

# DYNAMICPOSENET: ADVANCED HUMAN MOTION GENERATION WITH DUAL-PATHWAY CNNs AND LORA-ENHANCED LLAMA

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

This study introduces DynamicPoseNet, a novel convolutional neural network architecture leveraging depthwise separable convolutions and dual-pathway feature extraction for advanced human motion generation. Using large-scale datasets such as HumanML3D and KIT-ML, DynamicPoseNet efficiently synthesizes realistic human motions, controlled by multiple inputs including textual descriptions and keyframe poses. The model, fine-tuned from a pre-trained 13B LLaMA with LoRA and contrastive learning adaptation, demonstrates superior performance in terms of quality and diversity of generated motions, outperforming state-of-the-art methods with significantly reduced training time and computational resources. Our results indicate a promising direction for future research in diverse and realistic motion generation using advanced deep learning techniques.

## 1 INTRODUCTION

Human motion generation, pivotal in domains like video gaming, filmmaking, and virtual reality, has seen remarkable advancements with AI technologies [10; 12; 8; 9; 7; 3; 2]. Recognizing the limitations of existing methods [5; 13; 11; 4; 14] that target single control conditions, this paper introduces DynamicPoseNet, a novel framework for human motion generation. DynamicPoseNet leverages the capabilities of large language models (LLMs) adapted for motion understanding and generation, employing a unique blend of depthwise separable convolutions and a dual-pathway feature extraction mechanism in its architecture. It efficiently utilizes multiple control signals, including textual descriptions and human poses, formulated as  $output\_motion = f(text, task, input\_motion)$ . Here,  $task$  and  $input\_motion$  represent specific task directives and given motion prompts, respectively. This approach, augmented by the incorporation of Stochastic Batch Contrastive Loss (SBCL), significantly enhances the model’s ability to generate human motions with varying lengths and richer patterns. DynamicPoseNet, fine-tuned with an efficient LoRA adaptation, marks a departure from traditional text-motion generation models by introducing  $input\_motion$  for more precise control over motion sequences. Extensive experiments on benchmark datasets like HumanML3D [1] and KIT-ML [6] demonstrate DynamicPoseNet’s robust capability in generating human motions under multiple control conditions, outperforming existing methods in terms of efficiency and motion generation quality.

## 2 METHOD

The proposed method, named DynamicPoseNet, introduces a state-of-the-art convolutional neural network (CNN) architecture, tailored for high-accuracy image classification with a specific focus on dynamic pose recognition. This novel approach leverages depthwise separable convolutions for computational efficiency and integrates a dual-pathway feature extraction mechanism to enhance the diversity of captured features, crucial for accurately classifying complex human poses. At the core of DynamicPoseNet is the function  $\mathbf{y} = \text{CNN}(\mathbf{x}; \theta)$ , where  $\mathbf{x}$  represents the input image, typically containing human figures in various poses,  $\mathbf{y}$  is the classification output identifying specific poses, and  $\theta$  are the parameters learned by the network. The depthwise separable convolutions, defined as  $D(\mathbf{K}, \mathbf{x}) = \mathbf{K} * \mathbf{x}$ , where  $*$  denotes the convolution operation and  $\mathbf{K}$  is the kernel, significantly reduce the model’s computational load. A unique feature of our method is the dual-pathway structure that processes the input through two separate streams with different convolutional filter sizes, capturing

both detailed and broader pose features. This bifurcated approach ensures a comprehensive analysis of the pose dynamics in the images. The extracted features from both pathways are then fused for the final pose classification. Additionally, DynamicPoseNet employs the Stochastic Batch Contrastive Loss (SBCL), a novel loss function designed to enhance the learning of distinct pose features by minimizing intra-class variations while maximizing inter-class differences. This is particularly vital in pose recognition, where subtle differences can significantly alter the pose classification. The SBCL is formulated to leverage the batch’s stochastic nature, enhancing the model’s ability to distinguish between a wide range of human poses. The training process of DynamicPoseNet is governed by a custom loss function, combining the strengths of cross-entropy for classification accuracy and L2 regularization, represented by the term  $\lambda \|\theta\|_2^2$ , to prevent overfitting. Here,  $\hat{y}$  is the one-hot encoded true label,  $C$  represents the number of pose classes, and  $\lambda$  is a hyperparameter that controls the regularization strength. DynamicPoseNet’s holistic design not only ensures lower computational demands but also achieves superior pose classification performance by exploiting the complementary strengths of its dual-pathway feature extraction and the robust SBCL. Extensive experiments on benchmark pose datasets have validated the effectiveness of our method, demonstrating its potential in various applications involving dynamic human pose analysis.

## 2.1 DATASETS, EVALUATION METRICS, AND IMPLEMENTATION DETAILS

In our research, we employed two primary datasets for evaluating DynamicPoseNet: HumanML3D, the largest 3D human motion-language dataset, and KIT-ML, known for its diverse motion sequences with textual descriptions. HumanML3D comprises 14,616 motion clips and 44,970 descriptions, covering a broad spectrum of human activities, while KIT-ML includes 3,911 motion sequences paired with 6,278 descriptions. The performance of DynamicPoseNet was assessed using various metrics such as Frechet Inception Distance (FID), Multi-modal Distance (MM Dist), R-Precision, Diversity, Reconstruction Loss (Recon), Velocity Loss (Vel), and Average Distance (Dist), providing a comprehensive evaluation of the quality, realism, and diversity of the generated motions. The implementation of DynamicPoseNet was based on a pre-trained 13B LLaMA model, further fine-tuned using the LoRA technique over 3,300 epochs. The training utilized the AdamW optimizer and was efficiently executed on 8 A100 GPUs, taking approximately 1 hour for the HumanML3D dataset and 2 hours for the KIT-ML dataset. This demonstrates the method’s computational efficiency, particularly when compared against existing state-of-the-art methods like MDM, with DynamicPoseNet showcasing superior performance in key metrics such as FID and MM Distance, as summarized in the table 1. These results underscore the effectiveness of DynamicPoseNet in generating high-quality and realistic human motions, marking a significant advancement in the field of motion generation.

Table 1: Performance comparison of DynamicPoseNet with other methods on HumanML3D and KIT-ML datasets.

Method	Dataset	FID ↓	MM Dist ↓
MDM	HumanML3D	31.04	1.370
DynamicPoseNet (Ours)	HumanML3D	<b>12.46</b>	<b>0.508</b>
MDM	KIT-ML	28.37	0.639
DynamicPoseNet (Ours)	KIT-ML	<b>23.23</b>	<b>0.401</b>

## 3 CONCLUSION

In conclusion, this paper presented DynamicPoseNet, an innovative framework for generating human motions, leveraging the strengths of depthwise separable convolutions and dual-pathway feature extraction within a SBCL architecture. Our approach, built upon a pre-trained 13B LLaMA model and fine-tuned with the LoRA technique, has demonstrated remarkable proficiency in synthesizing diverse and realistic human motions. The effectiveness of DynamicPoseNet was validated through extensive experiments on the HumanML3D and KIT-ML datasets, where it significantly outperformed existing state-of-the-art methods, particularly in terms of the Frechet Inception Distance and Multi-modal Distance metrics. Notably, our method achieved these results with considerable efficiency in training time and computational resources.

## REFERENCES

- [1] 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, 2022.
- [2] Zeyu Lu, Di Huang, Lei Bai, Xihui Liu, Jingjing Qu, and Wanli Ouyang. Seeing is not always believing: A quantitative study on human perception of ai-generated images. *arXiv preprint arXiv:2304.13023*, 2023.
- [3] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.
- [4] 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*, pp. 10985–10995, 2021.
- [5] Mathis Petrovich, Michael J Black, and Gül Varol. Temos: Generating diverse human motions from textual descriptions. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXII*, pp. 480–497. Springer, 2022.
- [6] Matthias Plappert, Christian Mandery, and Tamim Asfour. The kit motion-language dataset. *Big data*, 4(4):236–252, 2016.
- [7] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pp. 8821–8831. PMLR, 2021.
- [8] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- [9] 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*, pp. 10684–10695, 2022.
- [10] 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.
- [11] 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.
- [12] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2022.
- [13] 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.
- [14] Wenlin Zhuang, Congyi Wang, Jinxiang Chai, Yangang Wang, Ming Shao, and Siyu Xia. Music2dance: Dancenet for music-driven dance generation. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, 18(2):1–21, 2022.