

LUMA: LOW-DIMENSION UNIFIED MOTION ALIGNMENT WITH DUAL-PATH ANCHORING FOR TEXT-TO-MOTION DIFFUSION MODEL

Anonymous authors

Paper under double-blind review

ABSTRACT

While current diffusion-based models, typically built on U-Net architectures, have shown promising results on the text-to-motion generation task, they still suffer from semantic misalignment and kinematic artifacts. Through analysis, we identify severe gradient attenuation in the deep layers of the network as a key bottleneck, leading to insufficient learning of high-level features. To address this issue, we propose **LUMA** (*Low-dimension Unified Motion Alignment*), a text-to-motion diffusion model that incorporates dual-path anchoring to enhance semantic alignment. The first path incorporates a lightweight MoCLIP model trained via contrastive learning without relying on external data, offering semantic supervision in the temporal domain. The second path introduces complementary alignment signals in the frequency domain, extracted from low-frequency DCT components known for their rich semantic content. These two anchors are adaptively fused through a temporal modulation mechanism, allowing the model to progressively transition from coarse alignment to fine-grained semantic refinement throughout the denoising process. Experimental results on HumanML3D and KIT-ML demonstrate that LUMA achieves state-of-the-art performance, with FID scores of 0.035 and 0.123, respectively. Furthermore, LUMA accelerates convergence by $1.4\times$ compared to the baseline, making it an efficient and scalable solution for high-fidelity text-to-motion generation.



Figure 1: LUMA generates high-fidelity 3D human motion from natural language input. It injects temporal semantic anchor from the lightweight MoCLIP encoder and frequency semantic anchor from low-frequency DCT directly into the diffusion backbone. This design overcomes gradient attenuation in deep layers and achieves state-of-the-art motion fidelity.

1 INTRODUCTION

Text-driven human motion generation has recently garnered considerable attention due to the semantic richness and intuitive nature of natural language descriptions. This technology holds great potential for a wide range of applications, including animation, film production, virtual/augmented reality (VR/AR), and robotics Chen et al. (2024; 2025); Pinyoanuntapong et al. (2024c;a;b). Re-

cent performance gains have largely stemmed from the development of increasingly sophisticated architectures. Notable examples include ReMoDiffuse Zhang et al. (2023b), MotionLCM Dai et al. (2024) and MMM Pinyoanuntapong et al. (2024c). While these approaches have achieved notable success, they often rely on progressively larger and more complex models, resulting in diminishing efficiency and limited practical improvement.

Current research in text-to-motion generation mainly follows two paradigms. The first employs VAE-based models that encode motion into discrete tokens, predicted via autoregressive (AR) Zhang et al. (2023a); Jiang et al. (2023) or non-autoregressive (NAR) Guo et al. (2023); Pinyoanuntapong et al. (2024c;b) frameworks. The second leverages diffusion-based models, which iteratively transform Gaussian noise into motion sequences under text guidance Tevet et al. (2022); Chen et al. (2023); Huang et al. (2024); Zhang et al. (2023b); Dai et al. (2024). While VQ-based discrete methods have gained popularity Guo et al. (2023); Pinyoanuntapong et al. (2024b), they often suffer from tokenization-induced information loss and limited motion diversity Zhang et al. (2023a); Guo et al. (2022b). In contrast, diffusion models offer finer motion details, greater diversity, and improved physical plausibility, making them a promising alternative Yuan et al. (2023b); Zhang et al. (2022); Tevet et al. (2022); Chen et al. (2025). Notably, recent state-of-the-art diffusion models such as StableMotion Huang et al. (2024) and MLD Chen et al. (2023) adopt U-Net architectures as the core component for intermediate feature extraction and motion reconstruction. However, these models are limited by their high computational cost during training and the relatively low semantic fidelity of the generated motions.

Through experimental analysis, we identify the root cause of this issue as gradient shifts in the deeper layers of the U-Net. Specifically, we observe that the gradients in the downsampling and bottleneck layers are significantly smaller, and in some cases, nearly vanish compared to those in the upsampling layers. This phenomenon leads to the network’s inability to effectively learn high-level abstract features from the input data Srivastava et al. (2015); He et al. (2016); Alzubaidi et al. (2021), resulting in slow or even stalled convergence during training. Consequently, the generated motion sequences often suffer from missing fine-grained details, incomplete structural representations, and inaccurate semantic alignment. These limitations ultimately hinder the overall performance of the model.

To address this issue, an intuitive approach is to introduce additional dense supervision signals at intermediate layers, rather than relying solely on sparse loss functions applied to the final output. Existing methods, such as Representational Alignment (REPA) Yu et al. (2025), significantly enhance semantic expressiveness by utilizing intermediate features from large-scale pre-trained models (e.g., DINOv2-g) Oquab et al. (2024) as additional supervisory signals. Nevertheless, REPA heavily relies on external pre-trained models, limiting its applicability in motion tasks lacking suitable pre-trained models or with limited high quality motion data.

While a wide range of options exists for selecting dense supervision signals Wang & He (2025), it is critical to ensure they effectively enhance the semantic representation capability of deep layers. To this end, we leverage supervision signals from both temporal and spatial-frequency dimensions. Consequently, we propose a dual-path unified semantic alignment framework that eliminates the reliance on large-scale pretrained model supervision. In the temporal path, we introduce *MoCLIP*, a lightweight contrastively trained text-motion encoder, to extract semantic representations directly from motion sequences without any external corpora or teachers. The second path introduces an orthogonal frequency domain alignment anchor which serves as a complementary counterpart to the temporal domain. Specifically, motion data is transformed into the frequency domain using Discrete Cosine Transform (DCT), and its low-frequency components, recognized for their rich semantic content, are used as stable and orthogonal supervisory signals. To effectively integrate these two complementary supervision signals, we design an adaptive FiLM Perez et al. (2017) modulation module that dynamically and progressively adjusts the strength of semantic feature injection throughout the network.

We summarize our contributions as follows:

- (i) We conduct the first systematic gradient analysis of diffusion-based text-to-motion models, revealing severe sparsity in down-sampling and bottleneck layers that hinders semantic learning and slows convergence.

- 108 (ii) We introduce *MoCLIP*, a new motion encoder for text–motion alignment. MoCLIP couples
109 a transformer based motion encoder with a CLIP text encoder under a contrastive objective,
110 and its token-wise text embeddings are fused into each U-Net block via cross-attention to
111 strengthen semantic grounding without relying on large pretrained teachers.
- 112 (iii) We propose *Dual-Path Timestep-Aware Semantic Anchoring*, which injects complementary
113 temporal and frequency semantic anchors via FiLM modulation to revitalize deep-layer
114 gradients, eliminating the need for large-scale pretrained teachers.
- 115 (iv) Extensive experiments on multiple datasets demonstrate that our LUMA framework
116 achieves state-of-the-art performance across various metrics, while converging significantly
117 faster than baselines.
- 118
119
120

121 2 RELATED WORK

122 2.1 TEXT-TO-MOTION GENERATION.

123
124
125 Text-to-motion generation has been primarily driven by two dominant paradigms: diffusion-
126 based denoising models and autoregressive token-based models. Diffusion frameworks such as
127 MDM Tevet et al. (2022), MotionDiffuse Zhang et al. (2024a), and PhysDiff Yuan et al. (2023a)
128 iteratively refine noise into motion using cross-attention and physics-guided denoisers, establishing
129 strong baselines for realism. Latent diffusion variants, exemplified by MLD Chen et al. (2023) and
130 MotionLCM Dai et al. (2024), compress motion trajectories into compact latent codes and achieve
131 real-time sampling. [Retrieval-based methods like ReMoDiffuse Zhang et al. \(2023b\) inject exem-](#)
132 [plar motions for style control, while control-oriented pipelines incorporate explicit constraints to](#)
133 [improve diversity, alignment, and user steering.](#)

134 On the other track, autoregressive approaches such as T2M-GPT Zhang et al. (2023a), MoMask Guo
135 et al. (2023), BAMB Pinyoanuntapong et al. (2024b), InfiniMotion Zhang et al. (2024b), and
136 TEACH Athanasiou et al. (2022) treat motion as token sequences and rely on transformer language
137 models to generate coherent long-range structure. [More recently, LaMP Li et al. \(2025\) revisits](#)
138 [autoregressive modeling through language–motion pretraining, improving cross-modal alignment.](#)
139 However, these models are still trained using single-scale and low-level supervision, typically frame-
140 wise MSE or token reconstruction. This weakens gradient propagation in deeper layers and leads
141 to semantic dilution and loss of detail on complex prompts. Multi-scale and cross-modal objectives
142 remain essential for addressing these limitations. Our work addresses this by introducing temporal
143 and frequency semantic anchors, combined with timestep-aware modulation to guide the denoising
144 process.

145 2.2 REPRESENTATION ALIGNMENT FOR DIFFUSION MODELS.

146
147
148 Recent studies have explored enhancing diffusion models by aligning their internal representations
149 with high-level semantic features to improve learning efficiency and generation quality. Representa-
150 tion Alignment (REPA) Yu et al. (2025) explicitly aligns intermediate hidden states of diffusion
151 transformers with semantically rich features from large-scale vision encoders (e.g., DINOv2, CLIP),
152 yielding over 17× faster convergence and improved generation quality. Extensions rapidly broaden
153 the scope: REPA-E Leng et al. (2025) unlocks end-to-end training by jointly tuning the VAE tok-
154 enizer alongside the diffusion backbone; U-REPA Tian et al. (2025) adapts the paradigm to U-Net
155 architectures through spatial up-sampling and a manifold loss; and HASTE Wang et al. (2025)
156 couples holistic (feature + attention) alignment with stage-wise termination to mitigate late-stage
157 capacity mismatch. However, these methods rely on large-scale pretrained teachers and abundant
158 paired vision data, which are unavailable for text-to-motion tasks. This makes it difficult to replicate
159 the strong teacher–student setup that REPA relies on. Our work differs from existing representa-
160 tion alignment methods by removing the reliance on pretrained teachers and introducing temporal
161 and frequency anchors that are learned jointly with the diffusion model through timestep-aware
modulation. This design improves deep-layer gradients and maintains motion detail under limited
supervision.

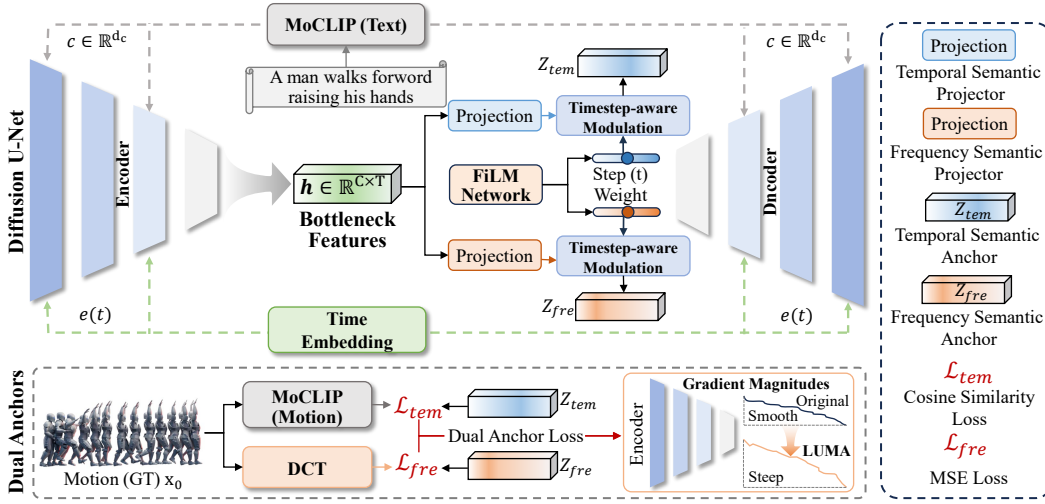


Figure 2: **Overview of the LUMA framework.** Text embeddings from MoCLIP are injected into each U-Net block via cross-attention (dashed lines). Intermediate features are projected into two paths: a temporal semantic branch aligned with MoCLIP features and a frequency semantic branch aligned with DCT coefficients. Both branches are modulated by timestep embeddings via FiLM and integrated with the main U-Net stream to guide denoising.

3 METHOD

3.1 PRELIMINARIES

3.1.1 TEXT-TO-MOTION DIFFUSION BASICS.

Let a natural-language prompt \mathcal{T} be embedded by a language encoder into a vector $\mathbf{c} = \text{CLIP}(\mathcal{T}) \in \mathbb{R}^{d_c}$. The ground-truth motion $\mathbf{x}_0 \in \mathbb{R}^{N \times d_m}$ contains N frames, each parameterised by d_m joints or 6-DoF pose coefficients. Following the DDPM formulation adopted by StableMoFusion Huang et al. (2024), we corrupt \mathbf{x}_0 with a variance schedule $\{\beta_t\}_{t=1}^T$:

$$\mathbf{x}_t = \sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad \alpha_t = \prod_{i=1}^t (1 - \beta_i) \quad (1)$$

A UNet-style denoiser G_θ is trained to predict \mathbf{x}_0 (or, equivalently, the added noise) from \mathbf{x}_t given the timestep t and the text embedding \mathbf{c} :

$$\min_{\theta} \mathbb{E}_{t, \mathbf{x}_0, \boldsymbol{\epsilon}, \mathbf{c}} \|G_\theta(\mathbf{x}_t, t, \mathbf{c}) - \mathbf{x}_0\|_2^2 \quad (2)$$

At inference time a sampler such as DPM-Solver++ Lu et al. (2023) iteratively denoises from pure Gaussian noise $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, while *classifier-free guidance* combines the conditional and unconditional scores:

$$G_{\text{CFG}}(\mathbf{x}_t, t, \mathbf{c}) = G_\theta(\mathbf{x}_t, t, \varnothing) + \omega [G_\theta(\mathbf{x}_t, t, \mathbf{c}) - G_\theta(\mathbf{x}_t, t, \varnothing)] \quad (3)$$

Here, \varnothing denotes the unconditional input, and ω is a scalar hyperparameter.

3.2 LUMA ARCHITECTURE

Building on the StableMoFusion Huang et al. (2024) backbone, we present the **Low-dimensional Unified Motion Alignment (LUMA)** framework, as illustrated in Figure 2. LUMA introduces dual semantic anchors at the final downsampling block of the UNet architecture, aiming to enhance semantic alignment and improve gradient flow in deep layers. These anchors deliver high-signal

gradients to otherwise deep, semantically sparse and low-gradient layers, thereby alleviating gradient divergence and thus achieving a more stable and efficient training process. Meanwhile, to better capture motion semantics and encode action descriptions in text, we propose a novel text-motion alignment model, **MoCLIP**, trained from scratch using contrastive learning (see details in Appendix B). Its text encoder provides token-wise embeddings that are injected into each UNet block via cross attention.

3.2.1 DUAL ANCHORS.

The low-dimensional feature $\mathbf{h} \in \mathbb{R}^{C \times T}$ is extracted from the final downsampling block of the UNet, which serves as the information bridge between the encoder and decoder. Here, C is the channel number and T is the number of motion frames. This bottleneck stage effectively aggregates global context while preserving essential temporal structures, making it an ideal location for extracting semantically rich and compact representations.

Accordingly, LUMA applies two lightweight multilayer perceptrons (MLPs) to project \mathbf{h} into distinct frequency semantic and temporal semantic representations: a frequency semantic projector $P_\theta : \mathbb{R}^{C \times T} \rightarrow \mathbb{R}^{d_m \times F}$, where d_m is the motion-feature dimension and F is the frequency semantic anchor dimension; and a temporal semantic projector $Q_\theta : \mathbb{R}^{C \times T} \rightarrow \mathbb{R}^{D_a}$, where D_a denotes the temporal semantic anchor dimension.

Timestep-aware modulation. To address varying gradient propagation across diffusion timesteps, where early stages have larger gradients and later stages smaller ones, we introduce Timestep-aware Modulation. For each projection, we employ Feature-wise Linear Modulation (FiLM) Perez et al. (2017) to adaptively scale and shift the features based on the current timestep. Let $\mathbf{e}(t) \in \mathbb{R}^{d_{\text{emb}}}$ denote the sinusoidal embedding of t , where d_{emb} is the embedding dimension. Two FiLM networks $F_\psi : \mathbb{R}^{d_{\text{emb}}} \rightarrow \mathbb{R}^{2d_m}$ and $F_{\psi'} : \mathbb{R}^{d_{\text{emb}}} \rightarrow \mathbb{R}^{2D_a}$ produce scaling and shifting vectors:

$$(\gamma_f(t), \beta_f(t)) = F_\psi(\mathbf{e}(t)), \quad (4)$$

$$(\gamma_t(t), \beta_t(t)) = F_{\psi'}(\mathbf{e}(t)), \quad (5)$$

where $\gamma_f(t), \beta_f(t) \in \mathbb{R}^{d_m}$ are for the frequency semantic anchor, and $\gamma_t(t), \beta_t(t) \in \mathbb{R}^{D_a}$ are for the temporal semantic anchor. The modulated features are:

$$\tilde{\mathbf{h}}_f = \left((1 + \gamma_f(t)) \odot P_\theta(\mathbf{h}) + \beta_f(t) \right)^\top, \quad (6)$$

$$\tilde{\mathbf{h}}_t = (1 + \gamma_t(t)) \odot Q_\theta(\mathbf{h}) + \beta_t(t). \quad (7)$$

The superscript \top permutes the last two dimensions, $[d_m, F] \rightarrow [F, d_m]$, not a matrix transpose. The operator \odot denotes element-wise multiplication, with broadcasting along F in the frequency branch and none in the temporal branch.

Frequency semantic anchor. The **frequency semantic anchor** $\mathbf{z}_{\text{fre}} = \tilde{\mathbf{h}}_f \in \mathbb{R}^{F \times d_m}$ is aligned with the first k DCT coefficients of the target motion $\mathbf{x}_0 \in \mathbb{R}^{N \times d_m}$ (where N is the number of frames and d_m is the motion feature dimension):

$$\mathcal{L}_{\text{fre}} = \|\mathbf{z}_{\text{fre}} - \text{DCT}_k(\mathbf{x}_0)\|_2^2, \quad (8)$$

where $\text{DCT}_k(\mathbf{x}_0) \in \mathbb{R}^{F \times d_m}$ contains the first k low-frequency DCT coefficients.

Temporal semantic anchor. The **semantic anchor** $\mathbf{z}_{\text{tem}} = \tilde{\mathbf{h}}_t \in \mathbb{R}^{D_a}$ is aligned with the output of a frozen MoCLIP motion encoder $f_{\text{tem}}(\mathbf{x}_0) \in \mathbb{R}^{D_a}$ by maximizing cosine similarity:

$$\mathcal{L}_{\text{tem}} = 1 - \cos(\mathbf{z}_{\text{tem}}, f_{\text{tem}}(\mathbf{x}_0)). \quad (9)$$

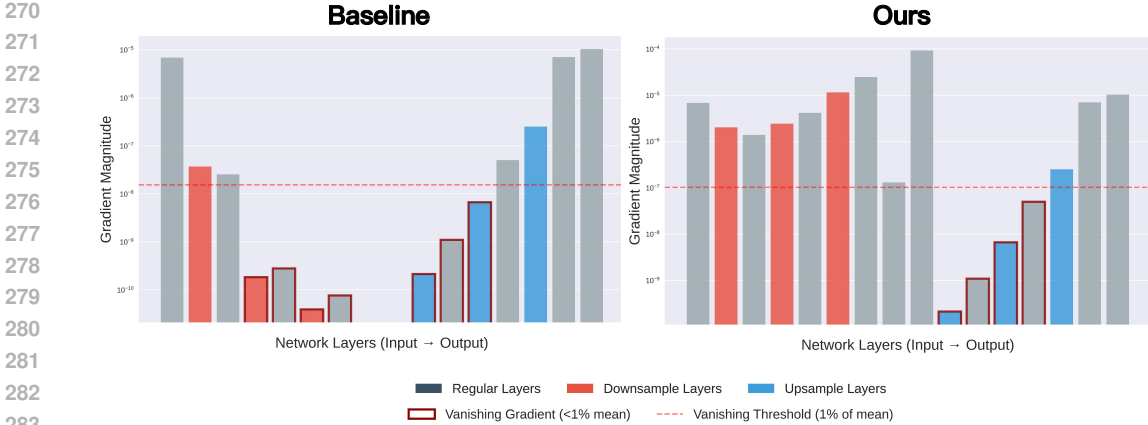


Figure 4: **Gradient magnitudes across network layers in LUMA.** Left: baseline without DAL; Right: with DAL. Red and blue bars denote downsampling and upsampling layers. The dashed line marks the vanishing gradient threshold (1% of mean). DAL boosts gradients in deep layers, mitigating vanishing issues and enhancing learning.

3.2.2 DYNAMIC ANCHOR WEIGHTING.

Inspired by prior work Wang et al. (2025) advocating attenuated auxiliary supervision in later training stages, we apply a cosine annealing schedule to progressively reduce the influence of semantic alignment losses. The overall training objective integrates the DDPM reconstruction loss with the **Dual Anchor Loss (DAL)**, which combines frequency and temporal alignment losses:

$$\mathcal{L} = \mathcal{L}_{\text{DDPM}} + \zeta(n) \cdot (\lambda_{\text{fre}} \mathcal{L}_{\text{fre}} + \lambda_{\text{tem}} \mathcal{L}_{\text{tem}}), \quad (10)$$

where $\mathcal{L}_{\text{DDPM}}$ denotes the standard denoising loss, λ_{fre} and λ_{tem} are weighting coefficients, and n is the current training step. The annealing factor $\zeta(n)$ is defined as:

$$\zeta(n) = \frac{1}{2} \left[1 + \cos \left(\pi \cdot \min \left(\frac{n}{N}, 1 \right) \right) \right], \quad (11)$$

with π denoting the mathematical constant and N the decay threshold. This schedule gradually suppresses DAL by step N , preventing over-regularization and allowing the model to refine motion generation autonomously.

4 EXPERIMENTS

We evaluate our approach on two standard motion-language benchmarks: HumanML3D Guo et al. (2022a) and KIT-ML Plappert et al. (2016). HumanML3D contains 14,616 motion sequences from AMASS Mahmood et al. (2019) and HumanAct12 Guo et al. (2020), each paired with three text descriptions (44,970 total), covering diverse actions like walking, exercising, and dancing. KIT-ML includes 3,911 motions and 6,278 descriptions, serving as a smaller-scale benchmark. We follow the standard evaluation protocol of StableMoFusion, including pose representation and mirroring-based augmentation. The data is split into training, validation, and test sets with a ratio of 0.8:0.15:0.05. More details about the datasets and evaluation metrics are available in **Appendix A**.

Experimental Setup. We adopt model architecture settings consistent with StableMoFusion Huang et al. (2024), closely following its training methodology with a total of 50,000 training steps. The diffusion process employs $T = 1000$ time steps with a linear beta schedule, and DPM-Solver is used for inference with 10 sampling steps. For DAL, λ_{fre} is set to 0.1, and λ_{tem} is actively used with a value of 0.5. The parameter k is set to 64. In our framework, anchor selection is performed at the Down-block3 stage, where FiLM modulation is also injected. During inference, the classifier-free guidance scale is set to 2.5. A conditional mask probability of 0.1 is applied during training to enable classifier-free guidance. The entire training procedure can be efficiently completed on a single RTX 4090 GPU with 24 GB of memory.

| Method | FID ↓ | R-Precision ↑ | | | Diversity → | MM. Dist. ↓ | MM. ↑ |
|-------------------------------|--------------------|--------------------|--------------------|---------------------|---------------------|--------------------|--------------------|
| | | Top1 | Top2 | Top3 | | | |
| HumanML3D | | | | | | | |
| Real | 0.002±.000 | 0.511±.003 | 0.703±.003 | 0.797±.002 | 9.503±.065 | - | - |
| MLD (Chen et al.) | 0.473±.013 | 0.481±.003 | 0.673±.003 | 0.772±.002 | 9.724±.082 | 3.196±.010 | <u>2.413</u> ±.079 |
| MDM (Tevet et al.) | 0.544±.044 | 0.320±.005 | 0.498±.004 | 0.611±.007 | <u>9.559</u> ±.086 | 5.556±.027 | 2.799 ±.072 |
| ReMoDiffuse (Zhang et al.) | 0.103±.004 | 0.510±.005 | 0.698±.006 | 0.795±.004 | 9.018±.075 | 3.025±.008 | 1.795±.043 |
| T2M-GPT (Zhang et al.) | 0.141±.005 | 0.492±.003 | 0.679±.002 | 0.775±.002 | 9.722±.082 | 3.121±.009 | 1.831±.048 |
| MotionGPT (Jiang et al.) | 0.232±.008 | 0.492±.003 | 0.681±.003 | 0.778±.002 | 9.528±.071 | 3.096±.008 | 2.008±.084 |
| MotionDiffuse (Zhang et al.) | 0.630±.001 | 0.491±.001 | 0.681±.001 | 0.782±.001 | 9.410±.049 | 3.113±.001 | 1.553±.042 |
| MotionLCM (Dai et al.) | 0.467±.012 | 0.502±.003 | 0.701±.002 | 0.803±.002 | 9.361±.660 | 3.012±.007 | 2.172±.082 |
| MMM (Pinyoanuntapong et al.) | 0.089±.002 | 0.515±.002 | 0.708±.002 | 0.804±.002 | 9.577±.050 | 2.926±.007 | 1.226±.035 |
| MoMask (Guo et al.) | <u>0.045</u> ±.002 | 0.521±.002 | 0.713±.002 | 0.807±.002 | - | 2.958±.008 | 1.241±.040 |
| StableMoFusion (Huang et al.) | 0.152±.004 | <u>0.546</u> ±.002 | <u>0.742</u> ±.002 | <u>0.835</u> ±.002 | 9.466 ±.002 | <u>2.781</u> ±.011 | 1.362±.062 |
| LUMA (Ours) | 0.035 ±.002 | 0.556 ±.003 | 0.750 ±.002 | 0.839 ±.0013 | 9.596±.083 | 2.763 ±.008 | 1.696±.049 |
| KIT-ML | | | | | | | |
| Real | 0.031±.004 | 0.424±.005 | 0.649±.006 | 0.779±.006 | 11.080±.097 | - | - |
| MLD (Chen et al.) | 0.404±.027 | 0.390±.008 | 0.609±.008 | 0.734±.007 | 10.800±.117 | 3.204±.027 | <u>2.192</u> ±.071 |
| MDM (Tevet et al.) | 0.497±.021 | 0.164±.004 | 0.291±.004 | 0.396±.004 | 10.847±.119 | 9.191±.022 | 1.907±.214 |
| ReMoDiffuse (Zhang et al.) | 0.155±.006 | 0.427±.014 | 0.641±.004 | 0.765±.055 | 10.800±.105 | 2.814±.028 | 1.239±.028 |
| T2M-GPT (Zhang et al.) | 0.514±.029 | 0.416±.006 | 0.627±.006 | 0.745±.006 | 10.921±.108 | 3.007±.023 | 1.570±.039 |
| MotionGPT (Jiang et al.) | 0.510±.016 | 0.366±.005 | 0.558±.004 | 0.680±.005 | 10.350±.084 | 3.527±.021 | 2.328 ±.117 |
| MotionDiffuse (Zhang et al.) | 1.954±.062 | 0.417±.004 | 0.621±.004 | 0.739±.004 | 11.100 ±.143 | 2.958±.005 | 0.730±.013 |
| MMM (Pinyoanuntapong et al.) | 0.316±.028 | 0.404±.005 | 0.621±.005 | 0.744±.004 | 10.910±.101 | 2.977±.019 | 1.232±.039 |
| MoMask (Guo et al.) | <u>0.204</u> ±.011 | 0.433±.007 | 0.656±.005 | 0.781±.005 | - | 2.779±.022 | 1.131±.043 |
| StableMoFusion (Huang et al.) | 0.258±.029 | <u>0.446</u> ±.006 | <u>0.660</u> ±.005 | <u>0.782</u> ±.004 | <u>10.936</u> ±.077 | <u>2.800</u> ±.018 | 1.362±.062 |
| LUMA (Ours) | 0.123 ±.014 | 0.454 ±.006 | 0.675 ±.005 | 0.796 ±.005 | 10.803±.076 | 2.711 ±.015 | 1.233±.055 |

Table 1: Experimental results on HumanML3D and KIT-ML datasets. ± indicates a 95% confidence interval. **Bold** and Underline indicate the best and the second-best results, respectively. The right arrow (→) denotes that a higher value is closer to real motion. Our method achieves state-of-the-art FID and R-Precision across both benchmarks.

4.1 COMPARISON TO STATE-OF-THE-ART APPROACHES

Quantitative comparisons. Following Zhang et al. (2023a); Guo et al. (2023), we report the average over 20 repeated generations with a 95% confidence interval. Table 1 presents evaluations on the HumanML3D Guo et al. (2022a) and KIT-ML Plappert et al. (2016) datasets, respectively, in comparison with state-of-the-art (SOTA) approaches.

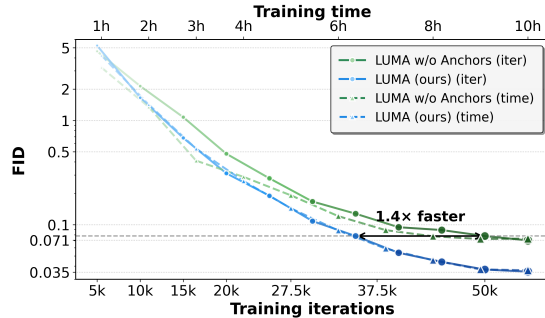


Figure 5: FID convergence curves for LUMA and the anchor-free baseline (vs. training steps and wall-clock time). The dashed line marks an FID threshold of 0.078. LUMA reaches this in **35k** iterations (**1.4×** faster), reducing training time from 8 to **7 hours**, demonstrating more efficient convergence.

among all evaluated methods and ranks first in R-Precision, demonstrating the broad adaptability of the LUMA architecture. Overall, these results underscore the capability of LUMA to significantly improve both the quality and diversity of text-driven motion generation.

In terms of empirical results, LUMA demonstrates strong performance across all key evaluation metrics. Our method achieves substantial improvements in FID (HumanML3D: **0.035**; KIT: **0.123**) and R-Precision (Top-3 HumanML3D: **0.839**; KIT: **0.796**), clearly illustrating the effectiveness of LUMA in enhancing the baseline model. Compared to state-of-the-art VQ-based approaches such as MMM Pinyoanuntapong et al. (2024c) and MoMask Guo et al. (2023), LUMA consistently achieves competitive or superior results on the HumanML3D dataset, outperforming other methods in terms of FID, R-Precision, and MultiModal Distance. This highlights LUMA’s ability to generate high-quality, semantically aligned, and diverse motion sequences. On the more challenging KIT-ML benchmark, LUMA attains the best FID

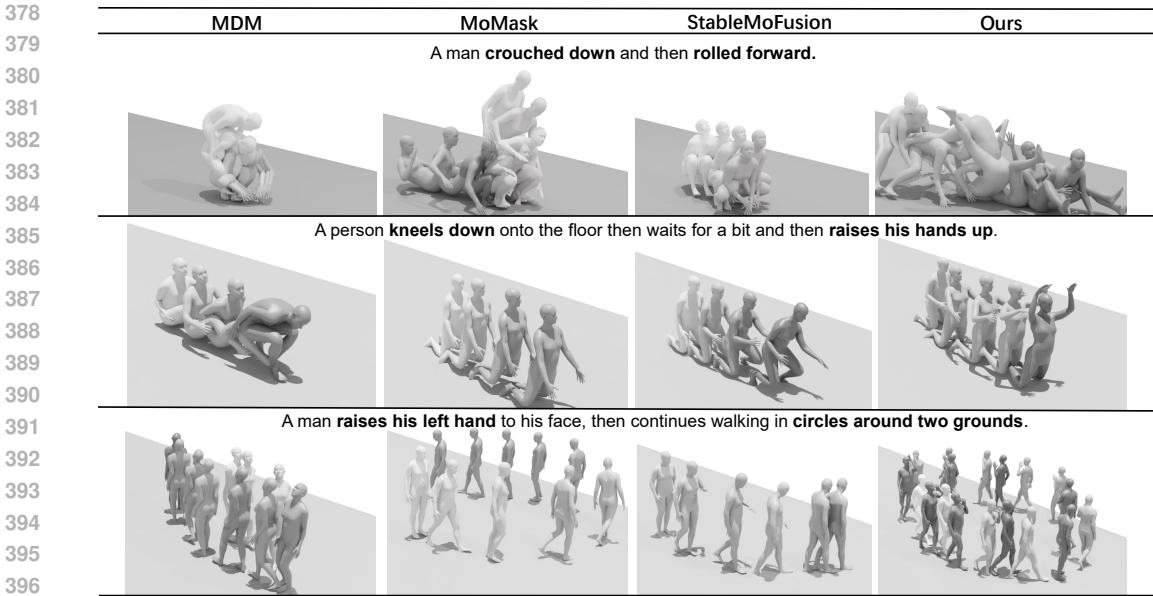


Figure 6: Visualization Comparison. We compare the visual results of LUMA with three state-of-the-art methods. In both examples, LUMA consistently produces more accurate, natural, and fine-grained motion than the other approaches.

Qualitative comparison. Figure 6 presents visual comparisons with MDM, StableMoFusion, and MoMask under three prompts. For "A man crouched down and then rolled forward," MDM and StableMoFusion only crouch, while MoMask fails to roll properly. Our model performs the full action accurately. For "A person kneels down onto the floor then waits and raises his hands up," MDM misses kneeling, and others raise hands poorly; our method captures both with correct timing. For the complex prompt involving hand-raising and circling, baselines fail to complete all sub-actions, while our model executes them with high fidelity. These results demonstrate superior semantic understanding, temporal consistency, and detail accuracy.

4.2 FASTER CONVERGENCE WITH DUAL ANCHORS

As shown in Figure 5, LUMA achieves significantly faster convergence compared to the anchor-free baseline. Our model surpasses the FID threshold of 0.078 with only **35k** training iterations, whereas the baseline requires **50k** steps, resulting in a **1.4x** speedup. This efficiency gain is also reflected in the wall-clock training time. As illustrated on the upper axis of the same figure, LUMA reaches the target FID in just **7 hours**, while the baseline takes **8 hours**. This demonstrates that our dual-anchor alignment not only improves motion quality but also makes training substantially more efficient, both in computation steps and real-world time.

4.3 GRADIENT ANALYSIS

As shown in Figure 4, removing the DAL causes the gradient magnitudes in the middle layers to drop sharply. The downsampling path is especially affected. Some layers fall below the vanishing threshold (1% of the mean). This indicates severe gradient attenuation in these regions. After introducing the DAL, the gradients in the deep downsampling layers increase by nearly two orders of magnitude and are brought to a level comparable with the upsampling path. This observation validates our earlier hypothesis that vanishing gradients in deep layers hinder effective learning of abstract features. The additional semantic and motion anchors introduced by DAL provide more informative supervisory signals to the encoder, substantially alleviating the vanishing gradient problem and accelerating effective representation learning in deeper layers.

432
433
434
435
436
437
438
439
440
441
442
443

| Design Choice | FID ↓ | R@3 ↑ |
|--------------------------------|----------------------------------|----------------------------------|
| DCT $k = 8$ | 0.044 \pm .003 | 0.835 \pm .003 |
| DCT $k = 16$ | 0.052 \pm .002 | 0.829 \pm .003 |
| DCT $k = 32$ | 0.052 \pm .004 | 0.832 \pm .002 |
| DCT $k = 64$ | 0.035\pm.002 | 0.839\pm.001 |
| DCT $k = 128$ | 0.063 \pm .004 | 0.831 \pm .003 |
| Down-block1 | 0.103 \pm .003 | 0.813 \pm .002 |
| Down-block2 | 0.076 \pm .001 | 0.827 \pm .001 |
| Down-block3 | 0.035\pm.002 | 0.839\pm.001 |
| Bottleneck | 0.092 \pm .002 | 0.838 \pm .003 |

Table 3: Ablation study on key design choices in LUMA. For each factor, the best configuration is highlighted in bold.

444
445
446
447
448
449
450
451
452

| Anchor Weighting Strategy | FID ↓ | R@3 ↑ |
|---------------------------|----------------------------------|----------------------------------|
| Static (fixed) | 0.047 \pm .001 | 0.837 \pm .002 |
| Learnable (global) | 0.051 \pm .002 | 0.834 \pm .002 |
| Dynamic (cosine) | 0.035\pm.002 | 0.839\pm.001 |

Table 4: Ablation on anchor weighting strategies. The dynamic schedule achieves the best performance.

453
454
455
456
457

4.4 ABLATION STUDY

458
459
460
461
462

To comprehensively assess the contribution of each component in our framework, we conduct a series of controlled ablation experiments. Following the evaluation protocol, we focus on two aspects: (i) the effect of each core module, (ii) the influence of key design choices. All experiments are performed under identical training settings for fair comparison.

463
464
465
466
467
468
469
470
471
472
473
474

Core Component Analysis. As shown in Table 2, each component of LUMA plays a critical and complementary role. The backbone alone achieves a baseline FID of 0.152, while removing both anchors but retaining the LUMA architecture already yields a substantial FID drop to 0.067, highlighting structural benefits. Introducing either the frequency semantic or temporal semantic anchor alone further reduces FID to 0.057 and improves R@3 to 0.840, indicating that both anchors independently enhance motion quality and text alignment. Their similar contributions suggest that each focuses on a different aspect of the generation process. One captures frequency-level structure, while the other improves semantic consistency. Removing MoCLIP leads to a slight drop in performance, confirming its importance in providing effective semantic supervision. Meanwhile, eliminating FiLM leads to a marginal drop in R@3 despite strong FID, showing its importance in temporal modulation. The full model outperforms all variants (FID 0.035, R@3 0.839), demonstrating that these modules act synergistically to improve fidelity, semantic alignment, and training stability.

475
476
477
478
479
480
481
482
483
484
485

Design Choice Exploration. We study two key hyperparameters in LUMA: the number of retained DCT coefficients (k) and the FiLM injection location. As shown in Table 3, an intermediate setting ($k = 64$) best balances global structure and noise suppression; smaller ($k = 8$) limits expressiveness, while larger ($k = 128$) introduces artifacts. For FiLM, injecting at *Down-block3* provides the best trade-off, capturing sufficient temporal context without losing fine details or adding notable compute.

| Method | FID ↓ | R@3 ↑ |
|-------------------------------------|----------------------------------|----------------------------------|
| StableMoFusion (backbone) | 0.152 \pm .004 | 0.835 \pm .002 |
| LUMA w/o DAL | 0.067 \pm .003 | 0.835 \pm .002 |
| LUMA w/o \mathcal{L}_{tem} | 0.057 \pm .002 | 0.840 \pm .002 |
| LUMA w/o \mathcal{L}_{fre} | 0.057 \pm .002 | 0.840 \pm .002 |
| LUMA w/o MoCLIP | 0.075 \pm .003 | 0.838 \pm .002 |
| LUMA w/o FiLM | 0.056 \pm .003 | 0.833 \pm .002 |
| LUMA (full) | 0.035\pm.002 | 0.839\pm.001 |

Table 2: Ablation results for core components of LUMA on the main benchmark. We report FID (↓) and R@3 (↑) with mean \pm 95% CI.

486 **Dynamic Anchor Weighting.** We compare
487 our dynamic cosine-annealed schedule with
488 (i) *static averaging* and (ii) *fixed learnable weights*. Table 4 shows the dynamic scheme consistently
489 achieves the best FID and R@3. It emphasizes temporal and frequency alignment early, then grad-
490 ually relaxes constraints to avoid over-regularization, enabling finer motion refinement and stronger
491 semantic consistency.

492

493 5 CONCLUSION

494

495 In conclusion, we identify severe gradient sparsity in deep layers of diffusion-based motion models,
496 which limits semantic and structural alignment. To address this, we propose **LUMA**, a dual-path
497 framework injecting temporal and frequency anchors into the diffusion backbone. This design en-
498 hances gradient flow, speeds up convergence, and improves motion fidelity. Experiments on Hu-
499 manML3D and KIT-ML show state-of-the-art FID and R-Precision, with strong diversity and mul-
500 timodality. Ablations and gradient analyses further validate the effectiveness of our approach and
501 support our hypothesis.

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593

ETHICS STATEMENT

This work focuses on advancing text-to-motion generation models. The proposed methods and experiments do not involve sensitive data, personally identifiable information, or human subjects. The potential societal impact lies primarily in creative applications such as animation, gaming, VR/AR, and robotics. We acknowledge the importance of responsible use and encourage future work to mitigate potential misuse in areas such as generating misleading or harmful content.

REPRODUCIBILITY STATEMENT

We have taken extensive steps to ensure reproducibility. Details of the model architecture, hyperparameters, datasets, and training procedure are provided in the main text and the Appendix. We use the publicly available HumanML3D and KIT-ML datasets. To further facilitate replication, we include the full source code, training scripts, and pre-trained models in the Supplementary Material, enabling direct reproduction of our results. Our experiments were performed with clearly specified hardware and software configurations, ensuring that other researchers can reliably replicate and extend our work.

REFERENCES

- 594
595
596 Laith Alzubaidi, Jinglan Zhang, Amjad J. Humaidi, Ayad Al-Dujaili, Ye Duan, Omran Al-Shamma,
597 J. Santamaría, Mohammed A. Fadhel, Muthana Al-Amidie, and Laith Farhan. Review of deep
598 learning: concepts, cnn architectures, challenges, applications, future directions. *Journal of Big*
599 *Data*, 8:53, 2021. doi: 10.1186/s40537-021-00444-8. URL [https://journalofbigdata.
600 springeropen.com/articles/10.1186/s40537-021-00444-8](https://journalofbigdata.springeropen.com/articles/10.1186/s40537-021-00444-8).
- 601 Nikos Athanasiou, Mathis Petrovich, Michael J. Black, and Gül Varol. Teach: Temporal action
602 composition for 3d humans, 2022. URL <https://arxiv.org/abs/2209.04066>.
- 603
604 Wenshuo Chen, Hongru Xiao, Erhang Zhang, Lijie Hu, Lei Wang, Mengyuan Liu, and Chen Chen.
605 Sato: Stable text-to-motion framework. In *Proceedings of the 32nd ACM International Con-*
606 *ference on Multimedia, MM '24*, pp. 6989–6997. ACM, October 2024. doi: 10.1145/3664647.
607 3681034. URL <http://dx.doi.org/10.1145/3664647.3681034>.
- 608 Wenshuo Chen, Haozhe Jia, Songning Lai, Keming Wu, Hongru Xiao, Lijie Hu, and Yutao Yue.
609 Free-t2m: Frequency enhanced text-to-motion diffusion model with consistency loss, 2025. URL
610 <https://arxiv.org/abs/2501.18232>.
- 611
612 Xin Chen, Biao Jiang, Wen Liu, Zilong Huang, Bin Fu, Tao Chen, Jingyi Yu, and Gang Yu. Execut-
613 ing your commands via motion diffusion in latent space, 2023. URL [https://arxiv.org/
614 abs/2212.04048](https://arxiv.org/abs/2212.04048).
- 615 Wenxun Dai, Ling-Hao Chen, Jingbo Wang, Jinpeng Liu, Bo Dai, and Yansong Tang. Motionlcm:
616 Real-time controllable motion generation via latent consistency model, 2024. URL [https:
617 //arxiv.org/abs/2404.19759](https://arxiv.org/abs/2404.19759).
- 618
619 Chuan Guo, Xinxin Zuo, Sen Wang, Shihao Zou, Qingyao Sun, Annan Deng, Minglun Gong, and
620 Li Cheng. Action2motion: Conditioned generation of 3d human motions. In *Proceedings of the*
621 *28th ACM International Conference on Multimedia*, pp. 2021–2029, 2020.
- 622 Chuan Guo, Shihao Zou, Xinxin Zuo, Sen Wang, Wei Ji, Xingyu Li, and Li Cheng. Generating
623 diverse and natural 3d human motions from text. In *Proceedings of the IEEE/CVF Conference on*
624 *Computer Vision and Pattern Recognition (CVPR)*, pp. 5152–5161, June 2022a.
- 625
626 Chuan Guo, Xinxin Zuo, Sen Wang, and Li Cheng. Tm2t: Stochastic and tokenized modeling for
627 the reciprocal generation of 3d human motions and texts, 2022b. URL [https://arxiv.org/
628 abs/2207.01696](https://arxiv.org/abs/2207.01696).
- 629 Chuan Guo, Yuxuan Mu, Muhammad Gohar Javed, Sen Wang, and Li Cheng. Momask: Generative
630 masked modeling of 3d human motions. 2023.
- 631
632 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
633 nition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*
634 *(CVPR)*, pp. 770–778, June 2016. doi: 10.1109/CVPR.2016.90.
- 635 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter.
636 Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Proceedings*
637 *of the 31st International Conference on Neural Information Processing Systems, NIPS'17*, pp.
638 6629–6640, Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964.
- 639
640 Yiheng Huang, Hui Yang, Chuanchen Luo, Yuxi Wang, Shibiao Xu, Zhaoxiang Zhang, Man Zhang,
641 and Junran Peng. Stablemofusion: Towards robust and efficient diffusion-based motion genera-
642 tion framework, 2024. URL <https://arxiv.org/abs/2405.05691>.
- 643 Biao Jiang, Xin Chen, Wen Liu, Jingyi Yu, Gang Yu, and Tao Chen. Motiongpt: Human motion as
644 a foreign language, 2023. URL <https://arxiv.org/abs/2306.14795>.
- 645
646 Xingjian Leng, Jaskirat Singh, Yunzhong Hou, Zhenchang Xing, Saining Xie, and Liang Zheng.
647 Repa-e: Unlocking vae for end-to-end tuning with latent diffusion transformers, 2025. URL
<https://arxiv.org/abs/2504.10483>.

- 648 Zhe Li, Weihao Yuan, Yisheng He, Lingteng Qiu, Shenhao Zhu, Xiaodong Gu, Weichao Shen, Yuan
649 Dong, Zilong Dong, and Laurence T. Yang. Lamp: Language-motion pretraining for motion gen-
650 eration, retrieval, and captioning, 2025. URL <https://arxiv.org/abs/2410.07093>.
- 651
- 652 Yihao Liao, Yiyu Fu, Ziming Cheng, and Jiangfeiyang Wang. Animationgpt:an aigc tool for gener-
653 ating game combat motion assets. <https://github.com/fyyakaxy/AnimationGPT>,
654 2024.
- 655 Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black.
656 SMPL: A skinned multi-person linear model. *ACM Trans. Graphics (Proc. SIGGRAPH Asia)*,
657 34(6):248:1–248:16, October 2015.
- 658
- 659 Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast
660 solver for guided sampling of diffusion probabilistic models, 2023. URL <https://arxiv.org/abs/2211.01095>.
- 661
- 662 Naureen Mahmood, Nima Ghorbani, Nikolaus F. Troje, Gerard Pons-Moll, and Michael J. Black.
663 AMASS: Archive of motion capture as surface shapes. In *International Conference on Computer*
664 *Vision*, pp. 5442–5451, October 2019.
- 665
- 666 Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov,
667 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud Assran, Nico-
668 las Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael
669 Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jegou, Julien Mairal, Patrick Labatut, Ar-
670 mand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision,
671 2024. URL <https://arxiv.org/abs/2304.07193>.
- 672
- 673 Ethan Perez, Florian Strub, Harm de Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual
674 reasoning with a general conditioning layer, 2017. URL <https://arxiv.org/abs/1709.07871>.
- 675
- 676 Ekkasit Pinyoanuntapong, Muhammad Usama Saleem, Korrawe Karunratanakul, Pu Wang, Hongfei
677 Xue, Chen Chen, Chuan Guo, Junli Cao, Jian Ren, and Sergey Tulyakov. Controlmm: Control-
678 lable masked motion generation, 2024a. URL <https://arxiv.org/abs/2410.10780>.
- 679
- 680 Ekkasit Pinyoanuntapong, Muhammad Usama Saleem, Pu Wang, Minwoo Lee, Srijan Das, and
681 Chen Chen. Bamm: Bidirectional autoregressive motion model, 2024b. URL <https://arxiv.org/abs/2403.19435>.
- 682
- 683 Ekkasit Pinyoanuntapong, Pu Wang, Minwoo Lee, and Chen Chen. Mmm: Generative masked
684 motion model, 2024c. URL <https://arxiv.org/abs/2312.03596>.
- 685
- 686 Matthias Plappert, Christian Mandery, and Tamim Asfour. The kit motion-language dataset. *Big*
687 *Data*, 4(4):236–252, December 2016. ISSN 2167-647X. doi: 10.1089/big.2016.0028. URL
688 <http://dx.doi.org/10.1089/big.2016.0028>.
- 689
- 690 Rupesh Kumar Srivastava, Klaus Greff, and Jürgen Schmidhuber. Highway networks. *CoRR*,
691 abs/1505.00387, 2015. URL <http://arxiv.org/abs/1505.00387>.
- 692
- 693 Guy Tevet, Sigal Raab, Brian Gordon, Yonatan Shafir, Daniel Cohen-Or, and Amit H Bermano.
694 Human motion diffusion model. *arXiv preprint arXiv:2209.14916*, 2022.
- 695
- 696 Yuchuan Tian, Hanting Chen, Mengyu Zheng, Yuchen Liang, Chao Xu, and Yunhe Wang. U-repa:
697 Aligning diffusion u-nets to vits, 2025. URL <https://arxiv.org/abs/2503.18414>.
- 698
- 699 Runqian Wang and Kaiming He. Diffuse and disperse: Image generation with representation regu-
700 larization, 2025. URL <https://arxiv.org/abs/2506.09027>.
- 701
- 702 Ziqiao Wang, Wangbo Zhao, Yuhao Zhou, Zekai Li, Zhiyuan Liang, Mingjia Shi, Xuanlei Zhao,
703 Pengfei Zhou, Kaipeng Zhang, Zhangyang Wang, Kai Wang, and Yang You. Repa works until it
704 doesn’t: Early-stopped, holistic alignment supercharges diffusion training, 2025. URL <https://arxiv.org/abs/2505.16792>.

702 Sihyun Yu, Sangkyung Kwak, Huiwon Jang, Jongheon Jeong, Jonathan Huang, Jinwoo Shin, and
703 Saining Xie. Representation alignment for generation: Training diffusion transformers is easier
704 than you think, 2025. URL <https://arxiv.org/abs/2410.06940>.
705

706 Ye Yuan, Jiaming Song, Umar Iqbal, Arash Vahdat, and Jan Kautz. Physdiff: Physics-guided human
707 motion diffusion model. In *Proceedings of the IEEE/CVF international conference on computer
708 vision*, pp. 16010–16021, 2023a.

709 Ye Yuan, Jiaming Song, Umar Iqbal, Arash Vahdat, and Jan Kautz. Physdiff: Physics-guided human
710 motion diffusion model, 2023b. URL <https://arxiv.org/abs/2212.02500>.
711

712 Jianrong Zhang, Yangsong Zhang, Xiaodong Cun, Shaoli Huang, Yong Zhang, Hongwei Zhao,
713 Hongtao Lu, and Xi Shen. T2m-gpt: Generating human motion from textual descriptions with
714 discrete representations, 2023a. URL <https://arxiv.org/abs/2301.06052>.

715 Mingyuan Zhang, Zhongang Cai, Liang Pan, Fangzhou Hong, Xinying Guo, Lei Yang, and Ziwei
716 Liu. Motiondiffuse: Text-driven human motion generation with diffusion model, 2022. URL
717 <https://arxiv.org/abs/2208.15001>.

718 Mingyuan Zhang, Xinying Guo, Liang Pan, Zhongang Cai, Fangzhou Hong, Huirong Li, Lei Yang,
719 and Ziwei Liu. Remodiffuse: Retrieval-augmented motion diffusion model. In *Proceedings of
720 the IEEE/CVF International Conference on Computer Vision*, pp. 364–373, 2023b.
721

722 Mingyuan Zhang, Zhongang Cai, Liang Pan, Fangzhou Hong, Xinying Guo, Lei Yang, and Ziwei
723 Liu. Motiondiffuse: Text-driven human motion generation with diffusion model. *IEEE transac-
724 tions on pattern analysis and machine intelligence*, 46(6):4115–4128, 2024a.

725 Zeyu Zhang, Akide Liu, Qi Chen, Feng Chen, Ian Reid, Richard Hartley, Bohan Zhuang, and Hao
726 Tang. Infinimotion: Mamba boosts memory in transformer for arbitrary long motion generation,
727 2024b. URL <https://arxiv.org/abs/2407.10061>.
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

APPENDIX

A MOTION DATA REPRESENTATION AND EVALUATION METRICS

To comprehensively evaluate text-to-motion generation, we utilize a suite of established quantitative metrics and adopt precise motion data representations. This enables rigorous analysis of both generation performance and data fidelity. We first introduce the evaluation metrics, followed by a detailed discussion of motion representations.

A.1 EVALUATION METRICS

We employ widely recognized quantitative metrics including Frechet Inception Distance (FID) Heusel et al. (2017), R-Precision, Diversity, Multimodal Distance (MM-Dist.), and Multimodality (MM.), consistent with the evaluation protocol of StableMoFusion Huang et al. (2024). Additionally, human evaluation is conducted to assess semantic accuracy and human preference, providing a supplementary subjective perspective.

- **Frechet Inception Distance (FID):** Measures the distributional difference between generated and real motion features, reflecting the overall motion quality.
- **R-Precision:** Evaluates motion-to-text retrieval by ranking Euclidean distances between each generated motion and 32 candidate text descriptions (1 ground-truth, 31 mismatched). Top-1, Top-2, and Top-3 retrieval accuracies are reported.
- **Diversity:** Quantifies motion diversity by averaging Euclidean distances between 300 randomly sampled pairs of generated motions.
- **Multimodal Distance (MM-Dist.):** Assesses the semantic alignment between generated motions and corresponding texts. Lower MM-Dist. indicates better cross-modal correspondence.
- **Multimodality (MM.):** Measures the diversity of motions generated from a single text prompt. For each prompt, 20 motions are generated to form 10 pairs; the mean Euclidean distance between pairs is computed and averaged across all prompts.

A.2 MOTION DATA REPRESENTATIONS

We analyze two predominant motion representation formats: the *HumanML3D Format* and the *SMPL-based Format* Loper et al. (2015). Both are widely adopted in prior works and offer complementary advantages in semantic expressiveness and biomechanical fidelity.

HumanML3D Format. The HumanML3D format encodes motion as a tuple of spatial and dynamic features:

$$x^i = \{r^a, r^x, r^z, r^y, j^P, j^v, j^r, \mathbf{f}_{\text{contact}}\}, \quad (12)$$

where r^a is the root angular velocity (Y-axis), r^x and r^z are root linear velocities (XZ-plane), r^y is root height, $j^P \in \mathbb{R}^{3N_j}$ denotes local joint positions, $j^v \in \mathbb{R}^{3N_j}$ and $j^r \in \mathbb{R}^{6N_j}$ represent joint velocities and joint rotations, and $\mathbf{f}_{\text{contact}} \in \mathbb{R}^4$ are binary foot-ground contact features. This comprehensive definition captures both high-level semantic structure and fine-grained motion detail.

SMPL-based Format. The SMPL model focuses on anatomical accuracy. Each motion is represented as:

$$x^i = \{r, \boldsymbol{\vartheta}, \beta\}, \quad (13)$$

where r denotes global translation, $\boldsymbol{\vartheta} \in \mathbb{R}^{3 \times 23 + 3}$ represents joint rotations for 23 joints and a root joint, and β encodes body shape parameters. The SMPL-based representation is well-suited for modeling biomechanically realistic motions.

Frequency-Domain Representation. In addition to these established formats, we further introduce a novel frequency-domain representation to enhance semantic planning and detailed motion refinement, bridging global temporal structure with local motion accuracy.

B DETAILS OF MOCLIP

Inspired by CLIP’s image-text alignment success MoCLIP extends this concept to human motion. It learns a shared embedding space mapping textual descriptions to corresponding motion sequences enabling cross-modal retrieval and understanding.

B.0.1 ARCHITECTURE

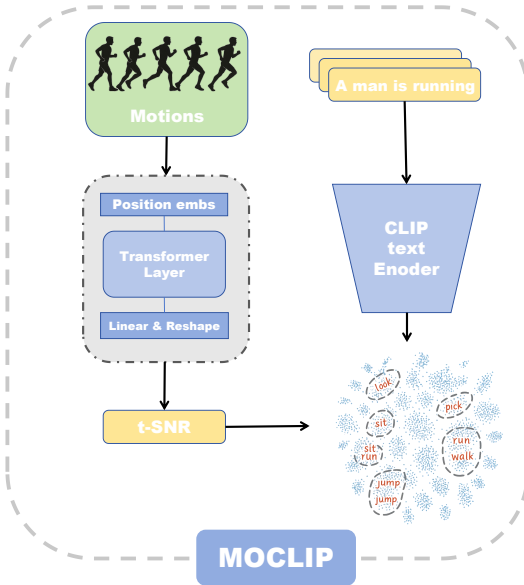


Figure 7: **Overview of MoCLIP.** MoCLIP aligns motion and text through contrastive learning, using a transformer-based motion encoder and a CLIP text encoder to build a shared semantic space.

As shown in Figure 7, MoCLIP employs a dual-encoder structure adapting the pretrained CLIP architecture. The text pathway uses a fine-tuned CLIP text encoder generating semantic embeddings $\mathbf{z}_{\text{text}} \in \mathbb{R}^{d_{\text{MoCLIP}}}$. The motion pathway features a dedicated MotionEncoder processing sequences $\mathbf{M} \in \mathbb{R}^{N_{\text{seq}} \times d_{\text{motion}}}$. This encoder applies linear projection \mathbf{W}_p adds sinusoidal positional encoding \mathbf{P} passes the result through L transformer layers handling variable lengths via mask \mathbf{M}_{mask} performs temporal average pooling and finally projects features via \mathbf{W}_o to obtain motion embeddings $\mathbf{z}_{\text{motion}} \in \mathbb{R}^{d_{\text{MoCLIP}}}$. The core operation within the transformer layers is multi-head self-attention:

$$\mathbf{s}_l = \text{TransformerLayer}(\mathbf{s}_{l-1}, \mathbf{M}_{\text{mask}}). \quad (14)$$

Contrastive learning in the shared d_{MoCLIP} -dimensional space aligns these embeddings scaled by temperature τ .

B.0.2 LOSS FUNCTION

MoCLIP utilizes a symmetric contrastive loss to align modalities:

$$\mathcal{L}_{\text{contrastive}} = \frac{1}{2} (\mathcal{L}_{\text{motion-to-text}} + \mathcal{L}_{\text{text-to-motion}}). \quad (15)$$

This loss averages the motion-to-text and text-to-motion cross-entropy terms. These terms are computed using cosine similarity between L2-normalized motion and text embeddings promoting high similarity for matched pairs and low similarity for mismatched pairs.

B.0.3 TRAINING STRATEGY

Training follows a two-stage strategy for progressive alignment. **Stage 1: Motion Encoder Pre-training.** The CLIP text encoder is frozen. Only the MotionEncoder components are trained optimizing the contrastive loss (Eq. 15). This initially aligns motion features to the fixed text embedding space. **Stage 2: Joint Fine-Tuning.** The final layers of the CLIP text encoder are unfrozen. The entire model is then fine-tuned jointly with a lower learning rate. This allows mutual refinement of both motion and text representations enhancing the joint embedding space. This approach facilitates stable learning and effective cross-modal integration.

B.1 PERFORMANCE OF MOCLIP

We evaluated MoCLIP’s core text-motion alignment capability on the HumanML3D Guo et al. (2022a) KIT Plappert et al. (2016) and CMP Liao et al. (2024) (Combat Motion Processed) datasets. Performance was measured using standard Top- k retrieval accuracy (Top-1 Top-2 Top-3). Table 5 shows MoCLIP significantly outperforms the baseline across all datasets. Notably on HumanML3D MoCLIP achieves 0.705 Top-1 accuracy versus the baseline’s 0.511. On CMP the improvement is also substantial reaching 0.748 Top-1 accuracy compared to 0.335. Consistent gains are observed on the KIT dataset. These results validate MoCLIP’s effectiveness in learning accurate text-motion semantic mappings.

864
865
866
867
868
869
870
871
872

| | Dataset | Top-1 | Top-2 | Top-3 |
|----------|-----------|-------|-------|-------|
| Baseline | Humanml3d | 0.511 | 0.703 | 0.797 |
| MoCLIP | Humanml3d | 0.705 | 0.856 | 0.913 |
| Baseline | KIT | 0.424 | 0.649 | 0.779 |
| MoCLIP | KIT | 0.469 | 0.676 | 0.788 |
| Baseline | CMP | 0.335 | 0.513 | 0.628 |
| MoCLIP | CMP | 0.748 | 0.891 | 0.942 |

Table 5: Top- k retrieval accuracy comparison between the baseline and MoCLIP on HumanML3D, KIT, and CMP datasets.873
874
875
876
877
878
879
880
881
882
883
884
885
886

| Method | FID ↓ | R-Precision ↑ | | |
|-------------------------|------------------|------------------|-------------------|------------------|
| | | top1 | top2 | top3 |
| MDM baseline | 0.544 \pm .044 | 0.320 \pm .005 | 0.498 \pm .004 | 0.611 \pm .007 |
| MDM MoCLIP | 0.527 \pm .034 | 0.514 \pm .003 | 0.719 \pm .001 | 0.820 \pm .001 |
| Momask baseline | 0.045 \pm .002 | 0.521 \pm .002 | 0.713 \pm .002 | 0.807 \pm .002 |
| Momask MoCLIP | 0.065 \pm .002 | 0.529 \pm .002 | 0.724 \pm 0.002 | 0.818 \pm .002 |
| StableMoFusion baseline | 0.152 \pm .004 | 0.546 \pm .002 | 0.742 \pm .002 | 0.835 \pm .002 |
| StableMoFusion MoCLIP | 0.067 \pm .003 | 0.549 \pm .002 | 0.742 \pm .002 | 0.835 \pm .002 |

Table 6: Evaluation results of MoCLIP integration with different models. The table shows the FID (lower is better) and R-Precision (higher is better) at top1, top2, and top3 for MDM, MoMask, and StableMoFusion models with and without MoCLIP. The results demonstrate the positive impact of MoCLIP on improving both FID and R-Precision scores in motion generation tasks.

887
888
889
890
891
892
893
894

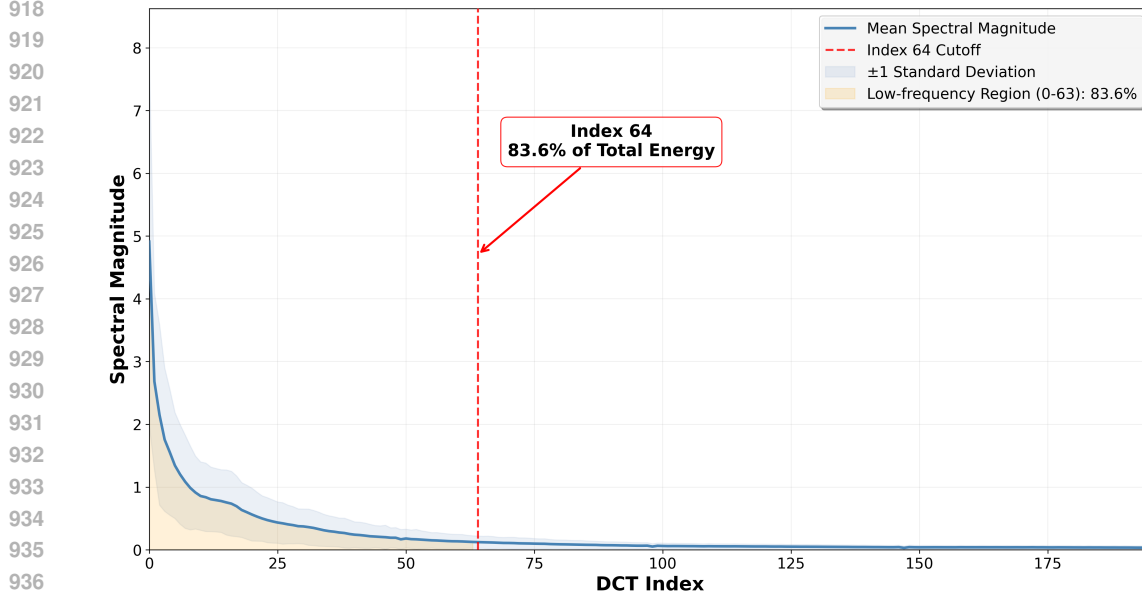
B.2 INTEGRATING MOCLIP FOR ENHANCED GENERATION

895
896
897
898
899
900
901
902
903
904

To evaluate the practical benefit of MoCLIP’s learned representations we integrated its fine-tuned text encoder $\mathcal{T}_{\text{MoCLIP}}$ into existing generation frameworks. Specifically we replaced the native text encoders of StableMoFusion MDM and MoMask with $\mathcal{T}_{\text{MoCLIP}}$. This modification supplies these generators with motion-aligned text embeddings $\mathbf{z}_{\text{text}}^{\text{MoCLIP}}$ leveraging the shared semantic space detailed in Section B. The resulting performance improvements detailed in Table 6 demonstrate the efficacy of this approach. Using $\mathcal{T}_{\text{MoCLIP}}$ consistently enhances generation quality across the tested models. This confirms that the MoCLIP encoder effectively extracts motion-relevant semantics transforming text descriptions into representations more conducive to high-fidelity motion synthesis.

905
906
907
908
909
910
911
912
913
914
915
916
917

To investigate the impact of our motion-text alignment training, we performed a comparative analysis of the attention mechanisms within the text encoders of MoCLIP and the original CLIP model. Our analysis focuses on the attention maps produced by the final Transformer layer. As shown in Figure 11, it is crucial to note that at this terminal stage of encoding, the initial, discrete token embeddings have been iteratively transformed into highly contextualized feature vectors, where each vector at a given position encapsulates rich, sequence-wide semantic information. Our methodology involved computing the total attention received by the feature vector at each sequence position from all other positions, a value which we then averaged across all attention heads and normalized. By differencing these final-layer attention distributions (MoCLIP vs. original CLIP), we uncovered a systematic and significant reallocation of attention. Specifically, MoCLIP learns to assign substantially greater attentional weight to the feature vectors located at positions that correspond to the initial action-oriented tokens (e.g., verbs like “crouched”, “rolled”, “waves” and adverbs of motion). This strategic focus indicates that our contrastive training paradigm has successfully guided the model to identify the semantic loci of action within the sequence and to prioritize these now-contextualized features when constructing the final, motion-aware text representation. This learned



938 Figure 8: Mean DCT spectrum of HumanML3D. The first 64 coefficients (shaded) account for
939 83.6% of the total energy.

940
941

942 ability to pinpoint and amplify motion-centric semantics is fundamental to MoCLIP’s superior per-
943 formance in translating textual descriptions into complex human motion.

944

945 C PSEUDO-CODE

946

947 Algorithm 1 shows LUMA’s Pseudo-code.

948

949 Algorithm 1: Concise LUMA Training Loop

950

951 **Input:** Dataset \mathcal{D} ; frozen MoCLIP encoders $f_{\text{text}}, f_{\text{tem}}$;
952 model params $\Theta = \{\phi, \theta, \psi, \psi'\}$; hyper-params $(k, \lambda_{\text{fre}}, \lambda_{\text{tem}}, N_{\text{decay}}, T)$

953

954 **Output:** Updated Θ

955

```

955 1 for  $n = 1$  to  $N_{\text{train}}$  do
956   // Diffuse a clean sample
957   2 Sample  $(\mathbf{x}_0, \mathbf{c}), t \sim \mathcal{U}[1, T], \epsilon \sim \mathcal{N}(0, I)$ ;
958   3  $\mathbf{x}_t \leftarrow \sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1 - \alpha_t} \epsilon$ ;
959   // Predict noise and extract bottleneck
960   4  $(\hat{\epsilon}, \mathbf{h}) \leftarrow \text{DENOISER}_{\phi}(\mathbf{x}_t, f_{\text{text}}(\mathbf{c}), t)$ ;
961   // Dual-anchor projection and FiLM modulation
962   5  $(\mathbf{z}_{\text{fre}}, \mathbf{z}_{\text{tem}}) \leftarrow \text{ANCHOR}_{\theta, \psi, \psi'}(\mathbf{h}, t)$ ;
963   // Compute total loss
964   6  $\zeta \leftarrow \frac{1}{2} \left[ 1 + \cos \left( \pi \min \left( \frac{n}{N_{\text{decay}}}, 1 \right) \right) \right]$ ;
965   7  $\mathcal{L}_{\text{total}} \leftarrow \underbrace{\|\hat{\epsilon} - \epsilon\|_2^2}_{\mathcal{L}_{\text{DDPM}}} + \zeta \left( \lambda_{\text{fre}} \|\mathbf{z}_{\text{fre}} - \text{DCT}_k(\mathbf{x}_0)\|_2^2 + \lambda_{\text{tem}} (1 - \cos(\mathbf{z}_{\text{tem}}, f_{\text{tem}}(\mathbf{x}_0))) \right)$ ;
966   8 // SGD / Adam optimization step
967    $\Theta \leftarrow \Theta - \eta \nabla_{\Theta} \mathcal{L}_{\text{total}}$ ;
968
969
970
971
```

971

D HYPERPARAMETER SEARCH RESULT

In our study, we conducted a grid search on the validation set to determine the optimal weights for the dual-anchor losses in LUMA, specifically λ_{fre} (frequency anchor) and λ_{tem} (temporal anchor). We selected five candidate values for each hyperparameter, resulting in 25 combinations, with ranges extended beyond prior work to more fully explore the parameter space. Each combination was evaluated five times to ensure the reliability of our results.

Our analysis revealed an inverse relationship between FID and top-k accuracy when either anchor weight was set to an extreme value: lower FID often came at the expense of top-k accuracy, and vice versa. However, for most mid-range values, performance remained robust and stable, indicating that LUMA is not overly sensitive to moderate changes in these hyperparameters.

Based on these results, we selected $\lambda_{\text{fre}} = 0.10$ and $\lambda_{\text{tem}} = 0.50$ as our final configuration, achieving an optimal trade-off between perceptual quality (FID) and semantic alignment (top-k metrics) for all subsequent experiments.

E DETAILS OF HUMAN EVALUATION

Task design. We utilize the Google Form platform to conduct a human evaluation study involving 20 independent participants. For each test prompt, we generate *four* motion clips: one from our model and three from baseline methods MDM, MoMask, StableMoFusion, denoted as $\mathcal{B}_1, \mathcal{B}_2, \mathcal{B}_3$. Each trial presents one baseline clip and the corresponding clip generated by our model, displayed side-by-side. The identity of the baseline is concealed, and the left/right order is randomized. Each participant is shown 100 such trials, covering different prompts and baseline combinations, and is asked to answer three questions:

Human-study questionnaire

Motion A (left) and **Motion B** (right) are two candidates for the same text prompt.

1. **Q1: Semantic accuracy.**

Is Motion A semantically accurate?

- (a) Yes
- (b) No

Is Motion B semantically accurate?

- (a) Yes
- (b) No

2. **Q2: Complex-action completion.**

Is the complex movement well completed in Motion A?

- (a) Yes
- (b) No

Is the complex movement well completed in Motion B?

- (a) Yes
- (b) No

3. **Q3: Preference.**

Which result do you prefer?

- (a) Motion A
- (b) Motion B

Metric Computation. We report three metrics for our model and each baseline, computed as follows:

- 1026 1. **Semantic Accuracy (%)**: The proportion of trials where the motion is rated as semantically
 1027 consistent with the input text (Q1).
 1028
 1029 2. **Complex-Movement Completion (%)**: The proportion of trials where the motion is judged to
 1030 have successfully completed the described complex action (Q2).
 1031
 1032 3. **User Preference (%)**: The proportion of pairwise comparisons in which our model is preferred
 1033 over a given baseline (Q3).

| Method | Sem. Acc. | Compl. Comp. | Pref. |
|----------------|--------------|--------------|-------|
| MDM | 55.3% | 52.7% | 84.1% |
| MoMask | 70.8% | 72.4% | 75.6% |
| StableMoFusion | 60.2% | 55.5% | 80.9% |
| Ours | 85.7% | 86.3% | – |

1040 Table 7: Human evaluation results (all values in %). “Sem. Acc.” is Semantic Accuracy, “Compl.
 1041 Comp.” is Complex-Movement Completion, and “Pref.” is User Preference.

1043 **Results.** Ours achieves the highest semantic accuracy (85.7%) and complex-action completion
 1044 (86.3%); in pairwise comparisons annotators preferred our results over each baseline in 75.6–84.1%
 1045 of trials, confirming both objective and subjective superiority while keeping the evaluation logic
 1046 transparent and reproducible.

1048 F DISCRETE COSINE TRANSFORM ANALYSIS

1050 F.0.1 DEFINITION

1051 For a discrete signal $\mathbf{v} = \{v[0], v[1], \dots, v[N_{\text{sig}} - 1]\}$, its type-II Discrete Cosine Transform (DCT)
 1052 is

$$1053 v_f[k] = A(k) \sum_{n=0}^{N_{\text{sig}}-1} v[n] \cos\left(\frac{\pi(2n+1)k}{2N_{\text{sig}}}\right), \quad k = 0, \dots, N_{\text{sig}} - 1, \quad (16)$$

1056 where the normalisation factor is

$$1057 A(k) = \begin{cases} \sqrt{\frac{1}{N_{\text{sig}}}}, & k = 0, \\ \sqrt{\frac{2}{N_{\text{sig}}}}, & k > 0. \end{cases} \quad (17)$$

1062 F.0.2 FREQUENCY ANALYSIS ON HUMANML3D

1063 Figure 8 shows the average spectral magnitude of HumanML3D motions after applying the
 1064 DCT equation 16. We observe that the signal energy is highly concentrated in the low-frequency
 1065 region: the first 64 DCT coefficients, comprising only one-third of the total frequency range, cap-
 1066 ture 83.6% of the total spectral energy. These low-frequency components primarily encode the
 1067 smooth, semantically meaningful motion trajectories, while higher frequencies contribute mainly to
 1068 fine-grained details and noise.

1069 Motivated by this empirical evidence, LUMA sets the frequency semantic anchor dimension to
 1070 $k = 64$, ensuring that the model retains the most informative, energy-rich motion patterns while
 1071 discarding redundant high-frequency noise and reducing computational cost.

1074 G GRADIENT-FLOW ANALYSIS

1075 To quantify optimisation dynamics we attach hooks to every UNet block and record gradient magni-
 1076 tudes after 40k training steps on HumanML3D. As Figure 10 shows, the baseline exhibits a spatial
 1077 profile *U-shaped*: the encoder and decoder peripheries maintain healthy gradients, while the bottle-
 1078 neck layers suffer from one to two orders of magnitude attenuation, sometimes dipping below the 1
 1079 % threshold, effectively failing to learn. Temporally, the profile is *non-monotonic*: gradients in deep

| Method | FID ↓ | R-Precision ↑ | | |
|----------------|-------------------------|-------------------------|-------------------------|-------------------------|
| | | top1 | top2 | top3 |
| StableMoFusion | 0.291 \pm .008 | 0.427 \pm .002 | 0.603 \pm .003 | 0.700 \pm .003 |
| LUMA | 0.135 \pm .010 | 0.430 \pm .003 | 0.606 \pm .004 | 0.703 \pm .004 |

Table 8: Comparison on synonym-perturbed text. LUMA achieves better FID and R-Precision, indicating higher robustness.

| $\lambda_{fre} \backslash \lambda_{tem}$ | 0.05 | 0.10 | 0.20 | 0.35 | 0.50 |
|--|--------------|--------------|--------------|--------------|---------------------|
| 0.05 | 0.5461/0.059 | 0.5524/0.042 | 0.5448/0.041 | 0.5379/0.036 | 0.5412/0.037 |
| 0.10 | 0.5397/0.054 | 0.5563/0.045 | 0.5481/0.049 | 0.5422/0.048 | 0.5556/0.035 |
| 0.20 | 0.5486/0.056 | 0.5619/0.044 | 0.5514/0.052 | 0.5457/0.050 | 0.5542/0.039 |
| 0.35 | 0.5449/0.068 | 0.5417/0.040 | 0.5485/0.047 | 0.5433/0.034 | 0.5396/0.036 |
| 0.50 | 0.5504/0.046 | 0.5557/0.049 | 0.5467/0.048 | 0.5479/0.038 | 0.5508/0.043 |

Table 9: LUMA dual-anchor grid search. Each cell shows Top-1 Accuracy ↑ / FID ↓. Anchor weights λ_{fre} and λ_{tem} vary in {0.05, 0.10, 0.20, 0.35, 0.50}.

layers are largest at early timesteps ($t=750$) when the network must infer global, abstract motion structure from noise, but they decay rapidly as $t \rightarrow 0$, indicating that late denoising relies chiefly on shallow layers to polish local details. This spatiotemporal imbalance both motivates our *Timestep-aware FiLM modulation*, which reinjects semantic signals proportional to t and justifies the design of *Dual Anchor Loss*, whose intermediate supervision revitalises the vanishing bottleneck gradients.

H ROBUSTNESS ANALYSIS

To evaluate the robustness of our LUMA model under real-world language variations, we conducted experiments on a synonym-perturbed dataset constructed following the protocol of SATO Chen et al. (2024), where the textual descriptions are randomly replaced with contextually appropriate synonyms. The quantitative results (Table 8) demonstrate that LUMA consistently outperforms the baseline StableMoFusion across all key metrics. Notably, LUMA achieves a substantially lower FID (0.135 vs. 0.291), indicating that the generated motions maintain higher fidelity to the ground truth even when the input text is perturbed. Similarly, LUMA yields higher R-Precision scores at top-1, top-2, and top-3 ranks, reflecting a stronger alignment between the generated motions and the intended semantic content of the perturbed prompts. These results confirm that LUMA is significantly more robust to synonym-level variations in natural language input, thereby improving the reliability and practicality of text-to-motion generation in real-world, linguistically diverse scenarios.

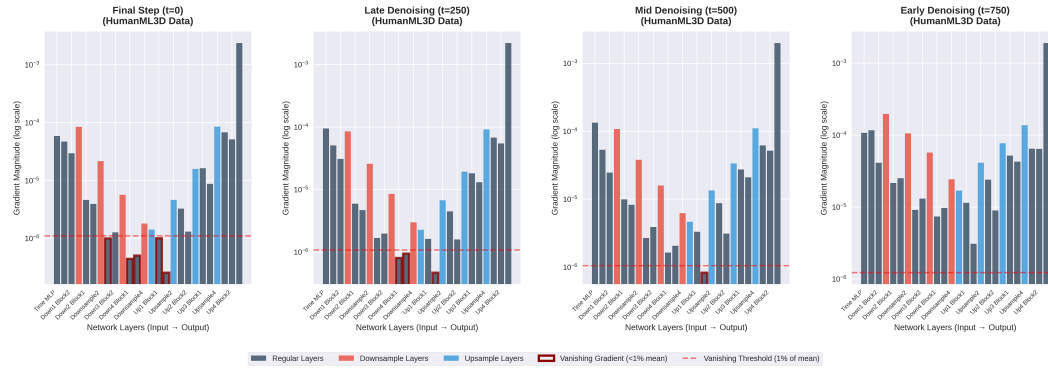


Figure 9: **Gradient magnitude across UNet layers (baseline, no DAL)**. At all timesteps ($t \in \{750, 500, 250, 0\}$), gradients form a U-shaped spatial profile, with severe attenuation in deep bottleneck layers. As t decreases, the gradient further weakens, suggesting spatio-temporal specialization during denoising.

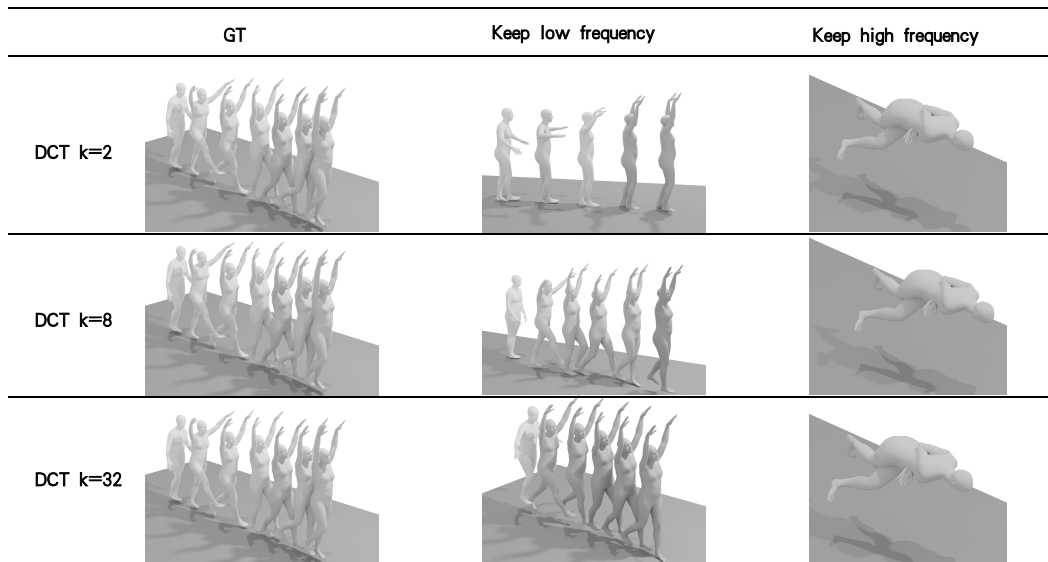


Figure 10: Visualization of motion reconstruction under DCT filtering for the text prompt “a man walks forward raising his hands”. Each row corresponds to a different cutoff value $K \in \{2, 8, 32\}$. The middle column (“Keep low frequency”) retains only the first K low-frequency DCT coefficients, while the right column (“Keep high frequency”) keeps only the last K high-frequency coefficients. Keeping low-frequency components preserves the global trajectory and the core semantic structure of the action, and fidelity improves as K increases. In contrast, retaining only high-frequency components removes semantic content and leads to fragmented, unstable motion. These results indicate that the linguistic and semantic meaning of human motion is primarily carried by low-frequency DCT bands.

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241



Figure 11: Comparison of final-layer attention between MoCLIP and original CLIP on motion descriptions. Red tokens receive more attention in MoCLIP, blue in original CLIP; red borders denote action words. MoCLIP consistently reallocates attention toward verbs and motion-centric tokens, highlighting its improved focus on the semantic core of actions.

1242 LLM USAGE STATEMENT
1243

1244 During the preparation of this manuscript, we strictly used a large language model (LLM) for lan-
1245 guage refinement purposes. Specifically, the LLM was used to improve the clarity, grammar, and
1246 readability of sentences originally written by the authors. The technical content, including the con-
1247 ceptual framework, methodology, mathematical formulations, experimental setup, and results, was
1248 entirely designed, developed, and validated by the authors without AI assistance.

1249 The LLM was not used to generate novel ideas, algorithms, experimental designs, analyses, or
1250 citations. All scientific contributions and claims in this work originate solely from the authors.
1251 We also performed thorough reviews to ensure that the final text accurately reflects our intended
1252 meaning and that no fabricated references or unsupported claims were introduced during the editing
1253 process.
1254

1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295