# SoPo: Text-to-Motion Generation Using Semi-Online Preference Optimization

3

Supplementary Material



Figure S1: Visual results on HumanML3D dataset. We integrate our SoPo into MDM [1] and MLD [2], respectively. Our SoPo improves the alignment between text and motion preferences. Here, the red text denotes descriptions inconsistent with the generated motion.

- This supplementary document contains the technical proofs of results and some additional experimental results. It is structured as follows. Sec. A provides the implementation and theoretical
  analysis of our SoPo. Sec. B gives the proofs of the main results, including Theorem 1, Theorem
- analysis of our SoPo. Sec. B gives the proofs of the main results, including Theorem 1, Theorem
   the objective function of DSoPo, the objective function of USoPo, and theorem of SoPo for
- text-to-motion generation. Then in Sec. C presents the additional experiment information, including
- additional experimental details (Sec. C.1 and C.2) and results (Sec. C.3).

# **10** A Details of SoPo for Text-to-Motion Generation

In this section, we first examine the objective function of SoPo and argue that it presents significant
challenges for optimization. Fortunately, we then discover and derive an equivalent form that is easier
to optimize (Sec. A.1). Finally, we design an algorithm to optimize it and finish discussing their
correspondence (Sec. A.2).

## 15 A.1 Equivalent form of SoPo

<sup>16</sup> In Eq. (15) and (16), the objective function of SoPo is defined as:  

$$\mathcal{L}_{SoPo}^{diff} = \mathcal{L}_{SoPo-vu}^{diff} + \mathcal{L}_{SoPo-bu}^{diff},$$

18

$$\mathcal{L}_{\text{SoPo-vu}}^{\text{diff}} = -\mathbb{E}_{t \sim \mathcal{U}(0,T), (x^{w},c) \sim \mathcal{D}, x_{\pi_{\theta}}^{1:K} \sim \pi_{\theta}^{vu*}(\cdot|c)} Z_{vu}(c) \\ \left[ \log \sigma \Big( -T\omega_{t} \big( \beta_{w}(x_{w}) \big( \mathcal{L}(\theta, \text{ref}, x_{t}^{w}) - \beta \mathcal{L}(\theta, \text{ref}, x_{t}^{l}) \big) \big) \Big] \right],$$

$$\mathcal{L}_{\text{SoPo-hu}}^{\text{diff}} = -\mathbb{E}_{t \sim \mathcal{U}(0,T), (x^{w},c) \sim \mathcal{D}} Z_{hu}(c) \left[ \log \sigma \Big( -T\omega_{t} \beta_{w}(x_{w}) \mathcal{L}(\theta, \text{ref}, x_{t}^{w}) \big) \right],$$
(S2)

(S1)

 $\mathcal{L}_{\text{SoPo-hu}} = -\mathbb{E}_{t \sim \mathcal{U}(0,T),(x^w,c) \sim \mathcal{D}} \mathcal{D}_{hu}(c) \left[ \log o \left( -T \omega_t \mathcal{D}_w(x_w) \mathcal{L}(0, \operatorname{ref}, x_t) \right) \right],$ However, these objectives can not be directly optimized, since the distribution  $\bar{\pi}_{\theta}^{vu*}$  and  $\bar{\pi}_{\theta}^{hu*}$  are not

<sup>19</sup> defined explicitly. To this end, we begin by inducing its equivalent form:  $C^{\text{diff}}(0) = \mathbb{R}$ 

$$\mathcal{L}_{\text{SoPo}}^{\text{soPo}}(\theta) = -\mathbb{E}_{t\sim\mathcal{U}(0,T),(x^{w},c)\sim\mathcal{D},x_{\bar{\pi}\theta}^{1:K}\sim\bar{\pi}_{\theta}(\cdot|c)} \\ \begin{cases} \log\sigma\Big(-T\omega_{t}\big(\beta_{w}(x_{w})\big(\mathcal{L}(\theta,\operatorname{ref},x_{t}^{w})-\beta\mathcal{L}(\theta,\operatorname{ref},x_{t}^{l})\big)\big)\Big), & \text{if } r(x^{l},c)<\tau, \\ \log\sigma\Big(-T\omega_{t}\beta_{w}(x_{w})\mathcal{L}(\theta,\operatorname{ref},x_{t}^{w})\Big), & \text{otherwise.} \end{cases} \end{cases}$$

$$(S3)$$

20 where  $x^{l} = \operatorname{argmin}_{\{x_{\pi_{\theta}}^{k}\}_{k=1}^{K} \sim \pi_{\theta}} r(x_{\pi_{\theta}}^{k}, c).$ 

*Proof.* Recall our definition of  $\mathcal{L}_{SoPo}^{diff}(\theta)$  in Eq. (15) and (16). Through algebraic maneuvers, we have:

$$\mathcal{L}_{SoPo}^{diff} = \mathcal{L}_{soPo-vu}^{diff} + \mathcal{L}_{SoPo-vu}^{diff} \sim \pi_{\theta}^{vus}(\cdot|c) Z_{vu}(c) \left[ \log \sigma \left( -T\omega_{t}(\beta_{w}(x_{w})(\mathcal{L}(\theta, \operatorname{ref}, x_{t}^{w}) - \beta\mathcal{L}(\theta, \operatorname{ref}, x_{t}^{l}))) \right) \right] - \mathbb{E}_{t\sim\mathcal{U}(0,T),(x^{w},c)\sim\mathcal{D}} Z_{hu}(c) \left[ \log \sigma \left( -T\omega_{t}\beta_{w}(x_{w})\mathcal{L}(\theta, \operatorname{ref}, x_{t}^{w}) \right) \right] = -\mathbb{E}_{t\sim\mathcal{U}(0,T),(x^{w},c)\sim\mathcal{D}} \mathbb{E}_{x_{\pi_{\theta}}^{1:K} \sim \pi_{\theta}^{vus}(\cdot|c)} Z_{vu}(c) \left[ \log \sigma \left( -T\omega_{t}(\beta_{w}(x_{w})(\mathcal{L}(\theta, \operatorname{ref}, x_{t}^{w}) - \beta\mathcal{L}(\theta, \operatorname{ref}, x_{t}^{l}))) \right) \right] - \mathbb{E}_{t\sim\mathcal{U}(0,T),(x^{w},c)\sim\mathcal{D}} \mathbb{E}_{x_{\pi_{\theta}}^{1:K} \sim \pi_{\theta}^{uus}(\cdot|c)} Z_{hu}(c) \left[ \log \sigma \left( -T\omega_{t}\beta_{w}(x_{w})\mathcal{L}(\theta, \operatorname{ref}, x_{t}^{w}) \right) \right] = -\mathbb{E}_{t\sim\mathcal{U}(0,T),(x^{w},c)\sim\mathcal{D}} \mathbb{E}_{x^{1:K} \sim \pi_{\theta}^{hus}(\cdot|c)} Z_{hu}(c) \left[ \log \sigma \left( -T\omega_{t}\beta_{w}(x_{w})\mathcal{L}(\theta, \operatorname{ref}, x_{t}^{w}) \right) \right] = -\mathbb{E}_{t\sim\mathcal{U}(0,T),(x^{w},c)\sim\mathcal{D}} \mathbb{E}_{x^{1:K} \sim \pi_{\theta}^{hus}(x^{1:K}|c)} \left[ \log \sigma \left( -T\omega_{t}\beta_{w}(x_{w})\mathcal{L}(\theta, \operatorname{ref}, x_{t}^{w}) \right) \right] = -\mathbb{E}_{t\sim\mathcal{U}(0,T),(x^{w},c)\sim\mathcal{D}} \mathbb{E}_{x^{1:K} \sim \pi_{\theta}^{hus}(x^{1:K}|c)} \left[ \log \sigma \left( -T\omega_{t}\beta_{w}(x_{w})\mathcal{L}(\theta, \operatorname{ref}, x_{t}^{w}) \right) \right] = -\mathbb{E}_{t\sim\mathcal{U}(0,T),(x^{w},c)\sim\mathcal{D}} \mathbb{E}_{x_{\pi_{\theta}}^{1:K} \sim \pi_{\theta}(\cdot|c)} \mathcal{P}_{\tau}(r(x_{\pi_{\theta}}^{l}, c) < \tau) \left[ \log \sigma \left( -T\omega_{t} \left( \beta_{w}(x_{w}) \left( \mathcal{L}(\theta, \operatorname{ref}, x_{t}^{w}) - \beta \mathcal{L}(\theta, \operatorname{ref}, x_{t}^{l}) \right) \right) \right] - \mathbb{E}_{t\sim\mathcal{U}(0,T),(x^{w},c)\sim\mathcal{D}} \mathbb{E}_{x_{\pi_{\theta}}^{1:K} \sim \pi_{\theta}(\cdot|c)} \mathcal{P}_{\tau}(r(x_{\pi_{\theta}}^{l}, c) \geq \tau) \left[ \log \sigma \left( -T\omega_{t}\beta_{w}(x_{w})\mathcal{L}(\theta, \operatorname{ref}, x_{t}^{w}) \right) \right] - \mathbb{E}_{t\sim\mathcal{U}(0,T),(x^{w},c)\sim\mathcal{D}} \mathbb{E}_{x_{\pi_{\theta}}^{1:K} \sim \pi_{\theta}(\cdot|c)} \mathcal{P}_{\tau}(r(x_{\pi_{\theta}}^{l}, c) \geq \tau) \left[ \log \sigma \left( -T\omega_{t}\beta_{w}(x_{w})\mathcal{L}(\theta, \operatorname{ref}, x_{t}^{w}) \right) \right] = -\mathbb{E}_{t\sim\mathcal{U}(0,T),(x^{w},c)\sim\mathcal{D}} \mathbb{E}_{x_{\pi_{\theta}}^{1:K} \sim \pi_{\theta}(\cdot|c)} \\ \left\{ \log \sigma \left( -T\omega_{t} \left( \beta_{w}(x_{w}) \mathcal{L}(\theta, \operatorname{ref}, x_{t}^{w}) - \beta \mathcal{L}(\theta, \operatorname{ref}, x_{t}^{l}) \right) \right\}, \quad \text{if } r(x^{l}, c) < \tau, \\ \log \sigma \left( -T\omega_{t} \beta_{w}(x_{w})\mathcal{L}(\theta, \operatorname{ref}, x_{t}^{w}) \right), \quad \text{otherwise.}$$

where  $\stackrel{(\square)}{=}$  holds since  $p_{\bar{\pi}_{\theta}^{vu*}}(\cdot) = \frac{p_{\bar{\pi}_{\theta}}^{vu*}(\cdot)}{Z_{vu}(c)}$  and  $p_{\bar{\pi}_{\theta}}^{vu}(x_{\bar{\pi}_{\theta}}^{1:K}|c) = p_{\bar{\pi}_{\theta}}(x_{\bar{\pi}_{\theta}}^{1:K}|c) \cdot p_{\tau}(r(x_{\bar{\pi}_{\theta}}^{l},c) \geq \tau)$ . The proof is completed.

#### 25 A.2 The process of SoPo for text-to-motion generation

Based on the equivalent form of SoPo in Eq. (S3), we can design an algorithm to directly optimize it, as shown in **Algorithm 1**.

## Algorithm 1 SoPo for text-to-motion generation

**Input:** Preference dataset  $\mathcal{D}$ ; diffusion steps T; iterations I; samples K; ref model  $\pi_{ref}$ ; policy  $\pi_{\theta}$ ; threshold  $\tau$ **Output:** Aligned model  $\pi_{\theta}$ 1: **for** i = 1 to *I* **do** for each  $(x^w, c) \in \mathcal{D}$  do 2: Sample  $t \sim \mathcal{U}(0,T)$ Sample  $x_{\overline{\pi}_{\theta}}^{1:K} \sim \overline{\pi}_{\theta}(\cdot|c)$ Compute  $S(x^w) = \min_k \cos(x^w, x_{\overline{\pi}_{\theta}}^k)$ 3: 4: 5:  $x^{l} = \arg\min_{k} r(x_{\pi_{\theta}}^{k}, c)$ 6: if  $r(x^l, c) < \tau$  then 7:  $\mathcal{L} = \log \sigma(-T\omega_t \beta_w(x^w)(\mathcal{L}(\theta, \mathrm{ref}, x_t^w) - \beta \mathcal{L}(\theta, \mathrm{ref}, x_t^l)))$ 8: 9: else  $\mathcal{L} = \log \sigma(-T\omega_t \beta_w(x^w) \mathcal{L}(\theta, \operatorname{ref}, x_t^w))$ 10: 11: end if Accumulate loss:  $\mathcal{L}_{SoPo}^{diff} = \mathcal{L}_{SoPo}^{diff} + \mathcal{L}$ 12: 13: end for Update  $\pi_{\theta}$  using  $\nabla_{\theta} \mathcal{L}_{SoPo}^{diff}$ 14: 15: end for 16: return  $\pi_{\theta}$ 

The SoPo optimizes a policy model  $\pi_{\theta}$  for text-to-motion generation through an iterative process 28 guided by a reward model. In each iteration, given a preferred motion  $x^w$  and a conditional code c, 29 a random diffusion step t is selected, and K candidate motions are generated by  $\pi_{\theta}$ . The motion 30 with the lowest preference score is then treated as the unpreferred motion. To determine the weight 31 of the preferred motion  $x^w$ , the similarities between all generated motions are computed, and the 32 lowest cosine similarity value is used to calculate its weight. Finally, the loss is calculated in two 33 ways, determined based on the preference scores of the unpreferred motion. If the preference score 34 of the selected unpreferred motion falls below a threshold  $\tau$ , it is identified as a valuable unpreferred 35 motion and used for training. Otherwise, it indicates that the motions generated by the policy model 36  $\pi_{\theta}$  are satisfactory. In such cases, the policy model is trained exclusively on high-quality preferred 37 motions, rather than on both preferred motions and relatively high-preference unpreferred motions. 38

<sup>39</sup> To further understand the objective function, we analyze the correspondence between the objective <sup>40</sup> function in Eq. (S3) and Algorithm 1:

# 41 **B** Theories

#### 42 B.1 Proof of Theorem 1

<sup>43</sup> *Proof.* The offline DPO based on Plackett-Luce model [3] can be denoted as:

$$\mathcal{L}_{\text{off}}(\theta) = -\mathbb{E}_{(x^{1:K},c)\sim\mathcal{D}}\Big[\log\prod_{k=1}^{K}\frac{\exp(\beta h_{\theta}(x^{k},c))}{\sum_{j=k}^{K}\exp(\beta h_{\theta}(x^{j},c))}\Big],\tag{S5}$$

44 where  $h_{ heta}(x,c) = \log rac{\pi_{ heta}(x|c)}{\pi_{ ext{ref}}(x|c)}$ . Then we have:

$$\begin{aligned} \mathcal{L}_{\text{off}}(\theta) &= -\mathbb{E}_{(x^{1:K},c)\sim\mathcal{D}} \bigg[ \log \prod_{k=1}^{K} \frac{\exp(\beta h_{\theta}(x^{k},c))}{\sum_{j=k}^{K} \exp(\beta h_{\theta}(x^{j},c))} \bigg] \\ &= -\mathbb{E}_{c\sim\mathcal{D},x^{1:K}} p_{\text{gt}}(x^{1:K}|c) \bigg[ \log \prod_{k=1}^{K} \frac{\exp(\beta h_{\theta}(x^{k},c))}{\sum_{j=k}^{K} \exp(\beta h_{\theta}(x^{j},c))} \bigg] \\ &= -\mathbb{E}_{c\sim\mathcal{D},x^{1:K}} p_{\text{gt}}(x^{1:K}|c) \bigg[ \log \prod_{k=1}^{K} \frac{\exp(\beta \log \frac{\pi_{\theta}(x^{k}|c)}{\pi_{\text{ref}}(x^{k}|c)})}{\sum_{j=k}^{K} \exp(\beta \log \frac{\pi_{\theta}(x^{j}|c)}{\pi_{\text{ref}}(x^{k}|c)})} \bigg] \\ &= -\mathbb{E}_{c\sim\mathcal{D},x^{1:K}} p_{\text{gt}}(x^{1:K}|c) \bigg[ \log \prod_{k=1}^{K} \frac{\exp\log(\frac{\pi_{\theta}(x^{k}|c)}{\pi_{\text{ref}}(x^{k}|c)})^{\beta}}{\sum_{j=k}^{K} \exp\log(\frac{\pi_{\theta}(x^{j}|c)}{\pi_{\text{ref}}(x^{k}|c)})^{\beta}} \bigg] \\ &= -\mathbb{E}_{c\sim\mathcal{D},x^{1:K}} p_{\text{gt}}(x^{1:K}|c) \bigg[ \log \prod_{k=1}^{K} \frac{(\frac{\pi_{\theta}(x^{k}|c)}{\pi_{\text{ref}}(x^{k}|c)})^{\beta}}{\sum_{p_{\theta}(x^{k}|c)}} \bigg] \\ &= -\mathbb{E}_{c\sim\mathcal{D},x^{1:K}} p_{\text{gt}}(x^{1:K}|c) \bigg[ \log \prod_{k=1}^{K} p_{\theta}(x^{k}|c) \bigg] \\ &= -\mathbb{E}_{c\sim\mathcal{D},x^{1:K}} p_{\text{gt}}(x^{1:K}|c) \bigg[ \log p_{\theta}(x^{1:K}|c) - \log p_{\text{gt}}(x^{1:K}|c) + \log p_{\text{gt}}(x^{1:K}|c) \bigg] \\ &= -\mathbb{E}_{c\sim\mathcal{D},x^{1:K}} p_{\text{gt}}(x^{1:K}|c) \bigg[ \log p_{\theta}(x^{1:K}|c) - \log p_{\text{gt}}(x^{1:K}|c) + \log p_{\text{gt}}(x^{1:K}|c) \bigg] \\ &= \mathbb{E}_{c\sim\mathcal{D},x^{1:K}} p_{\text{gt}}(x^{1:K}|c) \bigg[ \log p_{\theta}(x^{1:K}|c) - \log p_{\text{gt}}(x^{1:K}|c) \bigg] \\ &= \mathbb{E}_{c\sim\mathcal{D},x^{1:K}} p_{\text{gt}}(x^{1:K}|c) \bigg[ \log p_{\theta}(x^{1:K}|c) - \log p_{\text{gt}}(x^{1:K}|c) \bigg] \\ &= \mathbb{E}_{c\sim\mathcal{D},x^{1:K}} p_{\text{gt}}(x^{1:K}|c) \bigg[ \log p_{\theta}(x^{1:K}|c) - \log p_{\text{gt}}(x^{1:K}|c) \bigg] \end{aligned}$$

45 Therefore, we have:

$$\nabla_{\theta} \mathcal{L}_{\text{off}}(\theta) = \mathbb{E}_{c \sim \mathcal{D}, x^{1:K}} \nabla_{\theta} D_{KL}(p_{\text{gt}} || p_{\theta}).$$
(S7)

<sup>46</sup> The proof is completed.

# 47 **B.2** Proof of Theorem 2

*Proof.* Inspired by [4], we replace the one-hot vector in DPO with Plackett-Luce model [3], and then
 the online DPO can be expressed as

$$\mathcal{L}_{\text{DPO-On}}(\theta) = -\mathbb{E}_{c \sim \mathcal{D}, x^{1:K} \sim \bar{\pi}_{\theta}(\cdot|c)} \Big[ \sum_{k=1}^{K} p_r(x_k|c) \log \frac{\left(\frac{\pi_{\theta}(x^k|c)}{\pi_{\text{ref}}(x^k|c)}\right)^{\beta}}{\sum_{j=k}^{K} \left(\frac{\pi_{\theta}(x^j|c)}{\pi_{\text{ref}}(x^j|c)}\right)^{\beta}} \Big],$$
(S8)

so where  $p_r(x_{\pi_{\theta}}^k|c) = \frac{\exp r(x_{\pi_{\theta}}^k,c)}{\sum_{i=k}^K \exp r(x_{\pi_{\theta}}^i,c)}$ . Then we have:

$$\mathcal{L}_{\mathrm{on}}(\theta) = -\mathbb{E}_{c\sim\mathcal{D},x^{1:K}\sim\bar{\pi}_{\theta}(\cdot|c)} \left[ \sum_{k=1}^{K} p_{r}(x_{k}|c) \log \frac{\left(\frac{\pi_{\theta}(x^{k}|c)}{\pi_{\mathrm{ref}}(x^{k}|c)}\right)^{\beta}}{\sum_{j=k}^{K} \left(\frac{\pi_{\theta}(x^{j}|c)}{\pi_{\mathrm{ref}}(x^{k}|c)}\right)^{\beta}} \right]$$

$$= -\mathbb{E}_{c\sim\mathcal{D}} p_{\bar{\pi}_{\theta}}(x^{1:K}|c) \left[ \sum_{k=1}^{K} p_{r}(x_{k}|c) \log \frac{\left(\frac{\pi_{\theta}(x^{k}|c)}{\pi_{\mathrm{ref}}(x^{k}|c)}\right)^{\beta}}{\sum_{j=k}^{K} \left(\frac{\pi_{\theta}(x^{j}|c)}{\pi_{\mathrm{ref}}(x^{k}|c)}\right)^{\beta}} \right]$$

$$= -\mathbb{E}_{c\sim\mathcal{D}} p_{\bar{\pi}_{\theta}}(x^{1:K}|c) \left[ \sum_{k=1}^{K} p_{r}(x^{k}|c) \log \frac{\left(\frac{\pi_{\theta}(x^{k}|c)}{\pi_{\mathrm{ref}}(x^{k}|c)}\right)^{\beta}}{\sum_{j=k}^{K} \left(\frac{\pi_{\theta}(x^{j}|c)}{\pi_{\mathrm{ref}}(x^{j}|c)}\right)^{\beta}}} \right]$$

$$= -\mathbb{E}_{c\sim\mathcal{D}} p_{\bar{\pi}_{\theta}}(x^{1:K}|c) \left[ \sum_{k=1}^{K} p_{r}(x^{k}|c) \log p_{\theta}(x^{k}|c) \right]$$
(S9)

$$= -\mathbb{E}_{c \sim \mathcal{D}} p_{\bar{\pi}_{\theta}}(x^{1:K}|c) \Big[ \sum_{k=1}^{K} p_{r}(x^{k}|c) (\log p_{\theta}(x^{k}|c) - \log p_{r}(x^{k}|c) + \log p_{r}(x^{k}|c)) \Big]$$
$$= \mathbb{E}_{c \sim \mathcal{D}} p_{\bar{\pi}_{\theta}}(x^{1:K}|c) \Big[ D_{KL}(p_{r}|p_{\theta}) - p_{r}(x^{k}|c) \log p_{r}(x^{k}|c) \Big]$$

51 Therefore, we have:

$$\nabla_{\theta} \mathcal{L}_{\text{on}}(\theta) = \mathbb{E}_{c \sim \mathcal{D}} \nabla_{\theta} \ p_{\bar{\pi}_{\theta}}(x^{1:K}|c) D_{KL}(p_r||p_{\theta}).$$
(S10)

- 52 The proof is completed.
- Given a sample x with a tiny generative probability  $p_{\pi_{\theta}|c}(x) \to 0$ , and large reward value  $r(x,c) \to 1$ , we have  $\lim_{p_{\pi_{\theta}}(x|c)\to 0, r(x,c)\to 1} \nabla_{\theta} \mathcal{L}_{on} = \mathbf{0}$ .
- 55 *Proof.* Since x is contained in the sampled motion group  $x^{1:K}$ , we have:

$$\lim_{p_{\pi_{\theta}}(x|c)\to 0, r(x,c)\to 1} \nabla_{\theta} \mathcal{L}_{on}$$

$$= \lim_{p_{\pi_{\theta}}(x|c)\to 0, r(x,c)\to 1} \nabla_{\theta} p_{\bar{\pi}_{\theta}}(x^{1:K}|c) D_{KL}(p_{r}||p_{\theta})$$

$$\bigoplus_{p_{\pi_{\theta}}(x^{1:K}|c)\to 0, r(x,c)\to 1} \nabla_{\theta} p_{\bar{\pi}_{\theta}}(x^{1:K}|c) D_{KL}(p_{r}||p_{\theta})$$

$$= \mathbf{0},$$
(S11)

where ① holds since  $p_{\pi_{\theta}}(x^{1:K}|c) = p_{\pi_{\theta}}(x|c)p_{\pi_{\theta}}(x^{M}|c) \le p_{\pi_{\theta}}(x|c)$ , and  $x^{M}$  denotes a motion group obtained by removing the given motion x from the group  $x^{1:K}$ , i.e., satisfying that  $x^{M} = x^{1:K} - \{x\}$ . The proof is completed.

# 59 B.3 Proof of DSoPo

Proof. Eq. (10) suggests that DSoPo samples multiple unpreferred motion candidates instead of a
 single unpreferred motion. Thus, we should first extend Eq. (9) as:

$$\mathcal{L}_{\mathrm{DSoPo}}(\theta) = -\mathbb{E}_{(x^{w},c)\sim\mathcal{D}}\mathbb{E}_{x^{1:K}\sim\bar{\pi}_{\theta}(x|c)}\log\sigma\Big(\beta\mathcal{H}_{\theta}(x^{w},x^{l},c)\Big),\tag{S12}$$

62 where  $x^l = \operatorname{argmin}_{\{x_{\pi_{\theta}}^k\}_{k=1}^K \sim \pi_{\theta}} r(x_{\pi_{\theta}}^k, c)$ . Then, we have:

$$\mathcal{L}_{\text{DSoPo}}(\theta) = - \mathbb{E}_{(x^{w},c)\sim\mathcal{D}} \mathbb{E}_{x^{1:K}\sim\pi_{\theta}(x^{|c|}c)} \log \sigma \left(\beta \mathcal{H}_{\theta}(x^{w},x^{l},c)\right)$$

$$= - \mathbb{E}_{(x^{w},c)\sim\mathcal{D}} \mathbb{E}_{x^{1:K}} \underbrace{p_{\pi_{\theta}}(x^{1:K}|c)}_{\text{Substituting with (11)}} \log \sigma \left(\beta \mathcal{H}_{\theta}(x^{w},x^{l},c)\right)$$

$$= - \mathbb{E}_{(x^{w},c)\sim\mathcal{D}} \mathbb{E}_{x^{1:K}} \left(p_{\pi_{\theta}}(x^{\pi_{\theta}^{1:K}}|c)p_{\tau}(r(x^{l},c)\geq\tau) + p_{\pi_{\theta}}(x^{\pi_{\theta}^{1:K}}|c)p_{\tau}(r(x^{l},c)<\tau)\right) \log \sigma \left(\beta \mathcal{H}_{\theta}(x^{w},x^{l},c)\right)$$

$$= - \mathbb{E}_{(x^{w},c)\sim\mathcal{D}} \mathbb{E}_{x^{1:K}} \underbrace{p_{\pi_{\theta}}(x^{\pi_{\theta}^{1:K}}|c)p_{\tau}(r(x^{l},c)\geq\tau)}_{p_{\pi_{\theta}^{\mu}(x^{1:K}|c)}^{h_{\pi_{\theta}}(x^{1:K}|c)}} \log \sigma \left(\beta \mathcal{H}_{\theta}(x^{w},x^{l},c)\right)$$

$$= - \mathbb{E}_{(x^{w},c)\sim\mathcal{D}} \mathbb{E}_{x^{1:K}} \underbrace{p_{\pi_{\theta}}(x^{\pi_{\theta}^{1:K}}|c)p_{\tau}(r(x^{l},c)<\tau)}_{p_{\pi_{\theta}^{w}(x^{1:K}|c)}^{h_{\theta}(x^{w},x^{l},c)} \log \sigma \left(\beta \mathcal{H}_{\theta}(x^{w},x^{l},c)\right)$$

$$= - \mathbb{E}_{(x^{w},c)\sim\mathcal{D}} \mathbb{E}_{x^{1:K}} Z_{hu}(c) p_{\pi_{\theta}^{w}}^{h_{\theta}(x^{1:K}|c)} \log \sigma \left(\beta \mathcal{H}_{\theta}(x^{w},x^{l},c)\right)$$

$$= - \mathbb{E}_{(x^{w},c)\sim\mathcal{D}} \mathbb{E}_{x^{1:K}} Z_{hu}(c) p_{\pi_{\theta}^{w}}^{h_{\theta}(x^{1:K}|c)} \log \sigma \left(\beta \mathcal{H}_{\theta}(x^{w},x^{l},c)\right)$$

$$= - \mathbb{E}_{(x^{w},c)\sim\mathcal{D}} \mathbb{E}_{x^{1:K}} Z_{hu}(c) p_{\pi_{\theta}^{w}}^{h_{\theta}(x^{1:K}|c)} \log \sigma \left(\beta \mathcal{H}_{\theta}(x^{w},x^{l},c)\right)$$

$$= - \mathbb{E}_{(x^{w},c)\sim\mathcal{D}} \mathbb{E}_{x^{1:K}} \mathbb{E}_{\pi_{\theta}^{w}(x^{1:K}|c)} \log \sigma \left(\beta \mathcal{H}_{\theta}(x^{w},x^{l},c)\right)$$

$$= - \mathbb{E}_{(x^{w},c)\sim\mathcal{D}} \mathbb{Z}_{hu}(c) \mathbb{E}_{x^{1:K}} p_{\pi_{\theta}^{w}}^{h_{\theta}(x^{1:K}|c)} \log \sigma \left(\beta \mathcal{H}_{\theta}(x^{w},x^{l},c)\right)$$

$$= - \mathbb{E}_{(x^{w},c)\sim\mathcal{D}} \mathbb{Z}_{hu}(c) \mathbb{E}_{x^{1:K} \sim \pi_{\theta}^{w}} \log \sigma \left(\beta \mathcal{H}_{\theta}(x^{w},x^{l},c)\right)$$

$$= - \mathbb{E}_{(x^{w},c)\sim\mathcal{D}} \mathbb{Z}_{hu}(c) \mathbb{E}_{x^{1:K} \sim \pi_{\theta}^{w}} \log \sigma \left(\beta \mathcal{H}_{\theta}(x^{w},x^{l},c)\right)$$

$$= \mathcal{L}_{vu}(\theta) + \mathcal{L}_{hu}(\theta),$$

$$(S13)$$

where  $p_{\bar{\pi}_{\theta}^{vu*}}(\cdot) = \frac{p_{\pi_{\theta}}^{vu}(\cdot)}{Z_{vu}(c)}$  and  $p_{\bar{\pi}_{\theta}}^{hu*}(\cdot) = \frac{p_{\pi_{\theta}}^{hu}(\cdot)}{Z_{hu}(c)}$  respectively denote the distributions of valuable unpreferred and high-preference unpreferred motions. The proof is completed.

Accordingly, we rewrite  $\mathcal{L}_{hu}(\theta)$  and obtain the objective function of USoPo:

$$\mathcal{L}_{\text{USoPo-hu}}(\theta) = -\mathbb{E}_{(x^w,c)\sim\mathcal{D}}Z_{hu}(c)\log\sigma\Big(\beta h_{\theta}(x^w,c)\Big),$$
  
$$\mathcal{L}_{\text{USoPo}}(\theta) = \mathcal{L}_{\text{USoPo-hu}}(\theta) + \mathcal{L}_{\text{vu}}(\theta).$$
(S14)

**Implementation** Now, we discuss how to deal with the computation of  $Z_{vu}(c)$  and  $Z_{hu}(c)$  in our implementation. As discussed in Sec. A, directly optimizing the objective function  $\mathcal{L}_{SoPo}^{diff}(\theta)$  is challenging, and we used **Algorithm 1** optimized its equivalent form:

$$\mathcal{L}_{\text{SoPo}}^{\text{diff}}(\theta) = -\mathbb{E}_{t \sim \mathcal{U}(0,T), (x^w,c) \sim \mathcal{D}, x_{\pi_{\theta}}^{1:K} \sim \bar{\pi}_{\theta}(\cdot|c)} \\ \begin{cases} \log \sigma \Big( -T\omega_t \big( \beta_w(x_w) \big( \mathcal{L}(\theta, \text{ref}, x_t^w) - \beta \mathcal{L}(\theta, \text{ref}, x_t^l) \big) \big) \Big), & \text{if } r(x^l,c) < \tau, \ \text{(S15)} \\ \log \sigma \Big( -T\omega_t \beta_w(x_w) \mathcal{L}(\theta, \text{ref}, x_t^w) \Big), & \text{otherwise.} \end{cases}$$

Similarly, we can optimize the equivalent form of UDoPo to avoid the computation of  $Z_{vu}(c)$  and  $Z_{hu}(c)$ :

$$\mathcal{L}_{\text{USoPo}}(\theta) = -\mathbb{E}_{(x^w,c)\sim\mathcal{D}, x^{1:K}_{\bar{\pi}_{\theta}}\sim\bar{\pi}_{\theta}(\cdot|c)} \begin{cases} \log\sigma\Big(\beta\mathcal{H}_{\theta}(x^w,x^l,c)\Big), & \text{If } r(x^l,c)<\tau, \\ \log\sigma\Big(\beta h_{\theta}(x^w,c)\Big), & \text{Otherwise.} \end{cases}$$
(S16)

The proof of Eq. (S16) follows the same steps as the proof of Eq. (S15) in Sec. A.

# 72 B.4 Discussion of USoPo and DSoPo

<sup>73</sup> In this section, we discuss the relationship between USoPo and DSoPo and the difference between

their optimization. Here, USoPo and DSoPo are defined as:

$$\mathcal{L}_{\text{USoPo}}(\theta) = -\mathbb{E}_{(x^w,c)\sim\mathcal{D}} Z_{hu}(c) \log \sigma \Big(\beta h_\theta(x^w,c)\Big) + \mathcal{L}_{\text{vu}}(\theta).$$
(S17)

$$\mathcal{L}_{\rm DSoPo}(\theta) = \mathcal{L}_{\rm vu}(\theta) + \mathcal{L}_{\rm hu}(\theta), \tag{S18}$$

#### **Relationship between USoPo and DSoPo** We begin by analyzing the size relationship between 76 USoPo and DSoPo: 77

$$\mathcal{L}_{\mathrm{DSoPo}}(\theta) - \mathcal{L}_{\mathrm{USoPo}}(\theta)$$

$$= \mathcal{L}_{\mathrm{hu}}(\theta) + \mathbb{E}_{(x^{w},c)\sim\mathcal{D}}Z_{hu}(c)\log\sigma\Big(\beta h_{\theta}(x^{w},c)\Big)$$

$$= -\mathbb{E}_{(x^{w},c)\sim\mathcal{D}}Z_{hu}(c)\mathbb{E}_{x^{1:K}\sim\bar{\pi}_{\theta}^{hu*}}\log\sigma\Big(\beta \mathcal{H}_{\theta}(x^{w},x^{l},c)\Big) + \mathbb{E}_{(x^{w},c)\sim\mathcal{D}}Z_{hu}(c)\log\sigma\Big(\beta h_{\theta}(x^{w},c)\Big)$$

$$= -\mathbb{E}_{(x^{w},c)\sim\mathcal{D}}Z_{hu}(c)\mathbb{E}_{x^{1:K}\sim\bar{\pi}_{\theta}^{hu*}}\Big[\log\sigma\Big(\beta \mathcal{H}_{\theta}(x^{w},x^{l},c)\Big) - \log\sigma\Big(\beta h_{\theta}(x^{w},c)\Big)\Big].$$
(S19)

(S19) Considering that  $\mathcal{H}_{\theta}(x^w, x^l, c) = h_{\theta}(x^w, c) - h_{\theta}(x^l, c)$  and  $h_{\theta}(x, c) = \log \frac{\pi_{\theta}(x|c)}{\pi_{\text{ref}}(x|c)}$ , we have: 78

$$\mathcal{L}_{\text{DSoPo}}(\theta) - \mathcal{L}_{\text{USoPo}}(\theta)$$

$$= -\mathbb{E}_{(x^{w},c)\sim\mathcal{D}}Z_{hu}(c)\mathbb{E}_{x^{1:K}\sim\bar{\pi}_{\theta}^{hu*}}\left[\log\sigma\left(\beta\mathcal{H}_{\theta}(x^{w},x^{l},c)\right) - \log\sigma\left(\beta h_{\theta}(x^{w},c)\right)\right]$$

$$= -\mathbb{E}_{(x^{w},c)\sim\mathcal{D}}Z_{hu}(c)\mathbb{E}_{x^{1:K}\sim\bar{\pi}_{\theta}^{hu*}}\left[\log\frac{\exp\beta h_{\theta}(x^{w},c)}{\exp\beta h_{\theta}(x^{w},c) + \exp\beta h_{\theta}(x^{l},c)} - \log\frac{\exp\beta h_{\theta}(x^{w},c)}{\exp\beta h_{\theta}(x^{w},c) + 1}\right]$$

$$= -\mathbb{E}_{(x^{w},c)\sim\mathcal{D}}Z_{hu}(c)\mathbb{E}_{x^{1:K}\sim\bar{\pi}_{\theta}^{hu*}}\left[\log\frac{\exp\beta h_{\theta}(x^{w},c) + 1}{\exp\beta h_{\theta}(x^{w},c) + \exp\beta h_{\theta}(x^{l},c)}\right]$$

$$= -\mathbb{E}_{(x^{w},c)\sim\mathcal{D}}Z_{hu}(c)\mathbb{E}_{x^{1:K}\sim\bar{\pi}_{\theta}^{hu*}}\left[\log\frac{\left(\frac{\pi_{\theta}(x^{w}|c)}{\pi_{\text{ref}}(x^{w}|c)}\right)^{\beta} + 1}{\left(\frac{\pi_{\theta}(x^{w}|c)}{\pi_{\text{ref}}(x^{w}|c)}\right)^{\beta} + \left(\frac{\pi_{\theta}(x^{l}|c)}{\pi_{\text{ref}}(x^{l}|c)}\right)^{\beta}}\right].$$
(S20)

In general, DPO focuses on reducing the generative probability of loss samples (unpreferred motions). Consequently, the generative probability of the policy model  $\pi_{\theta}(x^l|c)$  will be lower than that of the reference model  $\pi_{\text{ref}}(x^l|c)$ , i.e.,  $\pi_{\theta}(x^l|c) \leq \pi_{\text{ref}}(x^l|c)$ , resulting in  $\frac{\pi_{\theta}(x^l|c)}{\pi_{\text{ref}}(x^l|c)} \leq 1$ . Hence, the 79 80 81 following relationship holds: 82

$$\frac{\pi_{\theta}(x^{l}|c)}{\pi_{\text{ref}}(x^{l}|c)} \leq 1$$

$$\Rightarrow \left(\frac{\pi_{\theta}(x^{w}|c)}{\pi_{\text{ref}}(x^{w}|c)}\right)^{\beta} + 1 \geq \left(\frac{\pi_{\theta}(x^{w}|c)}{\pi_{\text{ref}}(x^{w}|c)}\right)^{\beta} + \left(\frac{\pi_{\theta}(x^{l}|c)}{\pi_{\text{ref}}(x^{w}|c)}\right)^{\beta}$$

$$\Rightarrow \frac{\left(\frac{\pi_{\theta}(x^{w}|c)}{\pi_{\text{ref}}(x^{w}|c)}\right)^{\beta} + 1}{\left(\frac{\pi_{\theta}(x^{w}|c)}{\pi_{\text{ref}}(x^{w}|c)}\right)^{\beta} + \left(\frac{\pi_{\theta}(x^{l}|c)}{\pi_{\text{ref}}(x^{w}|c)}\right)^{\beta}} \geq 1$$

$$\Rightarrow \log \frac{\left(\frac{\pi_{\theta}(x^{w}|c)}{\pi_{\text{ref}}(x^{w}|c)}\right)^{\beta} + \left(\frac{\pi_{\theta}(x^{l}|c)}{\pi_{\text{ref}}(x^{w}|c)}\right)^{\beta}}{\left(\frac{\pi_{\theta}(x^{w}|c)}{\pi_{\text{ref}}(x^{w}|c)}\right)^{\beta} + \left(\frac{\pi_{\theta}(x^{l}|c)}{\pi_{\text{ref}}(x^{w}|c)}\right)^{\beta}} \geq 0$$

$$\Rightarrow -\mathbb{E}_{(x^{w},c)\sim\mathcal{D}}Z_{hu}(c)\mathbb{E}_{x^{1:K}\sim\pi_{\theta}^{hu*}}\left[\log\frac{\left(\frac{\pi_{\theta}(x^{w}|c)}{\pi_{\text{ref}}(x^{w}|c)}\right)^{\beta} + 1}{\left(\frac{\pi_{\theta}(x^{w}|c)}{\pi_{\text{ref}}(x^{w}|c)}\right)^{\beta} + \left(\frac{\pi_{\theta}(x^{l}|c)}{\pi_{\text{ref}}(x^{l}|c)}\right)^{\beta}}\right] \leq 0$$

$$\Rightarrow \mathcal{L}_{\text{DSoPo}}(\theta) \leq \mathcal{L}_{\text{USOPo}}(\theta).$$
(S21)

Eq. (S21) indicates that  $\mathcal{L}_{USoPo}$  is one of upper bounds of  $\mathcal{L}_{DSoPo}$ . 83

Difference between the optimization of USoPo and DSoPo The difference between the opti-84

mization of USoPo and DSoPo can be measured by that between their objective function. Let 85

7

 $\mathcal{L}_{d}(\theta) = \mathcal{L}_{USoPo}(\theta) - \mathcal{L}_{DSoPo}(\theta)$ , the difference between their objective function can be denoted as: 86  $\mathcal{L}_{\rm d}(\theta) = \mathcal{L}_{\rm USoPo}(\theta) - \mathcal{L}_{\rm DSoPo}(\theta)$ 

$$= \mathbb{E}_{(x^w,c)\sim\mathcal{D}} Z_{hu}(c) \mathbb{E}_{x^{1:K}\sim\bar{\pi}_{\theta}^{hu*}} \left[ \log \frac{\left(\frac{\pi_{\theta}(x^w|c)}{\pi_{\mathrm{ref}}(x^w|c)}\right)^{\beta} + 1}{\left(\frac{\pi_{\theta}(x^w|c)}{\pi_{\mathrm{ref}}(x^w|c)}\right)^{\beta} + \left(\frac{\pi_{\theta}(x^l|c)}{\pi_{\mathrm{ref}}(x^l|c)}\right)^{\beta}} \right] \stackrel{(\mathsf{S22})}{\stackrel{}{=} 0}$$

- where ① holds due to Eq. (S21). As discussed above, the generative probability of the policy model 87
- $\pi_{\theta}(x^{l}|c)$  will be lower than that of the reference model  $\pi_{ref}(x^{l}|c)$ , and thus  $\pi_{\theta}(x^{l}|c)$  falls in the range 88
- between 0 and  $\pi_{\text{ref}}(x^l|c)$ , i.e.,  $0 \le \pi_{\theta}(x^l|c) \le \pi_{\text{ref}}(x^l|c)$ . 89
- Assuming that the value of  $\pi_{\theta}(x^w|c)$  is fixed, the value of  $\mathcal{L}_{d}(\theta)$  is negatively correlated with  $\pi_{\theta}(x^l|c)$ , 90 since we have: 91

$$\nabla_{\theta} \mathcal{L}_{d}(\theta) = \nabla_{\theta} - \mathbb{E}_{(x^{w},c)\sim\mathcal{D}} Z_{hu}(c) \mathbb{E}_{x^{1:K}\sim\bar{\pi}_{\theta}^{hu*}} \left[ \log \frac{\left(\frac{\pi_{\theta}(x^{w}|c)}{\pi_{\text{ref}}(x^{w}|c)}\right)^{\beta} + 1}{\left(\frac{\pi_{\theta}(x^{w}|c)}{\pi_{\text{ref}}(x^{w}|c)}\right)^{\beta} + \left(\frac{\pi_{\theta}(x^{l}|c)}{\pi_{\text{ref}}(x^{l}|c)}\right)^{\beta}} \right]$$

$$= \mathbb{E}_{(x^{w},c)\sim\mathcal{D}} Z_{hu}(c) \mathbb{E}_{x^{1:K}\sim\bar{\pi}_{\theta}^{hu*}} \nabla_{\theta} - \log \left[ \left(\frac{\pi_{\theta}(x^{w}|c)}{\pi_{\text{ref}}(x^{w}|c)}\right)^{\beta} + \left(\frac{\pi_{\theta}(x^{l}|c)}{\pi_{\text{ref}}(x^{l}|c)}\right)^{\beta} \right]$$

$$= \mathbb{E}_{(x^{w},c)\sim\mathcal{D}} Z_{hu}(c) \mathbb{E}_{x^{1:K}\sim\bar{\pi}_{\theta}^{hu*}} \frac{1}{\left(\frac{\pi_{\theta}(x^{w}|c)}{\pi_{\text{ref}}(x^{w}|c)}\right)^{\beta} + \left(\frac{\pi_{\theta}(x^{l}|c)}{\pi_{\text{ref}}(x^{l}|c)}\right)^{\beta}} - \nabla_{\theta} \left(\frac{\pi_{\theta}(x^{l}|c)}{\pi_{\text{ref}}(x^{l}|c)}\right)^{\beta}$$

$$\stackrel{(S23)}{=} \nabla_{\theta} \left(\frac{\pi_{\theta}(x^{l}|c)}{\pi_{\text{ref}}(x^{l}|c)}\right)^{\beta}.$$

where () holds since  $\frac{1}{(\frac{\pi_{\theta}(x^w|c)}{\pi_{\mathrm{ref}}(x^w|c)})^{\beta} + (\frac{\pi_{\theta}(x^l|c)}{\pi_{\mathrm{ref}}(x^l|c)})^{\beta}} > 0.$ 92

Hence, when the generative probability of unpreferred motions  $\pi_{\theta}(x^l|c)$  is lower, the difference 93

between the optimization of USoPo and DSoPo is larger. However, the unpreferred motions are sampled from the relatively high-preference distribution  $\pi_{\bar{\theta}}^{hu*}$ , and thus should not be treated as unpreferred motions. Using  $\mathcal{L}_{\text{USoPo}}(\theta)$  to optimize policy model  $\pi_{\theta}$  instead of  $\mathcal{L}_{\text{DSoPo}}(\theta)$  can avoid 94

95

96

unnecessary optimization of these relatively high-preference unpreferred motion  $\mathcal{L}_{d}(\theta)$ . 97

#### **B.5 Proof of Eq. (16)** 98

- Before proving Eq. (16), we first present some useful lemmas from [5]. 99
- **Lemma 1.** [5] Given a winning sample  $x_w$  and a losing sample  $x_l$ , the DPO denoted as 100

$$\mathcal{L}_{\text{DPO}}(\theta) = \mathbb{E}_{(x^w, x^l, c) \sim \mathcal{D}} \left[ -\log \sigma \left( \beta \log \frac{\pi_{\theta}(x^w|c)}{\pi_{ref}(x^w|c)} - \beta \log \frac{\pi_{\theta}(x^l|c)}{\pi_{ref}(x^l|c)} \right) \right].$$
(S24)

Then the objective function for diffusion models can be denoted as: 101

$$\mathcal{L}_{DPO\text{-Diffusion}}(\theta) = -\mathbb{E}_{(x_0^w, x_0^l) \sim \mathcal{D}} \log \sigma(\beta \mathbb{E}_{x_{1:T}^w \sim \pi_{\theta}(x_{1:T}^w | x_0^w), x_{1:T}^l \sim \pi_{\theta}(x_{1:T}^l | x_0^l)} \\ [\log \frac{\pi_{\theta}(x_{0:T}^w)}{\pi_{\text{ref}}(x_{0:T}^w)} - \log \frac{\pi_{\theta}(x_{0:T}^l)}{\pi_{\text{ref}}(x_{0:T}^l)}]),$$
(S25)

- where  $x_t^*$  denoted the noised sample  $x^*$  for the t-th step. 102
- **Lemma 2.** [5] Given the objective function of diffusion-based DPO denoted as Eq. (S25), it has an 103 upper bound  $\mathcal{L}_{UB}(\theta)$ : 104

$$\mathcal{L}_{DPO\text{-Diffusion}}(\theta) \leq -\mathbb{E}_{(x_{0}^{w}, x_{0}^{l}) \sim \mathcal{D}, t \sim \mathcal{U}(0, T), x_{t-1, t}^{w} \sim \pi_{\theta}(x_{t-1, t}^{w} | x_{0}^{w}), x_{t-1, t}^{l} \sim \pi_{\theta}(x_{t-1, t}^{t} | x_{0}^{l}) \log \sigma} \\ \underbrace{\left(\beta T \log \frac{\pi_{\theta}(x_{t-1}^{w} | x_{t}^{w})}{\pi_{\mathrm{ref}}(x_{t-1}^{w} | x_{t}^{w})} - \beta T \log \frac{\pi_{\theta}(x_{t-1}^{l} | x_{t}^{l})}{\pi_{\mathrm{ref}}(x_{t-1}^{l} | x_{t}^{l})}\right)}_{\mathcal{L}_{\mathrm{UB}}(\theta)}, \quad (S26)$$

where T denotes the number of diffusion steps. 105

**Lemma 3.** [5] Given the objective function for diffusion model denoted as Eq. (S26), it can be rewritten as :

$$\mathcal{L}_{\rm UB}(\theta) = -\mathbb{E}_{(x_0^w, x_0^l) \sim \mathcal{D}, t \sim \mathcal{U}(0,T), x_t^w \sim q(x_t^w | x_0^w), x_t^l \sim q(x_t^l | x_0^l)} \log \sigma(-\beta T \omega_t) \\ (\|\epsilon - \epsilon_\theta(x_t^w, t)\|_2^2 - \|\epsilon - \epsilon_{\rm ref}(x_t^w, t)\|_2^2 - (\|\epsilon - \epsilon_\theta(x_t^l, t)\|_2^2 - \|\epsilon - \epsilon_{\rm ref}(x_t^l, t)\|_2^2))),$$
(S27)

108 where  $x_t^* = \alpha_t x_0^* + \sigma_t \epsilon$ ,  $\epsilon \sim \mathcal{N}(0, \mathbb{I})$  is a draw from the distribution of forward process  $q(x_t^* | x_0^*)$ .

Now, we proof Eq. (16) based on these lemmas.

*Proof.* This proof has three steps. In each step, we apply the three lemmas introduced above in succession. We begin with the loss function of SoPo for probability models:

$$\mathcal{L}_{\mathrm{SoPo}}(\theta) = \underbrace{-\mathbb{E}_{(x^{w},c)\sim\mathcal{D},x_{\overline{\pi}\theta}^{1:K}\sim\overline{\pi}_{\theta}^{vu*}(\cdot|c)} Z_{vu}(c) \Big[\log\sigma\Big(\beta_{w}(x^{w})h_{\theta}(x^{w},c) - \beta h_{\theta}(x^{l},c)\Big)\Big]}_{\mathcal{L}_{\mathrm{SoPo-vu}}(\theta)} - \underbrace{\mathbb{E}_{(x^{w},c)\sim\mathcal{D}} Z_{hu}(c)\log\sigma\Big(\beta_{w}(x^{w})h_{\theta}(x^{w},c)\Big)}_{\mathcal{L}_{\mathrm{SoPo-hu}}(\theta)}.$$
(S28)

<sup>112</sup> Based on Lemma 1, we can rewrite the objective function for diffusion models:

$$\begin{aligned} \mathcal{L}_{\mathrm{SoPo-Diffusion}}(\theta) &= \mathcal{L}_{\mathrm{SoPo-vu}}^{\mathrm{diff-ori}}(\theta) + \mathcal{L}_{\mathrm{SoPo-hu}}^{\mathrm{diff-ori}}(\theta) \\ \mathcal{L}_{\mathrm{SoPo-vu}}^{\mathrm{diff-ori}}(\theta) &= - \mathbb{E}_{(x_0^w,c)\sim\mathcal{D},x_0^{1:K}\sim\pi_\theta^w u^*(\cdot|c)} Z_{vu}(c) \\ \log \sigma \left( \mathbb{E}_{x_{1:T}^w \sim \pi_\theta(x_{1:T}^w|x_0^w), x_{1:T}^l \sim \pi_\theta(x_{1:T}^l|x_0^l)} [\beta_w(x_0^w) \log \frac{\pi_\theta(x_{0:T}^w)}{\pi_{\mathrm{ref}}(x_{0:T}^w)} - \beta \log \frac{\pi_\theta(x_{0:T}^l)}{\pi_{\mathrm{ref}}(x_{0:T}^l)}] \right), \\ \mathcal{L}_{\mathrm{SoPo-hu}}^{\mathrm{diff-ori}}(\theta) &= - \mathbb{E}_{(x_0^w,c)\sim\mathcal{D}} Z_{hu}(c) \log \sigma \left( \mathbb{E}_{x_{1:T}^w \sim \pi_\theta(x_{1:T}^w|x_0^w)} [\beta_w(x^w) \log \frac{\pi_\theta(x_{0:T}^w)}{\pi_{\mathrm{ref}}(x_{0:T}^w)}] \right), \end{aligned}$$

$$(S29)$$

where  $x_t^*$  denoted the noised sample  $x^*$  for the *t*-th step. According to Lemma 2, the upper bound of  $\mathcal{L}_{SoPo-vu}^{diff-ori}(\theta)$  and  $\mathcal{L}_{SoPo-hu}^{diff-ori}(\theta)$  can be denoted as:

$$\mathcal{L}_{\mathrm{SoPo-vu}}^{\mathrm{diff-ori}}(\theta) \leq -\mathbb{E}_{(x_{0}^{w},c)\sim\mathcal{D},x_{0}^{1:K}\sim\pi_{\theta}^{vu*}(\cdot|c),t\sim\mathcal{U}(0,T),x_{t-1,t}^{w}\sim\pi_{\theta}(x_{t-1,t}^{w}|x_{0}^{w}),x_{t-1,t}^{l}\sim\pi_{\theta}(x_{t-1,t}^{l}|x_{0}^{w})}{\log\sigma\left(\beta_{w}(x_{0}^{w})T\log\frac{\pi_{\theta}(x_{t-1}^{w}|x_{t}^{w})}{\pi_{\mathrm{ref}}(x_{t-1}^{w}|x_{t}^{w})}-\beta T\log\frac{\pi_{\theta}(x_{t-1,t}^{l}|x_{t}^{l})}{\pi_{\mathrm{ref}}(x_{t-1}^{l}|x_{t}^{l})}\right)},}{\mathcal{L}_{\mathrm{SoPo-vu}}^{\mathrm{diff-ori}}(\theta) \leq -\mathbb{E}_{(x_{0}^{w},c)\sim\mathcal{D},t\sim\mathcal{U}(0,T),x_{t-1,t}^{w}\sim\pi_{\theta}(x_{t-1,t}^{w}|x_{0}^{w})}\log\sigma\left(\beta_{w}(x_{0}^{w})T\log\frac{\pi_{\theta}(x_{t-1}^{u}|x_{t}^{w})}{\pi_{\mathrm{ref}}(x_{t-1}^{w}|x_{t}^{w})}\right)},}{\mathcal{L}_{\mathrm{SoPo-hu}}^{\mathrm{diff-ori}}(\theta)}$$
$$\mathcal{L}_{\mathrm{SoPo-biffusion}}(\theta) = \mathcal{L}_{\mathrm{SoPo-vu}}^{\mathrm{diff-ori}}(\theta) + \mathcal{L}_{\mathrm{SoPo-hu}}^{\mathrm{diff-ori}}(\theta) \leq \mathcal{L}_{\mathrm{SoPo-hu}}^{\mathrm{diff-ori}}(\theta) + \mathcal{L}_{\mathrm{SoPo-hu}}^{\mathrm{diff-ori}}(\theta) = \mathcal{L}_{\mathrm{SoPo-hu}}^{\mathrm{diff}}(\theta).$$
(S30)

115 Applying Lemma 3 to  $\mathcal{L}_{SoPo-vu}^{diff}(\theta)$  and  $\mathcal{L}_{SoPo-hu}^{diff}(\theta)$ , we have

$$\mathcal{L}_{\text{SoPo-vu}}^{\text{diff}}(\theta) = -\mathbb{E}_{(x_0^w, c)\sim\mathcal{D}, x_0^{1:K}\sim\bar{\pi}_{\theta}^{wu*}(\cdot|c), t\sim\mathcal{U}(0,T), x_t^w\sim q(x_t^w|x_0^w), x_t^l\sim q(x_t^l|x_0^l)} \log\sigma\bigg(-T\omega_t\Big(\beta_w(x_0^w)\big(\|\epsilon-\epsilon_\theta(x_t^w,t)\|_2^2 - \|\epsilon-\epsilon_{\text{ref}}(x_t^w,t)\|_2^2\Big) -\beta\big(\|\epsilon-\epsilon_\theta(x_t^l,t)\|_2^2 - \|\epsilon-\epsilon_{\text{ref}}(x_t^l,t)\|_2^2\big)\Big)\bigg),$$

$$\mathcal{L}_{\text{SoPo-bu}}^{\text{diff}}(\theta) = -\mathbb{E}_{(x_t^w,c)\sim\mathcal{D}, t\sim\mathcal{U}(0,T), x_t^w} \sim \pi_\theta(x_t^w, |x_t^w)$$
(S31)

$$\mathcal{L}_{\mathrm{SoPo-hu}}^{\mathrm{dmn}}(\theta) = -\mathbb{E}_{(x_0^w,c)\sim\mathcal{D},t\sim\mathcal{U}(0,T),x_{t-1,t}^w\sim\pi_\theta(x_{t-1,t}^w|x_0^w)} \\ \log\sigma\Big(-T\omega_t\beta_w(x_0^w)\big(\|\epsilon-\epsilon_\theta(x_t^w,t)\|_2^2 - \|\epsilon-\epsilon_{\mathrm{ref}}(x_t^w,t)\|_2^2\big)\Big), \\ \mathcal{L}_{\mathrm{SoPo}}^{\mathrm{diff}}(\theta) = \mathcal{L}_{\mathrm{SoPo-vu}}^{\mathrm{diff}}(\theta) + \mathcal{L}_{\mathrm{SoPo-hu}}^{\mathrm{diff}}(\theta)$$

<sup>116</sup> To simplify the symbolism, the objective functions can be rewritten as:

$$\mathcal{L}_{\text{SoPo-vu}}^{\text{diff}} = -\mathbb{E}_{t \sim \mathcal{U}(0,T), (x^{w},c) \sim \mathcal{D}, x_{\pi_{\theta}}^{1:K} \sim \bar{\pi}_{\theta}^{vu*}(\cdot|c)} Z_{vu}(c) \\ \left[ \log \sigma \Big( -T\omega_{t} \big( \beta_{w}(x_{w}) \big( \mathcal{L}(\theta, \operatorname{ref}, x_{t}^{w}) - \beta \mathcal{L}(\theta, \operatorname{ref}, x_{t}^{l}) \big) \big) \Big) \right],$$

$$\mathcal{L}_{\text{SoPo-hu}}^{\text{diff}} = -\mathbb{E}_{t \sim \mathcal{U}(0,T), (x^{w},c) \sim \mathcal{D}} Z_{hu}(c) \\ \left[ \log \sigma \Big( -T\omega_{t} \beta_{w}(x_{w}) \mathcal{L}(\theta, \operatorname{ref}, x_{t}^{w}) \Big) \Big],$$
(S32)

where  $\mathcal{L}(\theta, \operatorname{ref}, x_t) = \mathcal{L}(\theta, x_t) - \mathcal{L}(\operatorname{ref}, x_t)$ , and  $\mathcal{L}(\theta/\operatorname{ref}, x_t) = \|\epsilon_{\theta/\operatorname{ref}}(x_t, t) - \epsilon\|_2^2$  denotes the loss of the policy or reference model. The proof is completed.

## 119 C Experiment

#### 120 C.1 Details of Experiments on Synthetic Data

To simulate our preference optimization framework, we design a 2D synthetic setup with predefined generation and reward distributions. The generator distribution  $\pi_{\theta}$  is modeled as a Gaussian with mean [-2, 1] and covariance matrix diag(2.0, 2.0). The reward model is defined as a mixture of two

Gaussians with means [-3, 2] and [2, -2], covariances  $\begin{bmatrix} 1 & \pm 0.5 \\ \pm 0.5 & 1 \end{bmatrix}$ , and equal weights of 0.5.

For the offline dataset, preferred samples are randomly drawn from the reward distribution, while 125 unpreferred samples are sampled from a manually specified distribution dissimilar to the reward 126 model. These are used to fine-tune the generator via offline preference optimization. For the online 127 setting, we draw samples from the reference model and assign preference labels using the reward 128 model to distinguish preferred and unpreferred motions. This process is repeated iteratively to 129 optimize the model online. In SoPo, we combine offline preferred samples with online-generated 130 unpreferred ones to perform semi-online preference optimization, thereby leveraging the strengths of 131 both offline and online data. 132

## 133 C.2 Additional Experimental Datails

**Datasets & Evaluation** HumanML3D is derived from the AMASS [6] and HumanAct12 [7] datasets and contains 14,616 motions, each described by three textual annotations. All motion is split into train, test, and evaluate sets, composed of 23384, 1460, and 4380 motions, respectively. For both HumanML3D and KIT-ML datasets, we follow the official split and report the evaluated performance on the test set.

We evaluate our experimental results on two main aspects: alignment quality and generation quality. 139 Following prior research [8-10], we use motion retrieval precision (R-Precision) and multi-modal 140 distance (MM Dist) to evaluate alignment quality, while diversity and Fréchet Inception Distance 141 (FID) are employed to assess generation quality. (1) R-Precision evaluates the similarity between 142 generated motion and their corresponding text descriptions. Higher values indicate better alignment 143 quality. (2) MM Dist represents the average distance between the generated motion features and 144 their corresponding text embedding. (3) Diversity calculates the variation in generated samples. A 145 146 diversity close to real motions ensures that the model produces rich patterns rather than repetitive motions. (4) FID measures the distribution proximity between the generated and real samples in 147 latent space. Lower FID scores indicate higher generation quality. 148

**Implementation Details** For the preference alignment of MDM [1], we largely adopt the original implementation's settings. The model is trained using the AdamW optimizer [11] with a cosine decay learning rate scheduler and linear warm-up over the initial steps. We use a batch size of 64 and a learning rate of  $10^{-5}$ , with a guidance parameter of 2.5 during testing. Diffusion employs a cosine noise schedule with 50 steps, and an evaluation batch size of 32 ensures consistent metric computation. For fine-tuning MLD [2], we similarly follow its original parameter settings.

## 155 C.3 Additional Experimental Results

We visualize the generated motion for our SoPo. As shown in Fig. **S2**, our proposed approach helps text-to-motion models avoid frequent mistakes, such as incorrect movement direction and specific



(a) A person runs to their right and then curves to the left and continues to run then stops.

(b) A man jumps and brings both arms above his head as ... and then moves them back into the original position.

Figure S2: Visual results on HumanML3D dataset.

semantics. Additionally, we also present additional results generated by text-to-motion models with
SoPo, as illustrated in Fig. S1. Our proposed SoPo significantly enhances the ability of text-to-motion
models to comprehend text semantics. For instance, in Fig. S1 (j), a model integrated with SoPo can
successfully interpret the semantics of "zig-zag pattern", whereas a model without SoPo struggles to
do so.

# **163** References

- [1] Guy Tevet, Sigal Raab, Brian Gordon, Yoni Shafir, Daniel Cohen-or, and Amit Haim Bermano. Human motion diffusion model. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=SJ1kSy02jwu. 1, 10
- [2] Xin Chen, Biao Jiang, Wen Liu, Zilong Huang, Bin Fu, Tao Chen, and Gang Yu. Executing your commands
   via motion diffusion in latent space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18000–18010, 2023. 1, 10
- [3] R. L. Plackett. The analysis of permutations. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 24(2):193–202, 1975. ISSN 00359254, 14679876. URL http://www.jstor.org/stable/2346567. 3, 4
- [4] Cheng Lu Haozhe Ji, Pei Ke Yilin Niu, Jun Zhu Hongning Wang, and Minlie Huang Jie Tang. Towards
   efficient exact optimization of language model alignment. *The Forty-first International Conference on Machine Learning*, 2024. URL https://arxiv.org/abs/2402.00856. 4
- [5] Bram Wallace, Meihua Dang, Rafael Rafailov, Linqi Zhou, Aaron Lou, Senthil Purushwalkam, Stefano
   Ermon, Caiming Xiong, Shafiq Joty, and Nikhil Naik. Diffusion Model Alignment Using Direct Preference
   Optimization . In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8228–8238,
   Los Alamitos, CA, USA, June 2024. IEEE Computer Society. doi: 10.1109/CVPR52733.2024.00786.
   URL https://doi.ieeecomputersociety.org/10.1109/CVPR52733.2024.00786.
- [6] Naureen Mahmood, Nima Ghorbani, Nikolaus F. Troje, Gerard Pons-Moll, and Michael Black. AMASS:
   Archive of Motion Capture As Surface Shapes . In *IEEE/CVF International Conference on Computer Vision*,
   pages 5441–5450, Los Alamitos, CA, USA, 2019. IEEE Computer Society. doi: 10.1109/ICCV.2019.00554.
   URL https://doi.ieeecomputersociety.org/10.1109/ICCV.2019.00554.
- [7] Chuan Guo, Xinxin Zuo, Sen Wang, Shihao Zou, Qingyao Sun, Annan Deng, Minglun Gong, and Li Cheng. Action2motion: Conditioned generation of 3d human motions. In *Proceedings of the ACM International Conference on Multimedia*, page 2021–2029, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450379885. doi: 10.1145/3394171.3413635. URL https: //doi.org/10.1145/3394171.3413635. 10
- [8] Wenxun Dai, Ling-Hao Chen, Jingbo Wang, Jinpeng Liu, Bo Dai, and Yansong Tang. Motionlcm: Real time controllable motion generation via latent consistency model. In Aleš Leonardis, Elisa Ricci, Stefan
   Roth, Olga Russakovsky, Torsten Sattler, and Gül Varol, editors, *European Conference on Computer Vision*,
   pages 390–408, Cham, 2024. Springer Nature Switzerland. ISBN 978-3-031-72640-8. 10

- [9] Zeping Ren, Shaoli Huang, and Xiu Li. Realistic human motion generation with cross-diffusion models.
   *European Conference on Computer Vision*, 2024.
- 196 [10] Zeyu Zhang, Akide Liu, Ian Reid, Richard Hartley, Bohan Zhuang, and Hao Tang. Motion mamba: Efficient
- and long sequence motion generation. In Aleš Leonardis, Elisa Ricci, Stefan Roth, Olga Russakovsky,
   Torsten Sattler, and Gül Varol, editors, *European Conference on Computer Vision*, pages 265–282. Springer
- 199 Nature Switzerland, 2024. 10
- [11] I Loshchilov. Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101, 2017. 10