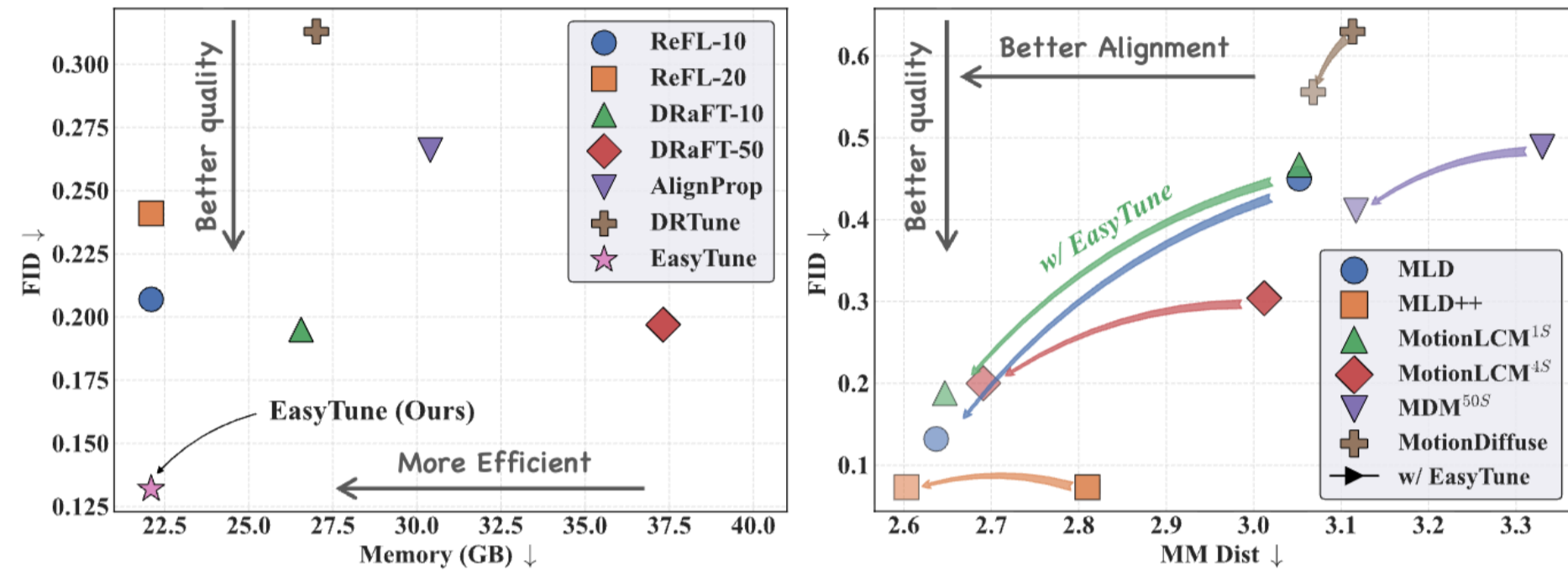


Introduction

Issue: Text-to-motion diffusion models trained on likelihood objectives misalign with downstream goals. While differentiable reward methods enable efficient preference fine-tuning without large-scale paired data, existing approaches suffer from (1) coarse-grained trajectory-level optimization, (2) O(T) memory overhead, and (3) vanishing gradients hindering early-step learning.



Key Observation: We theoretically and empirically identify that the root cause of these limitations is the recursive dependence between different steps in the denoising trajectory, leading to (1) O(T) memory complexity due to the entire computation graph, (2) delayed, sparse optimization signals, and (3) vanishing gradients.

TL;DR: We propose EasyTune, a reinforcement fine-tuning framework for diffusion models that decouples recursive dependencies and enables (1) dense and effective optimization, (2) memory-efficient training, and (3) fine-grained alignment.

Extension: We have released **MotionRFT**, a unified reinforcement fine-tuning for text-to-motion generation, with comprehensive code and documentation.

Motivation: Rethink differentiable reward

Recursive Dependence

Given a pre-trained motion diffusion model parameterized by ϵ_θ , the optimization objective is to fine-tune θ to maximize the reward value $\mathcal{R}_\phi(\mathbf{x}_0^t, c)$, with the loss defined as:

$$\mathcal{L}(\theta) = -\mathbb{E}_{c \sim \mathcal{D}_T, \mathbf{x}_0^t \sim \pi_\theta(\cdot|c)} [\mathcal{R}_\phi(\mathbf{x}_0^t, c)], \quad (1)$$

To compute the full gradient $\partial \mathcal{L}(\theta) / \partial \theta$, we unroll the $\partial \mathbf{x}_0^t / \partial \theta$ using Corollary 1 and substitute it into Eq. (3) resulting in (see proof in App. C.3):

$$\frac{\partial \mathcal{L}(\theta)}{\partial \theta} = -\mathbb{E}_{c \sim \mathcal{D}_T, \mathbf{x}_0^t \sim \pi_\theta(\cdot|c)} \left[\frac{\partial \mathcal{R}_\phi(\mathbf{x}_0^t)}{\partial \mathbf{x}_0^t} \cdot \prod_{s=1}^{t-1} \left(\frac{\partial \pi_\theta(\mathbf{x}_s, s, c)}{\partial \mathbf{x}_s^s} \right) \cdot \left(\frac{\partial \pi_\theta(\mathbf{x}_0^t, t, c)}{\partial \theta} \right) \right]. \quad (5)$$

tend to 0 when t is larger optimizing t-th step

Memory Complexity, Vanishing Gradients

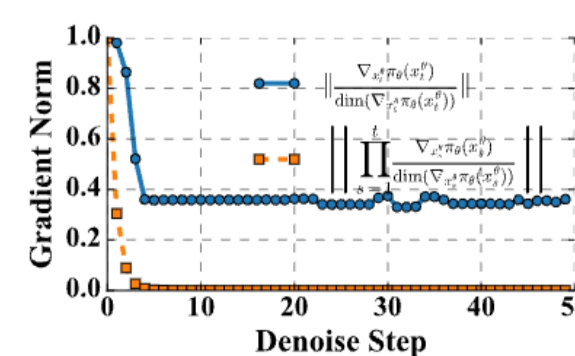


Fig. 4. Gradient norm with respect to denoising steps. Here, $\dim(\cdot)$ denotes the gradient dimension.

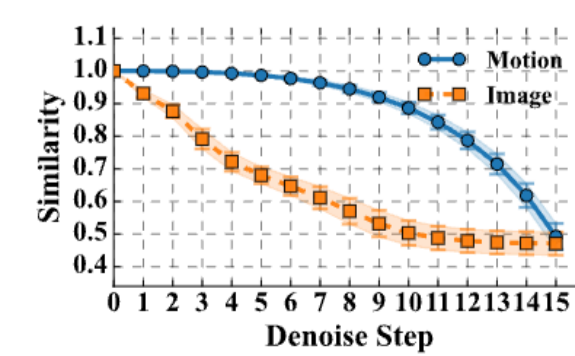


Fig. 5. Similarity between t-th step noised and clean motion. "w/o BP" indicates memory measured without back-propagation.

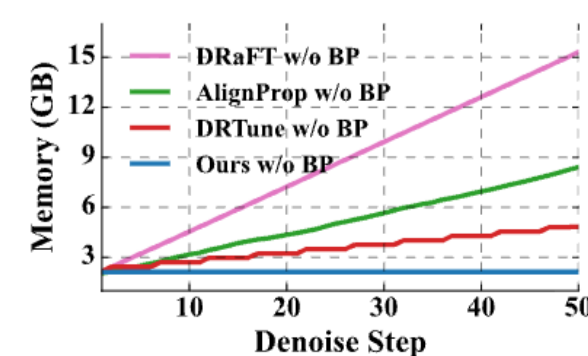
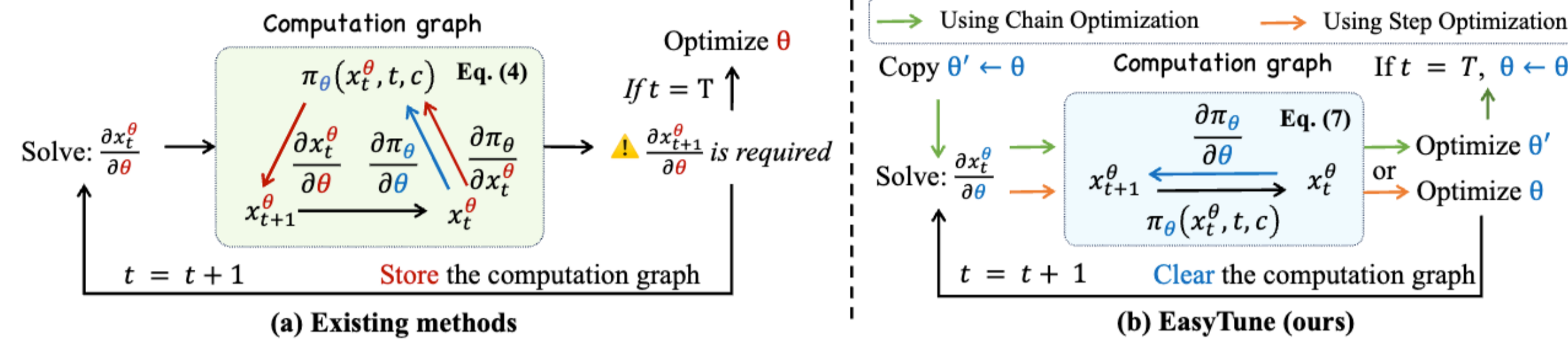


Fig. 6. Memory usage comparison. Here, "w/o BP" indicates memory measured without back-propagation.

Framework & Solution

Efficient Fine-Tune Method (EasyTune)



Step-aware Reward Model & Self-Preference Learning

Reward Model. Given a motion \mathbf{x} and a text description c , the reward value is computed based on the similarity between the motion features \mathbf{x} and text features c , denoted as:

$$\mathcal{R}_\phi(\mathbf{x}, c) = \mathcal{E}_M(\mathbf{x}) \cdot \mathcal{E}_T(c) \cdot \tau, \quad (11)$$

where \mathcal{E}_M and \mathcal{E}_T are the motion and text encoders from the pre-trained retrieval model (Weng et al., 2025a), and τ is a trainable temperature parameter.

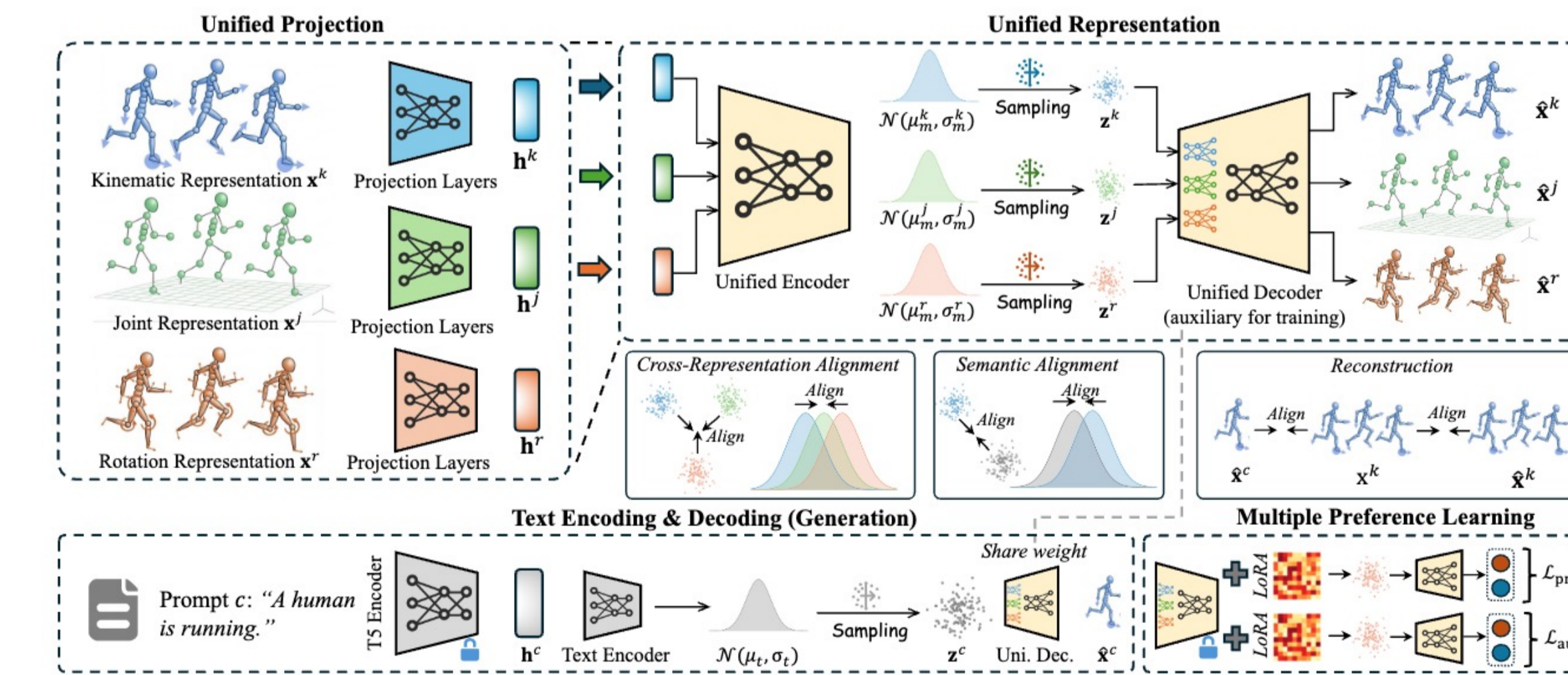
$$\mathcal{R}_\phi(\mathbf{x}_t, t, c) = \begin{cases} \mathcal{R}_\phi(\hat{\mathbf{x}}_0, 0, c), & \text{Only for ODE-based settings,} \\ \mathcal{R}_\phi(\mathbf{x}_t, t, c), & \text{For SDE- and ODE-based settings.} \end{cases} \quad (12)$$

To optimize the reward model ϕ , we minimize the KL divergence between them:

$$\mathcal{L}_{\text{SPL}}(\phi) = D_{\text{KL}}(\mathcal{Q} \parallel \mathcal{P}) = \sum_{\mathbf{x} \in \{\mathbf{x}^w, \mathbf{x}^l\}} \mathcal{Q}(\mathbf{x}, c) \log \frac{\mathcal{Q}(\mathbf{x}, c)}{\mathcal{P}(\mathbf{x}, c)}. \quad (17)$$

MotionReward (Extension)

Unified Motion Reward Model for Multi-Representation & Multi-Reward



Contact

Email: txf0620@gmail.com

WeChat: txf_06_20

The first author is seeking Ph.D opportunity in 2027 Spring/Fall



First Author



Project Page (Extension)



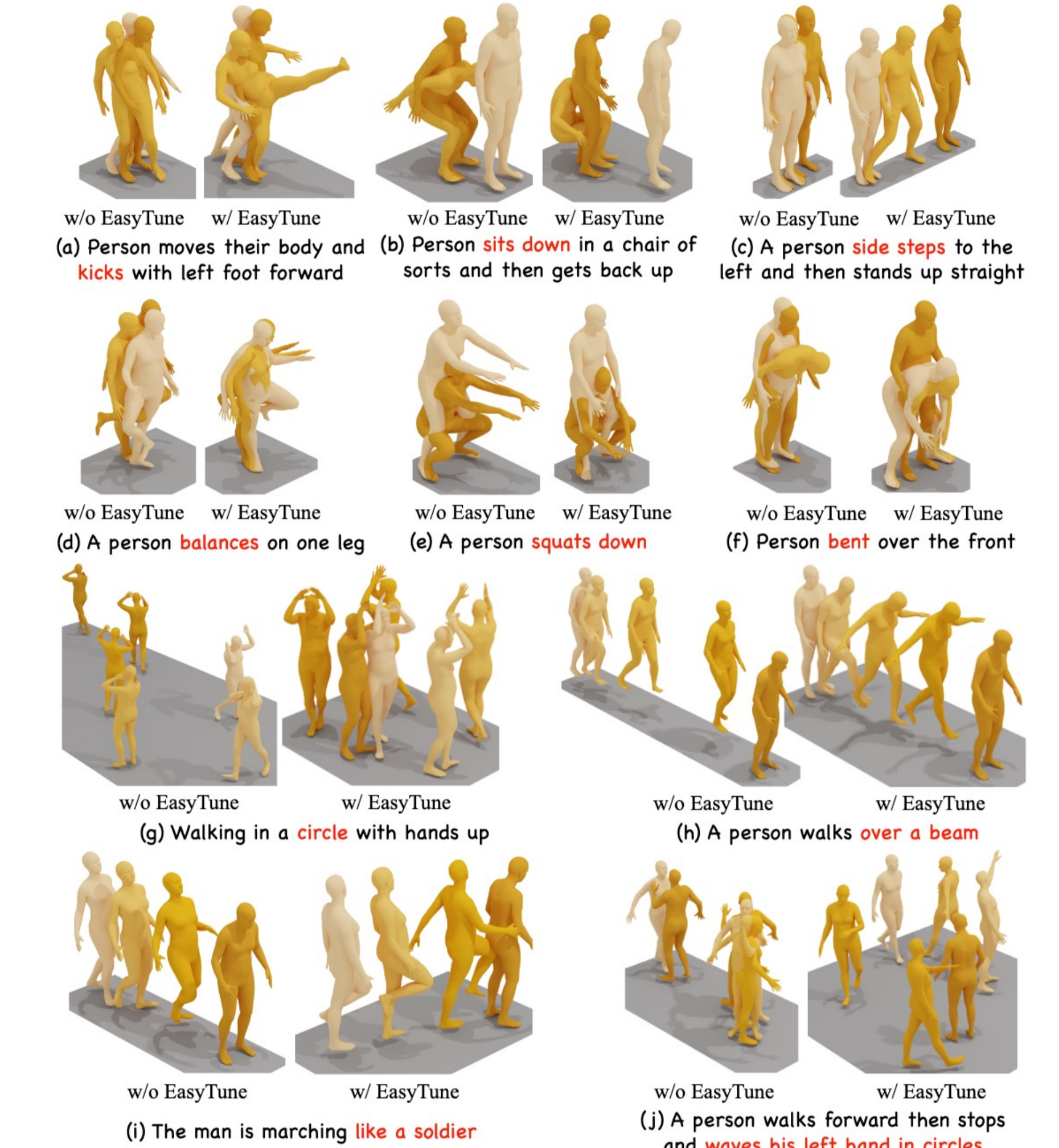
Code



Project Page (This work)

Experiments

Qualitative Results



Quantitative Results

Method	R Precision \uparrow			FID \downarrow	MM Dist \downarrow	Diversity \rightarrow	Memory (GB) \downarrow
	Top 1	Top 2	Top 3				
Real	0.511	0.703	0.797	0.002	2.974	9.503	-
MLD (Base Model)	0.504 \pm .002	0.698 \pm .003	0.796 \pm .002	0.450 \pm .011	3.052 \pm .009	9.634 \pm .064	15.21
w/ ReFL-10 (Clark et al., 2024)	0.533 \pm .58%	0.720 \pm 3.2%	0.821 \pm 3.1%	0.207 \pm 54.0%	2.852 \pm 6.6%	10.129 \pm .495	22.10 \pm 6.89
w/ ReFL-20 (Clark et al., 2024)	0.528 \pm 4.8%	0.718 \pm 2.9%	0.813 \pm 2.1%	0.241 \pm 46.4%	2.883 \pm 5.5%	10.189 \pm .555	22.10 \pm 6.89
w/ DRaFT-10 (Clark et al., 2024)	0.565 \pm 12.1%	0.757 \pm 8.5%	0.846 \pm 6.3%	0.195 \pm 56.7%	2.703 \pm 11.4%	9.851 \pm .217	26.56 \pm 11.35
w/ DRaFT-50 (Clark et al., 2024)	0.528 \pm 4.8%	0.724 \pm 3.7%	0.819 \pm 2.9%	0.197 \pm 56.2%	2.872 \pm 5.9%	9.641 \pm .007	37.32 \pm 22.11
w/ AlignProp (Prabhudesai et al., 2023)	0.560 \pm 11.1%	0.753 \pm 7.9%	0.841 \pm 5.7%	0.266 \pm 40.9%	2.739 \pm 10.3%	9.877 \pm .243	30.40 \pm 15.19
w/ DR Tune (Wu et al., 2025a)	0.549 \pm 8.9%	0.746 \pm 6.9%	0.836 \pm 5.0%	0.313 \pm 30.4%	2.795 \pm 8.4%	9.930 \pm .296	27.01 \pm 11.80
w/ EasyTune (Ours, Step Optimization)	0.581\pm15.3%	0.769\pm10.2%	0.855\pm7.4%	0.132\pm70.7%	2.637\pm13.6%	9.465\pm0.093	22.10\pm6.89
w/ EasyTune (Ours, Chain Optimization)	0.574 \pm 13.9%	0.766 \pm 9.7%	0.854 \pm 7.3%	0.172 \pm 61.8%	2.614\pm14.3%	9.348 \pm .024	24.21 \pm 9.00

Method	R Precision \uparrow			FID \downarrow	MM Dist \downarrow	Diversity \rightarrow
	Top 1	Top 2	Top 3			
Real	0.511	0.703	0.797	0.002	2.974	9.503
TM2T (Guo et al., 2022b)	0.424 \pm 0.003	0.618 \pm 0.003	0.729 \pm 0.002	1.501 \pm 0.017	3.467 \pm 0.011	8.589 \pm 0.076
T2M (Guo et al., 2022a)	0.455 \pm 0.002	0.636 \pm 0.003	0.736 \pm 0.003	1.087 \pm 0.002	3.347 \pm 0.008	9.175 \pm 0.002
MDM (Tevet et al., 2023)	0.455 \pm 0.006	0.645 \pm 0.007	0.749 \pm 0.006	0.489 \pm 0.047	3.330 \pm 0.025	9.920 \pm 0.083
T2M-GPT (Zhang et al., 2023a)	0.492 \pm 0.003	0.679 \pm 0.002	0.775 \pm 0.002	0.141 \pm 0.005	3.121 \pm 0.009	9.722 \pm 0.082
ReMoDiffuse (Zhang et al., 2023b)	0.510 \pm 0.005	0.698 \pm 0.006	0.795 \pm 0.004	0.103 \pm 0.004	2.974 \pm 0.016	9.018 \pm 0.075
Att2M (Zhang et al., 2023)	0.499 \pm 0.003	0.690 \pm 0.002	0.786 \pm 0.002	0.112 \pm 0.006	3.038 \pm 0.007	9.700 \pm 0.090
MotionDiffuse (Zhang et al., 2024a)	0.491 \pm 0.001	0.681 \pm 0.001	0.775 \pm 0.001	0.630 \pm 0.001	3.113 \pm 0.001	9.410 \pm 0.049
MotionLCM (Dai et al., 2024)	0.502 \pm 0.003	0.698 \pm 0.002	0.798 \pm 0.002	0.304 \pm 0.012	3.012 \pm 0.007	9.607 \pm 0.066
MotionMamba (Zhang et al., 2024b)	0.502 \pm 0.003	0.693 \pm 0.002	0.792 \pm 0.002	0.281 \pm 0.011	3.060 \pm 0.000	9.871 \pm 0.084
CoMo (Huang et al., 2024)	0.502 \pm 0.002	0.692 \pm 0.007	0.790 \pm 0.002	0.262 \pm 0.004	3.032 \pm 0.015	9.936 \pm 0.066
ParCo (Zou et al., 2024)	0.515 \pm 0.003	0.706 \pm 0.003	0.801 \pm 0.002	0.109 \pm 0.005	2.927 \pm 0.008	9.576 \pm 0.088
SoPo (Tan et al., 2025)	0.528 \pm 0.005	0.722 \pm 0.004	0.827 \pm 0.004	0.174 \pm 0.005	2.939 \pm 0.011	9.584 \pm 0.074
MLD (Chen et al., 2023) (Base Model)	0.504 \pm 0.002	0.698 \pm 0.003	0.796 \pm 0.002	0.450 \pm 0.011	3.052 \pm 0.009	9.634 \pm 0.064
w/ EasyTune (Ours)	0.581\pm15.3%	0.769\pm10.2%	0.855\pm7.4%	0.132\pm70.7%	2.637\pm13.6%	9.465\pm0.093
MLD++ (Dai et al., 2025) (Base Model)	0.548 \pm 0.003	0.738 \pm 0.003	0.829 \pm 0.002	0.073 \pm 0.003	2.810 \pm 0.008	9.658 \pm 0.089
w/ EasyTune (Ours)	0.591\pm10.04%	0.777\pm15.3%	0.859\pm10.02%	0.069\pm0.003	2.592\pm7.8%	9.705\pm0.086