# OMG: Towards Open-vocabulary Motion Generation via Mixture of Controllers

Han Liang[1]    Jiacheng Bao[1]    Ruichi Zhang[1]    Sihan Ren[1]    Yuecheng Xu[1]

Sibei Yang[1]    Xin Chen[2]    Jingyi Yu[1]    Lan Xu[1]

[1]ShanghaiTech University    [2]Tencent PCG

## Abstract

*We have recently seen tremendous progress in realistic text-to-motion generation. Yet, the existing methods often fail or produce implausible motions with unseen text inputs, which limits the applications. In this paper, we present OMG, a novel framework, which enables compelling motion generation from zero-shot open-vocabulary text prompts. Our key idea is to carefully tailor the pretrain-then-finetune paradigm into the text-to-motion generation. At the pre-training stage, our model improves the generation ability by learning the rich out-of-domain inherent motion traits. To this end, we scale up a large unconditional diffusion model up to 1B parameters, so as to utilize the massive unlabeled motion data up to over 20M motion instances. At the subsequent fine-tuning stage, we introduce motion ControlNet, which incorporates text prompts as conditioning information, through a trainable copy of the pre-trained model and the proposed novel Mixture-of-Controllers (MoC) block. MoC block adaptively recognizes various ranges of the sub-motions with a cross-attention mechanism and processes them separately with the text-token-specific experts. Such a design effectively aligns the CLIP token embeddings of text prompts to various ranges of compact and expressive motion features. Extensive experiments demonstrate that our OMG achieves significant improvements over the state-of-the-art methods on zero-shot text-to-motion generation. Project page:* **https://tr3e.github.io/omg-page**.

## 1. Introduction

Generating vivid motions of human characters is a longstanding task in computer vision, with numerous applications in movies, games, robotics, and VR/AR. The text-to-motion setting aims to democratize motion generation for novices and has recently received substantive attention.

Benefited from the existing text-annotated motion datasets [18, 47], recent advances generate diverse motions from text prompts, equipped with various generative models like VAEs [56, 76], autoregressive models [32, 87], and diffusion models [11, 12, 77, 89]. However, these meth-



Figure 1. Our Open-vocabulary Motion Generation (OMG) approach is capable of generating high-quality motions in response to unseen text prompts.

ods heavily rely on the paired text-motion data with limited text annotations, and hence fall short of generalizing to unseen open-vocabulary text inputs. Imagine generating "*Gollum breakdance footwork*", it requires out-of-domain generation ability with rich human knowledge to understand the motion traits in various words like characters (Gollum, ninja, etc.) or skills (breakdance, spinkick, etc.). To tackle the open-vocabulary texts, recent works [8, 27, 41] utilize the zero-shot text-image alignment ability of the large-scale pre-trained models, e.g., CLIP [57]. They adopt text-pose alignments and hence remove the reliance on pair-wise text and continuous motions. Yet, the discrete poses are insufficient to represent the rich motion traits, leading to unrealistic motion generation. On the other hand, recent advances [63, 88] enable impressive and controllable text-to-image generation, even from open-vocabulary text inputs. Revisiting their huge success, both the pretrain-then-finetune paradigm as well as scaling up the model have played an important role. Yet, adopting these strategies for

text-to-motion generation is challenging. First, there exists a huge imbalance between the quantities of unpaired motion data and paired text-motion data. Second, for open-vocabulary texts, various tokens correspond to various motion traits, constituting a complex many-to-many problem.

To tackle the challenges, we propose *OMG* – a novel scheme to generate motions of human characters from open-vocabulary texts (see Fig. 1). Our key idea is to carefully tailor the pretrain-then-finetune paradigm into text-to-motion generation. For pre-training, we adopt a minimalistic design to scale up a large unconditional diffusion model, so as to utilize the massive unlabeled motion data. For fine-tuning, we first freeze the pre-trained large model. We then adopt the trainable copy and treat the text prompts as extra conditions through a novel design, named mixture-of-controllers, to learn to predict the conditional residuals. It adaptively fuses the motion traits corresponding to various tokens from the input text prompts, so as to handle the ambiguity between the text and motion modalities.

Specifically, at the pre-training stage, we adopt the diffusion transformer [52] as the backbone with over-parameterized motion representation. We then scale up the models with parameters ranging from 88M to 1B, leveraging over 20M unlabeled motion instances from a diverse range of 13 publicly available datasets. During training, we also adopt a sliding random window strategy to crop various motion sequences, to improve the generation ability for arbitrary lengths of motion. To this end, the pre-trained model learns the rich inherent motion features with a large motion manifold to ensure the realism of the generated motions. For the subsequent stage, to incorporate the text prompts as conditioning information, we adopt a fine-tuning scheme called motion ControlNet, analogous to the famous ControlNet [88] for the text-image generation task. It includes a trainable copy of the large-scale pre-trained model, which serves as a strong backbone to retain the motion features, as well as a novel block called Mixture-of-Controllers (MoC). Such MoC block can effectively inject residual information into the pre-trained model, based on the motion features and the CLIP-based text embeddings. Our key design in the MoC block is the cross-attention mechanism between text and motion, as well as the text-token-specific Mixture-of-Experts [15, 67], which are to be detailed in later sections. Such a design effectively aligns the token embeddings of text prompts to various ranges of compact and expressive motion features, which are warmed up from a powerful pre-trained model. As a result, our OMG approach achieves compelling generation from open-vocabulary texts, significantly outperforming prior arts. In particular, only fine-tuned on the HumanML3D dataset [18], it can generate vivid motions of various human characters with diverse motion skills from the Mixamo dataset [1], as shown Fig. 1.

To summarize, our main contributions include:

- We propose a text-to-motion generation approach with a pretrain-then-finetune paradigm to scale up both data and model, achieving state-of-the-art zero-shot performance.
- We experimentally demonstrate that pre-training on large-scale unlabeled motion data improves the generation results from diverse and open-vocabulary texts.
- We propose a fine-tuning scheme for text conditioning, utilizing a mixture of controllers to effectively improve the alignment between text and motion.

## 2. Related Work

**Conditional Motion Synthesis.** Being able to generate realistic and contextually relevant motion sequences based on various types of conditions, conditional motion synthesis has received increasing attention in the field of motion generation. Common types of conditions include text [3, 4, 11, 19, 33, 40, 43, 56, 66, 76, 84], action [17, 55], music [2, 34, 36, 38], speech [5, 6, 20], control signals [54, 71, 73], action labels [17, 55, 81], incomplete motion sequences [13, 21], images [10, 24, 39, 60, 61] and scene [91]. The advent of Diffusion Models [25, 70] has given a big boost to text-driven motion synthesis. Kim et al. [33] develops a transformer-based diffusion model which could generate and edit motions well aligned with the given text. Motion Diffusion Model (MDM) [77] combines insights already well established in the motion generation domain and in each diffusion step predicts the sample rather than the noise. Motion Latent-based Diffusion model (MLD) [11] performs a diffusion process on the motion latent space. Besides the diffusion model, T2M-GPT [87] investigates a framework that combines VQ-VAE [79] and autoregressive model. However, since these models are trained only based on paired text-motion datasets such as HumanML3D [19], they cannot well handle unseen text prompts.

**Open-vocabulary Generation.** Instead of relying on pre-existing data, zero-shot text-driven generation leverages general knowledge learned during training to create novel content from text prompts. Reed et al. [59] uses DC-GAN architecture to steer for zero-shot text-to-image synthesis. CLIP [57], pre-trained on four hundred million image-text pairs, possesses the remarkable ability to effectively comprehend and generate meaningful associations between images and text. With CLIP's strong ability, many works are able to generate high-quality zero-shot text-driven images [16, 51, 82] or 3D objects [30, 31, 49, 53, 65, 80]. In the motion synthesis field, some works have attempted to investigate open-vocabulary text-to-motion generation and achieved good results. MotionCLIP [76] utilizes a decoder to decode the CLIP motion embeddings. AvatarCLIP [27] synthesizes a key pose and then retrieves the closest motion from the database. Lin et al. [41] pre-train a motion generator that learns to reconstruct the full motion from the
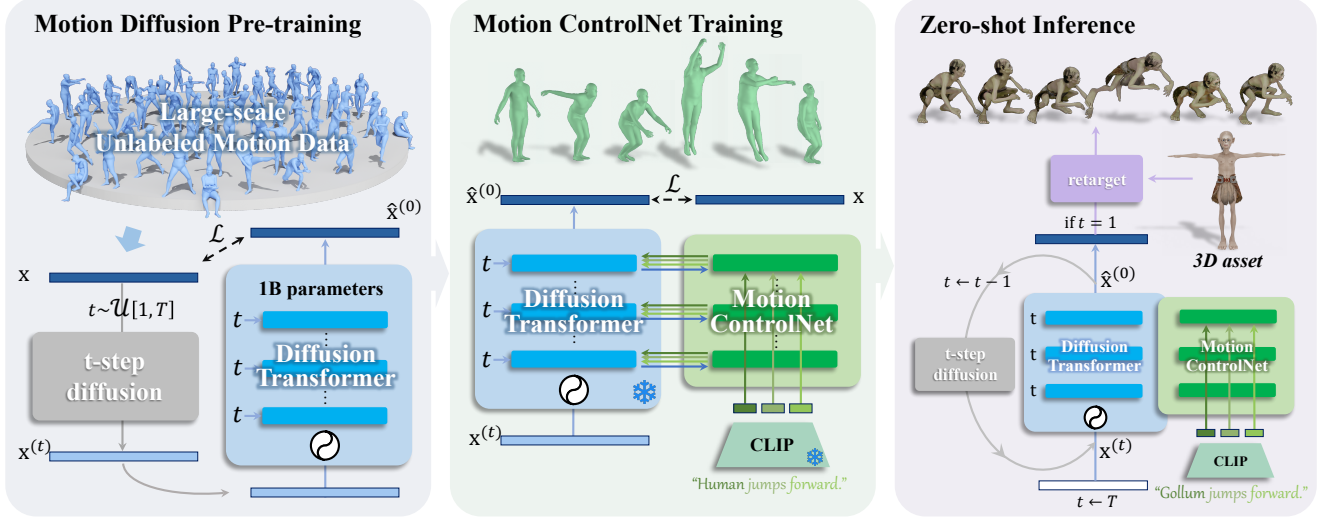
**Motion Diffusion Pre-training**

**Motion ControlNet Training**

**Zero-shot Inference**

Figure 2. **Method overview**. We train our OMG model in two stages. First, we leverage large-scale unlabeled motion data to pre-train an unconditional diffusion model with up to 1B parameters (Sec. 3.1). Then, we adopt a conditional fine-tuning scheme called motion ControlNet to condition the pre-trained diffusion model on text prompts (Sec. 3.2). During inference, the pre-trained unconditional denoiser and the fine-tuned conditional denoiser are combined with classifier-free guidance, generating realistic motions with zero-shot text inputs.

masked motion with the key pose. Make-An-Animation [8] pre-trains generative model with diverse in-the-wild text-pose pair. However, since these approaches are based on text-pose alignment, the generated motion is deficient in realism due to a lack of global rotation and translation. Recently, MotionGPT [32] pre-trains and fine-tunes the large language model with tokenized motion and text data. However, it still struggles to generate novel motion from open-vocabulary text prompts.

**Mixture-of-Experts.** In the deep learning field, Mixture-of-Experts (MoE) [15, 29, 67] in a neural network architecture divides specific model parameters into distinct groups, referred to as "experts". Eigen et al. [14] use a different gating network at each layer in a multi-layer network, increasing the number of effective experts by introducing an exponential number of paths through different combinations of experts at each layer. Other research that regards the entire model as an expert [37] also achieves good results. Moreover, endeavors are underway to enhance the training and inference speed within the MoE framework [23, 35]. The MoE paradigm is applicable not only to language processing tasks but also extends to visual models [62], Mixture-of-Modality-Experts in multi-modal transformers [68], as well as motion synthesis. Holden et al. [26] develop phase-functioned experts blended by pre-defined blending weights to control character at a specific phase, and Starke et al. [72, 73], Zhang et al. [85] uses a gating network to predict the blending weights, achieving impressive results of character control. Inspired by them, we devise text-token-specific experts that control the sub-motion for each text

token.

## 3. Methods

We aim to generate realistic and diverse human motions that are conditioned on text prompts, which capture complex and abstract motion characteristics with zero-shot open-vocabulary descriptions. To this end, we adopt pretrain-then-finetune paradigm to enhance the capability of our model. For the unconditional denoiser (Fig. 2 left), we leverage a large-scale unlabeled motion dataset for motion diffusion pre-training and scale the model size up (Sec. 3.1). For the conditional denoiser (Fig. 2 middle), we devise a specific conditional fine-tuning scheme called motion ControlNet, including a novel conditioning design called Mixture-of-Controllers, to exploit the zero-shot capability of the CLIP text embeddings (Sec. 3.2). During inference (Fig. 2 right), we further use classifier-free guidance to combine the unconditional denoiser and the conditional one.

### 3.1. Motion Diffusion Pre-training

**Large-scale Unlabeled Motion Dataset.** To facilitate large-scale pre-training in the motion domain, we collected a large amount of high-quality human motion data totaling over 20M frames from various sources, such as character animation datasets, marker-based, IMU-based, and multi-view markerless mocap datasets. We then standardize the frame rate and re-target them to the unified parametric skeleton of SMPL body model [45] using off-the-shelf re-targeting algorithms [86]. After that, following the previous successful works [18, 77] on motion generation, we

3

| Model Name | $n_{layers}$ | $d_{model}$ | $n_{heads}$ | $d_{head}$ | $n_{params}$ |
|---|---|---|---|---|---|
| OMG-Base | 8 | 1024 | 8 | 128 | 88M |
| OMG-Large | 12 | 1280 | 10 | 128 | 201M |
| OMG-Huge | 16 | 1664 | 13 | 128 | 405M |
| OMG-Giant | 24 | 2048 | 16 | 128 | 1B |

Table 1. Sizes and architectures of our 4 models.

enrich the SMPL skeletons with an over-parameterized motion representation that facilitates network learning.

**Model Scaling Up.** We design our pre-trained diffusion model with minimalism, only preserving the essential modules that are scalable and efficient to adapt to expanding data. To this end, we adopt the Diffusion Transformer (DiT) [52] architecture as our network backbone, due to its scalability and impressive performance with increasing input tokens and training data. The only difference is that we use rotary positional embedding [74] to encourage the model to capture the relative temporal relations among the frames. To study the dependence of performance on model size, we train 4 different sizes of our OMG model, ranging from 88 million parameters to 1 billion parameters. Tab. 1 shows the configuration details of our 4 models. Here $n_{layers}$ is the total number of layers, $d_{model}$ is the number of units in each bottleneck layer (we always have the feed-forward layer two times the size of the bottleneck layer, $d_{ff} = 2d_{model}$), and $d_{head}$ is the dimension of each attention head, and $n_{params}$ is the total number of trainable parameters.

**Training.** We train our unconditional denoiser $\mathcal{D}_u$ to predict the clean motion $\hat{\mathbf{x}}^{(0)}$ given diffusion timestep $t$ with the simple objective:

$$\mathcal{L}_{simple} = \mathbb{E}_{\mathbf{x},t,\epsilon}[\lambda_t||\mathbf{x} - \mathcal{D}_u(\mathbf{x}^{(t)}, t)||_2^2], \quad (1)$$

where $\mathbf{x}^{(t)} = \mathbf{x} + \sigma_t\epsilon$ is $t$-step noised motion, noise $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, $\lambda_t$ is the loss weighting factor. Furthermore, we introduce geometric losses, analogous to MDM [77], that contain velocity loss $\mathcal{L}_{vel}$, and foot contact loss $\mathcal{L}_{foot}$ to enforce physical properties and anti-sliding. Overall, the unconditional denoiser $\mathcal{D}_u$ is thus trained with the following objective:

$$\mathcal{L} = \mathcal{L}_{simple} + \lambda_{vel}\mathcal{L}_{vel} + \lambda_{foot}\mathcal{L}_{foot}. \quad (2)$$

In our experiments, $\lambda_{vel} = 30$ and $\lambda_{foot} = 30$. The number of diffusion timesteps $T$ is set to 1,000 with cosine noise level schedule [50]. We train all sizes of the model for 1M iterations with a batch size of 256 using the AdamW [46] optimizer with a weight decay of 0, a maximum learning rate of $10^{-4}$, and a cosine LR schedule with 10K linear warmup steps. The models are trained using PyTorch with ZeRO [58] memory redundancy strategy on 8
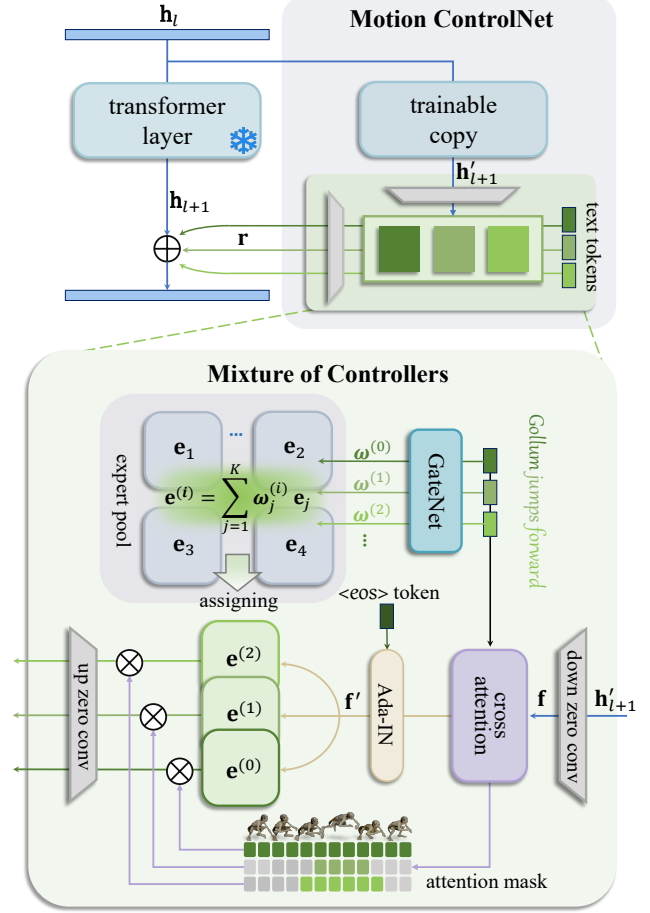


Figure 3. **Motion ControlNet (top)** freezes the parameters of the pre-trained transformer layer and combines a trainable copy of the layer with the proposed **Mixture-of-Controllers (bottom)** block. The MoC block first fuses the text features and motion features and simultaneously determines the sub-motion ranges for each text token with the cross-attention mechanism. Then it performs fine-grained control of sub-motions with text-token-specific experts.

NVIDIA A100 GPUs, and the largest model takes 1 week to train.

**Sliding Random Windows.**

During training, we propose sliding random windows that iterate over each motion state frame $x_i$ in the motion dataset to sample motion sequences starting from $x_i$. Specifically, for the $i^{th}$ frame, we randomly crop a subsequence $\mathbf{x}$ of length $l$ with a random window $[i, i + l)$, where $l \sim \mathcal{U}[1, L]$ is a uniform variable, and $L$ is the hyper-parameter that determines the maximum length of the generated motion sequence of our model. This pushes the model to learn the relations among both temporal frames and spatial features of a single motion state to process arbitrary lengths of motions even a single keyframe, to facilitate the following fine-grained control of sub-motions of arbitrary lengths. In our experiments, for balancing resolu-
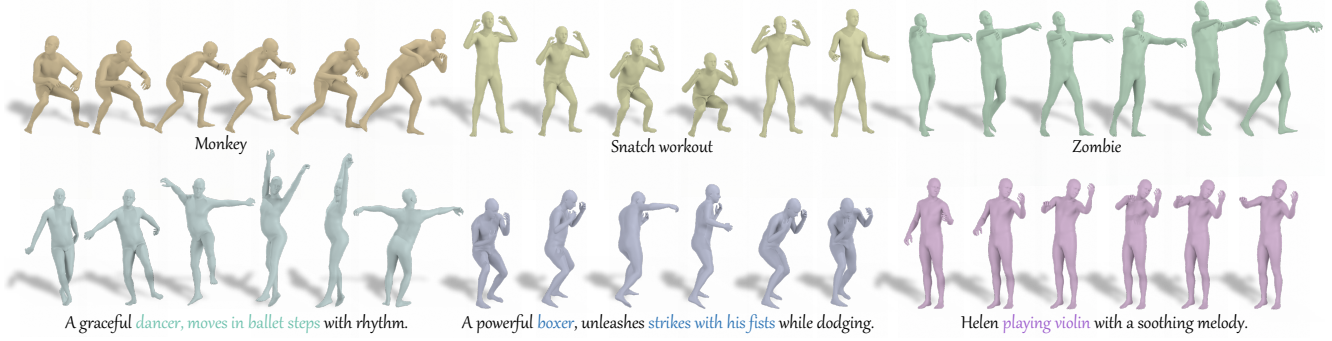
Figure 4. **Qualitative results** generated by our model given various text prompts. Our model effectively captures the motion characteristics from either a single phrase or longer natural sentences.

tion and duration, the maximum length $L$ of sliding random windows is set to 300 with the motion frame rate 30.

## 3.2. Motion ControlNet

To incorporate the text prompts $\mathbf{c} = \{\mathbf{c}_i\}_{i=1}^n$ into the pre-trained motion model, We adopt a fine-tuning scheme we call motion ControlNet to condition on text prompts, inspired by image ControlNet [88], including a trainable copy of original transformer layers and the proposed conditioning blocks called Mixture-of-Controllers (MoC).

As illustrated in Fig. 3 top, we freeze the parameters of the pre-trained transformer layer and combine a trainable copy of the layer with a MoC block that injects $n$ residuals corresponding to $n$ text tokens into the output of the original transformer layer. Specifically, the frozen parameters retain the pre-trained denoising ability, while the trainable copy reuses the large-scale pre-trained model as a strong initial backbone to learn to extract semantic motion features from intermediate representation $\mathbf{h}_l$ of the last layer. And based on that, the MoC block is employed to predict the $n$ conditional residuals $\mathbf{r}$ corresponding to $n$ text tokens, all of which are then added to the original output $\mathbf{h}_{l+1}$. To do so, a pre-trained CLIP text encoder $\mathcal{E}(\mathbf{c})$ is employed to extract the text token embeddings.

**Mixture-of-Controllers.** In the Mixture-of-Controllers block (Fig. 3 bottom), our key idea is to control the sub-motion sequences separately with different expert controllers in an MoE fashion [15, 67], so as to better align them to the corresponding text token embeddings in CLIP space. For instance, for the given text prompt *"Gollum jumps forward"*, the *"Gollum"* token corresponds to the entire sequence while *"jumps"* corresponds only to the sub-sequence when he is jumping.

To this end, we introduce two cooperative designs. The first one is a **cross-attention mechanism** to fuse the text features and motion features and simultaneously determine the sub-motion ranges for each text token. The second one is the **text-token-specific experts** selected from an expert pool to perform fine-grained control of sub-motions.

For the **cross-attention mechanism**, a sequence of motion features $\mathbf{f} \in \mathbb{R}^{l \times d_m}$ is input into a cross-attention layer along with text embeddings $\mathcal{E}(\mathbf{c}) \in \mathbb{R}^{n \times d_c}$, where $l$ denotes the length of the sequence, $d_m$ represents the dimension of the motion features, $n$ denotes the number of the text tokens, and $d_c$ represents the dimension of text embedding. The cross-attention layer projects $\mathbf{f}$ to a query matrix $\mathbf{Q} = \mathbf{f} \cdot \mathbf{W}_m^Q$, where $W_m^Q \in \mathbb{R}^{d_m \times d_m}$ is the motion projection matrix. And $\mathcal{E}(\mathbf{c})$ is projected to a key matrix $\mathbf{K} = \mathcal{E}(\mathbf{c})\mathbf{W}_c^K$ and a value matrix $\mathbf{V} = \mathcal{E}(\mathbf{c})\mathbf{W}_c^V$, where $W_c^K, W_c^V \in \mathbb{R}^{d_c \times d_m}$ are the text projection matrices. Then the text values $\mathbf{V}$ are distributed into the motion sequence with the attention mechanism:

$$\mathbf{f}' = \mathbf{f} + \text{softmax}(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}})\mathbf{V}. \tag{3}$$

Simultaneously an attention map between text tokens and motion sequence is produced, denoted by $\mathbf{A} = \text{softmax}(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}}) \in \mathbb{R}^{l \times n}$. Within this attention map, the $i^{th}$ column $\mathbf{A}_{*,i}$ indicates a correspondence between a sub-motion sequence and the $i^{th}$ text token. Besides, we employ the adaptive instance normalization (Ada-IN) conditioning on $\langle eos \rangle$ (end of sequence) token to normalize $\mathbf{f}'$ before sending them to the experts. Since CLIP $\langle eos \rangle$ token summarizes the entire text, we intuitively use it to unify the distribution of the entire motion sequence.

For the **text-token-specific experts**, we define an expert controller corresponding to $i^{th}$ text token $\mathcal{E}(\mathbf{c}_i)$ as a two-layer feed-forward network (denoted by $\mathcal{F}$) with parameters $\mathbf{e}^{(i)} = \{\mathbf{W}_0 \in \mathbb{R}^{d_m \times 2d_m}, \mathbf{W}_1 \in \mathbb{R}^{2d_m \times d_m}, \mathbf{b}_0 \in \mathbb{R}^{2d_m}, \mathbf{b}_1 \in \mathbb{R}^{d_m}\}$. Furthermore, the parameters $\mathbf{e}^{(i)}$ are generated by blending $K$ expert parameters $\{\mathbf{e}_1, ..., \mathbf{e}_K\}$ in an expert pool, each of which is in a form of the same parameter configuration:

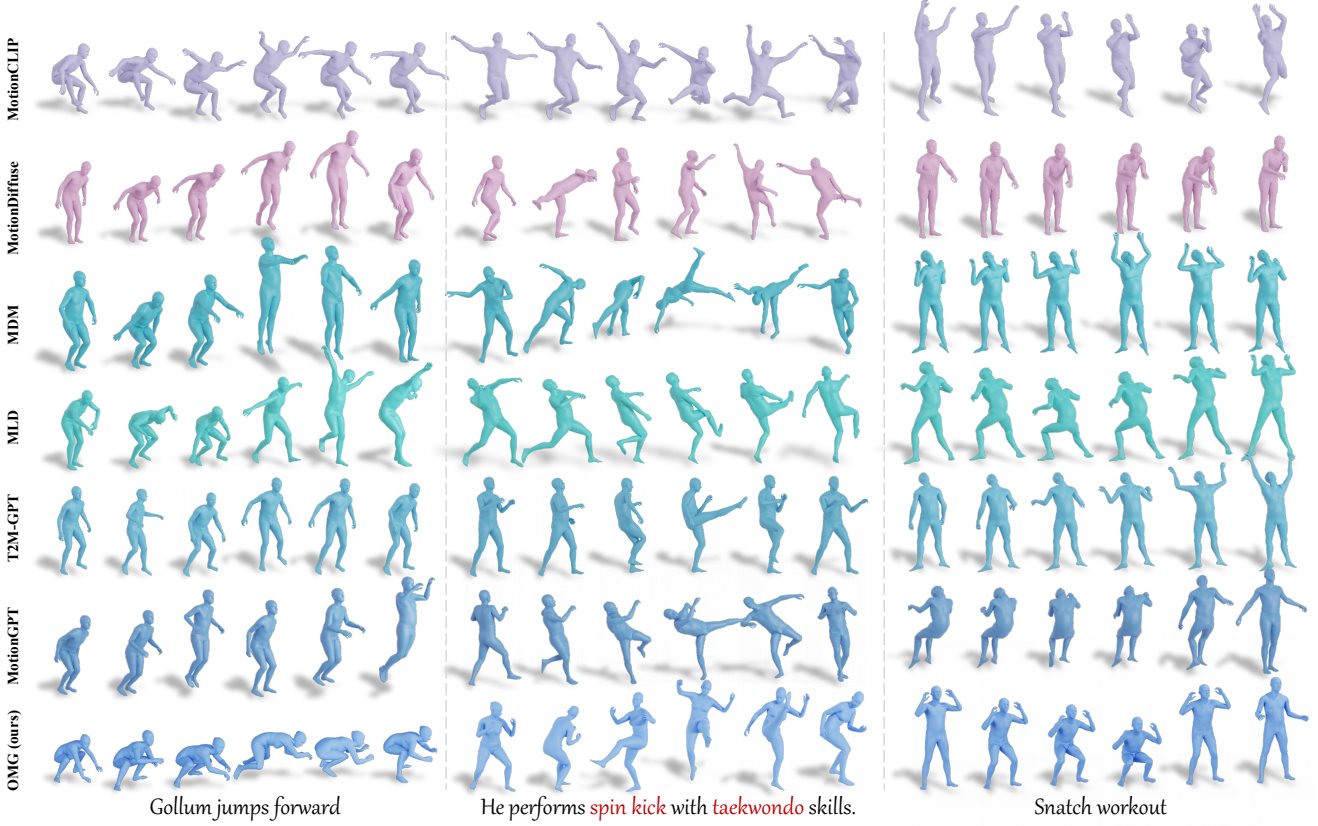$$\mathbf{e}^{(i)} = \sum_j^K \boldsymbol{\omega}_j^{(i)} \mathbf{e}_j, \tag{4}$$

Figure 5. **Qualitative comparison.** Our method can generate high-quality human motions that better align with text prompts than previous state-of-the-art methods.

where $j$ indicates the index of expert pool, $K$ is expert pool size that can be adjusted, and $\boldsymbol{\omega}^{(i)} = \{\boldsymbol{\omega}_j^{(i)}\}_{j=1}^K \in \mathbb{R}^K$ is the blending weights vector controlled by the text token $\mathbf{c}_i$ with a three-layer fully-connected gating network $\mathcal{G}$:

$$\boldsymbol{\omega}^{(i)} = \text{softmax}(\mathcal{G}(\mathcal{E}(\mathbf{c}_i))). \tag{5}$$

Then, the $n$ experts take motion features $\mathbf{f}'$ as input and output $n$ conditional masked residuals $\mathbf{r} = \{\mathbf{r}_i\}_{i=1}^n$:

$$\mathbf{r}_i = \mathbf{M}_{*,i} \circ \mathcal{F}(\mathbf{f}'|\mathbf{e}^{(i)}), \tag{6}$$

where $\mathbf{M}_{*,i} = \text{sigmoid}(\gamma(\mathbf{A}_{*,i} - \beta \max(\mathbf{A}_{*,i})))$ is an attention mask ranging from 0 to 1, and $\gamma$ and $\beta$ are hyperparameters to control the sharpness and threshold respectively.

Additionally, a down-projection and up-projection 1-D convolution pair is used to process input and output respectively to reduce the control latent dimension and thus the trainable parameters. Moreover, similar to image ControlNet [88], we use zero-initialized convolution parameters to protect the trainable copy from the harmful gradient noises in the initial training steps.

**Training and Inference.** The conditional denoiser $\mathcal{D}_c$ is also trained using the same objective $\mathcal{L}$ (Eq. (2)) with trainable parameters in motion ControlNet. In our experiments, we employ the frozen *CLIP-VIT-L/14* [57] text encoder to extract text embeddings at the final layer normalization with dimension $d_c = 768$ and we truncate the text tokens with the maximum token number 77. We set the down-projection latent dimension $d_m$ of MoC blocks to 256, the expert pool size $K = 12$, and attention mask sharpness $\gamma = 24$ and threshold $\beta = 0.25$. We train all the motion ControlNets for 500 epochs with a batch size of 64 using the AdamW optimizer with a weight decay of $1e - 5$, a maximum learning rate of $3 \times 10^{-5}$, and a cosine LR schedule with 1K linear warmup steps. In the training process, we randomly replace $\langle eos \rangle$ token with the empty token with 50% probability to increase the ability to capture global semantics in text tokens as a replacement.

During inference, classifier-free guidance is used to combine unconditional denoiser and conditional denoiser:

$$\hat{\mathbf{x}}^{(0)} = (1 - s) \cdot \mathcal{D}_u(\mathbf{x}^{(t)}, t) + s \cdot \mathcal{D}_c(\mathbf{x}^{(t)}, t, \mathbf{c}). \tag{7}$$

In our experiments, DDIM [69] sampling strategy with 200 timesteps is employed and the guidance strength $s = 4.5$.

| Methods | HumanML3D [19] | | | Mixamo [1] (zero-shot) | | |
|---|---|---|---|---|---|---|
| | FID↓ | R-Precision↑ | Diversity → | FID↓ | CLIP-score↑ | Diversity → |
| Real | $0.002^{\pm.000}$ | $0.797^{\pm.002}$ | $9.503^{\pm.065}$ | $0.106^{\pm.003}$ | $0.648^{\pm.001}$ | $2.665^{\pm.022}$ |
| MotionCLIP [76] | - | - | - | $2.542^{\pm.012}$ | $0.511^{\pm.004}$ | $2.205^{\pm.012}$ |
| MAA [8] | $0.774^{\pm.002}$ | $0.676^{\pm.001}$ | $8.230^{\pm.064}$ | - | - | - |
| MotionDiffuse [89] | $0.630^{\pm.001}$ | $\underline{0.782}^{\pm.001}$ | $9.410^{\pm.049}$ | $2.363^{\pm.010}$ | $0.505^{\pm.002}$ | $2.411^{\pm.016}$ |
| MDM [77] | $0.544^{\pm.044}$ | $0.611^{\pm.007}$ | $\underline{9.559}^{\pm.086}$ | $1.297^{\pm.004}$ | $0.536^{\pm.004}$ | $\underline{2.594}^{\pm.011}$ |
| MLD [11] | $0.473^{\pm.013}$ | $0.772^{\pm.002}$ | $9.724^{\pm.082}$ | $\underline{1.229}^{\pm.004}$ | $\underline{0.556}^{\pm.003}$ | $2.583^{\pm.018}$ |
| T2M-GPT [87] | $\mathbf{0.116}^{\pm.004}$ | $0.775^{\pm.002}$ | $9.844^{\pm.095}$ | $1.420^{\pm.003}$ | $0.541^{\pm.002}$ | $2.590^{\pm.022}$ |
| MotionGPT [32] | $\underline{0.232}^{\pm.008}$ | $0.778^{\pm.002}$ | $\mathbf{9.528}^{\pm.071}$ | $1.365^{\pm.003}$ | $0.552^{\pm.002}$ | $2.589^{\pm.018}$ |
| OMG (ours) | $0.381^{\pm.008}$ | $\mathbf{0.784}^{\pm.002}$ | $9.657^{\pm.085}$ | $\mathbf{1.164}^{\pm.009}$ | $\mathbf{0.588}^{\pm.002}$ | $\mathbf{2.632}^{\pm.021}$ |

Table 2. Comparison of text-to-motion generation on HumanML3D [19] and Mixamo [1] test set. We ran all the evaluations 20 times, with the average reported alongside a 95% confidence interval. The right arrow → means the closer to real motion the better. **Bold** and underline indicate the best and the second best result. The term (Zero-shot) implies that the dataset contains unseen open-vocabulary texts.



Figure 6. Quantitative evaluation on pre-training, model size, and expert pool size. (a) Models *w_pre-training* show consistently improved performance over *w/o_pre-training*, and *w_pre-training* models, which benefit from large-scale motion data, improve with increasing model size. (b) Larger expert pool sizes improve the performance.

# 4. Experiments

In this section, we first show our qualitative results in Fig. 4. Our method enables fine-grained control of complex and abstract motion trait descriptions. Then, we introduce dataset settings, and evaluation metrics (Sec. 4.1). We further compare our approach with various state-of-the-art methods on both in-domain and out-domain datasets (Sec. 4.2) and conduct extensive ablation studies on our model architecture (Sec. 4.3). We refer the reader to appreciate more qualitative results and details provided in supplements.

## 4.1. Datasets and Evaluation Metrics

**Training Datasets.** At the pre-training stage, we utilize various publicly available human motion datasets, such as artist-created datasets [21, 48], marker-based [7, 28, 44, 47, 75], IMU-based [39, 78] and multi-view markerless [9, 38, 40, 90] motion capture datasets, totaling over 20 million frames. In the subsequent conditional fine-tuning stage, we train our motion ControlNet using the text-motion HumanML3D [19] dataset, for fair comparisons with previous methods.

**Evaluation Datasets.** We test on two benchmarks, the HumanML3D [19] and the Mixamo [1] test set. The HumanML3D test set evaluates the in-domain performance. And the Mixamo dataset consists of abundant artist-created animations and human-annotated descriptions, offering a wide variety of diverse and dynamic motions, utilized to compare the zero-shot performance across the domains.

**Evaluation Metrics** are summarized as four parts. (1) Frechet Inception Distance (FID) is our principal metric to evaluate the feature distributions between the generated and real motions. (2) Motion-retrieval precision (R-Precision) calculates the text and motion Top 3 matching accuracy under feature space. (3) CLIP-score. Borrowing from text-to-image synthesis [64], we use CLIP-score to evaluate zero-
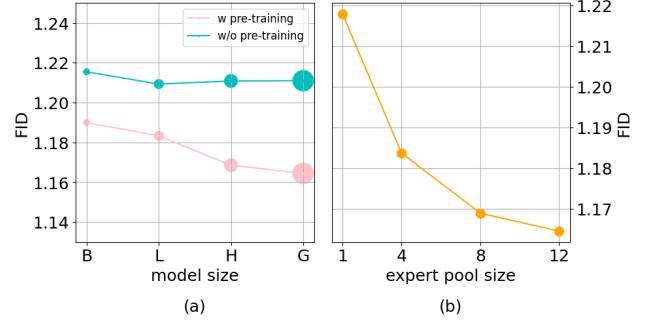
shot text-motion consistency by measuring cosine similarity in CLIP space. To do so, we train a motion encoder using the training set of HumanML3D and Mixamo to extract motion features that are aligned to text embeddings in CLIP space as same as [76]. (4) Diversity is assessed by calculating variance through features.

## 4.2. Comparison

We compare our approach with various state-of-the-art methods. Specifically, we apply conditional motion diffusion models, including MDM [77], MLD [11], MotionDiffuse [89], VAE-based model MotionCLIP [76], and autoregressive model T2M-GPT [87]. Besides, we also apply MAA [8] based on the text-pose alignment method, and MotionGPT [32] utilizing motion-language pre-training. Some results on HumanML3D are borrowed from their own paper.

The quantitative results are presented in Tab. 2. Our method demonstrates the best text-to-motion alignment (R-Precision) in the in-domain evaluation. Moreover, among all diffusion-based methods, our model achieves the best FID. In terms of zero-shot performance, our approach achieves the best FID and CLIP-score, outperforming previous methods. This suggests a superior capability for high-quality motion generation and effective matching with zero-shot text prompts. As depicted in Fig. 5 of the qualitative results, our model demonstrates the capability to generate motions that show more realism, and better alignment with each motion characteristic description, whether in the sentence or phrase.

## 4.3. Ablation Study

To examine the specific contributions of our novel OMG model architecture, we conduct a series of ablation studies focusing on the roles of pre-training, model scale, and

| Methods | FID ↓ | CLIP-score ↑ | Diversity → |
|---|---|---|---|
| Real | $0.106^{\pm.003}$ | $0.648^{\pm.001}$ | $2.665^{\pm.022}$ |
| (1) Cross-attn + FFN | $1.252^{\pm.006}$ | $0.552^{\pm.003}$ | $2.576^{\pm.025}$ |
| (2) w/o Zero Conv | $1.339^{\pm.005}$ | $0.535^{\pm.004}$ | $2.695^{\pm.014}$ |
| (3) w/o Attention Mask | $1.246^{\pm.008}$ | $0.557^{\pm.003}$ | $\mathbf{2.647}^{\pm.026}$ |
| ours (complete) | $\mathbf{1.164}^{\pm.009}$ | $\mathbf{0.588}^{\pm.002}$ | $2.632^{\pm.021}$ |

Table 3. Quantitative evaluation on MoC block. The damping performance of the three variants of our model highlights the effectiveness of our MoC block technical designs.
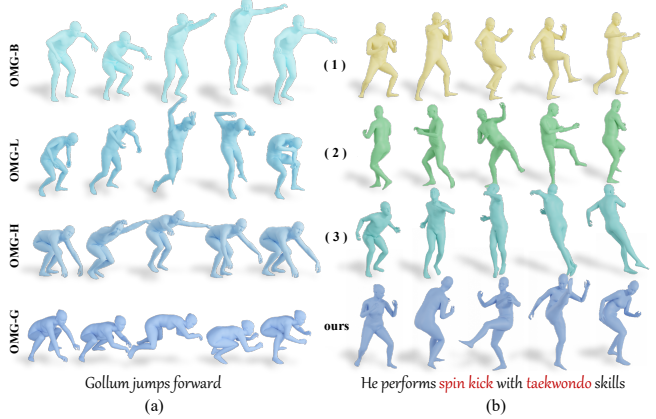


Figure 7. Qualitative evaluation on model sizes (a) and MoC block (b). Models with larger sizes effectively comprehend richer out-of-domain motion features to present better motion expressions. Besides, our technical designs effectively improve the alignment with the input texts.

multiple expert controllers. Additionally, we undertake an in-depth analysis of our MoC Block. These evaluations utilize the out-of-domain dataset Mixamo [1], providing a robust testbed to ascertain the effectiveness of the technical designs on the zero-shot performance of our model.

**Effect of Pre-training and Model Scale.** Our investigation into the impact of pre-training involves training several variant models without pre-training (*w/o_pre-training*) across four distinct scales (see Tab. 1). These models are then compared against counterparts with the pre-training process (*w_pre-training*). As illustrated in Fig. 6a, *w_pre-training* consistently outperforms *w/o_pre-training*, achieving lower FID across all sizes. This indicates that our model's performance on zero-shot motion generation is effectively enhanced through pre-training. Meanwhile, we observe that the performance of *w_pre-training* improves with increasing model size. Qualitative analysis, as shown in Fig. 7a, reveals that the *"Gollum"* characteristic of the generated motions becomes more pronounced with larger models. This suggests that larger model sizes enhance the overall quality and alignment of the generated motions.

**Effect of Expert Pool Size.** To systematically analyze the impact of multiple experts, we train four distinct variants of the OMG model, sweeping over different expert pool sizes ($K = 1, 4, 8$, and 12 respectively). The results, as depicted in Fig. 6b, demonstrate a discernible decrease in FID as expert pool size increases. This trend suggests that the larger expert hypothesis space enhances the experts' ability to align sub-motions with their respective text embeddings.

**Evaluation of the MoC Block.** In order to assess the effectiveness of our proposed MoC block, we explore its performance through several variant configurations: **(1)** Cross-attn + FFN, where the architecture is pruned to include only a cross-attention layer and a feed-forward network; **(2)** w/o Zero Conv, substituting zero convolutions with standard convolutional layers initialized using Gaussian distributions; and **(3)** w/o Attention Mask, omitting the multiplication of attention masks with the output of the expert layer. These variant models are rigorously compared against our original OMG model. Quantitative results, as shown in Tab. 3, demonstrate that the exclusion of specific

components within the MoC Block leads to a worse result in both the FID and CLIP-score. This observation underscores the integral contributions of these components to the motion generation process. Qualitative analysis, presented in Fig. 7b, indicates noticeable deficiencies in the motion generation of all three variants: variant **(1)** is not able to clearly depict *"spin"* motions; variant **(2)** exhibits less expressive motion features; and variant **(3)** inadequately captures the *"kick"* motion. In contrast, our OMG model demonstrates a superior ability to more accurately align motion with the text prompt.

## 5. Conclusion

In this paper, we present a novel text-to-motion generation framework, OMG, that combines the advantages of conditional generative models and text-pose alignment methods. It carefully tailors pretrain-then-finetune paradigm into text-to-motion generation. The pre-training stage leverages a large amount of unlabeled motion data to train a powerful unconditional diffusion model that ensures the realism and diversity of the generated motions. The fine-tuning stage introduces motion ControlNet including the proposed novel text conditioning block called Mixture-of-Controllers. With the cross-attention mechanism and text-token-specific experts, it adaptively aligns the sub-motion features to the text embeddings in the CLIP space in an MoE fashion. The extensive experiments demonstrate that OMG achieves state-of-the-art zero-shot performance on text-to-motion generation. We believe it is a significant step towards open-vocabulary motion generation of human characters, with wide potential applications in movies, games, robotics, and VR/AR.

# References

[1] Adobe. https://www.mixamo.com/. Accessed:2023-9-28. 2, 7, 8, 13

[2] Gunjan Aggarwal and Devi Parikh. Dance2music: Automatic dance-driven music generation. *arXiv preprint arXiv:2107.06252*, 2021. 2

[3] Hyemin Ahn, Timothy Ha, Yunho Choi, Hwiyeon Yoo, and Songhwai Oh. Text2action: Generative adversarial synthesis from language to action. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5915–5920. IEEE, 2018. 2

[4] Chaitanya Ahuja and Louis-Philippe Morency. Language2pose: Natural language grounded pose forecasting. In *2019 International Conference on 3D Vision (3DV)*, pages 719–728. IEEE, 2019. 2

[5] Tenglong Ao, Qingzhe Gao, Yuke Lou, Baoquan Chen, and Libin Liu. Rhythmic gesticulator: Rhythm-aware co-speech gesture synthesis with hierarchical neural embeddings. *ACM Transactions on Graphics (TOG)*, 41(6):1–19, 2022. 2

[6] Tenglong Ao, Zeyi Zhang, and Libin Liu. Gesturediffuclip: Gesture diffusion model with clip latents. *arXiv preprint arXiv:2303.14613*, 2023. 2

[7] Joao Pedro Araújo, Jiaman Li, Karthik Vetrivel, Rishi Agarwal, Jiajun Wu, Deepak Gopinath, Alexander William Clegg, and Karen Liu. Circle: Capture in rich contextual environments. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21211–21221, 2023. 7, 13

[8] Samaneh Azadi, Akbar Shah, Thomas Hayes, Devi Parikh, and Sonal Gupta. Make-an-animation: Large-scale text-conditional 3d human motion generation. *arXiv preprint arXiv:2305.09662*, 2023. 1, 3, 7

[9] Zhongang Cai, Daxuan Ren, Ailing Zeng, Zhengyu Lin, Tao Yu, Wenjia Wang, Xiangyu Fan, Yang Gao, Yifan Yu, Liang Pan, et al. Humman: Multi-modal 4d human dataset for versatile sensing and modeling. In *European Conference on Computer Vision*, pages 557–577. Springer, 2022. 7, 13

[10] Xin Chen, Zhuo Su, Lingbo Yang, Pei Cheng, Lan Xu, Bin Fu, and Gang Yu. Learning variational motion prior for video-based motion capture. *arXiv preprint arXiv:2210.15134*, 2022. 2

[11] Xin Chen, Biao Jiang, Wen Liu, Zilong Huang, Bin Fu, Tao Chen, and Gang Yu. Executing your commands via motion diffusion in latent space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18000–18010, 2023. 1, 2, 7

[12] Rishabh Dabral, Muhammad Hamza Mughal, Vladislav Golyanik, and Christian Theobalt. Mofusion: A framework for denoising-diffusion-based motion synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9760–9770, 2023. 1

[13] Yinglin Duan, Tianyang Shi, Zhengxia Zou, Yenan Lin, Zhehui Qian, Bohan Zhang, and Yi Yuan. Single-shot motion completion with transformer. *arXiv preprint arXiv:2103.00776*, 2021. 2

[14] David Eigen, Marc'Aurelio Ranzato, and Ilya Sutskever. Learning factored representations in a deep mixture of experts. *arXiv preprint arXiv:1312.4314*, 2013. 3

[15] William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *The Journal of Machine Learning Research*, 23(1):5232–5270, 2022. 2, 3, 5

[16] Kevin Frans, Lisa Soros, and Olaf Witkowski. Clipdraw: Exploring text-to-drawing synthesis through language-image encoders. *Advances in Neural Information Processing Systems*, 35:5207–5218, 2022. 2

[17] Chuan Guo, Xinxin Zuo, Sen Wang, Shihao Zou, Qingyao Sun, Annan Deng, Minglun Gong, and Li Cheng. Action2motion: Conditioned generation of 3d human motions. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 2021–2029, 2020. 2

[18] 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)*, pages 5152–5161, 2022. 1, 2, 3, 13

[19] 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*, pages 5152–5161, 2022. 2, 7, 13

[20] Ikhsanul Habibie, Mohamed Elgharib, Kripasindhu Sarkar, Ahsan Abdullah, Simbarashe Nyatsanga, Michael Neff, and Christian Theobalt. A motion matching-based framework for controllable gesture synthesis from speech. In *ACM SIGGRAPH 2022 Conference Proceedings*, pages 1–9, 2022. 2

[21] Félix G Harvey, Mike Yurick, Derek Nowrouzezahrai, and Christopher Pal. Robust motion in-betweening. *ACM Transactions on Graphics (TOG)*, 39(4):60–1, 2020. 2, 7, 13

[22] Mohamed Hassan, Yunrong Guo, Tingwu Wang, Michael Black, Sanja Fidler, and Xue Bin Peng. Synthesizing physical character-scene interactions. 2023. 18

[23] Jiaao He, Jiezhong Qiu, Aohan Zeng, Zhilin Yang, Jidong Zhai, and Jie Tang. Fastmoe: A fast mixture-of-expert training system. *arXiv preprint arXiv:2103.13262*, 2021. 3

[24] Yannan He, Anqi Pang, Xin Chen, Han Liang, Minye Wu, Yuexin Ma, and Lan Xu. Challencap: Monocular 3d capture of challenging human performances using multi-modal references. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11400–11411, 2021. 2

[25] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 2

[26] Daniel Holden, Taku Komura, and Jun Saito. Phase-functioned neural networks for character control. *ACM Transactions on Graphics (TOG)*, 36(4):1–13, 2017. 3

[27] Fangzhou Hong, Mingyuan Zhang, Liang Pan, Zhongang Cai, Lei Yang, and Ziwei Liu. Avatarclip: Zero-shot text-driven generation and animation of 3d avatars. *ACM Transactions on Graphics (TOG)*, 41(4):1–19, 2022. 1, 2

[28] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3. 6m: Large scale datasets and pre-

dictive methods for 3d human sensing in natural environments. *IEEE transactions on pattern analysis and machine intelligence*, 36(7):1325–1339, 2013. 7, 13

[29] Robert A Jacobs, Michael I Jordan, Steven J Nowlan, and Geoffrey E Hinton. Adaptive mixtures of local experts. *Neural computation*, 3(1):79–87, 1991. 3

[30] Ajay Jain, Ben Mildenhall, Jonathan T Barron, Pieter Abbeel, and Ben Poole. Zero-shot text-guided object generation with dream fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 867–876, 2022. 2

[31] Nikolay Jetchev. Clipmatrix: Text-controlled creation of 3d textured meshes. *arXiv preprint arXiv:2109.12922*, 2021. 2

[32] Biao Jiang, Xin Chen, Wen Liu, Jingyi Yu, Gang Yu, and Tao Chen. Motiongpt: Human motion as a foreign language. *arXiv preprint arXiv:2306.14795*, 2023. 1, 3, 7

[33] Jihoon Kim, Jiseob Kim, and Sungjoon Choi. Flame: Free-form language-based motion synthesis & editing. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 8255–8263, 2023. 2

[34] Hsin-Ying Lee, Xiaodong Yang, Ming-Yu Liu, Ting-Chun Wang, Yu-Ding Lu, Ming-Hsuan Yang, and Jan Kautz. Dancing to music. *Advances in neural information processing systems*, 32, 2019. 2

[35] Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang, Maxim Krikun, Noam Shazeer, and Zhifeng Chen. Gshard: Scaling giant models with conditional computation and automatic sharding. *arXiv preprint arXiv:2006.16668*, 2020. 3

[36] Buyu Li, Yongchi Zhao, Shi Zhelun, and Lu Sheng. Dance-former: Music conditioned 3d dance generation with parametric motion transformer. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 1272–1279, 2022. 2

[37] Margaret Li, Suchin Gururangan, Tim Dettmers, Mike Lewis, Tim Althoff, Noah A Smith, and Luke Zettlemoyer. Branch-train-merge: Embarrassingly parallel training of expert language models. *arXiv preprint arXiv:2208.03306*, 2022. 3

[38] Ruilong Li, Shan Yang, David A Ross, and Angjoo Kanazawa. Ai choreographer: Music conditioned 3d dance generation with aist++. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13401–13412, 2021. 2, 7, 13

[39] Han Liang, Yannan He, Chengfeng Zhao, Mutian Li, Jingya Wang, Jingyi Yu, and Lan Xu. Hybridcap: Inertia-aid monocular capture of challenging human motions. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 1539–1548, 2023. 2, 7, 13

[40] Han Liang, Wenqian Zhang, Wenxuan Li, Jingyi Yu, and Lan Xu. Intergen: Diffusion-based multi-human motion generation under complex interactions. *arXiv preprint arXiv:2304.05684*, 2023. 2, 7, 13

[41] Junfan Lin, Jianlong Chang, Lingbo Liu, Guanbin Li, Liang Lin, Qi Tian, and Chang-wen Chen. Being comes from not-being: Open-vocabulary text-to-motion generation with wordless training. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 23222–23231, 2023. 1, 2

[42] Jing Lin, Ailing Zeng, Shunlin Lu, Yuanhao Cai, Ruimao Zhang, Haoqian Wang, and Lei Zhang. Motion-x: A large-scale 3d expressive whole-body human motion dataset. *Advances in Neural Information Processing Systems*, 2023. 13

[43] Xiao Lin and Mohamed R Amer. Human motion modeling using dvgans. *arXiv preprint arXiv:1804.10652*, 2018. 2

[44] Haiyang Liu, Zihao Zhu, Naoya Iwamoto, Yichen Peng, Zhengqing Li, You Zhou, Elif Bozkurt, and Bo Zheng. Beat: A large-scale semantic and emotional multi-modal dataset for conversational gestures synthesis. In *European Conference on Computer Vision*, pages 612–630. Springer, 2022. 7, 13

[45] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. Smpl: A skinned multi-person linear model. In *Seminal Graphics Papers: Pushing the Boundaries, Volume 2*, pages 851–866. 2023. 3

[46] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 4

[47] Naureen Mahmood, Nima Ghorbani, Nikolaus F Troje, Gerard Pons-Moll, and Michael J Black. Amass: Archive of motion capture as surface shapes. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5442–5451, 2019. 1, 7, 13

[48] Ian Mason, Sebastian Starke, and Taku Komura. Real-time style modelling of human locomotion via feature-wise transformations and local motion phases. *Proceedings of the ACM on Computer Graphics and Interactive Techniques*, 5 (1), 2022. 7, 13

[49] Oscar Michel, Roi Bar-On, Richard Liu, Sagie Benaim, and Rana Hanocka. Text2mesh: Text-driven neural stylization for meshes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13492–13502, 2022. 2

[50] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021. 4

[51] Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. Styleclip: Text-driven manipulation of stylegan imagery. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2085–2094, 2021. 2

[52] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023. 2, 4

[53] Sida Peng, Junting Dong, Qianqian Wang, Shangzhan Zhang, Qing Shuai, Xiaowei Zhou, and Hujun Bao. Animatable neural radiance fields for modeling dynamic human bodies. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14314–14323, 2021. 2

[54] Xue Bin Peng, Ze Ma, Pieter Abbeel, Sergey Levine, and Angjoo Kanazawa. Amp: Adversarial motion priors for stylized physics-based character control. *ACM Transactions on Graphics (ToG)*, 40(4):1–20, 2021. 2

[55] Mathis Petrovich, Michael J Black, and Gül Varol. Action-conditioned 3d human motion synthesis with transformer

vae. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10985–10995, 2021. 2

[56] Mathis Petrovich, Michael J Black, and Gül Varol. Temos: Generating diverse human motions from textual descriptions. In *European Conference on Computer Vision*, pages 480–497. Springer, 2022. 1, 2

[57] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 1, 2, 6

[58] Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. Zero: Memory optimizations toward training trillion parameter models. In *SC20: International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–16. IEEE, 2020. 4

[59] Scott Reed, Zeynep Akata, Xinchen Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee. Generative adversarial text to image synthesis. In *International conference on machine learning*, pages 1060–1069. PMLR, 2016. 2

[60] Davis Rempe, Tolga Birdal, Aaron Hertzmann, Jimei Yang, Srinath Sridhar, and Leonidas J Guibas. Humor: 3d human motion model for robust pose estimation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 11488–11499, 2021. 2

[61] Yiming Ren, Chengfeng Zhao, Yannan He, Peishan Cong, Han Liang, Jingyi Yu, Lan Xu, and Yuexin Ma. Lidar-aid inertial poser: Large-scale human motion capture by sparse inertial and lidar sensors. *IEEE Transactions on Visualization and Computer Graphics*, 29(5):2337–2347, 2023. 2

[62] Carlos Riquelme, Joan Puigcerver, Basil Mustafa, Maxim Neumann, Rodolphe Jenatton, André Susano Pinto, Daniel Keysers, and Neil Houlsby. Scaling vision with sparse mixture of experts. *Advances in Neural Information Processing Systems*, 34:8583–8595, 2021. 3

[63] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 1

[64] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022. 7

[65] Aditya Sanghi, Hang Chu, Joseph G Lambourne, Ye Wang, Chin-Yi Cheng, Marco Fumero, and Kamal Rahimi Malekshan. Clip-forge: Towards zero-shot text-to-shape generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18603–18613, 2022. 2

[66] Yonatan Shafir, Guy Tevet, Roy Kapon, and Amit H Bermano. Human motion diffusion as a generative prior. *arXiv preprint arXiv:2303.01418*, 2023. 2

[67] Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. Outra-geously large neural networks: The sparsely-gated mixture-of-experts layer. *arXiv preprint arXiv:1701.06538*, 2017. 2, 3, 5

[68] Sheng Shen, Zhewei Yao, Chunyuan Li, Trevor Darrell, Kurt Keutzer, and Yuxiong He. Scaling vision-language models with sparse mixture of experts. *arXiv preprint arXiv:2303.07226*, 2023. 3

[69] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2020. 6

[70] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020. 2

[71] Sebastian Starke, He Zhang, Taku Komura, and Jun Saito. Neural state machine for character-scene interactions. *ACM Trans. Graph.*, 38(6):209–1, 2019. 2

[72] Sebastian Starke, Yiwei Zhao, Taku Komura, and Kazi Zaman. Local motion phases for learning multi-contact character movements. *ACM Transactions on Graphics (TOG)*, 39 (4):54–1, 2020. 3

[73] Sebastian Starke, Ian Mason, and Taku Komura. Deepphase: Periodic autoencoders for learning motion phase manifolds. *ACM Transactions on Graphics (TOG)*, 41(4):1–13, 2022. 2, 3

[74] Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *arXiv preprint arXiv:2104.09864*, 2021. 4

[75] Omid Taheri, Nima Ghorbani, Michael J Black, and Dimitrios Tzionas. Grab: A dataset of whole-body human grasping of objects. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IV 16*, pages 581–600. Springer, 2020. 7, 13

[76] Guy Tevet, Brian Gordon, Amir Hertz, Amit H Bermano, and Daniel Cohen-Or. Motionclip: Exposing human motion generation to clip space. In *European Conference on Computer Vision*, pages 358–374. Springer, 2022. 1, 2, 7

[77] Guy Tevet, Sigal Raab, Brian Gordon, Yonatan Shafir, Daniel Cohen-Or, and Amit H Bermano. Human motion diffusion model. *arXiv preprint arXiv:2209.14916*, 2022. 1, 2, 3, 4, 7

[78] Matthew Trumble, Andrew Gilbert, Charles Malleson, Adrian Hilton, and John Collomosse. Total capture: 3d human pose estimation fusing video and inertial sensors. In *Proceedings of 28th British Machine Vision Conference*, pages 1–13, 2017. 7, 13

[79] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017. 2

[80] Can Wang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. Clip-nerf: Text-and-image driven manipulation of neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3835–3844, 2022. 2

[81] Liang Xu, Ziyang Song, Dongliang Wang, Jing Su, Zhicheng Fang, Chenjing Ding, Weihao Gan, Yichao Yan, Xin Jin, Xiaokang Yang, et al. Actformer: A gan-based transformer

towards general action-conditioned 3d human motion generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2228–2238, 2023. 2

[82] Mengde Xu, Zheng Zhang, Fangyun Wei, Yutong Lin, Yue Cao, Han Hu, and Xiang Bai. A simple baseline for open-vocabulary semantic segmentation with pre-trained vision-language model. In *European Conference on Computer Vision*, pages 736–753. Springer, 2022. 2

[83] Heyuan Yao, Zhenhua Song, Yuyang Zhou, Tenglong Ao, Baoquan Chen, and Libin Liu. Moconvq: Unified physics-based motion control via scalable discrete representations. *arXiv preprint arXiv:2310.10198*, 2023. 18

[84] Ye Yuan, Jiaming Song, Umar Iqbal, Arash Vahdat, and Jan Kautz. Physdiff: Physics-guided human motion diffusion model. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16010–16021, 2023. 2, 18

[85] He Zhang, Sebastian Starke, Taku Komura, and Jun Saito. Mode-adaptive neural networks for quadruped motion control. *ACM Transactions on Graphics (TOG)*, 37(4):1–11, 2018. 3

[86] Jiaxu Zhang, Junwu Weng, Di Kang, Fang Zhao, Shaoli Huang, Xuefei Zhe, Linchao Bao, Ying Shan, Jue Wang, and Zhigang Tu. Skinned motion retargeting with residual perception of motion semantics & geometry. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13864–13872, 2023. 3

[87] 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. *arXiv preprint arXiv:2301.06052*, 2023. 1, 2, 7

[88] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847, 2023. 1, 2, 5, 6

[89] Mingyuan Zhang, Zhongang Cai, Liang Pan, Fangzhou Hong, Xinying Guo, Lei Yang, and Ziwei Liu. Motiondiffuse: Text-driven human motion generation with diffusion model. *arXiv preprint arXiv:2208.15001*, 2022. 1, 7

[90] Siwei Zhang, Qianli Ma, Yan Zhang, Zhiyin Qian, Taein Kwon, Marc Pollefeys, Federica Bogo, and Siyu Tang. Egobody: Human body shape and motion of interacting people from head-mounted devices. In *European Conference on Computer Vision*, pages 180–200. Springer, 2022. 7, 13

[91] Yan Zhang, Mohamed Hassan, Heiko Neumann, Michael J. Black, and Siyu Tang. Generating 3d people in scenes without people. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2

# Appendix

This appendix provides more qualitative results (Appendix A), dataset details (Appendix B), user study (Appendix C), and limitations (Appendix D).

## A. Qualitative Results

We show more qualitative results from both in-domain and out-of-domain text inputs of several text-motion datasets. First, we show the generated motions of our OMG model from the text inputs in the HumanML3D test set. As illustrated in Fig. B, our model enables realistic and diverse motion generation from complicated natural sentences. Then, we show the out-of-domain generation capability using text inputs of the Mixamo test set and the concurrent Motion-X [42] dataset. As illustrated in Fig. C and Fig. D, our model well-handles unseen high-level descriptions of motion traits, like "scary clown" or "imitating snake".

## B. Dataset Details

Here we provide the details of the motion-only datasets used at the pre-training stage, as illustrated in Tab. A. We utilize 13 publicly available human motion datasets captured from various motion modalities, such as artist-created datasets [21, 48], marker-based [7, 28, 44, 47, 75], IMU-based [39, 78] and multi-view markerless [9, 38, 40, 90] motion capture datasets, totaling over 22 million frames. Since the majority of motion data is in SMPL format, we apply the retargeting algorithm to standardize them to the SMPL skeleton with rotations and positions of 22 joints, and global translation.

Moreover, we utilize HumanML3D [19] training set to train our motion ControlNet for fair comparisons with previous methods. The dataset consists of 14616 motion clips with 44970 text annotations, totaling 3.1M motion instances, as illustated in Tab. B. We further introduce Mixamo [1] dataset, consisting of abundant artist-created animations and human-annotated descriptions. It is widely used in character animation applications, such as games and VR/AR. We employ it to benchmark the zero-shot performance due to its wide variety of diverse and dynamic motions and complicated and abstract motion trait descriptions.

## C. User Study

For the comparisons of the user study presented in Fig. A, we ask the users to "Rate the motion based on how realistic it is" and "Rate the match between motion and prompt". The provided motions are generated from 60 text descriptions, 30 of which are randomly generated from the HumanML3D [19] test set and 30 from Mixamo [1] test set. We invite 20 users, shuffle the order of results from the distinct compared methods, and ask them to complete the rat-

| Dataset | Duration (h) | Frame Number | Mocap Modality | Motion Format |
|---|---|---|---|---|
| HCM [39] | 2.9 | 0.3M | IMU | SMPL |
| AMASS [47] | 62.9 | 6.8M | Marker | SMPL |
| EgoBody [90] | 0.4 | 0.04M | RGB-D | SMPL |
| GRAB [75] | 3.8 | 0.4M | Marker | SMPL |
| AIST++ [38] | 4.0 | 0.4M | RGB | SMPL |
| HuMMan [9] | 0.9 | 0.1M | RGB-D | SMPL |
| InterHuman [40] | 13.1 | 1.4M | RGB | SMPL |
| CIRCLE [7] | 10.0 | 1.1M | Marker | SMPL |
| BEAT [44] | 76 | 8.2M | Marker | BVH |
| LaFan1 [21] | 4.6 | 0.5M | Marker | BVH |
| Human3.6M [28] | 5.0 | 0.5M | Marker | SMPL |
| Total Capture [78] | 0.8 | 0.09M | IMU | SMPL |
| 100style [48] | 22.1 | 2.4M | marker | BVH |
| Total | 206.5 | 22.3M | - | - |

Table A. The details of unlabeled motion datasets used at the pre-training stage.

| Dataset | Clip Number | Text Number | Duration (h) | Frame Number | Motion Format |
|---|---|---|---|---|---|
| HumanML3D [18] | 14616 | 44970 | 28.59 | 3.1M | SMPL |
| Mixamo [1] | 2254 | 2254 | 2.5 | 0.3M | FBX |

Table B. We use HumanML3D training set at the fine-tuning stage and HumanML3D and Mixamo test set for evaluation.
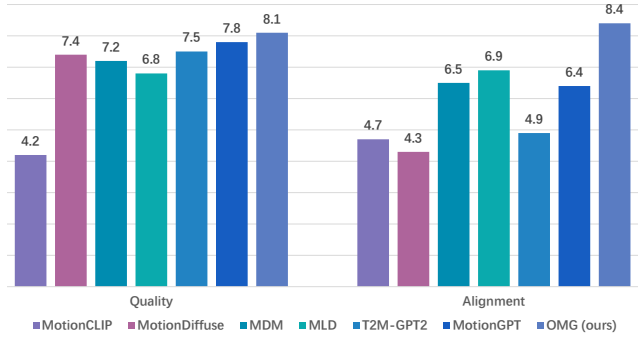


Figure A. **User Study.** We show the average quality rates and the average alignment rates of the compared methods, which indicate human evaluation of both motion quality and text-motion consistency respectively.

ing, as illustrated in Fig. E. As shown in Fig. A, our OMG was preferred over the other state-of-the-art methods in both motion quality and text-motion alignment.
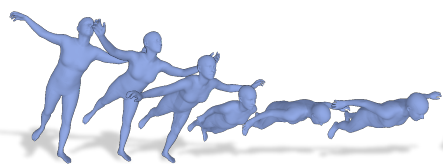
## D. Limitations

As the trial to explore realistic open-vocabulary motion generation, the proposed OMG still has limitations as follows.
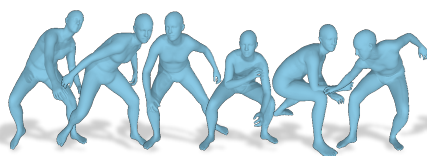
**Motion space.** Our method still relies on the training motion manifold and cannot generate motions that are beyond the scope of the training data, such as flying, yoga, or swimming.

**Precise control.** Our method does not explicitly model the temporal order and inclusion relations of sub-motions, which are unable to handle precise control, such as picking an object or reaching a goal.

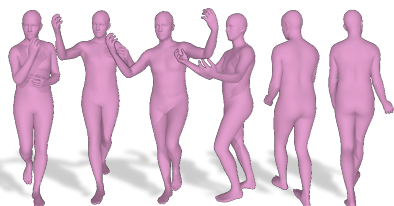**Physically implausible.** Our method does not explicitly

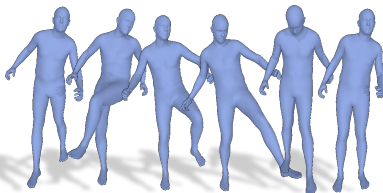The person is flying like a airplane.

Cross sideways step before crouching a bit sliding back and forth across the plane.

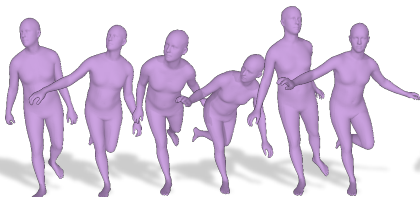A person walk forward, picks something up, then tosses it up.

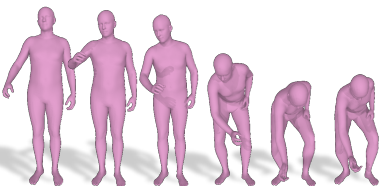Moves forward with arms moving dancing and then a turn then walks back.

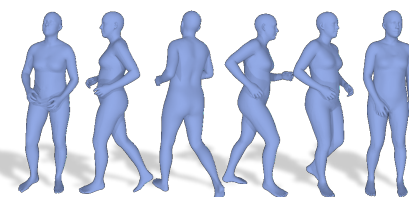The standing person kicks with their left foot before going back to their original stance.

A person bends their back to stretch.

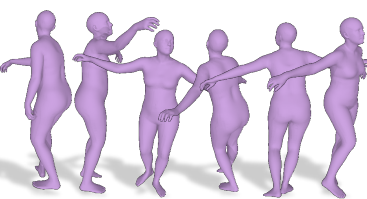Balancing on his right foot, touch down once with his left foot, then resume balancing.

Pokes their hand along the ground, like the might be planting seeds.

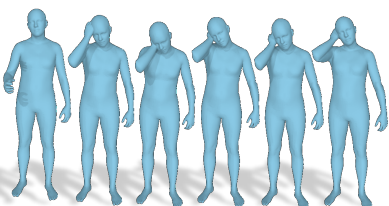Performing a right to left jogging move.
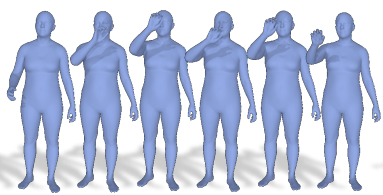
A person takes a wide swing with their left hand.

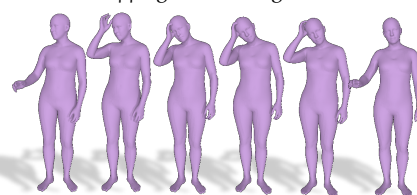Spinning dance where they turn around.

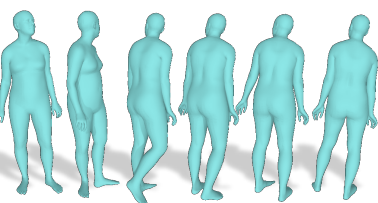Quickly waving arms above head and then clapping while looking around.

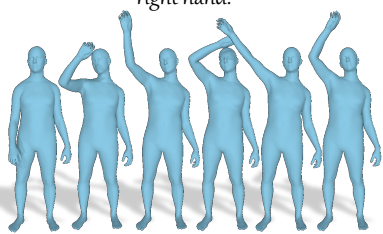Raises his right hand to talk on the phone.

Presses things in front of them with their right hand.
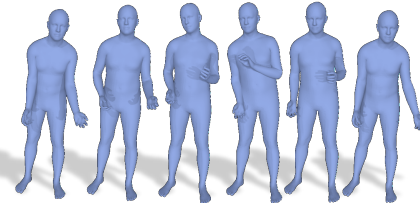
Scratch head with right hand.

A person turns around and looks to the left.

A person stands still and waves with right hand.

A person is pouring and serving drinks.

Figure B. Qualitative results on HumanML3D test set.

Air squat workout

Burpee exercise

Drunken walking backward right turn

Banging fist on table

Silly dance the cabbage patch

Pitching a baseball

Overhead bashing swing

Angry forward gesture

Juggling right kick of a soccer ball

Swing club and triple knee kicks

Left standing sneak behind cover

Jumping while strafing left

Male with headphones listening to music

Quick left handed punch

Nervously looking around

Female samba ijexa break

Hard floor stomp

Scary clown walk

Figure C. Qualitative results on Mixamo test set.

Bandaging

Arm leg lift to crunche

Baseball bunt

Imitating snake

Chainsaw tree

Climbing tree

Play basketball

Close-up dolchhaltung

Chest pain during walking

Butt kicks arm swing

Cross single-plank bridge

Dance krump arm swing

Dance middle hip hop

Eating a hamburger during walking

Move arrogant

Running side up punch

Dance ballet jazz entrelace

Ironing clothes

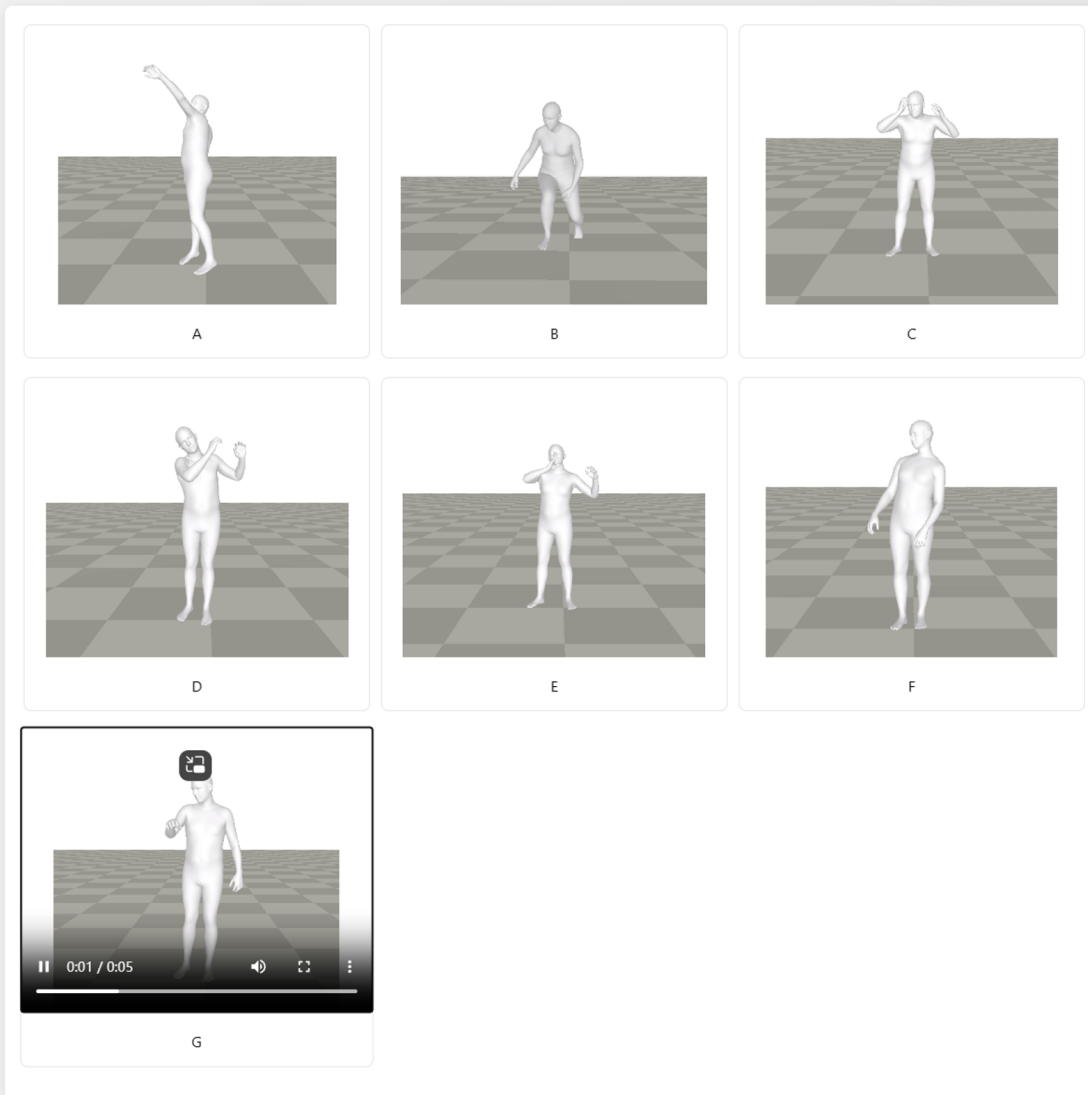Figure D. Qualitative results on Motion-X dataset.

Figure E. **User Study.** We ask 20 users to rate the motion quality and text-motion consistency of 60 results generated from each method. The rating range is from 1 to 10.

model physical dynamics, which leads to physically implausible motion generation. Recent physics-based motion control [22, 83, 84] approaches use reinforcement learning to control human characters in a physically simulated environment, achieving impressive motion quality. It's interesting to introduce physics into the conditional generative model pipeline.

**Maximum length.** Same as most motion generation methods, our method can generate arbitrary length results but still under the max-length in the dataset. It's interesting to model a non-stop human motion in temporal consistency.

**Full-body dynamics.** Our method focuses on articulated human bodies. How to model the full-body dynamics including the face, eyes, hands, and even toes, which enables complicated interactions with our complex physical world, remains a huge challenge.