guide} = \lambda \cdot \left[\frac{1}{2} \left(1 + \cos(\mu, \mu_{swap}) \right) + \eta \cdot \frac{1}{2} \left(1 - \cos(\ell, \ell_{swap}) \right) \right]$ 11: 12: Compute total loss: $\mathcal{L}_{total} = \mathcal{L}_{recon} + \overline{\mathcal{L}}_{KL} + \mathcal{L}_{guide}$ 1023 Update E, q_{ψ} , $F_{K_{\psi}}$ and r_{ϕ} by optimizing $\mathcal{L}_{\text{total}}$ 13: 1024 14: end while 1025

F POTENTIAL AND SOCIAL EFFECTS

We address the tendency of standard RLHF to bias rewards toward majority preferences by encoding a user preference from a small number of comparisons and conditioning the policy on a user-latent, yielding a personalized policy $\pi_{\theta}(y \mid x, z_k)$. We post-train the policy with

$$\max_{\substack{\pi_{\theta} \\ h \sim \mathbb{H}}} \mathbb{E} \left[\underset{\substack{x \sim \mathbb{D} \\ y \sim \pi_{\theta}(y|x) \\ \mathbf{z}_{k} \sim q_{\psi}(\mathbf{z}_{k}|\mathbb{D}_{h})}} \mathbb{E} \left[r_{\phi}(x, y, \mathbf{z}_{k}) \right] - \beta D_{\mathrm{KL}} \left[\pi_{\theta}(y \mid x, \mathbf{z}_{k}) \mid\mid \pi^{\mathrm{SFT}}(y \mid x) \right] \right],$$
(26)

which trains and deploys distinct behaviors conditioned on z_k inferred from the user's own choices. This differs from approaches that keep a single global policy and simply change the input x via a user context: here, the conditional policy itself learns to act differently under different latents, rather than relying on prompt-only adaptation (Dong et al., 2022).

The scheme in Eq.(26) naturally extends to implicit-reward objectives such as Direct Preference Optimization (DPO) (Rafailov et al., 2023) by conditioning the policy and implicit reward surrogate on z_k .

Plus, conditioning policy on user-latent is not limited to LLMs: Our swap-guided encoding and adaptive latent conditioning can be used when preferences are difficult to summarize, including in generative models or control settings (Poddar et al., 2024; Wang et al., 2025; Ng et al., 2025).

G NOTATIONS

Table 8: Notations

Notation	Meaning	Notation	Meaning	
Indices & counts		Embeddings, latents & contexts		
i n N j k K d	preference-pair index number of preference-pair in sample total number of preference-pair dimension index flow step total flow step latent dimension	$egin{array}{c} e_w \ e_l \ z \ \mu \ \sigma \ \ell \ \epsilon \ c \ c_d \ c_s \end{array}$	embedding of y_w embedding of y_l latent embedding mean standard deviation log-variance random noise shared context swap-reversal context swap-invariant context	
	Users & sets		Models & functions	
$egin{array}{c} h \ \mathbb{H} \ \mathbb{D} \ \mathbb{D}_h \ p \ \mathbb{P} \ \end{array}$	a user (annotator) user population full preference dataset user h's preference dataset a preference type set of preference types	$q(\cdot)$ $r(\cdot)$ $p(\cdot)$ $f(\cdot)$ $\mu_k(\cdot)$ $\sigma_k(\cdot)$	encoder / variational posterior decoder / reward function preference probability autoregressive transform shift function at step k scale function at step k	
	Prompt & response		Parameters & weights	
$ \begin{array}{c} x \\ y \\ y_w \\ y_l \end{array} $	prompt response chosen (winning) response rejected (losing) response	$egin{array}{c} \psi \ \phi \ eta \ \lambda \end{array}$	learnable params (encoder & flow) learnable params (decoder) KL-divergence weight guidance loss weight	

Norms Throughout, $\|\cdot\|$ denotes the Euclidean norm for vectors and the Frobenius norm for matrices. We use $\|\cdot\|_{\infty}$ for the entrywise max norm when needed.

THE USE OF LARGE LANGUAGE MODELS

We employed an LLM-assisted search to identify prior work on posterior collapse in VAEs and user-representation policies across domains. All retrieved items were manually reviewed by the authors to confirm their relevance before citation.