

MOTIONGPT3: HUMAN MOTION AS A SECOND MODALITY

Anonymous authors

Paper under double-blind review

ABSTRACT

With the rapid progress of large language models (LLMs), multimodal frameworks that unify understanding and generation have become promising, yet they face increasing complexity as the number of modalities and tasks grows. We observe that motion quantization introduces approximation errors that cap motion quality, while unifying discrete text and continuous motion within a single-stream backbone amplifies cross-modal interference. Motivated by recent multi-branch designs that separate signals from different modalities, we propose MotionGPT3, a bimodal motion–language model for both understanding and generation. MotionGPT3 encodes raw motion into a continuous latent space, thereby avoiding quantization-induced artifacts, while leveraging the semantic prior of pretrained language models. A dual-stream Transformer with shared attention preserves modality-specific routes while enabling controlled, bidirectional information flow, which reduces interference, stabilizing optimization, and empirically accelerates convergence without degrading fidelity. For multimodal joint training, a generate-then-align three-stage schedule further improves stability and limits cross-task interference. Experiments show that MotionGPT3 achieves 2× faster convergence in training loss and up to 4× faster convergence in validation, while maintaining state-of-the-art performance on standard motion understanding and motion generation benchmarks.

1 INTRODUCTION

Multimodal large language models (MLLMs) have recently achieved rapid progress in understanding and generation across text, images (Team, 2024; Wu et al., 2024a; Zhou et al., 2024), audio (Agostinelli et al., 2023; Copet et al., 2023; Liu et al., 2024), and video (Kondratyuk et al., 2023; Zhang et al., 2023a; 2024d). Built on the strong semantic priors and in-context learning capabilities of pretrained LLMs, these models capture long-range dependencies and compositional structure, enabling few-shot transfer and controllable citep across modalities Alayrac et al. (2022); Chowdhery et al. (2023); Dong et al. (2023); Li et al. (2023); Touvron et al. (2023). Toward Unified Motion–Language Modeling. While most prior work has focused mainly on text-driven motion synthesis (Shafir et al., 2023; Tevet et al., 2022a;b;c; Xin et al., 2023; Zhang et al., 2024a), unified motion–language models for both understanding and generation remain comparatively underexplored. Pursuing both tasks in a single model demands representations and training strategies that respect the distinctive statistics of human motion without sacrificing the reasoning benefits of language models.

Tokenizing motion into a fixed codebook, typically via VQ-based models, facilitates integration with Transformer-based LMs (Guo et al., 2022c; Zhang et al., 2023b; 2024c), however inevitably introduces quantization error, attenuating high-frequency components and degrading semantic-physical consistency. More importantly, treating motion as "language" (Jiang et al., 2023; Wang et al., 2023; Siyao et al., 2022) overlooks the gap between symbolic sequences and continuous trajectories (Wang et al., 2025b). Consequently, cross-modal alignment often remains at a symbolic level and struggles to capture the fine-grained kinematics demanded by nuanced linguistic semantics. In practice, limited codebook capacity and training coverage further constrain realism and controllability.

Recent MLLMs tend to process multiple modalities within a single backbone and attach modality-specific heads and supervision (Park et al., 2025; Team, 2024; Zhang et al., 2024b; Zhou et al., 2024). However, jointly optimizing multimodal objectives induces gradient interference and loss-scale mismatch, which increases hyperparameter sensitivity, destabilizes training, and can erode language competence (Driess et al., 2023; Kendall et al., 2018; Tsimpoukelli et al., 2021). Moreover, forcing distinct modalities into a shared space erodes modality-specific information and inductive biases, causing negative transfer. For **robust, controllable motion–language modeling**, a method is needed

[†]Corresponding authors.

that (i) adopts representations respect the continuous nature of human movement and (ii) explicitly balances multimodal, multi-objective training. Addressing these representational and optimization bottlenecks is, therefore, key to advancing unified motion understanding and generation.

Continuous Motion Latent Space Our approach first replace the motion tokens with a continuous, low-dimensional latent representation learned by a pretrained motion VAE (Xin et al., 2023). By ‘continuous’ we mean: (i) the latents are real-valued vectors rather than discrete code indices, and (ii) the VAE induces a smooth latent manifold in which nearby points correspond to gradually varying motions. Compared with VQ-based tokenization, this latent space is perceptually aligned with the original trajectories yet computationally compact, avoiding quantization artifacts and preserves high-frequency micro-dynamics for efficient and stable motion synthesis.

Diffusion Bridge within the LLM Framework Autoregressive generation and cross-entropy-based supervision in LLMs presumes discrete token targets and is therefore ill-suited to continuous motion latents. Conditioned on the LLM’s hidden states, we further attach a lightweight diffusion head that perform denoising directly in the motion latent space to predict motion VAE latents, which the motion decoder then converts into motion sequences. Operating in a low-dimensional latent domain with a relatively small expert, this diffusion scheme bridges the gap between LLM hidden states and motion latents while only brings little overhead in both training and inference.

Bimodal Architecture Following Mixture-of-Transformers (MoT) (Liang et al., 2024), we treat human motion as a second modality and introduce a motion branch symmetric to the language backbone. The two independent branches interact via shared attention layers, yet retain modality-specific embeddings and allow each module to be guided by its own objective. This bimodal design mitigates interference between modalities and preserves each modality’s structure, thereby enabling high-quality motion understanding and generation within a unified framework.

Three-Stage Training To effectively model motion branch under guidance of a pre-trained language model, we design a three-stage training scheme. First, we perform Uni-task Pretraining. with the text branch frozen, the motion branch is pre-trained on text-to-motion generation. Next, in Cross-Modal Alignment, motion-to-text and motion prediction objectives are introduced to align two branches. Finally, all parameters are optimized in Joint Fine-Tuning.

We summarize our contributions as follows: (i) Latent diffusion for motion. Unlike quantization-based pipelines (Zhang et al., 2023b; 2024c), we integrate latent diffusion (Rombach et al., 2022; Xin et al., 2023) into autoregressive backbone via a diffusion head, bridging the continuous motion motion with the next-token prediction framework for higher-fidelity and diverse synthesis. (ii) Architecture and training. We propose a bimodal motion-language framework with per-modal branches communicating through shared attention, reducing interference while preserving modality-specific intelligence. A three-stage generate-then-align scheme further stabilizes joint training and curbs negative transfer. (iii) Results and efficiency. Under comparable settings, MotionGPT3 achieves state-of-the-art performance on text-to-motion, motion-to-text, while reducing training time by approximately 2–3x.

2 RELATED WORK

Human Motion Modeling Early approaches leverage strong text encoders (Li et al., 2022; Radford et al., 2021; Raffel et al., 2020; Sanh et al., 2019) to develop motion–language understanding/retrieval via shared embeddings (Guo et al., 2022c; Tevet et al., 2022a; Yin et al., 2024) or contrastive learning (Chen et al., 2024; Petrovich et al., 2023). Recent methods (Athanasίου et al., 2024; Cohan et al., 2024; Shafir et al., 2023; Tevet et al., 2022b; Xin et al., 2023; Zhang et al., 2023c; 2024b) advance text-to-motion generation with diffusion backbones (Ho et al., 2020; Song et al., 2020), either operating directly on raw motion sequence or reconstructing a VAE latent. Working in a compressed latent space, LDMs (Rombach et al., 2022) keeps training computationally cheaper and inference faster while maintaining synthesis quality. In parallel, to fit next-token–prediction recipes in large language model (LLM), several works discretize motion with VQ-VAE (Esser et al., 2021; Van Den Oord et al., 2017) into token indices, enabling transformer-based generation (Zhang et al., 2023b; 2024c; 2025). However, quantization induces approximation error and a “symbolic–continuous mismatch” that attenuates fine-grained kinematics and limits controllability, while refined tokenization such as residual VQ (RVQ) (Guo et al., 2024) and post-training schemes (Wang et al., 2025b) only partially alleviate these issues and cannot fundamentally avoid the numerical and semantic discontinuities induced by tokenization. Accordingly, we adopt an approach that interface language models directly with unquantized VAE latent representations.

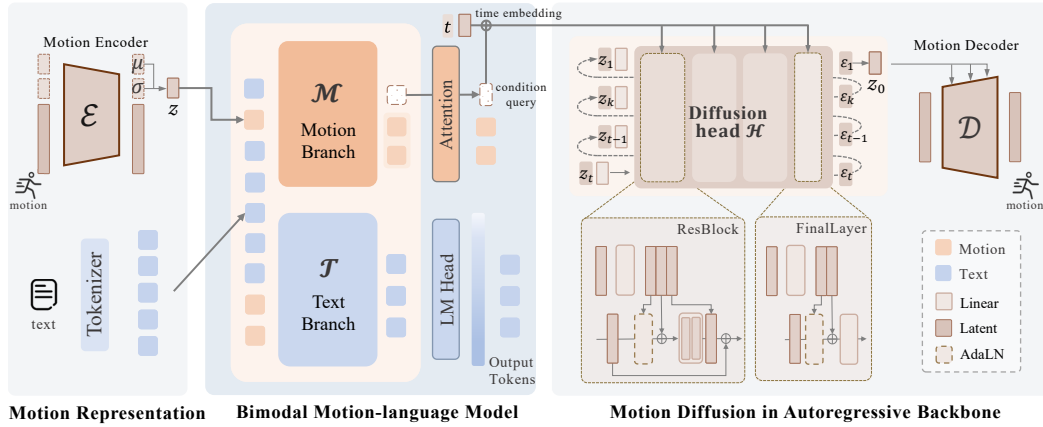


Figure 1: MotionGPT3 introduces hybrid motion-language model that takes motion as a second modality and processes the data through a new branch, with cross-modal attention mechanism to communicate with text branch (Sec. 3.2). We leverage a VAE network for continuous motion representation (Sec. 3.1), and design separate training objective for each modality (Sec. 3.3).

Human motion modeling has evolved from task-specific designs to **unified frameworks for multi-modal understanding and generation**. Recent motion models (Jiang et al., 2023; Park et al., 2025; Wang et al., 2024; Wu et al., 2025) adopt single-stream LM backbones with discretized motion tokens to support bidirectional text-motion mapping, and have been extended to other modalities such as music (Luo et al., 2024; You et al., 2024). In parallel, language-centric multimodal frameworks extend LLMs beyond text via lightweight adapters and cross-attention conditioning, offering a generic recipe for vision- and audio-grounded reasoning (Alayrac et al., 2022; Li et al., 2023; Copet et al., 2023; Liu et al., 2023a). Works such as NExT-GPT (Wu et al., 2024b) and Janus (Wu et al., 2024a; Ma et al., 2024) employ pretrained encoders/decoders to map inputs into modality-specific latent spaces and attach adapters for flexible multimodal generation. In vision-language, Chameleon (Team, 2024) and Transfusion (Zhou et al., 2024) discretize or encode images into token sequences to support interleaved text-image training, while Show-o (Xie et al., 2024) and Fuyu (Bavishi et al., 2023) employ masked or causal attention for joint reasoning. Despite their versatility, single-stream architectures often suffer from cross-modal interference, limiting scalability and robustness. Even with carefully tuned objectives, newly introduced modalities can disrupt existing representations, underscoring the challenge of preserving modality-specific capability while scaling to new domains.

Mixture-of-Experts and Multi-Stream Architecture address these limitations by routing inputs to modality-specific experts while maintaining a shared fusion interface (Alayrac et al., 2022; Li et al., 2023; Tsimpoukelli et al., 2021). This separation reduces gradient interference between modalities, and enables branch to be guided by its own objective (Liu et al., 2021; Sener & Koltun, 2018), and simplifies the introduction of new modalities. These insights motivate hybrid strategies (Cho et al., 2024; Shi et al., 2025; Wang et al., 2025a) that combine discrete and continuous representations and decouple modality-specific encoders with *minimal* modification on the LLM backbone, thereby enhancing alignment and expressiveness. Mixture-of-Transformers (MoT) (Liang et al., 2024) instantiates this idea with modality-specific Transformer experts coupled through shared attention, facilitating modular training and reducing interference when incorporating new modalities. Guided by these observations, we adopt a MoT-style architecture that isolates motion representation learning while leveraging the language competence of pretrained LLMs (Bai et al., 2023; Radford et al., 2019).

3 METHOD

To couple motion understanding and generation into language-centric LLMs, we observe that although discretization in prior unified systems (Jiang et al., 2023; Wang et al., 2024; Wu et al., 2025) facilitates reuse of text-style training and inference pipelines, it inevitably removes fine-grained details and complicates optimization. Moreover, single-stream backbones exacerbate cross-modal interference and yield imbalanced training. We circumvent these limitations by representing motion in a continuous, perceptually faithful VAE latent space (Sec. 3.1) and adopting a hybrid motion-text backbone that processes the two streams separately while permitting controlled interaction via

shared self-attention (Sec. 3.2). On top of this backbone, we attach a diffusion head conditioned on LLM hidden states to bridge language and motion latents, enabling bidirectional understanding and generation (Sec. 3.3). Finally, together with a three-stage training schedule (Sec. 3.4), our MotionGPT3 avoids quantization bottlenecks and improves training and inference efficiency while maintaining generation quality.

3.1 MOTION REPRESENTATION IN CONTINUOUS TOKENS

To align motion with the autoregressive generation paradigm of large language models (LLMs) (Bai et al., 2023; Radford et al., 2019; Raffel et al., 2020; Touvron et al., 2023), previous approaches typically *descretizes* motion with vector-quantized autoencoders (Esser et al., 2021; Van Den Oord et al., 2017), converting a N -length sequences into latents $y \in \mathbb{R}^{n \times d}$ and replace each vector $y^i \in \mathbb{R}^d$ by its nearest code e_k from a learned K -entry codebook. The corresponding indexes $k^{i \dots n}$ then serve as tokens for LLM training (Zhang et al., 2023b; Jiang et al., 2023; Wang et al., 2024). While they integrate cleanly with standard next-token prediction objectives in LLMs such as cross-entropy, quantization process inevitably introduces approximation error and disrupts motion continuity, weakening fine-grained dynamics and constraining controllability. Guo et al. (2024), equipped with residual vector quantization (RVQ), leverages multiple codebooks whose decoded latents are summed to reduce information loss. In parallel, Wang et al. (2025b) explores refined post-training tokenization. However, neither strategy fundamentally resolves the numerical and semantic discontinuities inherent to discretization.

In contrast, we adopt a continuous latent space learned by a motion VAE. Given a N frame motion sequence $m^{1:N}$, the encoder \mathcal{E} map m into a compact continuous latent vector $z \in \mathbb{R}^d$, and the decoder \mathcal{D} reconstructs $m^{1:M} = \mathcal{D}(z) = \mathcal{D}(\mathcal{E}(m^{1:M}))$. The VAE is trained once with a reconstruction term (optionally including kinematic losses on pose and velocity) and a KL regularizer (Kullback & Leibler, 1951) to prevent high-variance latents and promote a smooth manifold. This compressed, continuous representation learns the inherent structure of z , and preserves subtle variations and maintains numerical and semantic continuity, while providing a compact domain in which our downstream generator operates. Further details can be found in the supplement.

3.2 BIMODAL MOTION-LANGUAGE FRAMEWORK

To accommodate the distinct characteristics of language and motion while enabling efficient cross-modal interaction, we augment a decoder-only transformer backbone Radford et al. (2019) with a *parallel* motion branch. Unlike single-stream designs that merge all modalities into one pathway Jiang et al. (2023); Wu et al. (2025), our architecture preserves modality-specific routes: a text branch \mathcal{T} and a motion branch \mathcal{M} . Each branch maintains its own embeddings, feed-forward blocks, and normalization, and information exchange occurs only in shared self-attention layers (Alayrac et al., 2022; Shi et al., 2025). The motion branch is initialized from scratch and trained primarily under its own objective, thereby capturing motion-specific inductive biases and reducing cross-modal interference during multimodal training (Yu et al., 2020; Zhou et al., 2023).

Hybrid Sequence Route As illustrated in Fig. 1, given an input sequence $S = s^{1:k}$, each element is embedded either as a text embedding τ_i or as a motion latent z_i , with a routing indicator $\vartheta_i \in \{0, 1\}$ dispatches them to \mathcal{T} or \mathcal{M} . The branches compute hidden states h_t and h_m separately, which are then reassembled in input order for shared self-attention layers. This hybrid routing supports interleaved text–motion processing without collapsing modalities into a single embedding space, laying the foundation for high-quality, condition-aware generation.

Interfaces for Continuous Motion Latents Because that the continuous motion representation (Xin et al., 2023) does not rely on a tokenized vocabulary or codebook, the index-to-embedding lookup and softmax decoding employed for text cannot be reused for motion. We therefore introduce motion-specific interfaces that bridge continuous latents and transformer hidden states. First, we augment the text vocabulary with a small set of motion-boundary/holder tokens (i.e. `<som>` `<eom>`, `<motion_in>`, `<motion_out>`) to mark motion spans and I/O positions in interleaved sequences as in Zhou et al. (2024). Second, a Motion Understanding Head (MUH) linearly maps motion latents into the Transformer’s input embedding space for captioning and prediction. Finally, a lightweight Motion Generation Head (MGH) projects hidden states back to the VAE latent space via diffusion (Ho et al., 2020; Rombach et al., 2022; Xin et al., 2023).

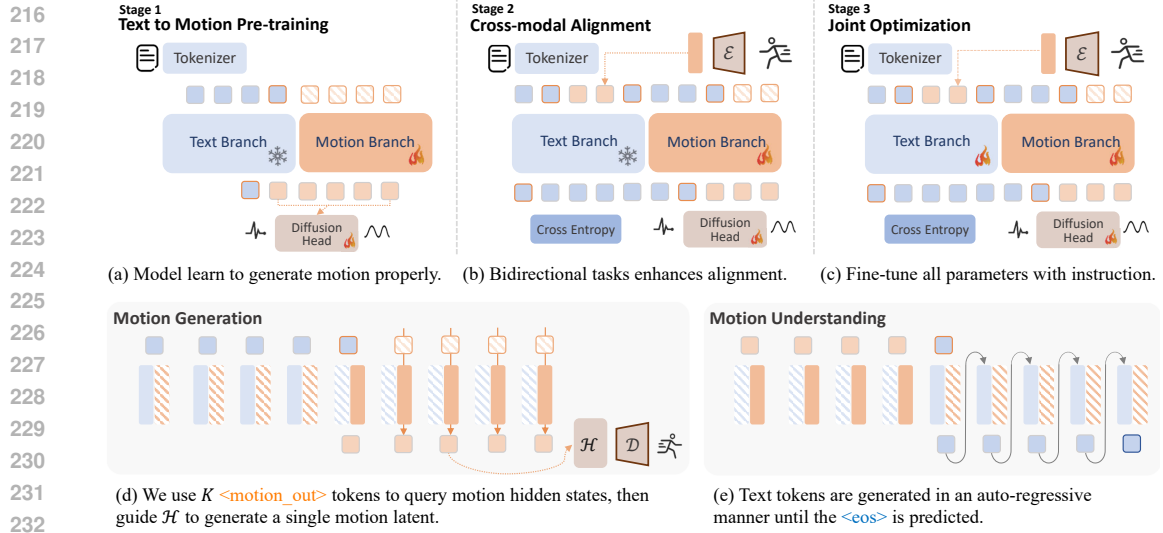


Figure 2: We propose a **three-stage alignment** strategy for our hybrid motion-language model: (a) The text branch is frozen, and only motion output is supervised. (b) Motion reasoning is introduced to further align the motion branch with language, with supervision on both modalities. (c) All modules are jointly fine-tuned with text branch unfrozen. (d)(e) shows inference time behavior of two branches which process data only with the same modality tags (differentiated by colors). Each rectangle block represents for a whole text/motion branch, while shadowed ones denote inactive modules. Modalities are color-coded: **blue for text** and **orange for motion**. Shadowed orange squares represent `<motion_out>`, and orange-outlined squares indicate boundary tokens `<som>` or `<eom>`.

3.3 MOTION DIFFUSION IN AUTOREGRESSIVE BACKBONE

Continuous representations are inherently misaligned with the discrete nature of token-based generation in LLMs, and requires more sophisticated modeling to support generation. Inspired by recent advances in diffusion-based generative modeling (Tevet et al., 2022b; Song et al., 2020; Rombach et al., 2022; Li et al., 2024a), we attach a lightweight diffusion module \mathcal{H} to bridge this gap. \mathcal{H} predicts motion latents directly from the backbone’s hidden states, enabling integration of continuous motion representation within an autoregressive framework.

Diffusion Process Given a ground-truth motion sequence x , we obtain the target latent $z_0 = \mathcal{E}(x) \in \mathbb{R}^d$ via the motion encoder \mathcal{E} . We adopt a fixed forward noising process over $t \in \{1, \dots, T\}$ with Gaussian perturbations: $z_t = \sqrt{\bar{\alpha}_t} z_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon$, where $\epsilon \sim \mathcal{N}(0, I)$ and $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$ is the cumulative product of noise scheduling coefficients. A time-aware denoiser \mathcal{H} is conditioned on the Transformer hidden states from the motion stream (denoted h_m) and learns to reverse the diffusion process. Conditioning is implemented via lightweight linear projections into the denoiser inputs, following recent practice (Ho & Salimans, 2022; Li et al., 2024a). We train \mathcal{H} with the standard DDPM objective (Ho et al., 2020; Song et al., 2020): $\mathcal{L}_{\text{diff}} = \mathbb{E}_{z_0, t, \epsilon} [\|\epsilon - \mathcal{H}(z_t, t, h_m)\|_2^2]$.

Inference The text branch autoregressively generates tokens until a motion start marker `<som>` is produced. We then insert K placeholder tokens (i.e., `<motion_out>`) to elicit the span-aligned hidden states $h_m^{i:i+K}$ in a single forward pass. The diffusion head runs the reverse process conditioned on $h_m^{i:i+K}$ to sample the noise-free motion latent \hat{z}_0 , which the VAE decoder \mathcal{D} finally decoded to the raw motion sequence. As in Team (2024), generation then resumes in the text stream with a concatenated `<eom>` until end-of-sequence. Operating on compact motion latents, this diffusion head adds only minimal overhead during training and inference.

3.4 TRAINING PROCEDURE

We first train a motion VAE to obtain a compact, continuous latent space, following prior works (Rombach et al., 2022; Xin et al., 2023). The bimodal backbone then adopts a pretrained decoder-only LLM (Radford et al., 2019) as the text branch \mathcal{T} . As described in Fig. 2, the motion branch \mathcal{M} is initialized from scratch and brought into alignment with \mathcal{T} through a three-stage schedule.

Stage I: Text-to-motion pretraining We begin by pretraining \mathcal{M} on text-to-motion, while freezing \mathcal{T} . This provides stable linguistic conditioning and biases the model toward motion-specific semantics. In this stage, \mathcal{M} conditions on the frozen language representations and is trained via diffusion, to synthesize VAE motion latents, where diverse text-motion pairs encourages a rich and flexible mapping from language to the latent space (Rombach et al., 2022; Xin et al., 2023).

Stage II: Cross-Modal Alignment Keeping \mathcal{T} frozen, we introduce additional objectives to couple understanding and generation. Concretely, training includes multiple tasks of text-to-motion (T2M), motion-to-text (M2T), and motion prediction. Following instruction-style formulations in Jiang et al. (2023), these tasks are further presented as prompts covering generation, captioning, prediction, and inbetweening. Multi-task optimization fosters bidirectional alignment without forcing a single shared representation and encourages motion representations that are semantically coherent with language features (Alayrac et al., 2022; Li et al., 2023).

Stage III: Joint Fine-Tuning Finally, we unfreeze \mathcal{T} and fine-tune all parameters via instruction tuning on a mixture of paired text-motion data and, optionally, text-only prompts (Dai et al., 2023; Liu et al., 2023b; Wei et al., 2021). Including text-only prompts can further improve language competence for downstream applications.

4 EXPERIMENTS

We empirically validate MotionGPT3, a dual-stream architecture for efficient, language-grounded multimodal motion understanding and generation, across motion-centric tasks. Dataset configurations, evaluation metrics, and implementation details are summarized in Sec. 4.1. Begin with analyzing optimization dynamics and inference efficiency via training loss and validation curves (Sec. 4.2), we then present controlled ablations that isolate the contributions of the continuous VAE motion representation and the bimodal design (Sec. 4.4). Next, We benchmark MotionGPT3 on text-to-motion generation and motion-to-text understanding, comparing against both specialized single-task methods and unified state-of-the-art systems (Sec. 4.3). Finally, we ablate the proposed three-stage training scheme (Sec. 4.4). Additional qualitative results are provided in the supplement.

4.1 EXPERIMENTAL SETUP

Datasets We train and evaluate our model on Guo et al. (2022b), a large-scale benchmark for text-motion generation and understanding. For comparison with prior works (Xin et al., 2023; Jiang et al., 2023), we adopt the 263-dim pose proposed in Guo et al. (2022b), which combines joint velocities/ positions/ rotations, and foot-contact signals, following the standard data split.

Evaluation Metrics We evaluate two tasks. For the *text-to-motion*, we follow the previous works (Guo et al., 2022c; Jiang et al., 2023; Xin et al., 2023; Zhang et al., 2023b) to report motion quality (FID), diversity (DIV and MM), and text-motion alignment (R-Precision and MMDist). For *motion-to-text*, we use both alignment metrics (R-Precision and MMDist) and linguistic metrics from NLP (Bleu (Papineni et al., 2002), Rouge-L (Lin, 2004), CIDEr (Vedantam et al., 2015), and BertScore (Zhang et al., 2019)). See Sec. D.3 for metric definitions and computation details.

Implementation Details Our framework comprises three main components: a motion VAE, a lightweight diffusion head, and a dual-stream backbone. We adopt the Transformer-based motion VAE of Xin et al. (2023), where both encoder and decoder consists of 9 layers and 4 heads with skip connections, producing a $1 \times 1 \times 256$ latent per motion sequence. Our Diffusion Head \mathcal{H} is implemented as a 3-layer MLP with ResBlock-style layers and hidden dimension 1024, following Li et al. (2024a). We train diffusion with a scaled linear noise schedule for 1000 denoising steps, while inference uses 100 steps by default. The text and motion branches share the GPT-2 base configuration but use disjoint parameters, both are decoder-only with 12 Transformer layers, model dimension 768, and MLP dimension 3072, unless stated otherwise. The text branch is initialized from a pretrained 124M GPT-2 checkpoint, while the motion branch from scratch, yielding total 238M parameters.

Training Protocol We use AdamW for all components, with a learning rate of 2×10^{-4} for the motion backbone and 1×10^{-4} for the diffusion head. Training uses a mini-batch size of 32 on 2 NVIDIA RTX 3090 GPUs, with identical training/inference settings on HumanML3D (Guo et al., 2022b). The motion VAE is trained with a learning rate of 1×10^{-4} , batch size 256, over 150K iterations. The motion-language backbone is trained for 100k iterations in text-to-motion pretraining, followed by 300k iterations for cross-modal alignment.

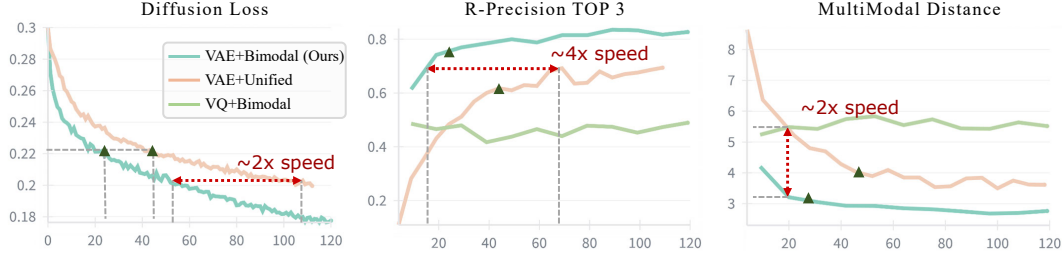


Figure 3: Training loss and validation curves on motion generation on HumanML3D for architecture variants of dual-stream and single-stream and representation variants of VAE and VQ latents. The right figures illustrate validation metrics of R-Precision TOP 3 ($R@3\uparrow$) and Multimodal Distance ($MMDist\downarrow$). Triangle markers indicate matched-loss checkpoints (~ 0.22). Our hybrid architecture with continuous motion representation helps accelerating convergence for about $2\times$, as well as achieves better quality especially in early training stage.

4.2 TRAINING EFFICIENCY WITH BIMODAL ARCHITECTURE AND CONTINUOUS LATENTS

To assess our two main design choices: (i) the dual-stream (bimodal) backbone and (ii) continuous representation from motion VAE, we analyze training dynamics on the text-to-motion task under identical settings on HumanML3D. Fig. 3 plots training loss and validation metrics over time for three variants of single-stream+VAE, dual-stream+VAE, and dual-stream+VQ.

We observe that i) **Discrete VQ latents plateau at a lower quality ceiling.** The VQ baseline (VQ + bimodal; green curve) reaches ~ 0.5 R-Precision Top 3 ($R@3$) early in training and then saturates, yielding substantially lower $R@3$ and higher MultiModal Distance (MMDist) than the VAE-based counterparts. This likely stems from quantization-induced information loss (Guo et al., 2024; Wang et al., 2025a) and tokenization that disrupts the semantic continuity of motion. ii) **Bimodal design accelerates optimization.** Compared with single-stream backbone, the dual-stream variant reduces diffusion loss roughly $2\times$ faster and sustains superior validation performance over the entire training trajectory, in terms of $R@3$ and MMDist. iii) **Superior quality at matched optimization states.** At comparable diffusion loss, the dual-stream model still leads. For instance, we mark by triangles a loss of ~ 0.22 , reached at ~ 20 epochs for dual-stream model and ~ 40 epochs for the single-stream model, the former achieves higher $R@3$ and lower MMDist.

We attribute these effects to **modality-aware optimization**: the motion branch learns motion-specific semantics while the language branch focuses on textual cues. Decoupling the streams and supervising them separately mitigates gradient interference and avoids the representational compromises common in single-stream models (Liu et al., 2021; Sener & Koltun, 2018; Yu et al., 2020). Although a shared space might be expected to bring modalities "closer", in practice single-stream coupling can entangle modality structure and yield counterintuitive outcomes. In summary, continuous VAE latents within a dual-stream backbone strike a favorable efficiency–fidelity trade-off, enabling high-quality motion synthesis with reduced training time.

4.3 COMPARISONS

By modeling human motion as a second modality alongside language, our bimodal motion-language model supports both text-to-motion (T2M) and motion-to-text (M2T). We report results for two settings: (i) single-task models trained specifically for target task (MotionGPT3 †), and (ii) a unified model (MotionGPT3) trained on both tasks with the three-stage scheme.

Text-to-Motion Generation The text-to-motion task involves generating realistic and diverse motion sequences conditioned on natural language descriptions. We train a single-task generator MotionGPT3 † for 200 epochs on HumanML3D and compare against recent methods (Guo et al., 2022b;c; Lou et al., 2023; Jiang et al., 2023; Xin et al., 2023; Zhang et al., 2023b;c; Guo et al., 2024; Wu et al., 2024c; Li et al., 2024b; Zhang et al., 2025). Following Guo et al. (2022b), each evaluation is repeated 20 times and reported with 95% confidence intervals. As shown in Tab. 1, MotionGPT3 † matches or exceeds generation-only baselines (Zhang et al., 2023b; Lou et al., 2023; Guo et al., 2024) on alignment metrics ($R@k$, MMDist), with competitive FID and diversity. The unified MotionGPT3 performs on par with or better than recent unified systems (Guo et al., 2022c; Jiang et al., 2023; Wu et al., 2024c). Sec. C.1 provides additional comparisons.

Table 1: Evaluation of text-guided motion generation on HumanML3D (Guo et al., 2022a). Rows are grouped by training tasks: *Gen. only* for generation-only and *Gen. & Und.* for both. *Real* is obtained by ground-truth motions, and \rightarrow indicate values closer to *Real* are desirable. \dagger marks our single-task model trained for 200 epochs, and MotionGPT3 is a three-stage model trained with unified tasks.

Types	Methods	R@1	R@2	R@3	FID \downarrow	MMDist \downarrow	Diversity \rightarrow	MModality \uparrow
	Real	0.511	0.703	0.797	0.002	2.974	9.503	-
Gen. only	T2M-GPT Zhang et al. (2023b)	0.491	0.68	0.775	0.116	3.118	9.761	1.856
	DiverseMotion Lou et al. (2023)	0.515	0.706	0.802	<u>0.072</u>	2.941	9.683	1.869
	MoMask Guo et al. (2024)	0.521	0.713	0.807	0.045	2.958	<u>9.620</u>	1.241
	MotionGPT3\dagger	<u>0.533</u>	0.731	<u>0.826</u>	0.239	<u>2.797</u>	9.688	1.560
Gen. & Und.	TM2T Guo et al. (2022c)	0.424	0.618	0.729	1.501	3.467	8.589	2.424
	MotionGPT Jiang et al. (2023)	0.492	0.681	0.733	0.232	3.096	9.528	2.00
	MoTe Wu et al. (2024c)	0.548	<u>0.737</u>	0.825	0.075	2.867	-	<u>2.399</u>
	MotionGPT3	0.553	0.747	0.837	0.208	2.725	9.700	1.018

Table 2: Comparison of motion captioning on HumanML3D (Guo et al., 2022a), evaluation follows (Guo et al., 2022c). MotionGPT3 \dagger denotes our single-task captioning model trained for 100 epochs, and MotionGPT3 is an unified model trained on both tasks with the three-stage scheme (Sec. 3.4). Both variants achieve R@k on par with recent state of the art, and surpass the GT metrics.

Methods	R@1	R@2	R@3	MMDist \downarrow	Bleu@1 \uparrow	Bleu@4 \uparrow	Rouge \uparrow	Cider \uparrow	BertScore \uparrow
Real	0.523	0.725	0.828	2.901	-	-	-	-	-
TM2T (Guo et al., 2022c)	0.516	-	0.823	2.935	48.9	7.00	38.1	16.8	32.2
MotionGPT (Jiang et al., 2023)	0.543	-	0.827	2.821	48.2	12.5	37.4	29.2	32.4
LaMPM2T (Li et al., 2024b)	0.547	-	0.831	2.808	47.8	13.04	37.1	28.9	32.7
MoTe (Wu et al., 2024c)	0.577	-	0.871	2.649	46.7	11.15	37.4	31.5	30.3
MotionGPT3\dagger	0.553	0.756	0.853	<u>2.524</u>	<u>56.363</u>	<u>17.661</u>	<u>44.997</u>	<u>30.980</u>	35.850
MotionGPT3	<u>0.573</u>	0.773	<u>0.864</u>	2.426	59.083	<u>19.412</u>	46.173	28.721	<u>35.231</u>

Motion-to-Text Understanding The motion-to-text task involves understanding motion sequences and generating semantically appropriate textual descriptions. We train a single-task captioner MotionGPT3 \dagger for 100 epochs and compare with recent SOTA (Guo et al., 2022c; Jiang et al., 2023; Li et al., 2024b; Wu et al., 2024c). Following Jiang et al. (2023), we evaluate on the raw ground truth texts using the TM2T protocol (Guo et al., 2022c). Results in Tab. 2 show that both MotionGPT3 \dagger and MotionGPT3 achieve strong retrieval performance and language metrics. Notably, we observe a marked reduction in Multimodal Distance (MMDist), indicating effective motion-language alignment under the dual-stream design.

4.4 ABLATION STUDIES

This section reports quantitative ablations. In contrast to training-curve analysis in Sec. 4.2, we evaluate final test-set performance on both text-to-motion (T2M) and motion-to-text (M2T). First, we assess the contributions of a dual-stream backbone and continuous VAE motion latents by varying one factor at a time. Then, we analyze the proposed three-stage training schedule and quantify its effects. We also examine the impact of hidden-state processing in the Diffusion Head \mathcal{H} and the use of classifier-free guidance (CFG). See Sec. C for more detailed experiments.

Model Design Tab. 3 summarized test-set results on HumanML3D (Guo et al., 2022a). We evaluate T2M and M2T separately and compare four variants obtained by crossing architecture (single- vs. dual-stream) with representation (discrete VQ vs. continuous VAE). Under the same evaluation protocol, replacing VAE with VQ or replacing a dual-stream backbone (Bimodal) with a single-stream one (Unified) consistently degrades performance on both tasks. Notably, changing the architecture change to Bimodal yields larger gains on M2T, whereas changing the representation to VAE yields larger gains on T2M. This task-dependent sensitivity is consistent with Sec. 4.2: decoupling streams mitigates cross-modal interference and benefits semantic-level alignment, while continuous latents reduce quantization loss and improve synthesis fidelity for motion generation.

Training Stage We ablate the three-stage schedule in Sec. 3.4, including 100k iters on text-to-motion pretraining (SI), 300k iters on cross-modal alignment (SII), and 50k iters on joint fine-tuning (SIII), and evaluate on both T2M and M2T. Results are summarized in Tab. 4. SI already yields strong generation and provides a motion-specialized initialization. Optimization in SII confers

Table 3: Component ablations on HumanML3D for representation choice and architecture design. *Unified* denotes a single-stream backbone, where one branch is shared by text and motion, as employed in Jiang et al. (2023), and *Bimodal* denotes a dual-stream backbone described in Sec. 3.2. VQ and VAE indicate discrete and continuous motion latents, respectively. For each configuration we train separate models for motion generation (T2M) and motion captioning (M2T) under the same protocol and report test-set metrics. All variants share the same GPT-2-style branch and hyperparameters, and training is run for 100 epochs on M2T and 200 epochs on T2M. Best and second-best results are highlighted in **bold** and underline.

Settings	Text-to-Motion				Motion-to-Text			
	R@1 ↑	R@3 ↑	MMDist ↓	FID ↓	R@1 ↑	R@3 ↑	MMDist ↓	BertScore ↑↑
Real	0.511 \pm 0.003	0.797 \pm 0.002	2.974 \pm 0.008	0.002 \pm 0	0.523	0.828	2.901	-
Unified+VQ	0.237 \pm 0.003	0.435 \pm 0.003	5.684 \pm 0.018	0.403 \pm 0.014	-	-	-	-
Unified+VAE	<u>0.501</u> \pm 0.003	<u>0.792</u> \pm 0.002	<u>2.841</u> \pm 0.011	0.489 \pm 0.017	0.234	0.426	5.976	16.197
Bimodal+VQ	0.300 \pm 0.005	0.532 \pm 0.02	4.937 \pm 0.077	0.454 \pm 0.078	<u>0.379</u>	<u>0.702</u>	<u>3.545</u>	<u>18.085</u>
Bimodal+VAE	0.533 \pm 0.002	0.826 \pm 0.003	2.797 \pm 0.007	0.239 \pm 0.008	0.553	0.853	2.524	35.850

T2M with *Bimodal+VQ* and *Unified+VQ* is extended to 400 epochs to approach convergence.

M2T results for *Unified+VQ* is not reported are omitted because performance remained unevaluable after more than 400 training epochs.

Table 4: Ablation on training-scheme. Enabled stages are marked with ✓, and colors encode the text branch **updated** or **frozen**. Best results are **bold** and second best are underline.

Stage I	Stage II	Stage III	Text-to-Motion			Motion-to-Text		
			R TOP3 ↑	FID ↓	MMDist ↓	R TOP1 ↑	Bleu@4 ↑	Bert ↑
✓	✗	✗	0.826	0.239	2.797	-	-	-
✓	✓	✗	0.831	0.215	2.755	0.571	18.328	33.993
✓	✓	✓	0.837	0.208	2.725	0.573	19.412	35.231
✗	✓	✓	0.772	0.325	3.108	0.573	18.277	35.546

M2T capability and, importantly, further improves T2M, by -0.10 on FID and -0.2 on MMD, indicating that explicit alignment benefits both directions. Without extra text-only supervision, SIII adds small additional gains on M2T while preserving T2M, serving as a light joint refinement rather than a substitute for S2. As shown in the last row of Tab. 4, a two-stage model that *omits S1* keeps M2T largely intact but markedly degrades T2M, underscoring the role of S1 in learning motion-specific features. Overall, the full three-stage schedule provides the best trade-off, delivering reliable generation and captioning with well-aligned motion-language representations. Additional variants are reported in Sec. C.6, including experiments with an unfrozen text branch.

5 DISCUSSION

To address quantization-induced degradation and cross-modal interference in multi-objective training, we present **MotionGPT3**, a dual-stream motion-language framework that unifies motion understanding and generation while preserving modality-specific inductive biases. By encoding motion as continuous VAE latents and generating in latent space with a lightweight diffusion head, the model avoids quantization artifacts and improves synthesis fidelity. The dual-stream Transformer with shared attention enables controlled bidirectional exchange, which strengthens text-motion alignment, reduces cross-modal interference, and empirically accelerates single-task convergence without degrading quality. For joint training of understanding and generation, a generate-then-align three-stage training schedule further stabilizes optimization and mitigates cross-task interference.

Limitations and Failure Cases Fine-grained control can fail on directional cues (e.g., left/right). Because the current VAE yields a single latent per sequence, segment-level composition and local semantic alignment for long motions are not explicitly supported. Generalization to out-of-domain descriptions is constrained by data coverage. Potential remedies include incorporating diverse text-only corpora in the final alignment stage, adopting stronger language backbones, and exploring hierarchical or segment-wise latent representations to enable compositional control. **Future Work** We will scale training to (i) larger, more diverse datasets, (ii) develop controllable motion with local semantic alignment and segment-level for long-horizon generation, and (iii) evaluate the framework with stronger language models and larger-scale training regimes to assess efficiency and robustness.

6 REPRODUCIBILITY STATEMENT

We have taken multiple steps to enable reproducibility. We document models, metrics, and training/evaluation settings in the paper/appendix and provide example code and demonstration materials in the supplementary package/website to facilitate replication and community auditing.

Model/algorithm details are specified in Sec. 3 and Figs. 1 and 11 with training/inference procedures in Secs. 3.4, D.1 and D.2, hyperparameters and optimizer settings are provided in Sec. 4.1. Datasets and all preprocessing steps are fully following previous methods Guo et al. (2022b); Jiang et al. (2023). Evaluation protocols, metrics, and test-time settings are described in Secs. 4.1 and D.3. Implementation details and scripts for training/evaluation are available in our supplemental code.

7 ETHICS STATEMENT

We have read and will adhere to the ICLR Code of Ethics.

This work does not involve the collection of new data from human subjects or interventions. All experiments use publicly available research datasets and evaluators for human-motion understanding/generation, primarily HumanML3D Guo et al. (2022b) and standard retrieval/captioning protocols reported in prior work Guo et al. (2022b); Jiang et al. (2023); Guo et al. (2022c); Petrovich et al. (2023).

Potential Risks and Responsible Use. Text-to-motion models can, in principle, be misused to synthesize deceptive behavior traces, to aid surveillance/individual tracking, or to control physical systems without adequate safety validation. Our release (code and research models) is intended for research/educational purposes only. We highlight known failure modes, e.g., errors on directional cues and limited long-horizon compositional control, which could lead to unsafe motions if outputs were executed on hardware without verification, underscoring the need for downstream safety safeguards. This work uses only publicly available research datasets with motion sequences; no new human-subjects data were collected, no interventions were conducted, and no personally identifiable information is used. All data are used under their licenses strictly for research and evaluation.

Fairness, Bias, and Inclusivity. Public motion/text datasets can reflect domain and cultural biases in action types and language. To mitigate risk, our study relies on task-standard, identity-agnostic motion representations and reports multiple alignment and retrieval metrics to monitor semantic fidelity and diversity. We encourage future audits with broader, culturally diverse data.

Privacy and Security. Our experiments rely on public benchmarks and evaluator code; we do not share any personally identifying information or private media. Any visualizations are derived from public benchmark samples or synthesis for illustration; we will not attempt to re-identify, link, or deanonymize any dataset items.

Legal Compliance and IP. We follow dataset licenses and applicable laws. Users who integrate additional third-party data must independently verify usage rights, privacy obligations, export/use restrictions, and local compliance.

Conflicts of Interest and Sponsorship. The authors have no undisclosed conflicts of interest or sponsorships related to this submission. Any support will be disclosed in the final camera-ready version, consistent with double-blind review policies.

Broader Impact. Our findings aim to improve motion–language alignment and robustness under standardized protocols and evaluators. We believe the above safeguards align with the principles of the ICLR Code of Ethics and we welcome community feedback and third-party audits on fairness, privacy, safety, and societal impact.

REFERENCES

- Andrea Agostinelli, Timo I Denk, Zalán Borsos, Jesse Engel, Mauro Verzetti, Antoine Caillon, Qingqing Huang, Aren Jansen, Adam Roberts, Marco Tagliasacchi, et al. Musiclm: Generating music from text. *arXiv preprint arXiv:2301.11325*, 2023.
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems*, 35:23716–23736, 2022.
- Nikos Athanasiou, Alpár Cseke, Markos Diomataris, Michael J Black, and Gül Varol. Motionfix: Text-driven 3d human motion editing. In *SIGGRAPH Asia 2024 Conference Papers*, pp. 1–11, 2024.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- Rohan Bavishi, Erich Elsen, Curtis Hawthorne, Maxwell Nye, Augustus Odena, Arushi Somani, and Sağnak Taşlılar. Introducing our multimodal models, 2023. URL <https://www.adept.ai/blog/fuyu-8b>.
- Ling-Hao Chen, Shunlin Lu, Ailing Zeng, Hao Zhang, Benyou Wang, Ruimao Zhang, and Lei Zhang. Motionllm: Understanding human behaviors from human motions and videos. *arXiv preprint arXiv:2405.20340*, 2024.
- Jungbin Cho, Junwan Kim, Jisoo Kim, Minseo Kim, Mingu Kang, Sungeun Hong, Tae-Hyun Oh, and Youngjae Yu. Discord: Discrete tokens to continuous motion via rectified flow decoding. *arXiv preprint arXiv:2411.19527*, 2024.
- Jungbin Cho, Junwan Kim, Jisoo Kim, Minseo Kim, Mingu Kang, Sungeun Hong, Tae-Hyun Oh, and Youngjae Yu. Discord: Discrete tokens to continuous motion via rectified flow decoding, 2025. URL <https://arxiv.org/abs/2411.19527>.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113, 2023.
- Setareh Cohan, Guy Tevet, Daniele Reda, Xue Bin Peng, and Michiel van de Panne. Flexible motion in-betweening with diffusion models. In *ACM SIGGRAPH 2024 Conference Papers*, pp. 1–9, 2024.
- Jade Copet, Felix Kreuk, Itai Gat, Tal Remez, David Kant, Gabriel Synnaeve, Yossi Adi, and Alexandre Défossez. Simple and controllable music generation. *Advances in Neural Information Processing Systems*, 36:47704–47720, 2023.
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale N Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. *Advances in neural information processing systems*, 36:49250–49267, 2023.
- Runpei Dong, Chunrui Han, Yuang Peng, Zekun Qi, Zheng Ge, Jinrong Yang, Liang Zhao, Jianjian Sun, Hongyu Zhou, Haoran Wei, et al. Dreamllm: Synergistic multimodal comprehension and creation. *arXiv preprint arXiv:2309.11499*, 2023.
- Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, Wenlong Huang, et al. Palm-e: An embodied multimodal language model. 2023.
- Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 12873–12883, 2021.
- 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 5152–5161, June 2022a.
- 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 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 5152–5161, 2022b.
- Chuan Guo, Xinxin Zuo, Sen Wang, and Li Cheng. Tm2t: Stochastic and tokenized modeling for the reciprocal generation of 3d human motions and texts. In *ECCV*, 2022c.

- 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*, pp. 1900–1910, 2024.
- Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.
- 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*, 36:20067–20079, 2023.
- Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 7482–7491, 2018.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Dan Kondratyuk, Lijun Yu, Xiuye Gu, José Lezama, Jonathan Huang, Grant Schindler, Rachel Hornung, Vighnesh Birodkar, Jimmy Yan, Ming-Chang Chiu, et al. Videopoet: A large language model for zero-shot video generation. *arXiv preprint arXiv:2312.14125*, 2023.
- Solomon Kullback and Richard A Leibler. On information and sufficiency. *The annals of mathematical statistics*, 22(1):79–86, 1951.
- Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pp. 12888–12900. PMLR, 2022.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023.
- Tianhong Li, Yonglong Tian, He Li, Mingyang Deng, and Kaiming He. Autoregressive image generation without vector quantization. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 56424–56445. Curran Associates, Inc., 2024a. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/66e226469f20625aaebddbe47f0ca997-Paper-Conference.pdf.
- Zhe Li, Weihao Yuan, Yisheng He, Lingteng Qiu, Shenhao Zhu, Xiaodong Gu, Weichao Shen, Yuan Dong, Zilong Dong, and Laurence T Yang. Lamp: Language-motion pretraining for motion generation, retrieval, and captioning. *arXiv preprint arXiv:2410.07093*, 2024b.
- Weixin Liang, Lili Yu, Liang Luo, Srinivasan Iyer, Ning Dong, Chunting Zhou, Gargi Ghosh, Mike Lewis, Wen-tau Yih, Luke Zettlemoyer, et al. Mixture-of-transformers: A sparse and scalable architecture for multi-modal foundation models. *arXiv preprint arXiv:2411.04996*, 2024.
- Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pp. 74–81, 2004.
- Bo Liu, Xingchao Liu, Xiaojie Jin, Peter Stone, and Qiang Liu. Conflict-averse gradient descent for multi-task learning. *Advances in Neural Information Processing Systems*, 34:18878–18890, 2021.
- Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo Liu, Danilo Mandic, Wenwu Wang, and Mark D Plumbley. Audioldm: Text-to-audio generation with latent diffusion models. *arXiv preprint arXiv:2301.12503*, 2023a.
- Haohe Liu, Yi Yuan, Xubo Liu, Xinhao Mei, Qiuqiang Kong, Qiao Tian, Yuping Wang, Wenwu Wang, Yuxuan Wang, and Mark D Plumbley. Audioldm 2: Learning holistic audio generation with self-supervised pretraining. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 32:2871–2883, 2024.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916, 2023b.
- Yunhong Lou, Linchao Zhu, Yaxiong Wang, Xiaohan Wang, and Yi Yang. Diversemotion: Towards diverse human motion generation via discrete diffusion. *arXiv preprint arXiv:2309.01372*, 2023.
- Mingshuang Luo, Ruibing Hou, Zhuo Li, Hong Chang, Zimo Liu, Yaowei Wang, and Shiguang Shan. M³gpt: An advanced multimodal, multitask framework for motion comprehension and generation. *arXiv preprint arXiv:2405.16273*, 2024.

- Yiyang Ma, Xingchao Liu, Xiaokang Chen, Wen Liu, Chengyue Wu, Zhiyu Wu, Zizheng Pan, Zhenda Xie, Haowei Zhang, Liang Zhao, et al. Janusflow: Harmonizing autoregression and rectified flow for unified multimodal understanding and generation. *arXiv preprint arXiv:2411.07975*, 2024.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pp. 311–318, 2002.
- Jeongeun Park, Sungjoon Choi, and Sangdoo Yun. A unified framework for motion reasoning and generation in human interaction, 2025. URL <https://arxiv.org/abs/2410.05628>.
- Mathis Petrovich, Michael J. Black, and Gül Varol. Action-conditioned 3D human motion synthesis with transformer VAE. In *International Conference on Computer Vision (ICCV)*, 2021.
- Mathis Petrovich, Michael J. Black, and Gül Varol. TEMOS: Generating diverse human motions from textual descriptions. In *European Conference on Computer Vision (ECCV)*, 2022.
- Mathis Petrovich, Michael J Black, and Gül Varol. Tmr: Text-to-motion retrieval using contrastive 3d human motion synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 9488–9497, 2023.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 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*, pp. 8748–8763. PMLR, 2021.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. URL <https://github.com/CompVis/latent-diffusionhttps://arxiv.org/abs/2112.10752>.
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pp. 234–241. Springer, 2015.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.
- Ozan Sener and Vladlen Koltun. Multi-task learning as multi-objective optimization. *Advances in neural information processing systems*, 31, 2018.
- Yonatan Shafir, Guy Tevet, Roy Kapon, and Amit H Bermano. Human motion diffusion as a generative prior. *arXiv preprint arXiv:2303.01418*, 2023.
- Weijia Shi, Xiaochuang Han, Chunting Zhou, Weixin Liang, Xi Victoria Lin, Luke Zettlemoyer, and Lili Yu. Lmfusion: Adapting pretrained language models for multimodal generation, 2025. URL <https://arxiv.org/abs/2412.15188>.
- Li Siyao, Weijiang Yu, Tianpei Gu, Chunze Lin, Quan Wang, Chen Qian, Chen Change Loy, and Ziwei Liu. Bailando: 3d dance generation by actor-critic gpt with choreographic memory. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11050–11059, 2022.
- Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020.
- Chameleon Team. Chameleon: Mixed-modal early-fusion foundation models. *arXiv preprint arXiv:2405.09818*, 2024.
- Guy Tevet, Brian Gordon, Amir Hertz, Amit H Bermano, and Daniel Cohen-Or. Motionclip: Exposing human motion generation to clip space. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXII*, pp. 358–374. Springer, 2022a.

- Guy Tevet, Sigal Raab, Brian Gordon, Yonatan Shafir, Amit H Bermano, and Daniel Cohen-Or. Human motion diffusion model. *arXiv preprint arXiv:2209.14916*, 2022b.
- Guy Tevet, Sigal Raab, Brian Gordon, Yonatan Shafir, Daniel Cohen-Or, and Amit H Bermano. Human motion diffusion model. *arXiv preprint arXiv:2209.14916*, 2022c.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- Maria Tsimpoukelli, Jacob L Menick, Serkan Cabi, SM Eslami, Oriol Vinyals, and Felix Hill. Multimodal few-shot learning with frozen language models. *Advances in Neural Information Processing Systems*, 34: 200–212, 2021.
- Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4566–4575, 2015.
- Wenhui Wang, Hangbo Bao, Li Dong, Johan Bjorck, Zhiliang Peng, Qiang Liu, Kriti Aggarwal, Owais Khan Mohammed, Saksham Singhal, Subhojit Som, et al. Image as a foreign language: Beit pretraining for vision and vision-language tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 19175–19186, 2023.
- Yuan Wang, Di Huang, Yaqi Zhang, Wanli Ouyang, Jile Jiao, Xuetao Feng, Yan Zhou, Pengfei Wan, Shixiang Tang, and Dan Xu. Motiongpt-2: A general-purpose motion-language model for motion generation and understanding. *arXiv preprint arXiv:2410.21747*, 2024.
- Yuqing Wang, Zhijie Lin, Yao Teng, Yuanzhi Zhu, Shuhuai Ren, Jiashi Feng, and Xihui Liu. Bridging continuous and discrete tokens for autoregressive visual generation. *arXiv preprint arXiv:2503.16430*, 2025a.
- Yuqing Wang, Zhijie Lin, Yao Teng, Yuanzhi Zhu, Shuhuai Ren, Jiashi Feng, and Xihui Liu. Bridging continuous and discrete tokens for autoregressive visual generation, 2025b. URL <https://arxiv.org/abs/2503.16430>.
- Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*, 2021.
- Bizhu Wu, Jinheng Xie, Keming Shen, Zhe Kong, Jianfeng Ren, Ruibin Bai, Rong Qu, and Linlin Shen. Mg-motionllm: A unified framework for motion comprehension and generation across multiple granularities. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 27849–27858, 2025.
- Chengyue Wu, Xiaokang Chen, Zhiyu Wu, Yiyang Ma, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu, Chong Ruan, et al. Janus: Decoupling visual encoding for unified multimodal understanding and generation. *arXiv preprint arXiv:2410.13848*, 2024a.
- Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. Next-gpt: Any-to-any multimodal llm. In *Forty-first International Conference on Machine Learning*, 2024b.
- Yiming Wu, Wei Ji, Kecheng Zheng, Zicheng Wang, and Dong Xu. Mote: Learning motion-text diffusion model for multiple generation tasks. *arXiv preprint arXiv:2411.19786*, 2024c.
- Jinheng Xie, Weijia Mao, Zechen Bai, David Junhao Zhang, Weihao Wang, Kevin Qinghong Lin, Yuchao Gu, Zhijie Chen, Zhenheng Yang, and Mike Zheng Shou. Show-o: One single transformer to unify multimodal understanding and generation. *arXiv preprint arXiv:2408.12528*, 2024.
- Chen Xin, Biao Jiang, Wen Liu, Zilong Huang, Bin Fu, Tao Chen, Jingyi Yu, 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 (CVPR)*, June 2023.
- Kangning Yin, Shihao Zou, Yuxuan Ge, and Zheng Tian. Tri-modal motion retrieval by learning a joint embedding space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1596–1605, June 2024.

- Fuming You, Minghui Fang, Li Tang, Rongjie Huang, Yongqi Wang, and Zhou Zhao. Momu-diffusion: On learning long-term motion-music synchronization and correspondence. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- Tianhe Yu, Saurabh Kumar, Abhishek Gupta, Sergey Levine, Karol Hausman, and Chelsea Finn. Gradient surgery for multi-task learning. *Advances in neural information processing systems*, 33:5824–5836, 2020.
- Hang Zhang, Xin Li, and Lidong Bing. Video-llama: An instruction-tuned audio-visual language model for video understanding. *arXiv preprint arXiv:2306.02858*, 2023a.
- Jianrong Zhang, Yangsong Zhang, Xiaodong Cun, Shaoli Huang, Yong Zhang, Hongwei Zhao, Hongtao Lu, and Xi Shen. T2m-gpt: Generating human motion from textual descriptions with discrete representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023b.
- 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*, pp. 364–373, 2023c.
- Mingyuan Zhang, Zhongang Cai, Liang Pan, Fangzhou Hong, Xinying Guo, Lei Yang, and Ziwei Liu. Motion-diffuse: Text-driven human motion generation with diffusion model. *IEEE transactions on pattern analysis and machine intelligence*, 46(6):4115–4128, 2024a.
- 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*, pp. 397–421. Springer, 2024b.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*, 2019.
- Yaqi Zhang, Di Huang, Bin Liu, Shixiang Tang, Yan Lu, Lu Chen, Lei Bai, Qi Chu, Nenghai Yu, and Wanli Ouyang. Motiongpt: Finetuned llms are general-purpose motion generators. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2024c.
- Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. Video instruction tuning with synthetic data. *arXiv preprint arXiv:2410.02713*, 2024d.
- Zeyu Zhang, Yiran Wang, Wei Mao, Danning Li, Rui Zhao, Biao Wu, Zirui Song, Bohan Zhuang, Ian Reid, and Richard Hartley. Motion anything: Any to motion generation. *arXiv preprint arXiv:2503.06955*, 2025.
- Chunting Zhou, Lili Yu, Arun Babu, Kushal Tirumala, Michihiro Yasunaga, Leonid Shamis, Jacob Kahn, Xuezhe Ma, Luke Zettlemoyer, and Omer Levy. Transfusion: Predict the next token and diffuse images with one multi-modal model. *arXiv preprint arXiv:2408.11039*, 2024.
- Yuwei Zhou, Xin Wang, Hong Chen, Xuguang Duan, and Wenwu Zhu. Intra-and inter-modal curriculum for multimodal learning. In *Proceedings of the 31st ACM International Conference on Multimedia*, pp. 3724–3735, 2023.