

# MoLA: Motion Generation and Editing with Latent Diffusion Enhanced by Adversarial Training

Kengo Uchida<sup>1</sup> Takashi Shibuya<sup>1</sup> Yuhta Takida<sup>1</sup> Naoki Murata<sup>1</sup>  
Julian Tanke<sup>1</sup> Shusuke Takahashi<sup>2</sup> Yuki Mitsufuji<sup>1,2</sup>

<sup>1</sup>Sony AI, <sup>2</sup>Sony Group Corporation

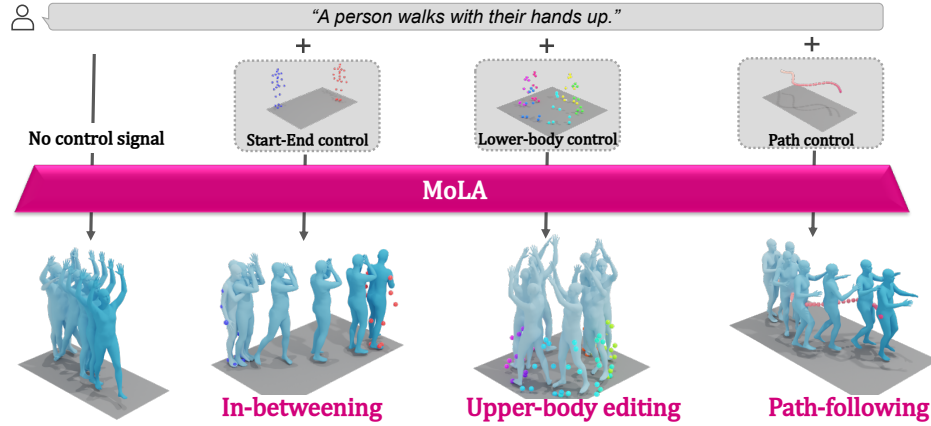


Figure 1. MoLA achieves fast and high-quality human motion generation given textual descriptions while enabling motion editing applications. With MoLA, we can deal with various types of motion editing tasks in a single framework.

## Abstract

In text-to-motion generation, controllability as well as generation quality and speed has become increasingly critical. The controllability challenges include generating a motion of a length that matches the given textual description and editing the generated motions according to control signals, such as the start-end positions and the pelvis trajectory. In this paper, we propose MoLA, which provides fast, high-quality, variable-length motion generation and can also deal with multiple editing tasks in a single framework. Our approach revisits the motion representation used as inputs and outputs in the model, incorporating an activation variable to enable variable-length motion generation. Additionally, we integrate a variational autoencoder and a latent diffusion model, further enhanced through adversarial training, to achieve high-quality and fast generation. Moreover, we apply a training-free guided generation framework to achieve various editing tasks with motion control inputs. We quantitatively show the effectiveness of adversarial learning in text-to-motion generation, and demonstrate the applicability of our editing framework to multiple

editing tasks in the motion domain. Our code is available at <https://github.com/sony/MoLA>.

## 1. Introduction

Human motion synthesis from text is an emerging task with highly relevant applications in fields such as multimedia production and computer animation. For example, an animator might wish to create or edit a motion prototype to verify their artistic intent before time-consuming animation or motion capture commences, generate smooth in-between motion between two motion capture clips, or generate specific motion that follows a predefined trajectory. In Figure 1, we demonstrate results of our approach on such motion generation and editing tasks. To be useful in those real-world applications, a method has to excel in three domains: (1) motion quality, which encompasses both the general motion quality and the adherence to the textual description; (2) fast inference time; (3) efficient motion editing.

Several recent works have attempted to address these desired properties: Methods based on vector quantization (VQ), such as T2M-GPT [44], MoMask [13], MMM [29],

ParCo [49], and BAMB [28], achieve impressive generation quality by compressing human motion into discrete tokens and then sampling those tokens to synthesize motion. Diffusion-based methods in data space, such as MDM [37], provide impressive flexibility with regard to quality and motion editability. To further improve motion quality and inference time, latent-space based methods such as MLD [3] or MotionMamba [47] operate in a learned continuous latent space.

However, no state-of-the-art method excels in all three domains necessary for real-world applications. For example, while latent-space based approaches such as VQ-based methods or MLD achieve impressive motion quality and fast inference time, they cannot edit a given motion sequence in a training-free manner. In contrast, data-space based methods such as MDM are able to edit motion in a training-free manner, which, however, comes at the cost of slow inference time and lower generation quality. Moreover, while these models excel at motion generation, they require users to manually specify the motion length instead of automatically determining it from the textual input. This often necessitates length estimation or iterative adjustment, which limits their flexibility. For instance, MoMask [13] uses a length estimator during inference to predict motion length based on text input. However, inaccurate predictions may result in significant motion drift [40].

To close this gap, we propose **MoLA**, **M**otion **G**eneration and **E**ditin**G** with **L**atent **D**iffusion **E**nhan**C**ed by **A**dversarial **T**raining. We revisit the representation of motion features used as inputs and outputs of the model and introduce an activation variable that characterizes the length of the motion. We utilize a variational auto-encoder (VAE) [20] for its continuous latent space, which allows training-free editing. This is in contrast to discrete latent spaces that do not allow for training-free motion editing. We also enhance the VAE training with adversarial learning, which has been shown to work well in the image [8, 16, 32] and audio [23] domains. We empirically show that our model achieves variable-length motion generation aligned with textual descriptions and high generation performance on a commonly used dataset [12]. We also demonstrate that training-free guided diffusion can enable multiple motion editing tasks such as path-following, in-betweening, and upper body editing. These experiments also highlight the significant improvement in speed and performance of our model compared to the performance of existing training-free motion editing models (see Figure 2).

The main contributions of this paper are threefold. First, we introduce an activation variable into the motion representation and show that our method can perform variable-length motion generation conditioned on text. Second, we propose a new continuous latent-based motion generation model that introduces adversarial training into the motion

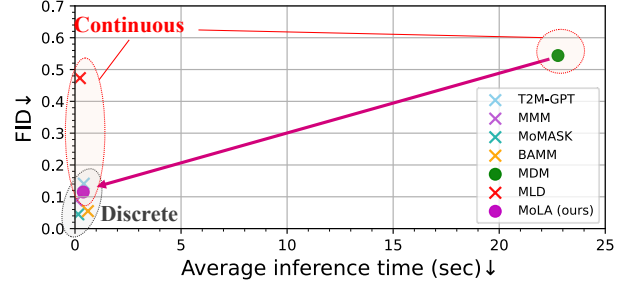


Figure 2. Comparison of inference cost, generation performance, and editability for text-to-motion methods on HumanML3D dataset (All tests are performed on NVIDIA A100 GPU). • means a method that can edit motion in a training-free manner, and × means a method that cannot edit motion in a training-free manner. The pink arrow in the figure indicates that MoLA significantly extends the performance boundaries (in terms of generation quality and speed) of methods categorized as enabling training-free editing.

VAE and quantitatively show that it significantly pushes the limits of existing continuous-based methods. Finally, we demonstrate that our model not only shows high generation performance but also can deal with various types of motion editing tasks in a training-free manner.

## 2. Related Work

### 2.1. Motion Generation

Text-to-motion generation technology has made rapid progress with the diffusion models and VQ-based models. MDM [37] and MotionDiffuse [46] adopt diffusion models for motion generation, which leads to better performance in terms of generation quality. However, these methods directly apply diffusion processes to raw motion sequences, thus resulting in slow generation. To address this issue, MLD [3] adopts a diffusion model in a low-dimensional latent space provided by a VAE trained on motion data, drawing inspiration from latent diffusion models [32]. VQ-based models have been studied as well, inspired by the success of VQ-based models in image generation [2, 10, 42]. In this approach, a VQ-VAE model is first trained on motion data to acquire discrete motion representations, and deep generative models are then applied to generate sequences of discrete representations (also called tokens). T2M-GPT [44], AttT2M [48], MotionGPT [17] and ParCo [49] utilize autoregressive (AR) models to generate motion tokens. However, AR models are slow in inference because motion tokens are generated sequentially. To address this issue, M2DM [21] and DiverseMotion [24] apply discrete diffusion models to motion tokens in a latent space, whereas MMM [29], MoMask [13] and BAMB [28] adopt a mask prediction model.

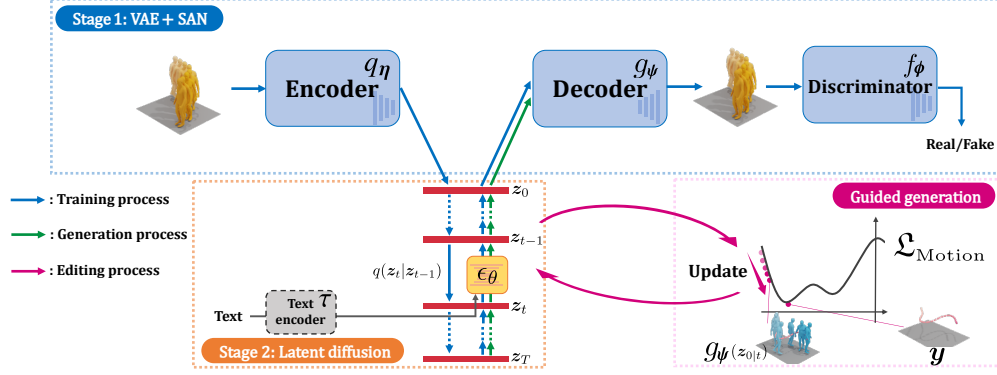


Figure 3. The overall framework of MoLA. Stage 1: A motion VAE enhanced by adversarial training learns a low-dimensional latent representation of diverse motion sequences. Stage 2: A text-conditioned latent diffusion model leverages this representation for fast and high-quality text-to-motion generation. Guided generation: During inference, a gradient-based method minimizes a loss function  $\mathcal{L}_{\text{Motion}}$  for each desired editing task, enabling multiple motion editing tasks within a unified framework.

## 2.2. Motion Editing

Motion editing has attracted much research interest as well. MDM [37] demonstrated upper body editing and motion in-betweening by applying diffusion inpainting to motion data in both the spatial and temporal domains. LGD [34] demonstrated path-following motion generation with a guided diffusion that utilizes multiple samples from a suitable distribution to reduce bias. GMD [18] guides the position of the root joint to control motion trajectories. OmniControl [41] controls any joints at any time by guiding a pre-trained motion diffusion model with an analytic function. DNO [19] optimizes the diffusion latent noise of a pre-trained text-to-motion model with user-provided criteria in the motion space, achieving multiple editing tasks. However, these methods employ data-space diffusion models similar to MDM and have not been demonstrated with a latent diffusion model. MMM [29] demonstrated motion editing by placing masked tokens where editing is needed and applying a mask prediction framework, and MotionLCM [6] proposed a fast controllable motion generation framework by introducing latent consistency distillation and the motion ControlNet [45] manipulation in the latent space.

## 3. Method

The goal of this study is to develop a framework for fast and high-quality text-guided motion generation and to deal with multiple control tasks in a training-free way. To achieve this, we propose the following training and inference tricks (I)-(IV): (I) We review the representation of motion features to achieve improved accuracy and variable length motion generation. We select the necessary features to reduce the burden on our model’s VAE encoder and add an activation variable to the representation to determine the motion length (Section 3.1). (II) We train a motion

VAE enhanced by adversarial training to achieve high quality generation performance (Section 3.2). (III) To reduce computational complexity while simultaneously enabling high-quality text-driven motion generation, we train a text-conditioned diffusion model on the low-dimensional latent space obtained by VAE model (Section 3.3). (IV) Guided generation: Adopting training-free guided generation in inference enables multiple editing functions required in motion generation without additional training (Section 3.4). The outline of our text-to-motion generation model, MoLA, is shown in Figure 3.

### 3.1. Motion representation

We build upon the motion representation used in [12, 30] and improve the representation to enable variable length generation and improve performance in our framework. The pose representation used in [12, 30] is expressed as  $\mathbf{m} \in \mathbb{R}^{(4+12N_j+4) \times L}$ , where  $N_j$  is the number of joints. A motion can be represented as a sequence of this representation. The  $i$ -th ( $1 \leq i \leq L$ ) pose is defined as follows:

$$\mathbf{m}^i = [\dot{r}_a^i, \dot{r}_x^i, \dot{r}_z^i, r_y^i, (\mathbf{j}_p^i)^\top, (\mathbf{j}_r^i)^\top, (\mathbf{j}_v^i)^\top, (\mathbf{c}^i)^\top]^\top \quad (1)$$

where  $\dot{r}_a^i \in \mathbb{R}$  is root angular velocity along the  $Y$ -axis,  $\dot{r}_x^i$  and  $\dot{r}_z^i \in \mathbb{R}$  are root linear velocities on  $XZ$ -plane,  $r_y^i$  is root height,  $\mathbf{j}_p^i \in \mathbb{R}^{3(N_j-1)}$  is local joints positions,  $\mathbf{j}_r^i \in \mathbb{R}^{6(N_j-1)}$  is rotations in root space,  $\mathbf{j}_v^i \in \mathbb{R}^{3N_j}$  is velocities and  $\mathbf{c}^i \in \mathbb{R}^4$  is binary foot-ground contact features by thresholding the heel and toe joint velocities. To achieve variable-length motion generation during stage 2 inference (explained in Section 3.3), we concatenate an activation variable  $a^i \in \mathbb{R}$  for the  $i$ -th pose to Equation (1) as  $[(\mathbf{m}^i)^\top, a^i]^\top$ . Additionally, to reduce the burden on the encoder in Section 3.2, we remove the substantially redundant information  $\mathbf{j}_v$  and  $\mathbf{c}$  from Equation (1) as  $\tilde{\mathbf{m}}^i = [\dot{r}_a^i, \dot{r}_x^i, \dot{r}_z^i, r_y^i, (\mathbf{j}_p^i)^\top, (\mathbf{c}^i)^\top]^\top$ , and then we set the vector

for the encoder input in Section 3.2 as follows:

$$\mathbf{x}^i = [(\tilde{\mathbf{m}}^i)^\top, \mathbf{a}^i]^\top = [\dot{r}_a^i, \dot{r}_x^i, \dot{r}_z^i, \dot{r}_y^i, (\mathbf{j}_p^i)^\top, (\mathbf{j}_r^i)^\top, \mathbf{a}^i]^\top. \quad (2)$$

We employ  $\mathbf{x} = [\mathbf{x}^1, \dots, \mathbf{x}^L] \in \mathbb{R}^{N \times L}$  as motion representation in our model training, where  $N = 4 + 9N_j + 1$ .<sup>1</sup>

### 3.2. Stage 1: Continuous Motion Latent Representation with Adversarial Training

#### 3.2.1. Learning a motion latent representation with VAE-GAN

We propose a motion variational autoencoder (VAE), enhanced by adversarial training, to learn a low-dimensional latent representation for diverse human motion sequences. Assume an observed motion  $\mathbf{x} \in \mathbb{R}^{N \times L}$ , where  $N$  and  $L$  denote raw motion data dimension per frame and motion length, respectively. To learn such a latent representation, we first define a latent variable  $\mathbf{z} \in \mathbb{R}^{d_z \times d_l}$ , which is assumed to generate data sample  $\mathbf{x}$ , where  $d_z, d_l \in \mathbb{N}$ . The generative process is modeled as  $\mathbf{x} \sim p_\psi(\mathbf{x}|\mathbf{z})$  with a prior  $p(\mathbf{z})$ . The prior is assumed to be a standard Gaussian distribution, i.e.,  $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$ . We model the conditional distribution as a Gaussian distribution:  $p_\psi(\mathbf{x}|\mathbf{z}) = \mathcal{N}(g_\psi(\mathbf{z}), \sigma^2 \mathbf{I})$  with  $g_\psi : \mathbb{R}^{d_z \times d_l} \rightarrow \mathbb{R}^{N \times L}$  and  $\sigma^2 \in \mathbb{R}_+$ , where  $\mathbb{R}_+$  indicates the set of all positive real numbers. As in a usual VAE, we introduce an approximated posterior, which is modeled by  $q_\eta(\mathbf{z}|\mathbf{x}) : \mathbb{R}^{N \times L} \rightarrow \mathbb{R}^{d_z \times d_l}$ . As a result, the VAE consists of an encoder and a decoder, parameterized by  $\eta$  and  $\psi$ , respectively. The objective function for the VAE is formulated as the negative evidence lower bound (negative ELBO) per sample  $\mathbf{x}$ , which is a weighted summation of mean squared error and KL regularization terms:

$$\mathcal{J}_{\text{VAE}}(\psi, \eta; \mathbf{x}) = \mathbb{E}_{q_\eta(\mathbf{z}|\mathbf{x})} [\|\mathbf{x} - g_\psi(\mathbf{z})\|_2^2] + \lambda_{\text{reg}} D_{\text{KL}}(q_\eta(\mathbf{z}|\mathbf{x}) \| p(\mathbf{z})), \quad (3)$$

where  $\lambda_{\text{reg}}$  is a hyperparameter for balancing the two terms. The encoder  $q_\eta$  can be trained to produce a low-dimensional latent representation, and the decoder  $g_\psi$  can also be trained to accurately reconstruct the input motion sequences from the latent representations. Besides, we separate the decoder output into two components, denoted as  $g_\psi(\mathbf{z}) = [(\mathbf{m}'(\mathbf{z}))^\top, \mathbf{a}'(\mathbf{z})]^\top$ , where  $\mathbf{m}'(\mathbf{z}) \in \mathbb{R}^{(N-1) \times L}$  and  $\mathbf{a}'(\mathbf{z}) \in \mathbb{R}^L$  correspond to the original motion and proposed activation variable. We adopt a binary cross entropy (BCE) loss for the activation variable. The following  $\mathcal{J}_{\text{MotionVAE}}$  is used as the loss function for the VAE part:

$$\mathcal{J}_{\text{MotionVAE}}(\psi, \eta; \mathbf{x}) = \mathbb{E}_{q_\eta(\mathbf{z}|\mathbf{x})} [\|\mathbf{m} - \mathbf{m}'(\mathbf{z})\|_2^2] + \lambda_{\text{act}} \mathcal{L}_{\text{BCE}}(\mathbf{a}, \mathbf{a}'(\mathbf{z})) + \lambda_{\text{reg}} D_{\text{KL}}(q_\eta(\mathbf{z}|\mathbf{x}) \| p(\mathbf{z})), \quad (4)$$

<sup>1</sup>In practice, during training, we pad zeros to  $\mathbf{m}$  or  $\tilde{\mathbf{m}}$  in inactive frames used to align sequence lengths. Then,  $\mathbf{a}^i = 0$  is assigned to padded frames, and  $\mathbf{a}^i = 1$  is assigned to frames containing motion data.

where  $\lambda_{\text{act}}$  is a hyperparameter for weighting the BCE loss term.

To achieve high-quality generation, we need to push the limits of compression. Hence, we propose incorporating adversarial training into the motion VAE (c.f. [23, 32]). More specifically, we introduce a discriminator, denoted as  $f_\phi : \mathbb{R}^{N \times L} \rightarrow \mathbb{R}$ , that aims to distinguish real and reconstructed motions. The adversarial training is formulated as a two-player optimization between the VAE and the discriminator. The discriminator is trained by the maximization of  $\mathbb{E}_{p(\mathbf{x})} \mathcal{L}_{\text{GAN}}(\phi; \psi, \eta, \mathbf{x})$  with respect to  $\phi$ , where

$$\mathcal{L}_{\text{GAN}}(\phi; \psi, \eta, \mathbf{x}) = \min\{0, -1 + f_\phi(\mathbf{x})\} + \mathbb{E}_{q_\eta(\mathbf{z}|\mathbf{x})} [\min\{0, -1 - f_\phi(g_\psi(\mathbf{z}))\}]. \quad (5)$$

We formulate the overall loss for the VAE as the sum of the negative ELBO and adversarial loss,

$$\min_{\phi, \psi} \mathbb{E}_{p(\mathbf{x})} [\mathcal{J}_{\text{MotionVAE}}(\psi, \eta; \mathbf{x}) + \lambda_{\text{adv}} \mathcal{J}_{\text{GAN}}(\psi, \eta; \phi, \mathbf{x})], \quad (6)$$

where  $\lambda_{\text{adv}}$  is a positive scalar that adjusts the balance between the two terms, and

$$\mathcal{J}_{\text{GAN}}(\psi, \eta; \phi, \mathbf{x}) = -\mathbb{E}_{q_\eta(\mathbf{z}|\mathbf{x})} [f_\phi(g_\psi(\mathbf{z}))]. \quad (7)$$

We train both the VAE and the discriminator with Equations (5) and Equation (6) in an alternating way.

#### 3.2.2. From GAN- to SAN-based discriminator

We apply the slicing adversarial network (SAN) framework [36] to further enhance the motion VAE, based on a prior report showing SAN-based models perform better than the GAN counterparts. The impact of this replacement on motion generation is discussed in Section 4.

#### 3.2.3. Architectures

We use a standard CNN-based architecture in the motion VAE encoder  $q_\eta$ , decoder  $g_\psi$  and discriminator  $f_\psi$ , consisting of 1D convolution, a residual block, and Leaky ReLU. For temporal downsampling and upsampling, we use stride 2 convolution and nearest interpolation, respectively. Specifically, the motion sequence  $\mathbf{x} \in \mathbb{R}^{N \times L}$  is encoded into a latent vector  $\mathbf{z} \in \mathbb{R}^{d_z \times d_l}$  with downsampling ratio of  $d_l = L/4$ . This architecture is inspired by [8, 44].

### 3.3. Stage 2: Motion Latent Diffusion

#### 3.3.1. Text-conditional motion generation

In this section, we train a text-conditioned diffusion model on the low-dimensional motion latent space obtained by the autoencoder learned in the stage 1 (Section 3.2). Using the trained model, we perform motion generation conditioned on text. First, we define a time-dependent sequence  $\mathbf{z}_0, \mathbf{z}_1, \dots, \mathbf{z}_t, \dots, \mathbf{z}_T \in \mathbb{R}^{d_z \times d_l}$  (starting from the



VAE encoder output  $z_0 = z \sim q_\eta(z|x)$ , which is derived from the following Markov diffusion process in the latent space:  $q(z_t|z_{t-1}) = \mathcal{N}(\sqrt{\alpha_t}z_{t-1}, (1 - \alpha_t)\mathbf{I})$ , where  $T > 0$  and the constant  $\alpha_t \in (0, 1)$  is a pre-defined noise-scheduling parameter that determines the forward process. The forward process allows for the sampling of  $z_t$  at an arbitrary time step  $t$  in a closed form:  $z_t = \sqrt{\bar{\alpha}_t}z_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon$ , where  $\bar{\alpha}_t := \prod_{s=1}^t \alpha_s$  and  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ .

As our goal is text-to-motion generation, our interest is in the conditional distribution  $p(z|c)$  given the text prompt  $c$ . Here, similar to many text-conditioned latent diffusion models [3, 32], we train the conditional model  $\epsilon_\theta(z_t, t, \tau(c))$  conditioned on the output of a text encoder  $\tau(c)$ , using the following objective function:

$$\mathcal{J}_{\text{CLDM}}(\theta) = \mathbb{E}_{\{z_0, c\}, \epsilon, t} [\|\epsilon - \epsilon_\theta(z_t, t, \tau(c))\|_2^2], \quad (8)$$

where  $z_0$  and  $c$  are drawn from the joint empirical distribution. In addition, as done in prior works, we adopt classifier-free guidance [15] and train the model unconditionally, i.e., without a text prompt, with a certain probability during training.

During inference, the trained diffusion model  $\epsilon_\theta(z_t, t, \tau(c))$  is used to generate  $z_0$  through a denoising process conditioned on the text prompt  $c$ . We adopt the sampling scheme of DDIM [33] with trailing sample steps [22], in which each sampling step is defined as:

$$\begin{aligned} z_{t-1} = & \sqrt{\bar{\alpha}_{t-1}} \left( \frac{z_t - \sqrt{1 - \bar{\alpha}_{t-1}}\epsilon_\theta(z_t, t, \tau(c))}{\sqrt{\bar{\alpha}_t}} \right) \\ & + \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2\epsilon_\theta(z_t, t, \tau(c))} + \sigma_t\epsilon, \end{aligned} \quad (9)$$

where  $\sigma_t > 0$  determines the stochasticity of the sampling process, and the sampling process becomes deterministic when  $\sigma_t = 0$ . The part  $(z_t - \sqrt{1 - \bar{\alpha}_{t-1}}\epsilon_\theta(z_t, t, \tau(c)))/\sqrt{\bar{\alpha}_t}$  in the first term corresponds to a direct estimate of the clean latent  $z_0$  from the noisy sample  $z_t$  using the diffusion model based on Tweedie’s formula [7]; this estimate is denoted as  $z_{0|t}$ . In actual inference, we use the estimated  $z_0$  and the VAE decoder trained in stage 1 to obtain  $g_\psi(z_0)$  as the generated motion sequences. Variable length motion generation is then achieved by clipping a part of  $g_\psi(z_0)$  where the activation variable  $a^i$  introduced in Section 3.1 satisfies  $a^i < \delta$  ( $0 < \delta < 1$ ). The effect of this approach is discussed in Section 4.1.

### 3.3.2. Architecture

We employ a diffusion transformer (DiT)-based architecture in our stage 2 model. The transformer used follows a standard structure of stacked blocks consisting of an attention layer and gated multilayer perceptrons (MLP) connected in series, with skip connections around each. Additionally, layer normalization is employed on the inputs of

both the attention layer and the MLP. At the input and output of the transformer, linear mapping is used to convert from the latent dimension of the stage 1 model to the embedded dimension of the transformer. Text CLIP [31] embedding is also added as input to the transformer, along with an embedding describing the current time step of the diffusion process. This architecture is inspired by [9], which established SOTA performance in audio generation.

### 3.4. Controllable Motion Generation on Latent Diffusion Sampling

In this section, we present a guided generation framework that leverages the pre-trained motion latent diffusion model (Section 3.3) for conditional motion generation and editing tasks without extra training. Major training-free methods (e.g. [4, 14, 43]) are based on the fact that the conditional score function can be decomposed into two additive terms: the unconditional score function and the log-likelihood term. Specifically, for a new condition  $y$ , we have  $\nabla_{z_t} \log p(z_t|c, y) = \nabla_{z_t} \log p(z_t|c) + \nabla_{z_t} \log p(y|z_t, c)$ , which is derived from Bayes’ rule. The conditional generation based on the aforementioned property can be regarded as a sequential procedure implemented as follows. First, a denoised sample  $z_{t-1}$  is obtained from  $z_t$  by the sampling step in Equation (9), without considering the given condition  $y$  (the first term in Equation (9)). Subsequently,  $z_{t-1}$  is further updated using the gradient of the log-likelihood term with respect to  $z_t$  (the second term).

The challenge here is that the log-likelihood term is computed based on noisy samples  $z_t$ . In the classifier-guided diffusion [35], time-dependent classifiers for  $z_t$  have to be trained, which requires additional training. In contrast, in training-free methods, this term is approximated with the current clean data estimate  $z_{0|t}$  and a loss function  $\mathcal{L}(x; y)$  defined for clean data. This loss function can be flexibly set depending on the task. For example, in inverse problems of the form  $y = \mathcal{A}(x)$ , where  $\mathcal{A}$  is a differentiable function with respect to  $x$ , the loss function can be set as  $\|y - \mathcal{A}(x_{0|t})\|_2^2$ , where  $x_{0|t} = g_\psi(z_{0|t})$  is a clean data estimate in the original data domain and obtained through the VAE decoder. Here, we adopt MPGD [14], a fast yet high-quality guidance method applicable to latent diffusion models. Following the denoising step, it updates the denoised sample based on the loss function  $\mathcal{L}(x; y)$  as follows:

$$z_{t-1} \leftarrow z_{t-1} - \rho_t \sqrt{\bar{\alpha}_{t-1}} \nabla_{z_{0|t}} \mathcal{L}(g_\psi(z_{0|t}); y), \quad (10)$$

where  $\rho_t$  is a time-dependent step size parameter.

More specifically, in editing motion tasks we generate motions to match given specific poses or trajectory control signals. To deal with various motion editing tasks in this framework, the loss function is designed as follows:

$$\mathcal{L}_{\text{Motion}}(g_\psi(z_{0|t}); y) = \sum_n \sum_l m_{nl} \|\Re(g_\psi(z_{0|t}))_{nl} - y_{nl}\|_2, \quad (11)$$

where  $n$  and  $l$  are indices of joint and frame, respectively, and  $m_{nl}$  is a binary value indicating whether the control position  $y_{nl}$  contains a valid value at frame  $l$  for joint  $n$ , and  $\mathcal{R}(\cdot)$  is a function that converts the motion features including the joint’s local positions to global absolute locations. We set  $\mathcal{L}_{\text{Motion}}$  for loss function  $\mathcal{L}$  in the guided generation to measure the distance between desired constraints  $\mathbf{y}$  and the joint locations of the generated motion. Target locations as constraint  $\mathbf{y}$  can be specified for any subset of joints in any subset of motion frames.<sup>2</sup> Editing a generated motion to match specific poses or follow a specific trajectory is achieved by minimizing  $\mathcal{L}_{\text{Motion}}$  using the update rule in Equation (10).

## 4. Experiments

We evaluate the performance of MoLA on two tasks: motion generation (Section 4.1) and motion editing (Section 4.2). Our results demonstrate that MoLA achieves its three key objectives: (1) fast and high-quality generation, (2) variable-length generation, and (3) multiple motion editing tasks in a training-free manner. To validate these objectives, we utilize the HumanML3D [12].<sup>3</sup>

### 4.1. Motion Generation

#### 4.1.1. Performance comparison with other text-to-motion models

We evaluate our proposed MoLA in comparison to current SOTA methods [3, 6, 13, 21, 24, 28, 29, 37, 38, 44, 46, 48, 49] using five metrics (R-Precision, Fréchet Inception Distance (FID), Multi-modal distance (MMDist), Diversity, and MultiModality (MModality)) proposed by Guo *et al.* [12]. For evaluation, we select the model that achieves the best FID, which is a metric that evaluates the overall motion quality, on the validation set and report its performance on the test set of HumanML3D. We show the results in Table 1. The methods are organized into three groups: i) those using VQ-based latent representations (**Discrete**), ii) those using data-space diffusion model (**Continuous (raw data)**), and iii) those using VAE-based latent representations (**Continuous (latent)**). The discrete approaches perform well in motion generation (e.g., [13, 24, 29]). However, those models cannot control an arbitrary set of joints in a training-free manner [29] as we have discussed so far. MDM [37], MotionDiffuse [46] and Fg-T2M [38] adopt

<sup>2</sup>As examples of motion editing, if  $\mathbf{y}$  is given as the start-end positions, we can handle the motion in-betweening. If  $\mathbf{y}$  is given as the lower body positions, we can edit the upper body corresponding to the lower one. If  $\mathbf{y}$  is given as the pelvis trajectory, it corresponds to the path-following task. We leave the task details to Section 4.2. Note that the guided generation framework in Equations (10) and Equation (11) has the potential to generate motion while dealing with a variety of time and spatial constraints not limited to these three task examples.

<sup>3</sup>This dataset contains 14,616 human motions from the AMASS [26] and HumanAct12 [11] datasets and 44,970 text descriptions.

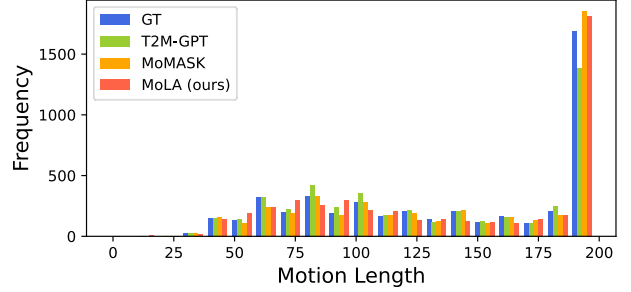


Figure 4. Comparison of motion length distributions between the HumanML3D test set and the generated samples. The Jensen-Shannon divergence (JSD) for each distribution is given by  $\text{JSD}(\text{GT}||\text{T2M-GPT}) = 0.041$ ,  $\text{JSD}(\text{GT}||\text{MoMASK}) = 0.040$ , and  $\text{JSD}(\text{GT}||\text{MoLA}) = 0.026$ . Similarly, the Earth Mover’s Distance (EMD) for each distribution is given by  $\mathcal{D}_{\text{EMD}}(\text{GT}, \text{T2M-GPT}) = 6.706$ ,  $\mathcal{D}_{\text{EMD}}(\text{GT}, \text{MoMASK}) = 3.673$ , and  $\mathcal{D}_{\text{EMD}}(\text{GT}, \text{MoLA}) = 3.538$  (a unit in this EMD means 1 frame).

data-space diffusion, while MLD [3], MotionLCM [6] and MoLA utilize a diffusion model in a lower-dimensional latent space. Therefore, the former are grouped in Continuous (raw data) and the latter in Continuous (latent) in the table. MoLA achieves the best performance in continuous methods, especially in the R-precision, FID, and MMDist metrics as shown in Table 1. Moreover, we also evaluate generation quality and inference cost in comparison to existing text-to-motion methods that have publicly available implementations [3, 13, 29, 37, 44]. Figure 2 shows that MoLA can generate much faster than MDM, which is the only existing method that performs training-free motion editing.

#### 4.1.2. Variable-length motion generation

Next, we discuss the impact of the activation variable in the motion representation introduced in Section 3.1. As shown in Equation (2), incorporating the activation variable into the motion features enables the Stage 2 model (Section 3.3) to generate variable-length motions. Figure 4 shows a comparison of the length distributions of samples generated by two existing methods (T2M-GPT and MoMask) and MoLA against the test set of HumanML3D. T2M-GPT represents a typical auto-regressive Text-to-Motion approach, while MoMask is a mask-prediction-based method that requires specifying the motion length. To address this limitation, MoMask introduces a length estimator that generates motion lengths conditioned on the input text. From Figure 4, we observe that our method successfully generates variable-length motions. Furthermore, the distribution of motion lengths generated by MoLA is closer to the actual motion distribution than those of the two existing methods, demonstrating the effectiveness of our approach.

Category	Method	R-Precision $\uparrow$			FID $\downarrow$	MMDist $\downarrow$	Diversity $\rightarrow$	MModality $\uparrow$
		Top-1	Top-2	Top-3				
N/A	Real motion data	0.511 $\pm$ .003	0.703 $\pm$ .003	0.797 $\pm$ .002	0.002 $\pm$ .000	2.974 $\pm$ .008	9.503 $\pm$ .065	-
Discrete	M2DM [21]	0.497 $\pm$ .003	0.682 $\pm$ .002	0.763 $\pm$ .003	0.352 $\pm$ .005	3.134 $\pm$ .010	9.926 $\pm$ .073	3.587 $\pm$ .072
	AttT2M [48]	0.499 $\pm$ .003	0.690 $\pm$ .002	0.786 $\pm$ .002	0.112 $\pm$ .006	3.038 $\pm$ .007	9.700 $\pm$ .090	2.452 $\pm$ .051
	T2M-GPT [44]	0.492 $\pm$ .003	0.679 $\pm$ .002	0.775 $\pm$ .002	0.141 $\pm$ .005	3.121 $\pm$ .009	9.722 $\pm$ .081	1.831 $\pm$ .048
	MoMask [13]	0.521 $\pm$ .002	0.713 $\pm$ .002	0.807 $\pm$ .002	0.045 $\pm$ .002	2.958 $\pm$ .008	-	1.241 $\pm$ .040
	DiverseMotion [24]	0.496 $\pm$ .004	0.687 $\pm$ .004	0.783 $\pm$ .003	0.070 $\pm$ .004	3.063 $\pm$ .011	9.551 $\pm$ .068	2.062 $\pm$ .079
	MMM [29]	0.504 $\pm$ .003	0.696 $\pm$ .003	0.794 $\pm$ .002	0.080 $\pm$ .003	2.998 $\pm$ .007	9.411 $\pm$ .058	1.164 $\pm$ .041
	ParCo [49]	0.515 $\pm$ .003	0.706 $\pm$ .003	0.801 $\pm$ .002	0.109 $\pm$ .003	2.927 $\pm$ .008	9.576 $\pm$ .088	1.382 $\pm$ .060
	BAMM [28]	0.525 $\pm$ .002	0.720 $\pm$ .003	0.814 $\pm$ .003	0.055 $\pm$ .002	2.919 $\pm$ .008	9.717 $\pm$ .089	1.687 $\pm$ .051
Continuous (raw data)	MotionDiffuse [46]	0.491 $\pm$ .001	0.681 $\pm$ .001	0.782 $\pm$ .001	0.630 $\pm$ .001	3.113 $\pm$ .001	9.410 $\pm$ .049	1.553 $\pm$ .042
	MDM [37]	0.320 $\pm$ .005	0.498 $\pm$ .004	0.611 $\pm$ .007	0.544 $\pm$ .044	5.566 $\pm$ .027	9.559 $\pm$ .086	2.799 $\pm$ .072
	Fg-T2M [38]	0.492 $\pm$ .002	0.683 $\pm$ .003	0.783 $\pm$ .002	0.243 $\pm$ .019	3.109 $\pm$ .007	9.278 $\pm$ .072	1.614 $\pm$ .049
Continuous (latent)	MLD [3]	0.481 $\pm$ .003	0.673 $\pm$ .003	0.772 $\pm$ .002	0.473 $\pm$ .013	3.196 $\pm$ .010	9.724 $\pm$ .082	<b>2.413<math>\pm</math>.079</b>
	MotionLCM [6]	0.502 $\pm$ .003	0.698 $\pm$ .002	0.798 $\pm$ .002	0.304 $\pm$ .012	3.012 $\pm$ .007	9.607 $\pm$ .066	2.259 $\pm$ .092
	MoLA (ours)	<b>0.516<math>\pm</math>.006</b>	<b>0.712<math>\pm</math>.005</b>	<b>0.805<math>\pm</math>.004</b>	<b>0.115<math>\pm</math>.004</b>	<b>3.008<math>\pm</math>.016</b>	9.885 $\pm$ .152	2.156 $\pm$ .157

Table 1. Comparison with state-of-the-art methods on HumanML3D dataset. Note that discrete representations do not allow for training-free motion editing; therefore, methods based on VQ-based latent representations (Discrete) are **grayed out**. The best scores for each metric in the methods using VAE-based latent representations (Continuous (latent)) are highlighted in **bold**.

Editing type	Methods	R-Precision Top-3 $\uparrow$	FID $\downarrow$	Diversity $\rightarrow$	Traj. err. $\downarrow$	Loc. err. $\downarrow$	Avg. err. $\downarrow$	AITS $\downarrow$
Training-based editing	OmniControl	0.688	0.192	9.533	0.065	0.007	0.053	74.4
	MotionLCM	0.759	0.501	9.293	0.237	0.054	0.164	0.02
Training-free editing	MoLA (ours)	0.761	0.486	9.322	0.271	0.051	0.159	1.04

Table 2. Comparison of motion editing (path-following task) on HumanML3D dataset

## 4.2. Motion Editing

Here, we demonstrate three types of editing tasks using a unified framework: path-following (motion guided by a specified trajectory), in-betweening (editing in the time direction), and upper-body editing (modifying specific joints). In particular, for path-following, we quantitatively compare our approach with existing models [6, 41].

**Path-following** is a task of giving a trajectory (often the position of the pelvis) and generating the motion that matches the given route. Controlling the trajectory of generated motion enables more motion variation, avoidance of obstacles, and the creation of motion that meets physical constraints. In this experimental case, the desired pelvis trajectory is set as the control signal  $\mathbf{y}$  in Equation (11). The upper row of Figure 5 shows the editing results with our model when different path controls are given as  $\mathbf{y}$  in the same text condition. In addition, we show a quantitative comparison with existing methods: OmniControl [41] and MotionLCM [6], where the former is Continuous (raw data) type method and the latter is Continuous (raw latent) one. We compare them in terms of FID, R-Precision, Diversity, Trajectory err (50cm), Location err (50cm), Average err, and Average inference time per sentence (AITS) in Table 2, following [6]. Table 2 demonstrates that MoLA achieves competitive editing performance across all compared to other editing methods. Notably, while MoLA shows lower performance in following control signals compared to OmniControl, it achieves significantly faster editing. Note that OmniControl and MotionLCM have been

trained for this specific editing task, while MoLA performs path-following without model fine-tuning. In other words, MoLA can perform different tasks without additional training as shown in the following sections, but OmniControl and MotionLCM require training a separate model for each task. MoLA is a more flexible and efficient framework.

**In-betweening** is an important editing task that interpolates or fills the gaps between keyframes or major motion joints to create smooth 3D motion animation. Our framework can coordinate and generate motion between past and future contexts without additional training. We only need to set the start-end positions or the motions of a few frames as the control signal  $\mathbf{y}$  in Equation (11). The middle row of Figure 5 depicts the motion in-betweening results with our model when different start-end controls are given as  $\mathbf{y}$  in the same text condition.

**Upper body editing** combines generated upper body parts with given lower body parts. Generating some joints while keeping the other body joints following a given control signal can be seen as the task of outpainting in the spatial dimension of motion. The control signal  $\mathbf{y}$  in Equation (11) is set to lower body positions that are not subject to editing.<sup>4</sup> The lower row of Figure 5 shows the upper body editing results with our model when different lower body controls are given as  $\mathbf{y}$  in the same text condition.

<sup>4</sup>Note that, although we are dealing with upper body editing in this experiment, it is in principle possible to specify a different joint subset.

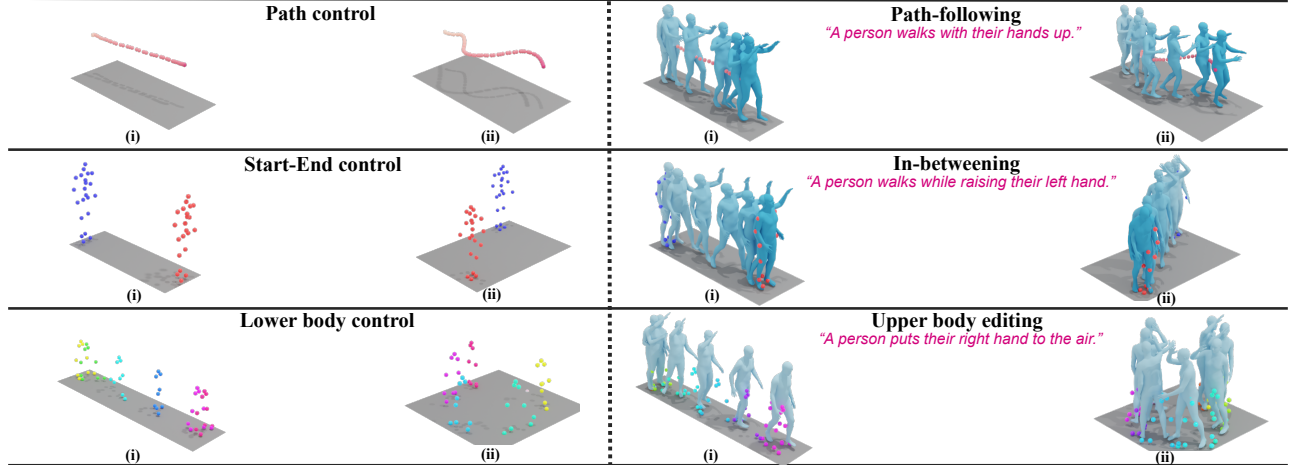


Figure 5. Qualitative results for the three editing tasks (path following, in-betweening, and upper-body editing). For these tasks, we treat each control signal (i) and (ii) in the left side of the figure as  $\mathbf{y}$  in Equation (10) and (11). The corresponding generated results using the same input text are shown on the right side of the figure as (i) and (ii), respectively.

Method	Reconstruction		Generation	
	rFID ↓	MPJPE ↓	FID ↓	MMDist ↓
Dimension of latent space				
MoLA ( $d_z = 8$ )	0.110	54.2	0.183	3.099
<b>MoLA</b> ( $d_z = 16$ )	0.030	29.3	<b>0.115</b>	<b>3.008</b>
MoLA ( $d_z = 32$ )	<b>0.028</b>	<b>26.8</b>	0.904	3.536
Adversarial training				
w/o GAN or SAN	0.038	29.3	0.126	3.053
w/ GAN instead of SAN	0.032	29.5	0.141	3.044
Input for the encoder $q_\eta$				
$[(\mathbf{m}^i)^\top, \mathbf{a}^i]^\top$ instead of Eq. (2)	0.029	31.2	<b>0.112</b>	3.024

Table 3. Analysis of motion reconstruction and generation performance on HumanML3D dataset

### 4.3. Ablation studies

We discuss the impact of latent space dimensionality, adversarial training, and motion representation on the stage 1 model. The dimensionality of the latent space may affect not only reconstruction quality but also influence the difficulty of training in stage 2, ultimately impacting the quality of the generated outputs. To investigate this, we evaluate the performance for cases with  $d_z = \{8, 16, 32\}$ . Table 3 shows that although the case of  $d_z = 16$  does not achieve the best reconstruction quality compared to the other cases, it performs best in terms of FID and MMDist. Therefore, we adopted  $d_z = 16$  for MoLA. The results for VAE combined with GAN/SAN and modifying the encoder inputs are also shown in Table 3. As shown in Table 3, the adversarial training is effective in the motion reconstruction task. It not only improves the reconstruction performance in stage 1 but also contributes to enhanced generation performance in stage 2. In particular, we can improve the performance of the stage 1 model by adopting the SAN framework instead of the conventional GAN. A better rFID is directly related to the upper bound of the overall performance of a text-to-motion model.

Method	Traj. err. ↓	Loc. err. ↓	Avg. err. ↓
	Input for the encoder $q_\eta$		
$[(\mathbf{m}^i)^\top, \mathbf{a}^i]^\top$ instead of Eq. (2)	0.281	0.068	0.174
<b>MoLA</b> ( $d_z = 16$ )	<b>0.271</b>	<b>0.051</b>	<b>0.159</b>

Table 4. Analysis of motion editing (path-following task) performance on HumanML3D dataset

Furthermore, as shown in Tables 3 and 4, modifying the encoder input Equation (2) improves reconstruction quality in terms of MPJPE without significantly compromising generation quality. As a result, this modification has a beneficial effect on the performance of the motion editing task, which requires fitting motion sequences into given control signals. This improvement can be attributed to the quality of the decoder  $g_\psi$ , which contributes to the performance of the editing task, as indicated by Equations (10) and (11). Thus, we adopted VAE trained with the SAN framework [36] and the motion representation Equation (2) for MoLA.

## 5. Conclusion

We proposed MoLA, a text-to-motion model that achieves fast, high-quality generation with multiple control tasks in a single framework. We rethought the motion representation and introduced an activation variable that characterizes the length of the motion. In addition, we integrated latent diffusion, adversarial training, and a guided generation framework. Our experiments demonstrated MoLA’s ability to generate variable-length motions with distributions close to real motions and perform diverse motion editing tasks, significantly extending the performance boundaries of methods categorized as enabling training-free editing.



## References

- [1] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *Proc. International Conference on Machine Learning (ICML)*, 2017. 11
- [2] Huiwen Chang, Han Zhang, Jarred Barber, Aaron Maschinot, Jose Lezama, Lu Jiang, Ming-Hsuan Yang, Kevin Patrick Murphy, William T. Freeman, Michael Rubinstein, Yuanzhen Li, and Dilip Krishnan. Muse: Text-to-image generation via masked generative transformers. In *Proc. International Conference on Machine Learning (ICML)*, 2023. 2
- [3] 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 *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 2, 5, 6, 7, 13
- [4] Hyungjin Chung, Jeongsol Kim, Michael Thompson McCann, Marc Louis Klasky, and Jong Chul Ye. Diffusion posterior sampling for general noisy inverse problems. In *Proc. International Conference on Learning Representation (ICLR)*, 2023. 5
- [5] CMU. CMU Graphics Lab Motion Capture Database., 2003. <http://mocap.cs.cmu.edu/>. (Accessed: 2024-02-08.). 13
- [6] 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 *Proc. European Conference on Computer Vision (ECCV)*, 2024. 3, 6, 7, 13
- [7] Bradley Efron. Tweedie’s formula and selection bias. *Journal of the American Statistical Association*, 106(496):1602–1614, 2011. 5
- [8] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2, 4
- [9] Zach Evans, Julian D Parker, CJ Carr, Zack Zukowski, Josiah Taylor, and Jordi Pons. Long-form music generation with latent diffusion. *arXiv preprint arXiv:2404.10301*, 2024. 5
- [10] Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and Baining Guo. Vector quantized diffusion model for text-to-image synthesis. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [11] 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 *Proc. ACM International Conference on Multimedia (ACMMM)*, 2020. 6
- [12] Chuan Guo, Shihao Zou, Xinxin Zuo, Sen Wang, Wei Ji, Xingyu Li, and Li Cheng. Generating diverse and natural 3d human motions from text. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 3, 6, 11
- [13] Chuan Guo, Yuxuan Mu, Muhammad Gohar Javed, Sen Wang, and Li Cheng. Momask: Generative masked modeling of 3d human motions. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 1, 2, 6, 7, 13
- [14] Yutong He, Naoki Murata, Chieh-Hsin Lai, Yuhta Takida, Toshimitsu Uesaka, Dongjun Kim, Wei-Hsiang Liao, Yuki Mitsufuji, J Zico Kolter, Ruslan Salakhutdinov, and Stefano Ermon. Manifold preserving guided diffusion. In *Proc. International Conference on Learning Representation (ICLR)*, 2024. 5, 11
- [15] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 5, 13
- [16] Vladimir Iashin and Esa Rahtu. Taming visually guided sound generation. In *Proc. British Machine Vision Conference (BMVC)*, 2021. 2
- [17] Biao Jiang, Xin Chen, Wen Liu, Jingyi Yu, Gang Yu, and Tao Chen. Motiongpt: Human motion as a foreign language. In *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, 2023. 2
- [18] Korrawe Karunratanakul, Konpat Preechakul, Supasorn Suwajanakorn, and Siyu Tang. Guided motion diffusion for controllable human motion synthesis. In *Proc. IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023. 3
- [19] Korrawe Karunratanakul, Konpat Preechakul, Emre Aksan, Thabo Beeler, Supasorn Suwajanakorn, and Siyu Tang. Optimizing diffusion noise can serve as universal motion priors. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 3
- [20] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *Proc. International Conference on Learning Representation (ICLR)*, 2013. 2
- [21] Hanyang Kong, Kehong Gong, Dongze Lian, Michael Bi Mi, and Xinchao Wang. Priority-centric human motion generation in discrete latent space. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 2, 6, 7, 13
- [22] Shanchuan Lin, Bingchen Liu, Jiashi Li, and Xiao Yang. Common diffusion noise schedules and sample steps are flawed. In *Proc. Winter Conference on Applications of Computer Vision (WACV)*, 2024. 5
- [23] Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo Liu, Danilo Mandic, Wenwu Wang, and Mark D Plumbley. AudioLDM: Text-to-audio generation with latent diffusion models. In *Proc. International Conference on Machine Learning (ICML)*, 2023. 2, 4
- [24] Yunhong Lou, Linchao Zhu, Yaxiong Wang, Xiaohan Wang, and Yi Yang. Diversemotion: Towards diverse human motion generation via discrete diffusion. *arXiv preprint arXiv:2309.01372*, 2023. 2, 6, 7, 13
- [25] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. RePaint: Inpainting using denoising diffusion probabilistic models. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 11
- [26] Naureen Mahmood, Nima Ghorbani, Nikolaus F Troje, Gerard Pons-Moll, and Michael J Black. Amass: Archive of motion capture as surface shapes. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 6

- [27] Christian Mandery, Ömer Terlemez, Martin Do, Nikolaus Vahrenkamp, and Tamim Asfour. The kit whole-body human motion database. In *Proc. IEEE International Conference on Advanced Robotics (ICAR)*, 2015. 13
- [28] Ekkasit Pinyoanuntapong, Muhammad Usama Saleem, Pu Wang, Minwoo Lee, Srijan Das, and Chen Chen. Bamm: bidirectional autoregressive motion model. In *Proc. European Conference on Computer Vision (ECCV)*, 2024. 2, 6, 7, 13
- [29] Ekkasit Pinyoanuntapong, Pu Wang, Minwoo Lee, and Chen Chen. Mmm: Generative masked motion model. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 1, 2, 3, 6, 7, 13
- [30] Matthias Plappert, Christian Mandery, and Tamim Asfour. The kit motion-language dataset. *Big data*, 4(4):236–252, 2016. 3, 11, 13
- [31] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *Proc. International Conference on Machine Learning (ICML)*, 2021. 5
- [32] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 4, 5
- [33] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *Proc. International Conference on Learning Representation (ICLR)*, 2021. 5
- [34] Jiaming Song, Qingsheng Zhang, Hongxu Yin, Morteza Mardani, Ming-Yu Liu, Jan Kautz, Yongxin Chen, and Arash Vahdat. Loss-guided diffusion models for plug-and-play controllable generation. In *Proc. International Conference on Machine Learning (ICML)*, 2023. 3
- [35] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *Proc. International Conference on Learning Representation (ICLR)*, 2021. 5
- [36] Yuhta Takida, Masaaki Imaizumi, Takashi Shibuya, Chieh-Hsin Lai, Toshimitsu Uesaka, Naoki Murata, and Yuki Mitsufuji. SAN: Inducing metrizable of gan with discriminative normalized linear layer. In *Proc. International Conference on Learning Representation (ICLR)*, 2024. 4, 8, 11
- [37] Guy Tevet, Sigal Raab, Brian Gordon, Yoni Shafir, Daniel Cohen-or, and Amit Haim Bermano. Human motion diffusion model. In *Proc. International Conference on Learning Representation (ICLR)*, 2023. 2, 3, 6, 7, 13
- [38] Yin Wang, Zhiying Leng, Frederick WB Li, Shun-Cheng Wu, and Xiaohui Liang. Fg-t2m: Fine-grained text-driven human motion generation via diffusion model. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 6, 7, 13
- [39] Yinhuai Wang, Jiwen Yu, and Jian Zhang. Zero-shot image restoration using denoising diffusion null-space model. In *Proc. International Conference on Learning Representation (ICLR)*, 2023. 11
- [40] Qi Wu, Yubo Zhao, Yifan Wang, Xinhang Liu, Yu-Wing Tai, and Chi-Keung Tang. Motion-agent: A conversational framework for human motion generation with llms. In *arXiv preprint arXiv:2405.17013*, 2024. 2
- [41] Yiming Xie, Varun Jampani, Lei Zhong, Deqing Sun, and Huaizu Jiang. Omnicontrol: Control any joint at any time for human motion generation. In *Proc. International Conference on Learning Representation (ICLR)*, 2024. 3, 7
- [42] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, Ben Hutchinson, Wei Han, Zarana Parekh, Xin Li, Han Zhang, Jason Baldridge, and Yonghui Wu. Scaling autoregressive models for content-rich text-to-image generation. *Transactions on Machine Learning Research*, 2022. 2
- [43] Jiwen Yu, Yinhuai Wang, Chen Zhao, Bernard Ghanem, and Jian Zhang. Freedom: Training-free energy-guided conditional diffusion model. In *Proc. IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023. 5, 11
- [44] Jianrong Zhang, Yangsong Zhang, Xiaodong Cun, Shaoli Huang, Yong Zhang, Hongwei Zhao, Hongtao Lu, and Xi Shen. T2m-gpt: Generating human motion from textual descriptions with discrete representations. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 1, 2, 4, 6, 7, 11, 13
- [45] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proc. IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023. 3
- [46] Mingyuan Zhang, Zhongang Cai, Liang Pan, Fangzhou Hong, Xinying Guo, Lei Yang, and Ziwei Liu. Motiondiffuse: Text-driven human motion generation with diffusion model. In *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024. 2, 6, 7, 13
- [47] Zeyu Zhang, Akide Liu, Ian Reid, Richard Hartley, Bohan Zhuang, and Hao Tang. Motion mamba: Efficient and long sequence motion generation. In *Proc. European Conference on Computer Vision (ECCV)*, 2025. 2
- [48] Chongyang Zhong, Lei Hu, Zihao Zhang, and Shihong Xia. Attt2m: Text-driven human motion generation with multi-perspective attention mechanism. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 2, 6, 7, 13
- [49] Qiran Zou, Shangyuan Yuan, Shian Du, Yu Wang, Chang Liu, Yi Xu, Jie Chen, and Xiangyang Ji. Parco: Part-coordinating text-to-motion synthesis. In *Proc. European Conference on Computer Vision (ECCV)*, 2024. 2, 6, 7, 13

This appendix provides the following supplementary information:

- Section A: A detailed explanation of the SAN-based discriminator as an alternative to GAN mentioned in Section 3.2
- Section B: An in-depth description of the editing method introduced in Section 3.4
- Section C: Details of the experimental settings used in Sections 4
- Section D: Definitions of evaluation metrics for the text-to-motion task
- Section E: Additional experiments on the performance of the proposed method

## A. SAN-based discriminator

Takida et al. [36] incorporated a sliced optimal transport perspective into a general GAN and established a new framework, slicing adversarial network (SAN). Following this work, we first decompose the discriminator  $f_\phi$  into the last linear layer and the remaining neural part, denoted as  $w \in \mathbb{R}^{d_w}$  and  $h_\varphi : \mathbb{R}^{N \times L} \rightarrow \mathbb{R}^{d_w}$  with  $d_w \in \mathbb{N}$ , respectively. This decomposition provides an interpretation of the discriminator: the neural function  $h_\varphi$  maps motion sequences to non-linear features, and then the linear layer  $w$  projects them into scalars. As a preliminary step for applying SAN to the adversarial training in Section 3.2.1, we normalize  $w$  using its norm, resulting in  $\omega = w/\|w\|_2$ , which indicates the direction of the projecting. Now, the discriminator is represented in the form of an inner product as  $f_\phi(x) = \omega^\top h_\varphi(x)$ , and its parameter is  $\phi = \{\varphi, \omega\}$ .

The prior work has shown that the discriminator obtained from the optimal solution of  $\omega$  in the hinge loss (5) does not guarantee gradients that make the generated distribution close to the data distribution. To address this issue, we adopt the SAN maximization problem instead of Equation (5). Specifically, we optimize the neural part  $\varphi$  with the original hinge loss, while applying the Wasserstein GAN loss [1] to the direction  $\omega$ . The modified maximization objective is formulated as

$$\mathcal{L}_{\text{SAN}}(\phi; \psi, \eta, x) = \mathcal{L}_{\text{GAN}}(\{\varphi, \omega^-\}; \psi, \eta, x) + \omega^\top (h_{\varphi^-}(x) - \mathbb{E}_{q_\eta(z|x)}[h_{\varphi^-}(g_\psi(z))]), \quad (12)$$

where  $(\cdot)^-$  indicates a stop gradient operator. The second term in Equation (12) induces the direction that best discriminates between the real and generated sample sets in the feature space. We employ the same minimization objective as defined in Equation (6).

## B. Editing procedure on guided generation

In Algorithm 1, we show the specific procedure for guided generation described in Section 3.4. Note that we apply a time-travel technique to the update rule in Section 3.4 as in [14, 25, 39, 43] to achieve better editing results. This technique adds noise after each gradient descent step to implicitly perform a multi-step optimization of minimizing the distance measuring function  $\mathcal{L}(\cdot, y)$  and lead to an improvement in the editing quality.

## C. Implementation Details

**Training setup for stage 1 model:** During the training, motion sequences are segmented into lengths of  $L = 64$ . We use AdamW optimizer and batch size of 128. The hyperparameters in Eq. (4) and (6) are set as  $\lambda_{\text{act}} = 1.0$ ,  $\lambda_{\text{reg}} = 1.0 \times 10^{-4}$ , and  $\lambda_{\text{adv}} = 1.0 \times 10^{-3}$ . Our model  $\{q_\eta, g_\psi, f_\psi\}$  are trained with a simple multi-step learning late, the first 10,000 iterations with a learning rate of  $2.0 \times 10^{-4}$ , and latter 5,000 with a learning rate of  $2.0 \times 10^{-5}$  in the case of HumanML3D dataset [12]. Note that in Section E, we also conducted training and evaluation on a different dataset, the KIT-ML dataset [30]. For this dataset, the model was trained for the first 10,000 iterations with a learning rate of  $5.0 \times 10^{-5}$ , followed by an additional 5,000 iterations with a reduced learning rate of  $5.0 \times 10^{-6}$ . In addition, we replace the reconstruction loss in (4) with smooth  $\ell_1$  loss (known as the special case of Huber loss) function and introduce a position enhancement term, following the technique in [44].

**Training setup for stage 2 model:** During the training, we use the AdamW optimizer with a batch size of 64. The model  $\epsilon_\theta$  is trained for 10,000 epochs with a cosine annealing learning rate schedule starting at  $1.0 \times 10^{-4}$ , including a warm-up phase. For inference, we adopt DDIM with 50 sampling steps and employ a trailing strategy. The classifier-free guidance scale  $s$  in Eq. (17) is set to  $s = 11$  for HumanML3D and  $s = 7$  for the KIT-ML dataset. Notably, we observed in Section E that the performance is significantly influenced by the scale  $s$ .

## D. Evaluation metrics details

We provide more details of evaluation metrics in Section 4.1. We use five metrics to quantitatively evaluate text-to-motion models. These metrics are calculated based on motion and text features extracted with pre-trained networks. More specifically, we utilize the motion encoder and text encoder provided in [12]. We denote ground-truth motion features, generated

---

**Algorithm 1** Guided Generation for Motion Editing

---

**Require:** motion control signal  $\mathbf{y}$ , output of text encoder  $\tau(c)$ , estimator in latent diffusion  $\epsilon_\theta(\cdot, t, \tau(c))$ , distance measuring function  $\mathcal{L}(\cdot; \mathbf{y})$ , pre-defined parameter  $\bar{\alpha}_t$ , time-dependent step-size  $\rho_t$ , and the repeat times of time-travel of each step  $\{r_1, \dots, r_T\}$ .

```
1:  $\mathbf{z}_T \sim \mathcal{N}(0, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:   for  $i = r_t, \dots, 1$  do
4:      $\epsilon_t \sim \mathcal{N}(0, \mathbf{I})$ 
5:      $\mathbf{z}_{0|t} = \frac{1}{\sqrt{\bar{\alpha}_t}}(\mathbf{z}_t - \sqrt{1 - \bar{\alpha}_t}\epsilon_\theta(\mathbf{z}_t, t, \tau(c)))$ 
6:      $\mathbf{z}_{0|t} = \mathbf{z}_{0|t} - \rho_t \nabla_{\mathbf{z}_{0|t}} \mathcal{L}_{\text{Motion}}(g_\psi(\mathbf{z}_{0|t}); \mathbf{y})$ 
7:      $\mathbf{z}_{t-1} = \sqrt{\bar{\alpha}_t}\mathbf{z}_{0|t} + \sqrt{1 - \bar{\alpha}_t - \sigma_t^2}\epsilon_\theta(\mathbf{z}_t, t, \tau(c)) + \sigma_t\epsilon_t$ 
8:     if  $i > 1$  then
9:        $\epsilon'_t \sim \mathcal{N}(0, \mathbf{I})$ 
10:       $\mathbf{z}_t = \sqrt{\alpha_t}\mathbf{z}_{t-1} + \sqrt{1 - \alpha_t}\epsilon'_t$ 
11:     end if
12:   end for
13: end for
14: return  $g_\psi(\mathbf{z}_0)$  ▷ edited motion generation sample
```

---

motion features, and text features as  $f_{\text{gt}} = \mathcal{E}_{\text{M}}(x_{\text{gt}}) \in \mathbb{R}^{512}$ ,  $f_{\text{pred}} = \mathcal{E}_{\text{M}}(x_{\text{pred}}) \in \mathbb{R}^{512}$ , and  $f_{\text{text}} = \mathcal{E}_{\text{T}}(c) \in \mathbb{R}^{512}$ , where  $\mathcal{E}_{\text{M}}(\cdot)$  and  $\mathcal{E}_{\text{T}}(\cdot)$  represent the motion encoder and the text encoder, respectively. We explain the five metrics below.

**R-Precision:** R-Precision evaluates semantic alignment between input text and generated motion in a sample-wise manner. Given one motion sequence and 32 text descriptions (one ground truth and thirty-one randomly selected mismatched descriptions), we rank the Euclidean distances between the motion and text embeddings. Then, we compute the average accuracy at top-1, top-2, and top-3 places.

**Fréchet Inception Distance (FID):** FID evaluates the overall motion quality by measuring the distributional difference between the motion features of the generated motions and those of real motions. We obtain FID by

$$\text{FID} := \|\mu_{\text{gt}} - \mu_{\text{pred}}\|_2^2 - \text{tr}(\Sigma_{\text{gt}} + \Sigma_{\text{pred}} - (\Sigma_{\text{gt}}\Sigma_{\text{pred}})^{\frac{1}{2}}), \quad (13)$$

where  $\mu_{\text{gt}}$  and  $\mu_{\text{pred}}$  are the mean of  $f_{\text{gt}}$  and  $f_{\text{pred}}$ , respectively,  $\Sigma_{\text{gt}}$  and  $\Sigma_{\text{pred}}$  are their corresponding covariance matrices, and  $\text{tr}$  denotes the trace of a matrix.

**Multi-modal distance (MMDist):** MMDist gauges sample-wise semantic alignment between input text and generated motion as well. MMDist is computed as the average Euclidean distance between the motion feature of each generated motion and the text feature of its corresponding description in the test set. Precisely, given  $N$  randomly generated samples, it computes the average Euclidean distance between each text feature and its corresponding generated motion feature:

$$\text{MMDist} := \frac{1}{N} \sum_{i=1}^N \|f_{\text{pred},i} - f_{\text{text},i}\|_2^2, \quad (14)$$

where  $f_{\text{pred},i}$  and  $f_{\text{text},i}$  are the features of the  $i$ -th text-motion pair.

**Diversity:** Diversity measures the variance of the generated motions across the test set. We randomly sample  $2S_d$  motions and extract motion features  $\{f_{\text{pred},1}, \dots, f_{\text{pred},2S_d}\}$  from the motions. Then, those motion features are grouped into the same size of two subsets,  $\{v_1, \dots, v_{S_d}\}$  and  $\{v'_1, \dots, v'_{S_d}\}$ . The Diversity for this set of motions is defined as

$$\text{Diversity} := \frac{1}{S_d} \sum_{i=1}^{S_d} \|v_i - v'_i\|_2^2. \quad (15)$$

$S_d = 300$  is used in our experiments, following previous works.



**Multi-modality (MModality):** MModality measures how much the generated motions diversify within each text description. Given a set of motions with  $K$  text descriptions, we randomly sample  $2S_l$  motions corresponding to the  $k$ -th text description and extract motion features  $\{f_{\text{pred},k,1}, \dots, f_{\text{pred},k,2S_l}\}$ . Then, those motion features are grouped into the same size of two subsets,  $\{v_{k,1}, \dots, v_{k,S_l}\}$  and  $\{v'_{k,1}, \dots, v'_{k,S_l}\}$ . The MModality for this motion set is formalized as

$$\text{MModality} = \frac{1}{K \cdot S_l} \sum_{k=1}^K \sum_{i=1}^{S_l} \|v_{k,i} - v'_{k,i}\|_2^2. \quad (16)$$

$S_l = 10$  is used in our experiments, following previous works.

## E. Additional experiments

### E.1. Performance comparison on the KIT-ML dataset

We evaluate our method not only on the HumanML3D dataset presented in Section 4.1 but also on the KIT-ML dataset [30], comparing its performance against other text-to-motion generation methods. The KIT-ML dataset includes 3,911 human motion sequences and 6,278 text descriptions derived from the KIT [27] and CMU [5] datasets. These data are split into training, validation, and test sets with proportions of 80%, 5%, and 15%, respectively. Consistent with Section 4.1, the evaluated methods are categorized into three groups: (i) those using VQ-based latent representations (discrete), (ii) those employing data-space diffusion models (continuous, raw data), and (iii) those utilizing VAE-based latent representations (continuous, latent). As shown in Table 5, MoLA achieves the best performance in terms of R-Precision and MMDist, particularly among continuous-based methods. However, it does not outperform methods like MLD in terms of FID. Comparing this with Table 1, we observe a discrepancy in the FID trends across the two datasets, suggesting there is room for improvement in this aspect of MoLA.

Category	Method	R-Precision $\uparrow$			FID $\downarrow$	MMDist $\downarrow$	Diversity $\rightarrow$	MModality $\uparrow$
		Top-1	Top-2	Top-3				
N/A	Real motion data	0.424 $\pm$ .005	0.649 $\pm$ .006	0.779 $\pm$ .006	0.031 $\pm$ .004	2.788 $\pm$ .012	11.08 $\pm$ .097	-
Discrete	M2DM [21]	0.416 $\pm$ .004	0.628 $\pm$ .004	0.743 $\pm$ .004	0.515 $\pm$ .029	3.015 $\pm$ .017	11.417 $\pm$ .97	3.325 $\pm$ .37
	AttT2M [48]	0.413 $\pm$ .006	0.632 $\pm$ .006	0.751 $\pm$ .006	0.870 $\pm$ .039	3.039 $\pm$ .021	10.96 $\pm$ .123	2.281 $\pm$ .047
	T2M-GPT [44]	0.416 $\pm$ .006	0.627 $\pm$ .006	0.745 $\pm$ .006	0.514 $\pm$ .029	3.007 $\pm$ .029	10.921 $\pm$ .108	1.570 $\pm$ .039
	MoMask [13]	0.433 $\pm$ .007	0.656 $\pm$ .005	0.781 $\pm$ .005	0.204 $\pm$ .011	2.779 $\pm$ .022	-	1.131 $\pm$ .043
	DiverseMotion [24]	0.416 $\pm$ .005	0.637 $\pm$ .008	0.760 $\pm$ .011	0.468 $\pm$ .098	2.892 $\pm$ .041	10.873 $\pm$ .101	2.062 $\pm$ .079
	MMM [29]	0.381 $\pm$ .005	0.590 $\pm$ .006	0.718 $\pm$ .005	0.429 $\pm$ .019	3.146 $\pm$ .019	10.633 $\pm$ .097	1.105 $\pm$ .026
	ParCO [49]	0.430 $\pm$ .004	0.649 $\pm$ .007	0.772 $\pm$ .008	0.453 $\pm$ .027	2.820 $\pm$ .028	10.95 $\pm$ .094	1.245 $\pm$ .022
	BAMM [28]	0.438 $\pm$ .009	0.661 $\pm$ .009	0.788 $\pm$ .005	0.183 $\pm$ .013	2.723 $\pm$ .026	11.008 $\pm$ .094	1.609 $\pm$ .065
Continuous (raw data)	MotionDiffuse [46]	0.417 $\pm$ .004	0.621 $\pm$ .004	0.739 $\pm$ .004	1.954 $\pm$ .062	2.958 $\pm$ .005	11.10 $\pm$ .143	0.730 $\pm$ .013
	MDM [37]	0.164 $\pm$ .004	0.291 $\pm$ .004	0.396 $\pm$ .004	0.497 $\pm$ .021	9.191 $\pm$ .022	10.847 $\pm$ .109	1.907 $\pm$ .214
	Fg-T2M [38]	0.418 $\pm$ .005	0.626 $\pm$ .004	0.745 $\pm$ .004	0.571 $\pm$ .047	3.114 $\pm$ .015	10.93 $\pm$ .083	1.019 $\pm$ .029
Continuous (latent)	MLD [3]	0.390 $\pm$ .008	0.609 $\pm$ .008	0.734 $\pm$ .007	<b>0.404<math>\pm</math>.027</b>	3.204 $\pm$ .027	10.80 $\pm$ .117	<b>2.192<math>\pm</math>.071</b>
	MotionLCM [6]	-	-	-	-	-	-	-
	MoLA (ours)	<b>0.432<math>\pm</math>.008</b>	<b>0.655<math>\pm</math>.008</b>	<b>0.770<math>\pm</math>.004</b>	0.529 $\pm$ .056	<b>2.942<math>\pm</math>.053</b>	11.129 $\pm$ .158	1.789 $\pm$ .174

Table 5. Comparison with state-of-the-art methods on KIT-ML dataset. Note that discrete representations do not allow for training-free motion editing; therefore, methods based on VQ-based latent representations (Discrete) are **grayed out**. The best scores for each metric in the methods using VAE-based latent representations (Continuous (latent)) are highlighted in **bold**.

### E.2. Ablation of classifier-free diffusion guidance on stage 2

As explained in Section 3.3, we employ a classifier-free diffusion guidance technique [15]. In training, we randomly drop the condition  $\tau(\mathbf{c})$  with a probability of 10% and train both the conditional model  $\epsilon_\theta(z_t, t, \tau(\mathbf{c}))$  and the unconditional model  $\epsilon_\theta(z_t, t, \emptyset)$ . In inference, the two predictions are linearly combined as follows:

$$\epsilon_\theta(z_t, t, \mathbf{c}) = s \cdot \epsilon_\theta(z_t, t, \tau(\mathbf{c})) + (1 - s) \cdot \epsilon_\theta(z_t, t, \emptyset), \quad (17)$$

where  $s$  is the guidance scale, and  $s > 1$  can amplify the effect of the guidance. We investigate how the performance of our method changes with different values of the hyperparameter  $s$ . From Figure 6, we observe that the best performance in FID is achieved at  $s = 11$  on the HumanML3D dataset. Consequently, we adopt  $s = 11$  in Section 4.1.

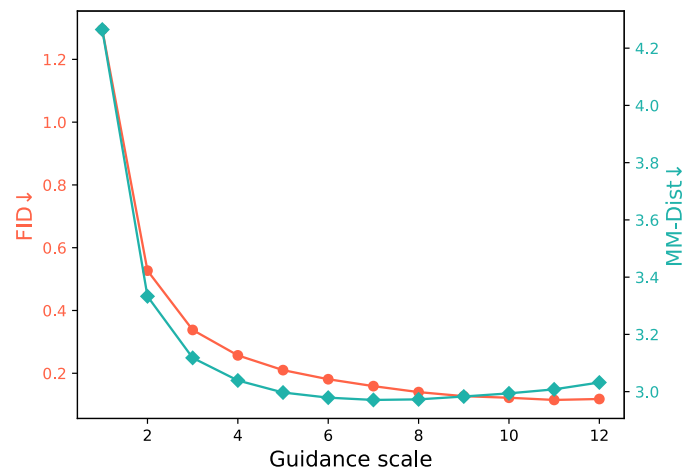


Figure 6. Evaluation sweep over guidance scale  $s$