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

Xiaofeng Tan[1,2], Hongsong Wang[1,2*], Xin Geng[1,2], Pan Zhou[3]

[1] Southeast University, [2] Key Laboratory of New Generation Artificial Intelligence Technology, [3] Singapore Management University

## Introduction

**Issue:** Existing text-to-motion methods struggle to generate semantically consistent motions.



(a) Someone with difficulty falls and kneels
(b) A man kneels down then stands back up
(c) The person is dribbling a basketball backwards
(d) Person is sprinting to their left
(e) A person, standing, raises his left hand ...
(f) A person ... and lost his balance, ...caught himself
(g) A person stands left leg forward in a karate pose, leans forward and makes a high karate kick with their right leg and returns to the karate pose

**Key Observation:** We observe that tasks of motion understanding and discrimination are generally less complex and demonstrate superior performance compared to motion generation.

**Question:** How effectively can discriminative models improve motion generation quality without any additional inference cost?

**Contributions:** We propose SoPo, a semi-online preference optimization method, combining the strengths of online and offline direct preference optimization to overcome their individual shortcomings, delivering enhanced motion generation quality and preference alignment.

## Motivation: Rethink Off-/Online DPO

**Offline DPO:** overfitting due to limited unpreferred motions.

**Theorem 1.** *Given a preference motion dataset $\mathcal{D}$, a reference model $\pi_{ref}$, and ground-truth preference distribution $p_{gt}$, the gradient of $\nabla_\theta \mathcal{L}_{off}$ can be written as:*

$$\nabla_\theta \mathcal{L}_{off}(\theta) = \mathbb{E}_{c \sim \mathcal{D}, x^{1:K}} \nabla_\theta D_{KL}(p_{gt} \| p_\theta). \quad (4)$$

*Here $p_\theta(x^{1:K}|c) = \prod_{k=1}^{K} p_\theta(x^k|c)$ represents the likelihood that policy model generates motions $x^{1:K}$ matching their rankings, where $p_\theta(x^k|c) = \frac{(\exp h_\theta(x^k,c))^\beta}{\sum_{j=k}^{K}(\exp h_\theta(x^j,c))^\beta}$.*

**Online DPO:** biased sampling, resulting in even high-preference samples being incorrectly categorized as low-preference motions.

**Theorem 2.** *Given a reward model $r$ and a reference model $\pi_{ref}$, for the online DPO loss $\mathcal{L}_{on}$, its gradient is:*

$$\nabla_\theta \mathcal{L}_{on}(\theta) = \mathbb{E}_{c \sim \mathcal{D}, x^{1:K}} \nabla_\theta \, p_{\pi_\theta}(x^{1:K}|c) D_{KL}(p_r \| p_\theta), \quad (6)$$

*where $p_{\pi_\theta}(x^{1:K}|c) = \prod_{k=1}^{K} p_{\pi_\theta}(x^k|c)$ with $p_{\pi_\theta}(x^k|c)$ being the generative probability of policy model to generate $x^k$ conditioned on c, and $p_\theta(x^k) = \frac{(\exp h_\theta(x_k,c))^\beta}{\sum_{j=k}^{K}(\exp h_\theta(x_j,c))^\beta}$ denotes the likelihood that policy model generates motion $x_k$ with the k-th largest probability.*

## Contact

Email: txf0620@gmail.com

WeChat: txf_06_20
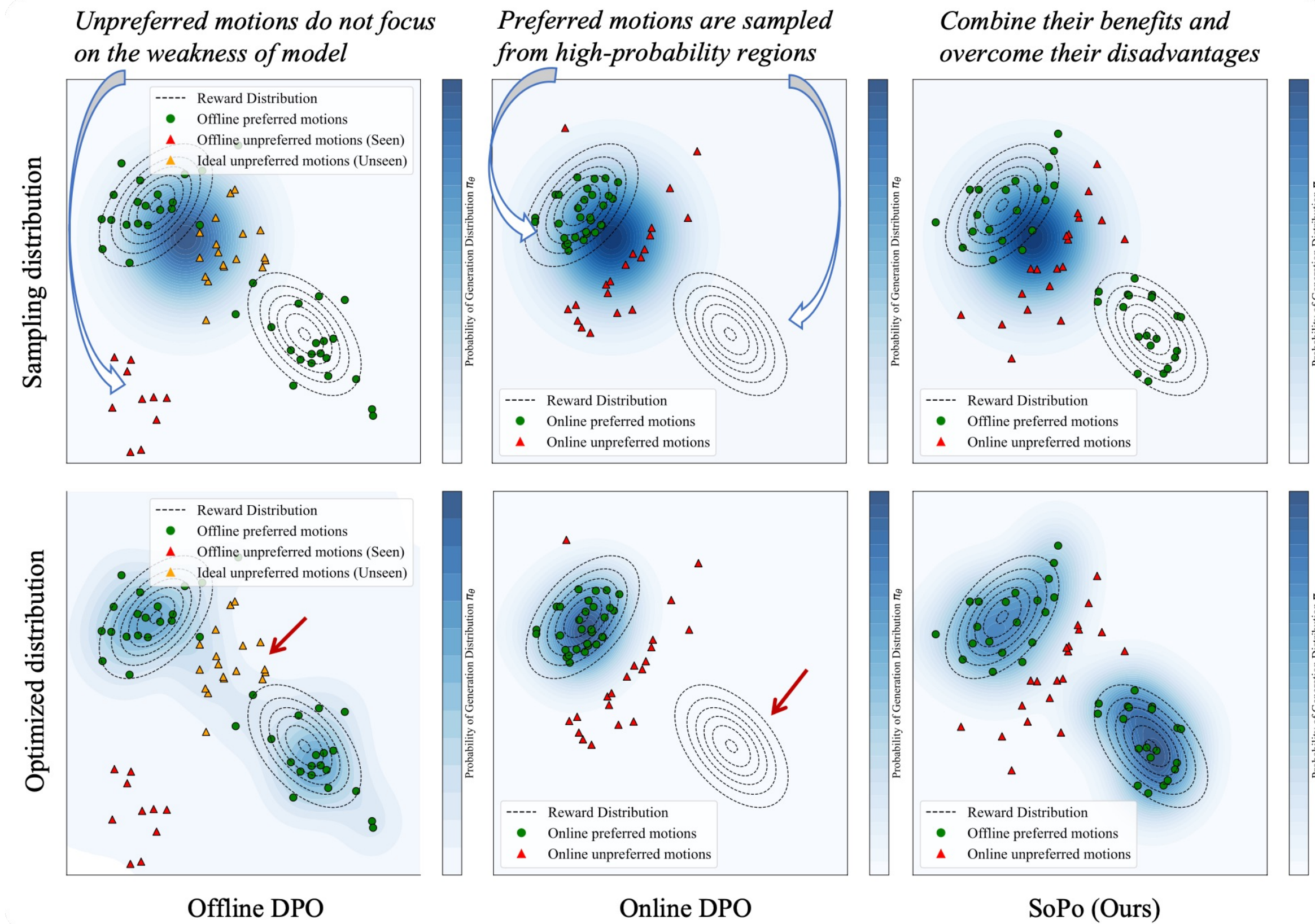
X: XiaofengTan85815

First Author | Project Page | Code | Paper

## Toy Example: Off-/Online DPO, and SoPo



*Unpreferred motions do not focus on the weakness of model*

*Preferred motions are sampled from high-probability regions*

*Combine their benefits and overcome their disadvantages*

Offline DPO | Online DPO | SoPo (Ours)

## Semi-Online Preference Optimization

### Insight for Unpreferred Motion Sampling

**Case 1:** The group $\{x_{\bar\pi_\theta}^k\}_{k=1}^K$ contains a low-preference unpreferred motion $x_{\bar\pi_\theta}^l$. Then we select these unpreferred motions iteratively which ensure diversity due to randomness of online generations and address the diversity lacking issue in offline DPO.

**Case 2:** The group contains no low-preference unpreferred motion $x_{\bar\pi_\theta}^l$, meaning all sampled motions are of high preference and should not be treated as unpreferred. This suggests the model performs well under condition c, so training should focus on high-quality preferred motions from offline data to further enhance generation quality.

### 1. Distribution Separation

$$p_{\bar\pi_\theta}(x_{\bar\pi_\theta}^{1:K}|c) = \underbrace{p_{\bar\pi_\theta}(x_{\bar\pi_\theta}^{1:K}|c)p_\tau(r(x_{\bar\pi_\theta}^l, c) \geq \tau)}_{\text{relatively high-preference unpreferred motions } \bar\pi_\theta^{hu}} + \underbrace{p_{\bar\pi_\theta}(x_{\bar\pi_\theta}^{1:K}|c)p_\tau(r(x_{\bar\pi_\theta}^l, c) < \tau)}_{\text{valuable unpreferred motions } \bar\pi_\theta^{vu}},$$

$$\mathcal{L}_{vu} = -\mathbb{E}_{(x^w,c)\sim\mathcal{D}} Z_{vu}(c) \mathbb{E}_{x_{\bar\pi_\theta}^{1:K} \sim \bar\pi_\theta^{vu*}(\cdot|c)} \log \sigma(\beta \mathcal{H}_\theta(x^w, x_{\bar\pi_\theta}^l, c)),$$

$$\mathcal{L}_{hu} = -\mathbb{E}_{(x^w,c)\sim\mathcal{D}} Z_{hu}(c) \mathbb{E}_{x_{\bar\pi_\theta}^{1:K} \sim \bar\pi_\theta^{hu*}(\cdot|c)} \log \sigma(\beta \mathcal{H}_\theta(x^w, x_{\bar\pi_\theta}^l, c)),$$

### 2. Training loss amendment

Accordingly, we rewrite the loss $\mathcal{L}_{hu}(\theta)$ into $\mathcal{L}_{USoPo-hu}(\theta)$ for filtering them:

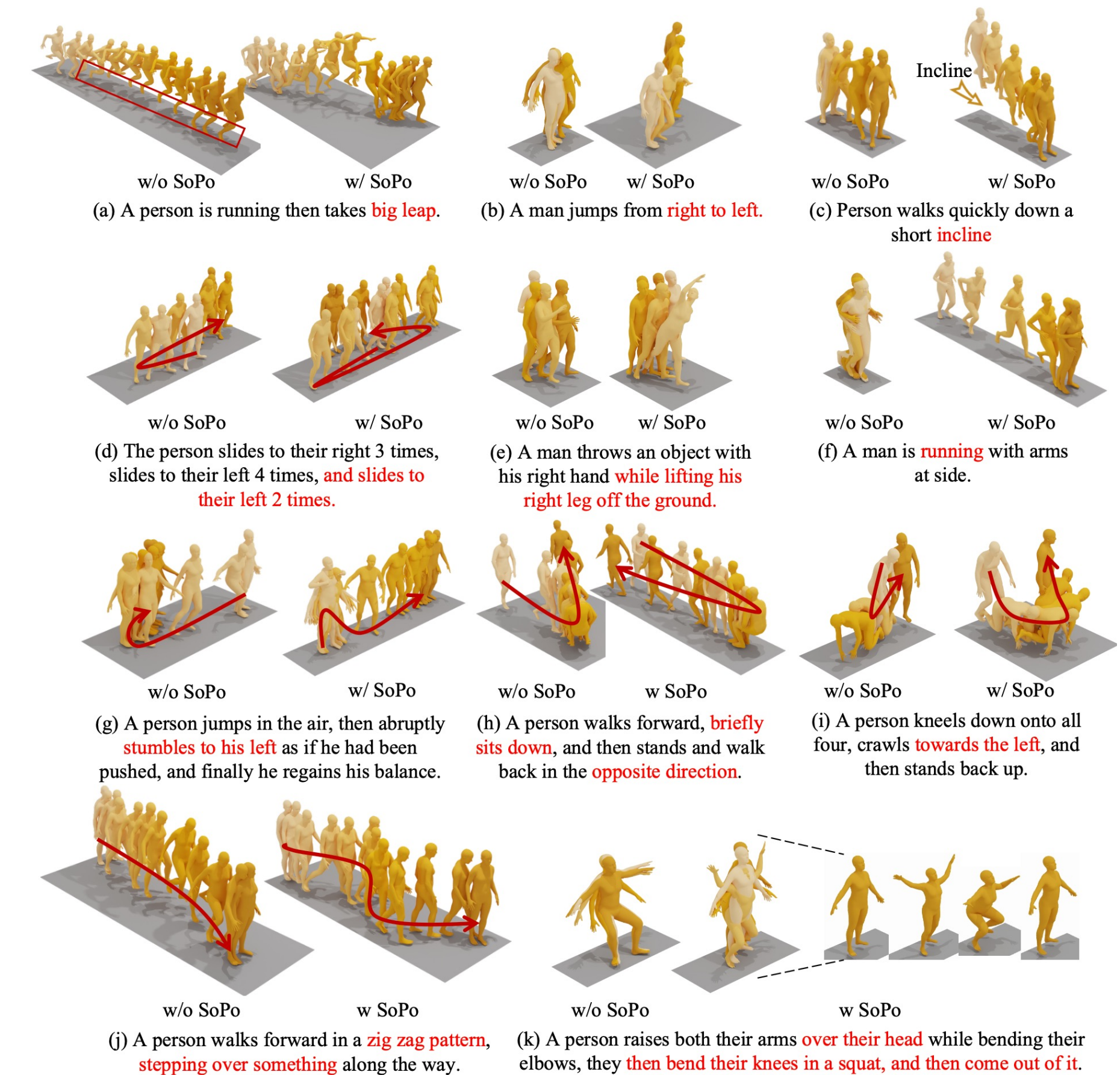$$\mathcal{L}_{USoPo-hu}(\theta) = -\mathbb{E}_{(x^w,c)\sim\mathcal{D}} Z_{hu}(c) \log \sigma(\beta h_\theta(x^w, c)), \quad \mathcal{L}_{USoPo}(\theta) = \mathcal{L}_{USoPo-hu}(\theta) + \mathcal{L}_{vu}(\theta)$$

### 3. SoPo for Diffusion

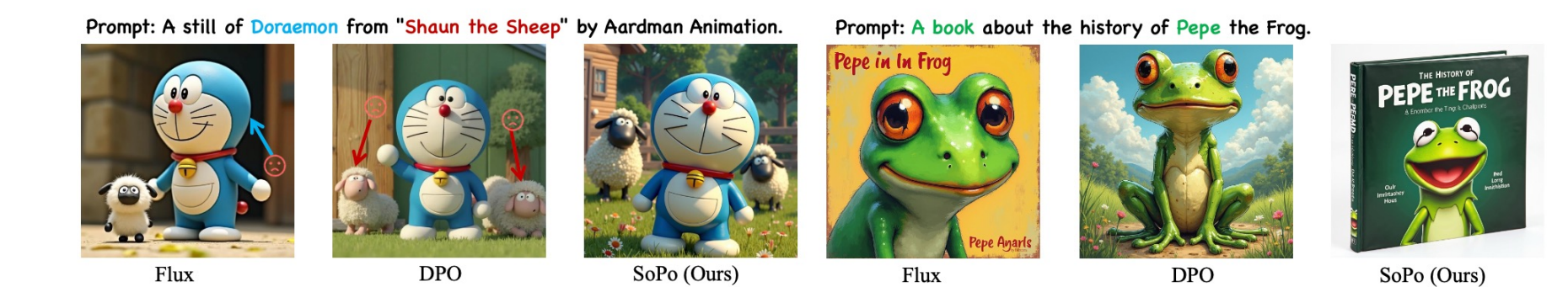$$\mathcal{L}_{SoPo}^{diff}(\theta) = -\mathbb{E}_{t\sim\mathcal{U}(0,T),(x^w,c)\sim\mathcal{D}, x_{\bar\pi_\theta}^{1:K}\sim\bar\pi_\theta(\cdot|c)} \begin{cases} \log \sigma(-T\omega_t(\beta_w(x_w)\mathcal{L}(\theta,ref,x_t^w) - \beta\mathcal{L}(\theta,ref,x_t^l))), & \text{if } r(x^l,c) < \tau, \\ \log \sigma(-T\omega_t\beta_w(x_w)\mathcal{L}(\theta,ref,x_t^w)), & \text{otherwise}. \end{cases} \quad (17)$$

## Experiments

### Text-to-Motion Qualitative Results



(a) A person is running then takes big leap.
(b) A man jumps from right to left.
(c) Person walks quickly down a short incline
(d) The person slides to their right 3 times, slides to their left 4 times, and slides to their left 2 times.
(e) A man throws an object with his right hand while lifting his right leg off the ground.
(f) A man is running with arms at side.
(g) A person jumps in the air, then abruptly stumbles to his left as if he had been pushed, and finally he regains his balance.
(h) A person walks forward, briefly sits down, and then stands and walk back in the opposite direction.
(i) A person kneels down onto all four, crawls towards the left, and then stands back up.
(j) A person walks forward in a zig zag pattern, stepping over something along the way.
(k) A person raises both their arms over their head while bending their elbows, they then bend their knees in a squat, and then come out of it.

### Text-to-Image Qualitative Results



Prompt: A still of Doraemon from "Shaun the Sheep" by Aardman Animation.
Flux | DPO | SoPo (Ours)

Prompt: A book about the history of Pepe the Frog.
Flux | DPO | SoPo (Ours)

### Text-to-Motion Quantitative Results

| Methods | Time* | R-Precision↑ | | | MM Dist↓ | Diversity → | FID↓ |
|---|---|---|---|---|---|---|---|
| | | Top 1 | Top 2 | Top 3 | | | |
| Real | - | 0.511±0.003 | 0.703±0.003 | 0.797±0.002 | 2.974±0.008 | 9.503±0.065 | 0.002±0.000 |
| MLD [1] | +0 X | 0.453±0.003 | 0.679±0.003 | 0.755±0.002 | 3.292±0.010 | 9.793±0.072 | 0.459±0.01 |
| +MoDiPO-T [9] | +121K X | 0.455±0.002 | 0.682±0.003 | 0.758±0.002 +0.40% | 3.267±0.010 +0.76% | 9.188±0.002 +33.9% | 0.303±0.031 +33.9% |
| +MoDiPO-X [9] | +121K X | 0.452±0.003 | 0.678±0.003 | 0.755±0.003 -0.26% | 3.294±0.010 +0.01% | 9.702±0.075 +1.02% | 0.281±0.013 +38.8% |
| +MoDiPO-O [9] | +121K X | 0.406±0.003 | 0.609±0.003 | 0.690±0.003 | 3.701±0.013 -12.4% | 9.241±0.079 +2.07% | 0.276±0.047 +39.9% |
| +SoPo (Ours) | +20 X | 0.463±0.003 -2.21% | 0.682±0.003 +1.23% | 0.763±0.002 -1.06% | 3.185±0.012 +3.25% | 9.525±0.065 +0.28% | 0.374±0.007 +18.5% |
| MDM [13] | +0 X | 0.418±0.005 | 0.604±0.005 | 0.703±0.005 | 3.658±0.025 | 9.546±0.096 | 0.501±0.037 |
| +MoDiPO-T [9] | +121K X | 0.424±0.006 | 0.613±0.005 | 0.706±0.004 +0.42% | 3.654±0.026 | 9.531±0.073 +0.16% | 0.490±0.035 |
| +MoDiPO-X [9] | +121K X | 0.420±0.006 | 0.632±0.005 | 0.704±0.001 +0.14% | 3.641±0.025 +0.46% | 9.495±0.071 +0.53% | 0.486±0.031 +2.99% |
| MDM (fast) [13] | +0 X | 0.455±0.006 | 0.645±0.007 | 0.728±0.005 | 3.304±0.023 | 9.948±0.084 | 0.534±0.025 |
| +SoPo (Ours) | +60 X | 0.479±0.006 +5.27% | 0.674±0.005 +4.50% | 0.770±0.006 +2.80% | 3.208±0.025 +2.91% | 9.906±0.091 +0.042 | 0.480±0.016 +10.1% |

| Methods | Year | R-Precision↑ | | | | MM Dist↓ | Diversity → | Multimodal↑ | FID↓ |
|---|---|---|---|---|---|---|---|---|---|
| | | Top 1 | Top 2 | Top 3 | Avg. | | | | |
| Real | - | 0.511±0.003 | 0.703±0.003 | 0.797±0.002 | 0.670 | 2.794±0.008 | 9.503±0.065 | - | 0.002±0.000 |
| TEMOS [40] | 2022 | 0.424±0.002 | 0.612±0.002 | 0.722±0.002 | 0.586 | 3.703±0.008 | 8.973±0.071 | 0.368±0.018 | 3.734±0.028 |
| T2M [3] | 2022 | 0.457±0.002 | 0.639±0.003 | 0.740±0.003 | 0.612 | 3.340±0.008 | 9.188±0.002 | 2.090±0.083 | 1.067±0.002 |
| MDM [13] | 2022 | 0.418±0.005 | 0.604±0.005 | 0.703±0.005 | 0.575 | 3.658±0.025 | 9.546±0.096 | 2.799±0.072 | 0.501±0.037 |
| MLD [1] | 2023 | 0.481±0.003 | 0.673±0.003 | 0.772±0.002 | 0.642 | 3.196±0.016 | 9.724±0.082 | 2.413±0.079 | 0.473±0.013 |
| MotionGPT [42] | 2023 | 0.492±0.003 | 0.681±0.003 | 0.778±0.002 | 0.650 | 3.096±0.008 | 9.528±0.071 | 2.008±0.084 | 0.232±0.008 |
| MotionDiffuse [14] | 2024 | 0.491±0.001 | 0.681±0.001 | 0.782±0.001 | 0.651 | 3.113±0.001 | 9.410±0.049 | 1.553±0.042 | 0.630±0.001 |
| OMG [43] | 2024 | - | - | 0.784±0.002 | - | - | 9.657±0.085 | - | 0.381±0.008 |
| Wang et. al. [6] | 2024 | 0.433±0.007 | 0.629±0.007 | 0.733±0.006 | 0.598 | 3.430±0.061 | 9.546±0.096 | - | 0.352±0.109 |
| MoDiPO-T [9] | 2024 | 0.424±0.006 | 0.613±0.005 | 0.706±0.004 | - | 3.267±0.010 | 9.747±0.073 | 2.663±0.111 | 0.490±0.035 |
| PriorMDM [12] | 2024 | 0.481±0.001 | - | - | - | 5.610±0.023 | 9.620±0.074 | - | 0.590±0.053 |
| LMM-T [41] | 2024 | 0.496±0.002 | 0.685±0.002 | 0.785±0.002 | 0.655 | 3.087±0.012 | 9.176±0.074 | 1.465±0.048 | 0.415±0.002 |
| CrossDiff [11] | 2024 | - | - | 0.730±0.005 | - | 3.358±0.011 | 9.577±0.082 | - | 0.281±0.013 |
| Motion Mamba [5] | 2024 | 0.502±0.003 | 0.693±0.002 | 0.792±0.002 | 0.662 | 3.060±0.009 | 9.871±0.084 | 2.294±0.058 | 0.281±0.009 |
| MLD* [1, 2] | 2024 | 0.504±0.001 | 0.698±0.003 | 0.796±0.002 | 0.666 | 3.052±0.009 | 9.634±0.064 | 2.267±0.082 | 0.450±0.011 |
| MLD* [2], SoPo | 2025 | 0.528 +4.76% | 0.722 +3.44% | 0.827 +3.89% | 0.692 +3.90% | 2.939 +3.70% | 9.584 +38.1% | 2.301±0.076 | 0.174 +61.3% |