

# Bilingual Text-to-Motion Generation via Step-Aware Reward-Guided Alignment

Anonymous ICCV submission

Paper ID 2603

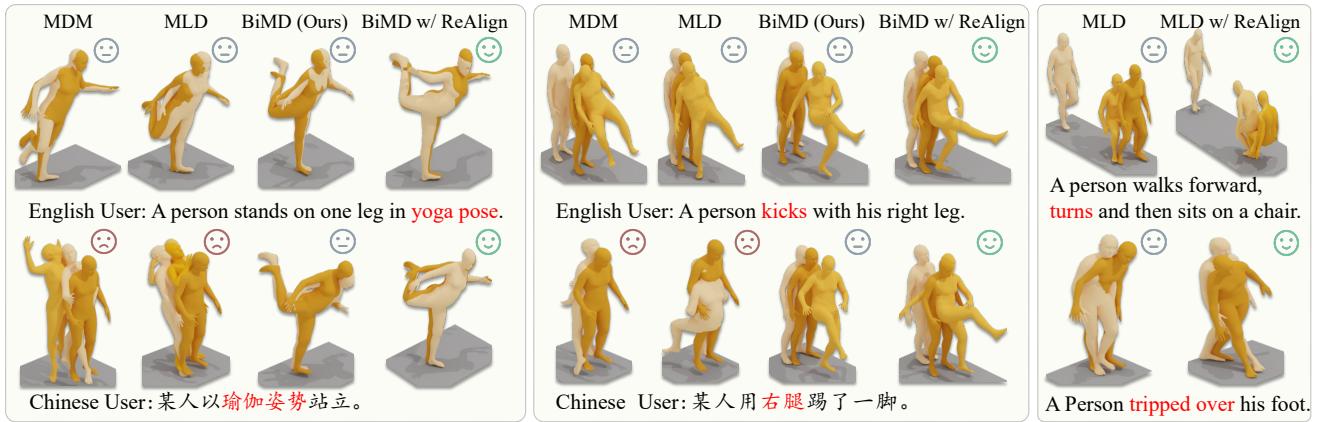


Figure 1. Visual results of bilingual (left and middle) and monolingual (right) text-driven motion generation. This figure presents motions generated by existing methods, such as MDM [27] and MLD [2], alongside our Bilingual Motion Diffusion model (BiMD) and plug-and-play Reward-guided Alignment (ReAlign), across both tasks. Observations reveal that: (1) in bilingual motion generation, MDM and MLD exhibit limitations in processing bilingual inputs; (2) in monolingual motion generation, generated motions demonstrate persistent misalignment with input texts. The left figure shows that our BiMD successfully generates motion from both English and Chinese inputs. Furthermore, the right figure highlights that BiMD, integrated with our ReAlign, successfully mitigates the misalignment issue.

## Abstract

001 *Bilingual text-to-motion generation, which synthesizes 3D*  
 002 *human motions from bilingual text inputs, holds immense*  
 003 *potential for cross-linguistic applications in gaming, film,*  
 004 *and robotics. However, this task faces critical challenges:*  
 005 *the absence of bilingual motion-language datasets and the*  
 006 *misalignment between text and motion distributions in dif-*  
 007 *fusion models, leading to semantically inconsistent or low-*  
 008 *quality motions. To address these challenges, we propose*  
 009 *BiHumanML3D, a novel bilingual human motion dataset,*  
 010 *which establishes a crucial benchmark for bilingual text-*  
 011 *to-motion generation models. Furthermore, we propose a*  
 012 *Bilingual Motion Diffusion model (**BiMD**), which leverages*  
 013 *cross-lingual aligned representations to capture semantics,*  
 014 *thereby achieving a unified bilingual model. Building upon*  
 015 *this, we propose **Reward-guided sampling Alignment (Re-***  
 016 ***Align**) method, comprising a step-aware reward model to*  
 017 *assess alignment quality during sampling and a reward-*  
 018 *guided strategy that directs the diffusion process toward*  
 019 *an optimally aligned distribution. This reward model inte-*  
 020 *grates step-aware tokens and combines a text-aligned mod-*

ule for semantic consistency and a motion-aligned module for realism, refining noisy motions at each timestep to balance probability density and alignment. Experiments demonstrate that our approach significantly improves text-motion alignment and motion quality compared to existing state-of-the-art methods.

## 1. Introduction

With the increasing demand for realistic and diverse 3D motion in gaming, filmmaking, and robotics [2, 6, 7, 12], text-to-motion generation has emerged as a key research topic, offering intuitive text-based control. Particularly, bilingual text-to-motion generation which aims to generate motion from bilingual text descriptions, holds significant potential in cross-linguistic applications [14, 41]. However, this task remains largely unexplored due to several key challenges.

Firstly, the scarcity of bilingual text-motion datasets presents a significant challenge to the development of bilingual text-to-motion generation. While many large-scale English motion datasets [7, 20] have been established and widely utilized, there are no publicly available motion

021  
022  
023  
024  
025  
026

027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040

041 datasets annotated in other languages, let alone bilingual  
042 datasets. This data scarcity greatly restricts research and  
043 progress in bilingual text-to-motion generation methods.  
044

045 Secondly, diffusion-based models [27, 34] struggle with  
046 text-motion alignment due to their reliance on text embeddings  
047 encoded by CLIP [22] which is trained on text-image  
048 pairs rather than text-motion pairs. Consequently, these  
049 models often fail to capture the semantic alignment between  
050 text and motion, resulting in synthesized motions that lack  
coherence with the input descriptions as shown in Fig. 1.

051 To the best of our knowledge, this is the first study to ex-  
052 plore bilingual text-to-motion generation. While research  
053 in bilingual generative models is gaining traction, most ef-  
054 forts remain domain-specific. For instance, Zuo et al. [41]  
055 developed a bilingual 3D sign language avatar generator,  
056 and Li et al. [14] introduced a bilingual image diffusion  
057 transformer. However, these works are designed for spe-  
058 cific applications, and do not generalize to motion synthesis,  
059 leaving bilingual text-to-motion generation an open chal-  
060 lenge. Furthermore, prior methods aiming to improve text-  
061 to-motion alignment, such as reinforcement learning with  
062 reward functions [9, 16, 25], primarily focus on fine-tuning  
063 generative models to enhance motion quality. However,  
064 these approaches do not explicitly address misalignment  
065 between text and motion. In contrast, we propose a plug-and-  
066 play reward model that can be seamlessly integrated into  
067 any diffusion model without additional fine-tuning, ensur-  
068 ing enhanced text-motion alignment.

069 **Contributions.** Our first contribution lies in the introduc-  
070 tion of a pioneering bilingual text-to-motion dataset, BiHu-  
071 manML3D, accompanied by a corresponding bilingual text-  
072 to-motion method, **Bilingual Motion Diffusion (BiMD)**.  
073 To address the scarcity of bilingual text-motion datasets,  
074 we extend the widely used text-to-motion dataset, Hu-  
075 manML3D [7], by introducing its bilingual version, BiHu-  
076 manML3D. Specifically, a multi-stage translation pipeline  
077 based on large language models and manual correction is  
078 designed to ensure high-quality annotations and accurate  
079 semantic translations. Furthermore, a unified bilingual mo-  
080 tion diffusion model, BiMD, is trained to efficiently handle  
081 bilingual text-to-motion generation. This is enabled by  
082 harmonizing semantics across languages via cross-lingual  
083 alignment.

084 Additionally, we propose a novel **Reward-guided sam-  
085 pling Alignment strategy (ReAlign)** to enhance text-motion  
086 alignment quality with the guidance of a well-aligned re-  
087 ward distribution. To this end, we derive the reward dis-  
088 tribution from a step-aware reward comprising two mod-  
089 ules: a text-aligned module to ensure semantic consistency,  
090 and a motion-aligned module to assess realism. Together,  
091 these modules adapt to noisy motions and variations across  
092 timesteps, guiding diffusion model toward a distribution  
093 that not only maximizes probability density but also main-

094 tains strong text-motion alignment. By explicitly address-  
095 ing both semantic misalignment and motion quality degra-  
096 dation, this approach improves the coherence and realism  
097 of the generated motion.

098 Finally, extensive experiments show that BiMD, driven  
099 by reward-guided sampling, generates high-quality motions  
100 aligned with both English and Chinese semantics. The  
101 step-aware reward model enhances diffusion-based models  
102 with plug-and-play versatility. On BiHumanML3D, BiMD  
103 with cross-lingual alignment outperforms baselines, includ-  
104 ing language-specific models, proving its strength in bilin-  
105 gual motion generation. On monolingual HumanML3D,  
106 our reward model enhances both our BiMD and the pre-  
107 vious SoTA MLD++ [6], e.g., improving BiMD by **55.2%**  
108 and MLD++ by **24.7%** in terms of FID, without any addi-  
109 tional training. These results highlight the superiority of our  
110 BiMD and ReAlign strategy.

## 2. Related Works

111 **Text-to-Motion Generation.** Text-to-motion generation  
112 represents a critical task in computer vision, exhibiting  
113 rapid advancements in recent years [26, 30, 32, 33, 36, 39].  
114 Specifically, Tevet et al.[27] and Zhang et al.[34] first pro-  
115 posed diffusion models to address text-driven motion gen-  
116 eration, laying the groundwork for the following innova-  
117 tions. Subsequently, Dai et al. [6] presented MotionLCM, a  
118 real-time controllable model that refines motion-latent dif-  
119 fusion, enabling precise spatiotemporal control via few-step  
120 inference. Zhang et al.[38] introduced motion mamba, a  
121 state-space framework that leverages hierarchical temporal  
122 and bidirectional spatial modules to improve efficiency and  
123 long-sequence modeling. However, these works are still  
124 confined to monolingual settings. Bilingual text-to-motion  
125 generation remains a challenging area due to the lack of  
126 datasets and misalignment between different languages.  
127

128 **Alignment in Motion Generation.** Alignment represents a  
129 versatile technique widely employed across the domains of  
130 language modeling [23], image generation [28], and policy  
131 optimization [1]. Recently, researchers [9, 16, 25] have ex-  
132 plored human preference alignment in text-to-motion gen-  
133 eration. Han et al. [9] introduced ReinDiffuse, a diffusion-  
134 based model refined through reinforcement learning, to en-  
135 hance the physical plausibility of generated motions. Liu et  
136 al. [16] investigated aligning human preferences using the  
137 proposed multi-reward reinforcement learning framework.  
138 Tan et al. [25] proposed SoPo, a novel semi-online prefer-  
139 ence optimization method that refines text-to-motion mod-  
140 els. However, these methods focus on fine-tuning genera-  
141 tive models to align preferences or enhance motion quality  
142 without explicitly addressing text-motion misalignment. In  
143 contrast, we tackle this issue with a plug-and-play reward  
144 model in the inference process.

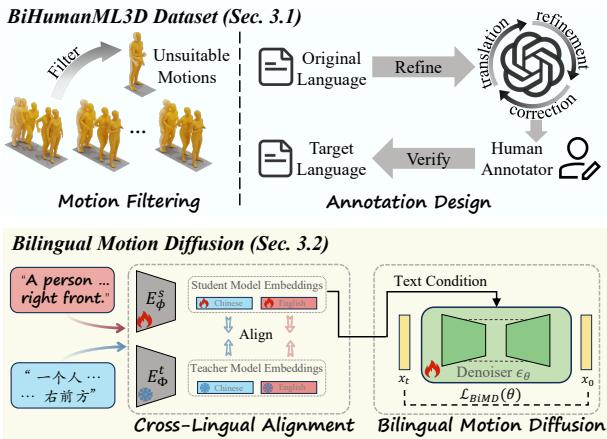


Figure 2. Pipeline for constructing our bilingual HumanML3D dataset (top) and training the bilingual motion diffusion model (bottom). We align English and Chinese text semantics in a shared latent space by freezing the teacher model  $E_\Phi^t$  and fine-tuning the student model  $E_\Phi^s$  with the cross-lingual alignment loss  $\mathcal{L}_{CLA}$  in Eq. (1). The aligned student model  $E_\Phi^s$  is then provided text conditions for training the diffusion model  $\epsilon_\theta$ , enabling bilingual motion generation while minimizing  $\mathcal{L}_{BiMD}$  in Eq. (2).

145

### 3. BiHumanML3D Dataset & Bilingual Model

146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156

We address the challenges of bilingual motion generation by introducing a new dataset and a unified bilingual motion diffusion model. In Sec. 3.1, we present a bilingual dataset to bridge the gap in multilingual motion generation. Instead of training separate models for each language—an inefficient approach that fails to leverage shared semantics—we propose a cross-lingually aligned diffusion model in Sec. 3.2, ensuring consistent motion generation across languages. Lastly, to tackle misalignment between text and motion, we introduce reward-guided sampling in Sec. 4, refining motion quality and semantic coherence.

157

#### 3.1. Bilingual HumanML3D Dataset

158  
159  
160  
161  
162  
163  
164

Despite the progress in text-to-motion models, their reliance on English-only datasets limits their usability in multilingual applications, reducing accessibility for non-English speakers and restricting cultural diversity in global industries like animation and robotics. So we construct a bilingual motion dataset for multilingual motion synthesis via two stages. See more construction details in Appendix B.

165  
166  
167  
168  
169  
170  
171

**Data Collection & Filtering.** Since large-scale English text-to-motion datasets are already well-established, we extend HumanML3D [7] into a bilingual version. Following previous work [2, 27], we filter out excessively short or long motion clips which often lack meaningful semantic signals, ensuring high-quality motion-caption pairs for annotation.

**Annotation Design.** To generate bilingual annotations

while preserving motion semantics, we develop an LLM-assisted translation pipeline as shown in Fig. 2. This process involves three stages: (1) initial translation of English captions into Chinese using an LLM, (2) refinement to enhance translation quality, and (3) validation via both automated and human review to ensure linguistic accuracy.

By introducing the Bilingual HumanML3D Dataset that contains 13,312 bilingual motions, we take a crucial step toward removing language barriers in text-to-motion generation. This dataset enhances the adaptability of motion synthesis models across languages while promoting fairness and inclusivity in AI-driven animation and robotics.

#### 3.2. Bilingual Motion Diffusion

While training separate diffusion models for each language on our bilingual dataset is a straightforward approach, it is computationally expensive and overlooks semantic similarities between languages. To address this, we propose a unified bilingual motion diffusion model that leverages cross-lingual alignment and bilingual diffusion training to reduce costs and enhance motion diversity and quality across languages, as illustrated in Fig. 2.

**Cross-Lingual Alignment.** Inspired by AltCLIP [3], we enhance cross-lingual text understanding by fine-tuning the pretrained language model XLM [5] using knowledge distillation. The goal is to align sentence semantics across languages, ensuring that motion descriptions in different languages share a consistent latent representation, allowing the motion diffusion model to interpret them interchangeably. Specifically, we align the text embeddings of a student model  $E_\Phi^s$  from XLM [5] with those of a teacher model  $E_\Phi^t$  from OpenCLIP [4] by optimizing:

$$\mathcal{L}_{CLA}(\phi) = D_{KL}(F_{en}^M | F_{cn}^\phi) + D_{KL}(F_{cn}^\phi | F_{en}^M), \quad (1)$$

where  $F_{en}^M$  and  $F_{cn}^M$  represent the English and Chinese text embeddings encoded by model  $M$ . Minimizing this loss ensures that the student model  $E_\Phi^s$  learns to interpret text descriptions in both languages interchangeably. This alignment is crucial for maintaining motion fidelity and coherence without relying on separate models for each language.

**Bilingual Training.** For robust bilingual motion generation, we use a motion diffusion model [2] as our generative backbone, conditioning it on cross-lingually aligned text embeddings. Instead of training two distinct models, we introduce a language-agnostic training strategy: each motion sample is randomly paired with either an English or Chinese text description. This encourages the model to capture shared motion patterns while remaining sensitive to language-specific nuances. Our training loss is as:

$$\mathcal{L}_{BiMD}(\theta) = \mathbb{E}_{\epsilon, t, c} \left[ \|\epsilon - \epsilon_\theta(x_t, t, E_\Phi^s(c_s))\|_2^2 \right], \quad (2)$$

where  $c_s$  is the randomly selected English or Chinese description of motion  $x$ , and  $E_\Phi^s(c_s)$  represents its cross-

222 lingually aligned feature. This approach enables the model  
 223 to synthesize high-quality motion sequences from multilingual  
 224 text prompts without requiring separate models.

225 By integrating cross-lingual alignment with bilingual  
 226 diffusion training, our method overcomes the English-only  
 227 limitation of existing text-to-motion models, making text-  
 228 driven motion generation more accessible and effective  
 229 across languages. This paves the way for broader applica-  
 230 tions like virtual animation and multilingual interactions.

## 231 4. Step-Aware Reward-Guided Alignment

### 232 4.1. Motivation & Framework

233 **Preliminaries.** Existing diffusion-based motion generation  
 234 methods [2, 27] operate via a forward process and a reverse  
 235 process. The forward process gradually adds noise into the  
 236 real motion distribution  $p_{\text{data}}(\cdot)$  over timestep, and can be  
 237 modeled as a stochastic differential equation (SDE) [24]:

$$238 \quad dx = f(x, t)dt + g(t)dw, \quad (3)$$

239 where  $t$  is timestep,  $f(\cdot, \cdot)$  and  $g(\cdot)$  are the drift and diffusion  
 240 coefficients, and  $w$  is the Wiener process. For reverse pro-  
 241 cess, motions  $x$  are generated via trajectory sampling [24]:

$$242 \quad dx = [f(x, t) - g(t)^2 \nabla \log p_t(x)]dt + g(t)dw, \quad (4)$$

243 where  $\nabla \log p_t(x)$  is the score function of  $p_t(x)$ , directing  
 244 sampling toward higher-density regions.

245 **Motivation.** While bilingual diffusion models enable motion  
 246 generation across languages, they often fail to produce motions that accurately align with textual descriptions.  
 247 For example, as illustrated in Fig. 3, the diffusion model  
 248 prompted to generate a person walking forward to the right  
 249 may instead veer left. This misalignment arises as the sam-  
 250 pling distribution  $p_t(x)$ , learned from the diffusion, priori-  
 251 tizes high-probability regions over semantic fidelity.

252 Upon analyzing the diffusion sampling process (Fig. 3),  
 253 we identify a key issue: sampled motions  $x_t$  (stars) are  
 254 guided by gradient descent toward high-density regions  
 255  $p_t(\cdot)$  but consistently diverge from text embeddings  $c$  (tri-  
 256 angles). This bias prioritizes probability density over semantic  
 257 alignment, largely due to the reliance on CLIP [22] as the  
 258 text encoder. While aligning text with static images, CLIP  
 259 struggles with the temporal dynamics of motion, hindering  
 260 the diffusion model’s ability to learn a semantically coher-  
 261 ent sampling distribution.

262 A direct solution is to learn a latent space that aligns  
 263 motion-text pairs and then train the diffusion model accord-  
 264 ingly. However, the scarcity of motion-text pairs makes  
 265 it difficult to train a generalized text encoder for motion,  
 266 reducing the diffusion model’s generalization ability. In-  
 267 stead, we propose a more effective approach: leveraging  
 268 an already well-aligned distribution to guide the misaligned

270 sampling process. Accordingly, we first estimate a re-  
 271 ward distribution  $p_t^r(x)$  from text-motion pairs, capturing  
 272 semantic alignment. We then integrate this reward distri-  
 273 bution with the vanilla sampling distribution to construct  
 274 an ideal distribution  $p_t^I(x)$ . Crucially, our method is inde-  
 275 pendent of the diffusion training process, allowing seam-  
 276 less integration into any diffusion model without retraining.  
 277 As shown in Fig. 3, sampling from this ideal distribution  
 278 ensures both high-probability density and strong semantic  
 279 alignment, overcoming previous limitations.

280 **Overall Framework.** Our framework enhances diffusion-  
 281 based motion generation by constructing an ideal sam-  
 282 pling distribution that balances motion probability with  
 283 text-motion alignment. This section describes how we inte-  
 284 grate the reward distribution into the diffusion process and  
 285 sample from the resulting ideal distribution. Sec. 4.2 details  
 286 the estimation of the reward distribution, while Sec. 4.3 de-  
 287 scribes the motion sampling process.

288 Formally, assume a reward distribution  $p_t^r(x|c)$  has been  
 289 estimated. Then we define the ideal distribution as:

$$290 \quad p_t^I(x|c) = p_t(x|c)p_t^r(x|c)/Z(c), \quad (5)$$

291 where  $Z(c) = \int p_t(x|c)p_t^r(x|c)dx$  is a normalizing con-  
 292 stant. This formulation integrates the original sampling dis-  
 293 tribution  $p_t(x|c)$  with the reward distribution  $p_t^r(x|c)$ , bal-  
 294 ancing both probability density and text-motion alignment.

295 Using this ideal distribution, we modify the reverse  
 296 process for trading-off semantic alignment and high-  
 297 probability sampling as stated in the following theorem.

298 **Theorem 1.** When using the ideal sampling distribution  
 299  $p_t^I(x|c)$  in Eq. (5) to replace the vanilla sampling distri-  
 300 bution  $p_t(x|c)$ , the reverse SDE becomes:

$$301 \quad dx = [f(x, t) - g(t)^2 \nabla (\log p_t(x|c) + \log p_t^r(x|c))]dt + g(t)dw. \quad (6)$$

302 See its proof in App. D.1. Theorem 1 shows that the gra-  
 303 dient of the ideal sampling distribution decomposes into the  
 304 gradients of  $p_t(x|c)$  and  $p_t^r(x|c)$ . Since  $p_t(x|c)$  is already  
 305 known, the estimated reward distribution can directly guide  
 306 the sampling process toward the ideal distribution. Next, we  
 307 detail the estimation of the reward distribution (Sec. 4.2)  
 308 and outline the motion sampling procedure (Sec. 4.3).

### 309 4.2. Step-Aware Alignment for Reward Distribution

310 A core challenge in estimating the reward distribution  
 311  $p_t^r(x|c)$  is achieving precise motion-text alignment under  
 312 varying noise levels in the diffusion process. Existing meth-  
 313 ods [19] assume clean and noise-free motion sequences, and  
 314 overlook timestep-dependent distortions, resulting in coarse  
 315 and inconsistent alignments. This misalignment hinders accu-  
 316 rate reward estimation, which is critical for guiding sam-  
 317 pling toward semantically faithful motion generation. To  
 318 address this, we introduce a step-aware reward model for

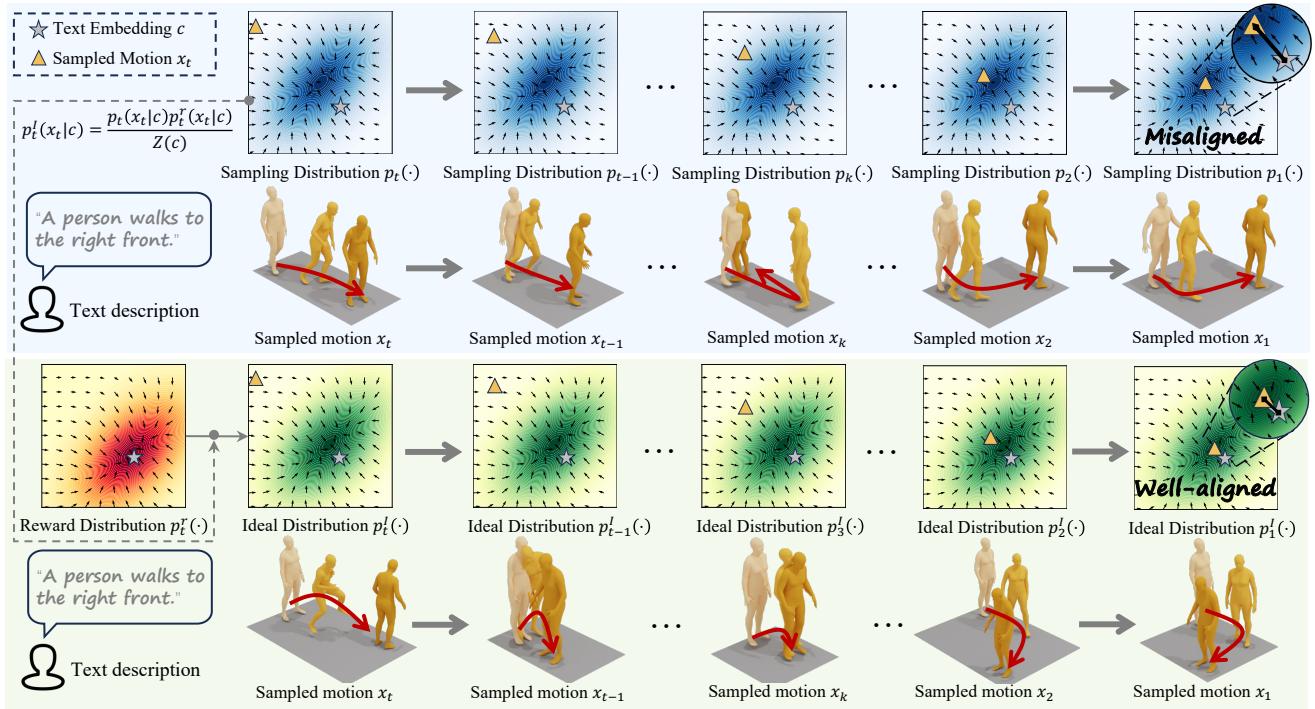


Figure 3. Illustration of the sampling process in diffusion-based motion generation frameworks. The blue region represents the sampling distribution  $p_t(\cdot)$  learned by the diffusion model, while the green region depicts the ideal sampling distribution  $p_t^I(\cdot)$  achieved by incorporating our proposed reward-guided sampling strategy with the sampling distribution  $p_t(\cdot)$ .

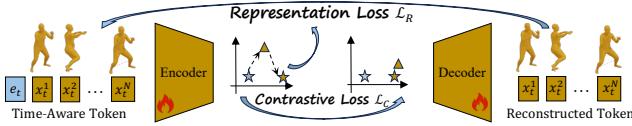


Figure 4. Framework of step-aware reward model. During this process, time-aware tokens, consisting of timestep embedding  $t$  and motion embeddings  $x_t^k$ , are aligned with text embedding  $c$  in the latent space and reconstructed via the decoder, with the encoder and decoder jointly optimized by contrastive loss  $\mathcal{L}_C$  [19] and representation loss  $\mathcal{L}_R$  [18].

noise-adaptive alignment and a motion-to-motion reward to ensure consistency with real-world motion patterns implied by text. These components are integrated into a unified reward distribution to enhance alignment and motion quality.

**Step-Aware Reward Model.** To mitigate timestep-dependent misalignment, we introduce a step-aware reward model  $R(\cdot)_\varphi$  illustrated in Fig. 4, which explicitly accounts for noise variations across diffusion timesteps. Unlike conventional alignment models [19], our approach incorporates a timestep token  $[e_t]$  into the motion representation, allowing the model to learn noise-dependent alignment patterns. Given an  $N$ -frame motion sequence  $[x_t^1, x_t^2, \dots, x_t^N]$ , we augment it with the timestep token to form the enriched representation  $[e_t, x_t^1, x_t^2, \dots, x_t^N]$ . This enables

the transformer-based encoder to process motion dynamics while adapting to different noise levels.

During training, noise is added to motion at timestep  $t$ , and the step-aware reward model  $R_\varphi(\mathbf{x}_t, c)$  is optimized by two complementary losses: a representation loss  $\mathcal{L}_R$  [18] to learn meaningful motion embeddings, and a contrastive loss  $\mathcal{L}_C$  [19] to ensure accurate text-motion retrieval:

$$\mathcal{L}_{RM}(\varphi; \mathbf{x}_t, c) = \mathcal{L}_C(\varphi; \mathbf{x}_t, c) + \mathcal{L}_R(\varphi; \mathbf{x}_t, c). \quad (7)$$

See more details to train this model in App. C.3.

Once trained, the step-aware reward model establishes a well-aligned latent space. Given a motion  $\mathbf{x}$  and text condition  $c$ , it evaluates their semantic alignment as:

$$R_\varphi(\mathbf{x}, c) = \cos(\mathbf{z}_x, \mathbf{z}_c), \quad (8)$$

where  $\mathbf{z}_x$  and  $\mathbf{z}_c$  are the respective motion and text embeddings in the learned latent space.

**Motion-to-Motion Reward.** While text-to-motion alignment is essential, text descriptions often exhibit ambiguity, leading to inconsistencies in generated motions. To mitigate this, we introduce a motion-to-motion reward, which evaluates alignment by comparing the generated motion  $\mathbf{x}_t$  with a reference motion  $\mathbf{x}^c$  retrieved from the training set  $\mathcal{D}_{tr}$ . The step-aware reward model is used to select  $\mathbf{x}^c$  as

333  
334  
335  
336  
337  
338  
339

341  
342  
343  
344

346  
347  
348  
349  
350  
351  
352  
353  
354

355 the closest match to the text condition  $c$ :

$$356 \quad \mathbf{x}^c = \arg \max_{\mathbf{x} \in \mathcal{D}_{tr}} R_\varphi(\mathbf{x}, c). \quad (9)$$

357 This retrieved motion  $\mathbf{x}^c$  acts as a dynamic anchor, ensuring  
358 that generated motions remain faithful to real-world motion  
359 patterns implied by the text. Accordingly, The motion-  
360 aligned reward is then computed as:

$$361 \quad R_m(\mathbf{x}_t, c) = \cos(\mathbf{z}_x, \mathbf{z}_{x^c}), \quad (10)$$

362 where  $\mathbf{z}_x$  and  $\mathbf{z}_{x^c}$  are the embeddings of the generated and  
363 retrieved motions, respectively. This ensures generated motions  
364 adhere to real-world motion patterns while maintaining  
365 semantic consistency.

366 **Reward Distribution.** With both the step-aware reward  
367 model and the motion-to-motion reward, we define the dual-  
368 alignment reward as:

$$369 \quad R(\mathbf{x}_t, c) = \mu R_\varphi(\mathbf{x}_t, c) + \eta R_m(\mathbf{x}_t, c), \quad (11)$$

370 where  $\mu$  and  $\eta$  control the contributions of text-based and  
371 motion-based alignment. This reward formulation defines  
372 the reward distribution over noised motion as:

$$373 \quad p_t^r(\mathbf{x}_t | c) = \exp(R(\mathbf{x}_t, c)) / Z^r(c). \quad (12)$$

374 Here,  $Z^r(c) = \int \exp(R_\varphi(\mathbf{x}, c)) d\mathbf{x}$  is for normalization.

375 By integrating text-motion and motion-motion alignment,  
376 our approach constructs a robust reward signal that  
377 captures both semantic consistency and motion coherence.  
378 This enables more precise guidance of the diffusion sam-  
379 pling process, ensuring that generated motions are not only  
380 probable but also faithful to their textual descriptions.

### 381 4.3. Reward-Guided Sampling

382 Building on the dual-alignment reward  $R(\mathbf{x}_t, c)$  and its as-  
383 sociated distribution  $p_t^r(\mathbf{x}_t | c)$ , we now integrate them into  
384 the reverse SDE to refine motion generation. The following  
385 theorem establishes how this reward distribution enhances  
386 sampling for precise text-conditioned synthesis.

387 **Theorem 2.** Given the reward distribution  $p_t^r(\mathbf{x}|c)$  defined  
388 in Eq. (12), the reverse SDE can be rewritten as:

$$389 \quad d\mathbf{x} = \left[ \mathbf{f}(\mathbf{x}, t) - g(t)^2 \nabla (\log p_t(\mathbf{x}|c) + R(\mathbf{x}_t, c)) \right] dt + g(t) d\mathbf{w}. \quad (13)$$

390 See its proof in App. D.2. Theorem 2 reveals that the  
391 reward gradient  $\nabla R(\mathbf{x}_t, c)$ , derived from both text-aligned  
392 and motion-aligned reward components, directly influences  
393 the sampling trajectory. Integrating these gradients into the  
394 reverse SDE can dynamically steer the sampling toward a  
395 distribution that better aligns with both textual conditions  
396 and realistic structures.

397 Building upon this continuous-time formulation, for  
398 practical motion generation we then derive its discrete ap-  
399 proximation within the DDPM [10] framework in the fol-  
400 lowing theorem. See proof in App. D.4.

---

### Algorithm 1 Reward-Guided Denoise Process

---

**Input:** diffusion model  $\epsilon_\theta$ , reward model  $R$ , training set  
 $\mathcal{D}_{tr}$ , condition  $c$ , timestep  $T$ .

**Output:** generated motion  $\mathbf{x}_0$

- 1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
  - 2:  $\mathbf{x}^c = \arg \max_{\mathbf{x} \in \mathcal{D}_{tr}} R_\varphi(\mathbf{x}, c)$
  - 3: **for**  $t = T, \dots, 1$  **do**
  - 4:   use  $\mathbf{x}^c$  to obtain reward score
  - 5:    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  **if**  $t > 1$  **else**  $\epsilon = \mathbf{0}$
  - 6:   use Eq. (15) to generate  $\mathbf{x}_{t-1}$
  - 7: **end for**
  - 8: **return**  $\mathbf{x}_0$
- 

401 **Theorem 3.** Given a reverse SDE defined in Eq. (13),  
402 adopting standard DDPM settings [10, 24] where  $\mathbf{f}(\mathbf{x}, t) =$   
403  $-\frac{1}{2}\bar{\beta}_{t+\Delta t}\mathbf{x}_t$ ,  $g(t) = \sqrt{\beta_{t+\Delta t}}$ , and  $\bar{\beta}_t = \frac{\beta_{t+\Delta t}}{\Delta t}$ , with time  
404 steps  $N \rightarrow \infty$  and step size  $\Delta t = \frac{1}{N}$ , the reward-guided  
405 denoising process is given by:

$$406 \quad \mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \bar{\mathbf{x}}_{t-1} + \sqrt{\beta_t} \epsilon \right) + \frac{\beta_t}{\sqrt{\alpha_t}} \nabla R(\mathbf{x}_t, c), \quad (14)$$

407 where  $\bar{\mathbf{x}}_{t-1} = \mathbf{x}_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t, c)$ ,  $\beta_t$  and  $\alpha_t$  are  
408 the noise schedule parameters,  $\epsilon_\theta(\cdot)$  represents the diffu-  
409 sion model network, and  $\epsilon$  is Gaussian noise sampled from  
410  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ .

411 Theorem 3 demonstrates that the reward gradient  
412  $\nabla R(\mathbf{x}_t, c)$ , weighted by  $\frac{\beta_t}{\sqrt{\alpha_t}}$ , progressively influences the  
413 denoising process, adapting the sampling trajectory toward  
414 a distribution that reflects the intended motion semantics.  
415 To ensure the sampling stability (see detailed discussion in  
416 App. D.4), we remove the weight  $\frac{\beta_t}{\sqrt{\alpha_t}}$  on the reward term,  
417 leading to a revised denoising process:

$$418 \quad \mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \bar{\mathbf{x}}_{t-1} + \sqrt{\beta_t} \epsilon \right) + \nabla R(\mathbf{x}_t, c). \quad (15)$$

419 Based on this theoretical framework, we introduce Algo-  
420 rithm 1 that integrates the step-aware reward model into the  
421 diffusion-based generation process.

## 5. Experiment

### 5.1. Experiment Setting

424 **Datasets & Evaluation.** We employ two widely used text-  
425 to-motion datasets, HumanML3D [7] and KIT-ML [21],  
426 alongside the proposed BiHumanML3D, for evaluation pur-  
427 poses. Consistent with the majority of prior studies [8, 13],  
428 we adopt R-Precision for Top  $k$ , Fréchet Inception Dis-  
429 tance (FID), Multi-Modal Distance (MM Dist), and Diver-  
430 sity as evaluation metrics to assess the generation quality  
431 and alignment accuracy of our model.

432 **Implementation Details.** Our bilingual motion diffusion  
433 model adopts MLD [2] as the backbone, adhering to its

Method	R Precision ↑			FID ↓	MM Dist ↓	Diversity →
	Top 1	Top 2	Top 3			
Real	0.511	0.703	0.797	0.002	2.974	9.503
T2M (2022) [7]	0.455±0.002	0.636±0.003	0.736±0.003	1.087±0.002	3.347±0.008	9.175±0.002
MDM (2023) [27]	0.455±0.006	0.645±0.007	0.749±0.006	0.489±0.047	3.330±0.25	9.920±0.083
T2M-GPT (2023) [31]	0.492±0.003	0.679±0.002	0.775±0.002	0.141±0.005	3.121±0.009	9.722±0.082
MLD (2023) [2]	0.481±0.003	0.673±0.003	0.772±0.002	0.473±0.013	3.196±0.010	9.724±0.082
Mo.Diffuse (2024) [34]	0.491±0.001	0.681±0.001	0.775±0.001	0.630±0.001	3.113±0.001	9.410±0.49
OMG (2024) [15]	-	-	0.784±0.002	0.381±0.008	-	9.657±0.085
MotionLCM (2025) [6]	0.502±0.003	0.698±0.002	0.798±0.002	0.304±0.012	3.012±0.007	9.607±0.066
LMM-T <sup>2</sup> (2024) [35]	0.496±0.002	0.685±0.002	0.785±0.002	0.415±0.002	3.087±0.012	9.176±0.074
Mo.Mamba (2025) [37]	0.502±0.003	0.693±0.002	0.792±0.002	0.281±0.011	3.060±0.000	9.871±0.084
CoMo (2024) [11]	0.502±0.002	0.692±0.007	0.790±0.002	0.262±0.004	3.032±0.015	9.936±0.066
ParCo (2025) [40]	0.515±0.003	0.706±0.003	0.801±0.002	0.109±0.005	2.927±0.008	9.576±0.088
BiMD (Ours)	0.499±0.002	0.691±0.002	0.789±0.003	0.397±0.011	3.105±0.009	9.635±0.089
w/ ReAlign (Ours)	0.566±0.003 (+13.4%)	0.759±0.002 (+9.7%)	0.847±0.002 (+7.4%)	0.178±0.006 (+55.2%)	2.714±0.007 (+12.6%)	9.573±0.068 (+47.0%)
MLD++ [6] (Baseline)	0.548±0.003	0.738±0.003	0.829±0.002	0.073±0.003	2.810±0.008	9.658±0.089
w/ ReAlign (Ours)	<b>0.572±0.002 (+4.4%)</b>	<b>0.764±0.002 (+3.5%)</b>	<b>0.852±0.001 (+2.8%)</b>	<b>0.055±0.003 (+24.7%)</b>	<b>2.648±0.008 (+5.8%)</b>	<b>9.478±0.055 (+83.9%)</b>

Table 1. Comparison of text-to-motion generation performance on the HumanML3D dataset. The arrows ↑, ↓, and → indicate higher, lower, and closer-to-real-motion values are better, respectively. **Bold** highlights the best results. Percentages in subscripts indicate improvements over respective baselines. Our BiMD adopts a similar backbone of MLD’s [2], and surpasses it on all metrics.

Method	R Precision ↑			FID ↓	MM Dist ↓	Diversity →
	Top 1	Top 2	Top 3			
Real	0.424	0.649	0.779	0.031	2.788	11.08
T2M (2022) [7]	0.361	0.559	0.681	3.022	2.052	10.72
MLD (2023) [2]	0.390	0.609	0.734	0.404	3.204	10.80
T2M-GPT (2023) [31]	0.416	0.627	0.745	0.514	3.007	10.86
CoMo (2024) [11]	0.422	0.638	0.765	0.332	2.873	10.95
Mo.Mamba (2025) [37]	0.419	0.645	0.765	0.307	3.021	11.02
ParCo (2025) [40]	0.430	0.649	0.772	0.453	2.820	10.95
Baseline [34]	0.417	0.621	0.739	1.954	2.958	<b>11.10</b>
w/ ReAlign (Ours)	0.419	0.639	0.764	0.805	2.801	10.66
Baseline (MDM) [27]	0.403	0.606	0.731	0.497	3.096	10.74
w/ ReAlign (Ours)	<b>0.451</b>	<b>0.664</b>	<b>0.784</b>	<b>0.276</b>	<b>2.775</b>	10.76

Table 2. Comparison of text-to-motion generation performance on the KIT-ML dataset. **Bold** highlights the best results. Since the models MLD [2] and MLD++ [6] for the KIT-ML dataset have not been released, we use the widely used MDM [27] as the baseline.

training configuration. For the cross-lingual alignment, we employ OpenCLIP [4] as the teacher model and use XLM-B [5] as the backbone of the student model, optimizing with  $\mathcal{L}_{CLA}$  defined in Eq. (1). The step-aware reward model is built based on SkipTransformer [2]. The architecture consists of a transformer encoder for text and motion and a decoder for motion, all with 9 layers and 4 heads, with a latent space dimension of 256. During training, the probability of noisy motion augmentation is set to 0.5, and the model is aware of a maximum timestep of 1000. We use AdamW [17] as the optimizer with a learning rate of  $10^{-4}$ , while other hyperparameter settings follow the TMR [19]. More details about experiment settings and results are provided in App. C.

## 5.2. Main Results

**Text-to-Motion Generation.** As shown in Tab. 1, our reward-guided sampling significantly enhances performance when integrated with state-of-the-art text-to-motion models. Specifically, our BiMD achieves performance com-

Method	CLA	Lang.	R Precision ↑			FID ↓	MM Dist ↓	Diversity →
			Top 1	Top 2	Top 3			
Real	-	CN	0.543	0.732	0.821	0.002	3.338	10.750
	-	EN	0.511	0.703	0.797	0.002	2.974	9.503
Mo.Diffuse [34]	X	CN	0.478	0.680	0.783	1.024	3.512	11.586
	✓	CN	0.502	0.696	0.791	0.643	3.356	11.064
	X	EN	0.491	0.681	0.782	0.630	3.113	9.410
MDM [27]	X	CN	0.481	0.673	0.774	0.908	3.482	11.674
	✓	CN	0.497	0.693	0.791	0.627	3.347	11.612
	X	EN	0.455	0.645	0.749	0.489	3.330	9.920
MLD [2]	X	CN	0.482	0.671	0.769	0.789	3.557	11.204
	X	EN	0.481	0.673	0.772	0.473	3.196	9.724
BiMD (Ours)	✓	CN	<b>0.505</b>	<b>0.696</b>	<b>0.792</b>	<b>0.528</b>	<b>3.338</b>	<b>10.741</b>
	✓	EN	<b>0.499</b>	<b>0.691</b>	<b>0.789</b>	<b>0.397</b>	<b>3.105</b>	<b>9.635</b>

Table 3. Comparison of text-to-motion generation performance on the BiHumanML3D dataset. “Lang.” indicates the evaluated language. English and Chinese results are assessed from the original evaluator [7] and our proposed evaluator, respectively. **Bold** highlights the best results. Given the absence of bilingual motion generation methods in current literature, our study adapts established monolingual frameworks [2, 27, 34], training separate, language-specific models to advance the field. The symbol “✓” at CLA” denotes methods that employ our cross-lingual alignment representation, indicating the use of a unified model for generation, whereas the symbol “X” means methods trained specifically for individual languages.

parable to previous SoTA methods and surpasses them by integrating the proposed reward-guided sampling. Notably, by using our proposed reward-guided sampling, MLD++ [6] achieves new SoTA results, with an R Precision@3 of 85.2% (+2.8%), alongside a reduction in FID of 0.055 (+24.7%) and an MM Dist to 2.648 (+5.8%). Furthermore, our reward-guided sampling significantly enhances the performance of MDM [2], yielding SoTA results on the KIT-ML dataset, with an R Precision@3 of 78.4% (+7.3%), alongside a reduction in FID of 0.276 (+44.5%) and an MM Dist to 2.775 (+10.4%). These consistent improvements over the baseline without reward-guided sampling demon-

434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464

Method	R Precision ↑			FID ↓	MM Dist ↓	Diversity →
	Top 1	Top 2	Top 3			
Real	0.511	0.703	0.797	0.002	2.974	9.503
MDiff [34]	0.491	0.681	0.775	0.630	3.113	9.410
w/ ReAlign	0.534	<sup>+8.8%</sup> 0.733	<sup>+7.6%</sup> 0.829	<sup>+7.0%</sup> 0.370	<sup>+41.3%</sup> 2.807	<sup>+9.9%</sup> 9.372
						<sup>-0.04</sup> <sub>-40.9%</sub>
MDM [27]	0.455	0.645	0.749	0.489	3.330	9.920
w/ ReAlign	0.470	<sup>+3.3%</sup> 0.677	<sup>+5.0%</sup> 0.789	<sup>+5.3%</sup> 0.325	<sup>+33.5%</sup> 3.129	<sup>+6.0%</sup> 9.355
						<sup>+0.27</sup> <sub>+64.5%</sub>
MLD [2]	0.481	0.673	0.772	0.473	3.196	9.724
w/ ReAlign	0.567	<sup>+17.9%</sup> 0.759	<sup>+12.8%</sup> 0.848	<sup>+9.8%</sup> 0.195	<sup>+58.8%</sup> 2.704	<sup>+15.4%</sup> 9.474
						<sup>+0.19</sup> <sub>+86.9%</sub>
MLCM <sup>1</sup> [6]	0.546	0.743	0.837	0.072	2.767	9.577
w/ ReAlign	0.555	<sup>+1.7%</sup> 0.751	<sup>+1.1%</sup> 0.841	<sup>+0.5%</sup> 0.088	<sup>-22.2%</sup> 2.726	<sup>+1.5%</sup> 9.541
						<sup>+0.04</sup> <sub>+48.6%</sub>
MLCM <sup>4</sup> [6]	0.502	0.698	0.798	0.304	3.012	9.607
w/ ReAlign	0.540	<sup>+7.6%</sup> 0.739	<sup>+5.9%</sup> 0.833	<sup>+4.4%</sup> 0.273	<sup>+10.2%</sup> 2.797	<sup>+7.1%</sup> 9.683
						<sup>-0.08</sup> <sub>-73.1%</sub>
MLD++ [6]	0.548	0.738	0.829	0.073	2.810	9.658
w/ ReAlign	0.572	<sup>+4.4%</sup> 0.764	<sup>+3.5%</sup> 0.852	<sup>+2.8%</sup> 0.055	<sup>+24.7%</sup> 2.648	<sup>+5.8%</sup> 9.478
						<sup>+0.13</sup> <sub>+83.9%</sub>
BiMD	0.499	0.691	0.789	0.397	3.105	9.635
w/ ReAlign	0.566	<sup>+13.4%</sup> 0.759	<sup>+9.7%</sup> 0.847	<sup>+7.4%</sup> 0.178	<sup>+55.2%</sup> 2.714	<sup>+12.6%</sup> 9.573
						<sup>+0.06</sup> <sub>+47.0%</sub>

Table 4. **Performance enhancement of motion generation methods with plug-and-play step-aware reward guidance.** Results are evaluated on the HumanML3D dataset, with improvements reported relative to baseline methods. Here, MLCM<sup>1</sup> and MLCM<sup>4</sup> denote the 1-step and 4-step model in MotionLCM [6]. MDiff is an abbreviation of MotionDiffuse [34].

strate the effectiveness of our reward-guided sampling in enhancing text-motion alignment quality.

**Bilingual Text-to-Motion Generation.** Table 3 reports bilingual text-to-motion generation results on our BiHumanML3D dataset. For a fair evaluation, we adapt state-of-the-art models, Mo.Diffuse, MLD, and MDM, to the Chinese context using the corresponding BiHumanML3D subset, comparing our cross-lingual alignment (CLA) approach against a language-specific baseline. Since these models are designed for English text-to-motion generation, we use ChineseCLIP [29] to replace their text encoder, adapting to Chinese settings. To extend these models for bilingual settings, we further substitute ChineseCLIP [29] with our CLA encoder. For all baselines, CLA integration yields notable gains; for example, in MDM, R-Precision rises from 77.4% to 79.1%, and FID drops from 0.908 to 0.627 compared to non-CLA models. This underscores the critical role of cross-lingual alignment in bilingual text-to-motion generation. Moreover, our BiMD model surpasses all baselines, achieving superior performance in both languages. These results highlight the synergy between BiMD and CLA, as well as the substantial advancement our approach offers over existing state-of-the-art methods.

**Plug-and-Play Functionality of ReAlign.** To demonstrate the plug-and-play capability and generalizability of our step-aware reward-guided alignment, we integrate it into various baseline models for text-to-motion generation, as shown in Tab. 4. Across methods such as Mo.Diffuse [34], MDM [27], MLD [2], MotionLCM [6], MLD++ [6], and our BiMD, our ReAlign consistently enhances performance. Notably, it achieves substantial improvements in

T2M	M2M	SA	R Precision ↑			FID ↓	MM Dist ↓	Diversity →
			Top 1	Top 2	Top 3			
✗	✗	✗	0.499	0.691	0.789	0.397	3.105	9.635
✓	✗	✗	0.557	0.749	0.840	0.216	2.760	<b>9.513</b>
✗	✓	✗	0.522	0.715	0.809	0.188	2.932	0.455
✓	✓	✗	0.557	0.751	0.841	<u>0.182</u>	2.748	9.530
✓	✗	✓	<b>0.567</b>	<b>0.760</b>	<b>0.849</b>	0.196	<b>2.721</b>	9.598
✗	✓	✓	0.521	0.711	0.806	0.210	2.961	<b>9.526</b>
✓	✓	✓	<u>0.566</u>	<u>0.759</u>	<u>0.847</u>	<b>0.178</b>	<b>2.714</b>	9.573

Table 5. **Ablation study of the text-to-motion on HumanML3D dataset.** Here, “T2M”, “M2M” and “SA” denote whether the text-to-motion reward, motion-to-motion reward and step-aware training is used, respectively.

alignment quality and motion realism, with relative gains of up to 17.9% in R Precision Top 1 and 66.1% in FID for BiMD. The approach also refines multimodal alignment while maintaining diversity close to real motions. While diversity slightly decreases in some cases, this is expected and beneficial. Better diversity does not always indicate better quality, as it simply reflects motion variety. ReAlign prioritizes well-aligned motions over misaligned ones, leading to significant gains in other metrics without compromising generation quality. These results underscore the plug-and-play versatility of this module, effectively elevating the efficacy of diverse motion generation frameworks.

### 5.3. Ablation Study

**Effectiveness of Reward Model.** We assess the dual reward mechanism, including T2M and M2M alignment rewards, along with the step-aware strategy in text-to-motion generation. As shown in Tab. 5, results indicate that the T2M reward significantly improves the alignment between the generated motions and text descriptions, as well as the realism of the motions. While the M2M reward alone has limited impact, combining it with the step-aware strategy further enhances motion realism. The combination of T2M and step-aware already yields the best result, with the addition of M2M providing a further boost to realism.

## 6. Conclusion

This paper introduces bilingual text-to-motion generation with BiHumanML3D, the first bilingual dataset, and BiMD, a unified diffusion model utilizing cross-lingual alignment for efficient motion generation. To tackle text-motion misalignment, we propose a step-aware reward model plugged into a pretrained diffusion model, enabling reward-guided sampling without further training.

**Limitation discussion.** Our BiMD aims to sample motion from the pretrained motion representation space. However, the scarcity of motion data limits the ability of pretrained models, such as variational autoencoders, to effectively extract rich semantic features, thereby constraining the generalizability of our model.

534

## References

535

- [1] Huayu Chen, Cheng Lu, Zhengyi Wang, Hang Su, and Jun Zhu. Score regularized policy optimization through diffusion behavior. *arXiv preprint arXiv:2310.07297*, 2023. 2
- [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, 3, 4, 6, 7, 8
- [3] Zhongzhi Chen, Guang Liu, Bo-Wen Zhang, Fulong Ye, Qinghong Yang, and Ledell Wu. Altclip: Altering the language encoder in clip for extended language capabilities. *arXiv preprint arXiv:2211.06679*, 2022. 3
- [4] Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. Reproducible scaling laws for contrastive language-image learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2818–2829, 2023. 3, 7
- [5] Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. Unsupervised cross-lingual representation learning at scale. *arXiv preprint arXiv:1911.02116*, 2019. 3, 7
- [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 *European Conference on Computer Vision*, pages 390–408, Cham, 2024. Springer Nature Switzerland. 1, 2, 7, 8
- [7] 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 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5142–5151, 2022. 1, 2, 3, 6, 7
- [8] Chuan Guo, Yuxuan Mu, Muhammad Gohar Javed, Sen Wang, and Li Cheng. Momask: Generative masked modeling of 3d human motions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1900–1910, 2024. 6
- [9] Gaoge Han, Mingjiang Liang, Jinglei Tang, Yongkang Cheng, Wei Liu, and Shaoli Huang. Reindiffuse: Crafting physically plausible motions with reinforced diffusion model. *arXiv preprint arXiv:2410.07296*, 2024. 2
- [10] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 6
- [11] Yiming Huang, Weilin Wan, Yue Yang, Chris Callison-Burch, Mark Yatskar, and Lingjie Liu. Como: Controllable motion generation through language guided pose code editing. In *Computer Vision – ECCV 2024: 18th European Conference, Milan, Italy, September 29–October 4, 2024, Proceedings, Part XXIX*, page 180–196, Berlin, Heidelberg, 2024. Springer-Verlag. 7
- [12] Biao Jiang, Xin Chen, Wen Liu, Jingyi Yu, Gang Yu, and Tao Chen. Motiongpt: Human motion as a foreign language. *Advances in Neural Information Processing Systems*, 36, 2024. 1

534

- [13] Zhe Li, Weihao Yuan, Yisheng He, Lingteng Qiu, Shen-hao Zhu, Xiaodong Gu, Weichao Shen, Yuan Dong, Zilong Dong, and Laurence T. Yang. Lamp: Language-motion pretraining for motion generation, retrieval, and captioning. 2024. 6
- [14] Zhimin Li, Jianwei Zhang, Qin Lin, Jiangfeng Xiong, Yanxin Long, Xinch Deng, Yingfang Zhang, Xingchao Liu, Minbin Huang, Zedong Xiao, et al. Hunyuan-dit: A powerful multi-resolution diffusion transformer with fine-grained chinese understanding. *arXiv preprint arXiv:2405.08748*, 2024. 1, 2
- [15] Han Liang, Jiacheng Bao, Ruichi Zhang, Sihan Ren, Yuecheng Xu, Sibei Yang, Xin Chen, Jingyi Yu, and Lan Xu. Omg: Towards open-vocabulary motion generation via mixture of controllers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 482–493, 2024. 7
- [16] Xiaoyang Liu, Yunyao Mao, Wengang Zhou, and Houqiang Li. Motionrl: Align text-to-motion generation to human preferences with multi-reward reinforcement learning. *arXiv preprint arXiv:2410.06513*, 2024. 2
- [17] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 7
- [18] Mathis Petrovich, Michael J Black, and Gü̈l Varol. Temos: Generating diverse human motions from textual descriptions. In *European Conference on Computer Vision*, pages 480–497. Springer, 2022. 5
- [19] Mathis Petrovich, Michael J. Black, and Gü̈l Varol. TMR: Text-to-motion retrieval using contrastive 3D human motion synthesis. In *International Conference on Computer Vision (ICCV)*, 2023. 4, 5, 7
- [20] Matthias Plappert, Christian Mandery, and Tamim Asfour. The kit motion-language dataset. *Big Data*, 4(4):236–252, 2016. PMID: 27992262. 1
- [21] Matthias Plappert, Christian Mandery, and Tamim Asfour. The kit motion-language dataset. *Big data*, 4(4):236–252, 2016. 6
- [22] 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 *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 2, 4
- [23] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36:53728–53741, 2023. 2
- [24] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020. 4, 6
- [25] Xiaofeng Tan, Hongsong Wang, Xin Geng, and Pan Zhou. Sopo: Text-to-motion generation using semi-online preference optimization. *arXiv preprint arXiv:2405.14734*, 2024. 2
- [26] Guy Tevet, Brian Gordon, Amir Hertz, Amit H Bermano, and Daniel Cohen-Or. Motionclip: Exposing human motion 1

- 649 generation to clip space. In *European Conference on Computer Vision*, pages 358–374. Springer, 2022. 2
- 650
- 651 [27] 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. 1, 2, 3, 4, 7, 8
- 652
- 653 [28] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8228–8238, 2024. 2
- 654
- 655 [29] An Yang, Junshu Pan, Junyang Lin, Rui Men, Yichang Zhang, Jingren Zhou, and Chang Zhou. Chinese clip: Contrastive vision-language pretraining in chinese. *arXiv preprint arXiv:2211.01335*, 2022. 8
- 656
- 657 [30] Weihao Yuan, Yisheng He, Weichao Shen, Yuan Dong, Xiaodong Gu, Zilong Dong, Liefeng Bo, and Qixing Huang. Mogents: Motion generation based on spatial-temporal joint modeling. *Advances in Neural Information Processing Systems*, 37:130739–130763, 2025. 2
- 658
- 659 [31] Jianrong Zhang, Yangsong Zhang, Xiaodong Cun, Yong Zhang, Hongwei Zhao, Hongtao Lu, Xi Shen, and Ying Shan. Generating human motion from textual descriptions with discrete representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14730–14740, 2023. 7
- 660
- 661 [32] Mingyuan Zhang, Xinying Guo, Liang Pan, Zhongang Cai, Fangzhou Hong, Huirong Li, Lei Yang, and Ziwei Liu. Remodiffuse: Retrieval-augmented motion diffusion model. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 364–373, 2023. 2
- 662
- 663 [33] Mingyuan Zhang, Huirong Li, Zhongang Cai, Jiawei Ren, Lei Yang, and Ziwei Liu. Finemogen: Fine-grained spatio-temporal motion generation and editing. *Advances in Neural Information Processing Systems*, 36:13981–13992, 2023. 2
- 664
- 665 [34] Mingyuan Zhang, Zhongang Cai, Liang Pan, Fangzhou Hong, Xinying Guo, Lei Yang, and Ziwei Liu. Motondifuse: Text-driven human motion generation with diffusion model. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(6):4115–4128, 2024. 2, 7, 8
- 666
- 667 [35] Mingyuan Zhang, Daisheng Jin, Chenyang Gu, Fangzhou Hong, Zhongang Cai, Jingfang Huang, Chongzhi Zhang, Xinying Guo, Lei Yang, Ying He, and Ziwei Liu. Large motion model for unified multi-modal motion generation. In *European Conference on Computer Vision*, page 397–421. Springer, 2024. 7
- 668
- 669 [36] Mingyuan Zhang, Daisheng Jin, Chenyang Gu, Fangzhou Hong, Zhongang Cai, Jingfang Huang, Chongzhi Zhang, Xinying Guo, Lei Yang, Ying He, et al. Large motion model for unified multi-modal motion generation. In *European Conference on Computer Vision*, pages 397–421. Springer, 2025. 2
- 670
- 671 [37] Zeyu Zhang, Akide Liu, Ian Reid, Richard Hartley, Bohan Zhuang, and Hao Tang. Motion mamba: Efficient and long sequence motion generation. In *European Conference on Computer Vision*, pages 265–282. Springer Nature Switzerland, 2024. 7
- 672
- 673
- 674
- 675
- 676
- 677
- 678
- 679
- 680
- 681
- 682
- 683
- 684
- 685
- 686
- 687
- 688
- 689
- 690
- 691
- 692
- 693
- 694
- 695
- 696
- 697
- 698
- 699
- 700
- 701
- 702
- 703
- 704
- 705
- 706
- 707 [38] Zeyu Zhang, Akide Liu, Ian Reid, Richard Hartley, Bohan Zhuang, and Hao Tang. Motion mamba: Efficient and long sequence motion generation. In *European Conference on Computer Vision*, pages 265–282. Springer, 2024. 2
- 708
- 709
- 710
- 711 [39] Chongyang Zhong, Lei Hu, Zihao Zhang, and Shihong Xia. Attt2m: Text-driven human motion generation with multi-perspective attention mechanism. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 509–519, 2023. 2
- 712
- 713
- 714
- 715
- 716 [40] Qiran Zou, Shanyuan Yuan, Shian Du, Yu Wang, Chang Liu, Yi Xu, Jie Chen, and Xiangyang Ji. Parco: Part-coordinating text-to-motion synthesis. In *Computer Vision – ECCV 2024*, pages 126–143, Cham, 2025. Springer Nature Switzerland. 7
- 717
- 718
- 719
- 720
- 721 [41] Ronglai Zuo, Rolandos Alexandros Potamias, Evangelos Ververas, Jiankang Deng, and Stefanos Zafeiriou. Signs as tokens: An autoregressive multilingual sign language generator. *arXiv preprint arXiv:2411.17799*, 2024. 1, 2
- 722
- 723
- 724