

# Taming Diffusion Models for Music-driven Conducting Motion Generation

Zhuoran Zhao<sup>1\*</sup>, Jinbin Bai<sup>1\*</sup>, Delong Chen<sup>2</sup>, Debang Wang<sup>1</sup>, Yubo Pan<sup>1</sup>

<sup>1</sup> Dept. of Computer Science, School of Computing, National University of Singapore

<sup>2</sup> Xiaobing.AI

## Abstract

Generating the motion of orchestral conductors from a given piece of symphony music is a challenging task since it requires a model to learn semantic music features and capture the underlying distribution of real conducting motion. Prior works have applied Generative Adversarial Networks (GAN) to this task, but the promising diffusion model, which recently showed its advantages in terms of both training stability and output quality, has not been exploited in this context. This paper presents *Diffusion-Conductor*, a novel DDIM-based approach for music-driven conducting motion generation, which integrates the diffusion model to a two-stage learning framework. We further propose a random masking strategy to improve the feature robustness, and use a pair of geometric loss functions to impose additional regularizations and increase motion diversity. We also design several novel metrics, including Fréchet Gesture Distance (FGD) and Beat Consistency Score (BC) for a more comprehensive evaluation of the generated motion. Experimental results demonstrate the advantages of our model. The code will be released.

## Introduction

Human conductors have the remarkable ability to translate their rich comprehension of music contents into sequences of precise yet graceful conducting motion. Advancements in AIGC technologies for human motion (Mourot et al. 2022) have addressed the generation of various human motions such as speech gestures, dance movements, and instrumental motions over recent years, and researchers are now pivoting toward building AI conductors. Pioneered works of VirtualConductor (Chen et al. 2021) and M<sup>2</sup>S-GAN (Liu et al. 2022) demonstrated the promising possibilities of building such systems. These works leverage Generative Adversarial Network (GAN) (Goodfellow et al. 2020) to learn the probabilistic distribution of real conducting motion from a large-scale paired music-motion dataset. However, GAN-based models typically suffer from notorious issues such as mode collapse and unstable training, which impede the generation of plausible conducting motions.

Recently, diffusion models (Ho, Jain, and Abbeel 2020; Ho and Salimans 2022) have emerged as the new state-of-the-art family of deep generative models. Representa-

tive models such as GLIDE (Nichol et al. 2021), DALL-E 2 (Ramesh et al. 2022), Latent Diffusion (Rombach et al. 2021), ImageGen (Saharia et al. 2022), and Stable Diffusion (Rombach et al. 2022), yields impressive performance on conditional image generation, surpassing those GAN-based methods which dominated the field for the past few years. We hypothesize that such an advantage can be extended to the task of music-driven conducting motion generation, and in this paper, we introduce our *Diffusion-Conductor* model, which is the first diffusion-based AI conductor model.

Our learning framework comprises two consecutive stages, namely the contrastive learning stage and the generative learning stage. The first stage builds a two-tower structure and performs music-motion contrastive pre-training to learn rich music features, those learned features are subsequently transferred to the second stage with a random masking strategy. We incorporate a DDIM-based model to learn the conditional generation of conducting motion, and we modify the supervision signal from  $\epsilon$  to  $x_0$  for better generation performance. Furthermore, we impose additional supervision on the model via two geometric regularization losses, namely velocity loss and elbow loss, to enhance the consistency and diversity of generated motions.

We use a broad array of metrics, including Mean Squared Error (MSE), Fréchet Gesture Distance (FGD), Beat Consistency Score (BC), and Diversity, to evaluate the motion produced by *Diffusion-Conductor*. Thorough comparisons demonstrated that our model outperforms the previous GAN-based method (Liu et al. 2022).

In summary, our main contributions are as follows:

- Our method is the first work to use diffusion model for music-driven conducting motion generation.
- We modify the supervision signal from  $\epsilon$  to  $x_0$  to achieve the better performance on generating conducting motions, which will inspire later research.
- Extensive experiments demonstrate the superiority of our method with quantitative comparison like FGD, BC, and Diversity.

## Methods

In this section, we will explain our task definition, provide an overview of our approach, and illustrate the Training Objective applied at the various stages.

\*These authors contributed equally.

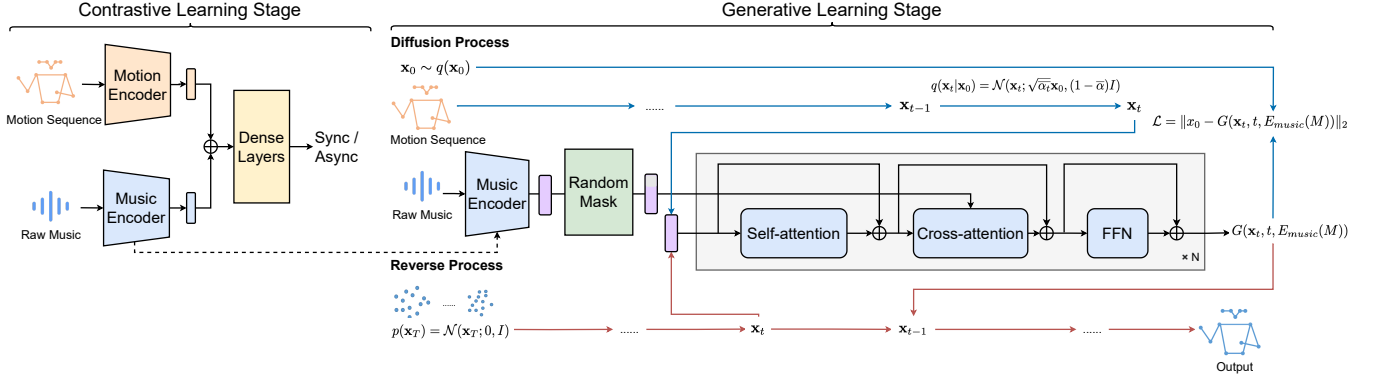


Figure 1: Overview of the proposed framework. The colors of the arrows in Generative Learning Stage represent different stages: **blue** for training, **red** for inference, and **black** for both training and inference.

## Task Definition

Conditional motion generation is a series of tasks that generates realistic and plausible human body motions  $\mathbf{Y}$  with specified actions in response to a given prompt. The structure of  $\mathbf{Y}$  comprises an array of poses  $[\mathbf{y}_i]$ , where each element  $\mathbf{y}_i \in \mathbb{R}^D$  denotes the pose state at the  $i$ -th frame, with  $F$  being the total number of frames. In the case of music-conditioned motion generation, the specified prompt is the music feature  $E_{music}(\mathbf{M})$  extracted from the raw music  $\mathbf{M}$ . Our objective is to learn a diffusion model  $G$ , which can generate a motion sequence  $\mathbf{Y}$  corresponding to the given  $E_{music}(\mathbf{M})$ .

## Overview of Our Approach

Our proposed architecture is illustrated in Fig. 1. In the contrastive learning stage, a contrastive pre-training network composed of a motion encoder  $E_{motion}$  and a music encoder  $E_{music}$  is used to learn music representations that are correctly aligned to their corresponding motion representations. Subsequently, a generation network  $G$  is employed during the generative learning stage to generate a motion sequence based on the music embeddings outputted by the pre-trained  $E_{music}$ . To further facilitate motion generation while undergoing the denoising process, we make use of the denoising diffusion implicit model (DDIM) (Song, Meng, and Ermon 2020) and introduce a Cross-Modality Linear Transformer. During inference, a Gaussian distribution noise is sampled according to the given random seed and fed into the denoising process with cross-attention between the music features. Finally, music-driven conducting motions will be generated. Detailed descriptions of our methods are presented in the following sections.

**Contrastive Pre-training** The contrastive pre-training network comprises three components: a Motion Encoder  $E_{motion}$ , a Music Encoder  $E_{music}$ , and a set of Dense Layers  $D$ . The motion embeddings and music embeddings generated by  $E_{motion}$  and  $E_{music}$  are concatenated and then passed to  $D$ , after which a binary cross-entropy loss is applied to assess whether music and motion are appropriately paired. Specifically, the Music Encoder  $E_{music}$  is used to generate music features from raw music and consists of three groups of layers, with each layer comprised of three residual

layers and a max-pooling layer. Meanwhile, the Motion Encoder  $E_{motion}$  is employed to generate motion features for the conducting motion sequence. To analyze the conducting motion both spatially and temporally, we make use of the Spatial-Temporal Graph Convolutional Network (ST-GCN) (Yan, Xiong, and Lin 2018), which has been used extensively in human pose estimation tasks.

**Diffusion Model for Motion Generation** Diffusion models involve a diffusion process and a reverse process. The diffusion process adds Gaussian noise to the motion sequence data in accordance with the Markov chain rule to approximate the posterior  $q(\mathbf{x}_{1:T}|\mathbf{x}_0)$ . Upon completion of the diffusion process, the data distribution  $\mathbf{x}_T$  should be equivalent to an isotropic Gaussian distribution:

$$q(\mathbf{x}_{1:T}|\mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1}) \quad (1)$$

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t I) \quad (2)$$

Using reparameterization trick, we can sample  $\mathbf{x}_t$  at any arbitrary time step  $t$  in a closed form:

$$q(\mathbf{x}_t|\mathbf{x}_0) = \sqrt{\alpha_t}\mathbf{x}_0 + \epsilon\sqrt{1 - \alpha_t}, \epsilon \sim \mathcal{N}(0, I) \quad (3)$$

In order to run the reverse process, we need to learn a model  $p_\theta$  to approximate  $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$  since  $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$  is intractable:

$$p_\theta(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) \quad (4)$$

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \Sigma_\theta(\mathbf{x}_t, t)) \quad (5)$$

Most prior works (Ho, Jain, and Abbeel 2020; Nichol et al. 2021; Zhu et al. 2023) train the model to predict the noise  $\epsilon_\theta(\mathbf{x}_t, t, E_{music}(\mathbf{M}))$  and then calculate the mean square error between  $\epsilon$  and  $\epsilon_\theta(\mathbf{x}_t, t, E_{music}(\mathbf{M}))$  to optimize the model. Here, we instead follow (Ramesh et al. 2022; Tevet et al. 2022) by directly predicting the motion  $\mathbf{x}_0$  and using the mean square error on this prediction which yields better generation performance. Subsequently, the reverse process can be employed to denoise the motion sequence step by step and generate a clean motion sequence conditioned on the given music embeddings.

**Cross-Modality Linear Transformer** To serve as the denoising model, we make use of a Transformer (Vaswani et al. 2017). We initially utilize a music encoder to extract the music embeddings, the pre-training of which during the contrastive learning stage can facilitate the generation process. Subsequently, a self-attention module is employed to enable motion features from different times to interact with each other. Additionally, a cross-attention module is utilized to fuse the music embeddings and motion sequence together while a feed-forward network is utilized to generate motion as  $\epsilon_\theta(\mathbf{x}_t, t, E_{music}(M))$ .

**Random Mask** Inspired by masked language modeling and masked image modeling, we incorporated a random mask (Zhong et al. 2020; Tevet et al. 2022) block after the music encoder to train the diffusion model with both music-conditional and unconditional elements. This can potentially allow us to trade off between diversity and quality for improved generalization performance.

### Training Objective

At the first stage, we adopt a binary cross-entropy loss to learn the representation of music under the supervision of motion, which can be formulated as:

$$\mathcal{L}_{bce} = \sum_{i,j=1}^N \left( c_{ij} \log_2(f[E_{music}(\mathbf{M}_i) \oplus E_{motion}(\mathbf{Y}_j)]) + (1 - c_{ij}) \log_2(1 - f[E_{music}(\mathbf{M}_i) \oplus E_{motion}(\mathbf{Y}_j)]) \right) \quad (6)$$

where  $c_{ij}$  is defined by

$$c_{ij} = \begin{cases} 1, & i = j \\ 0, & \text{otherwise} \end{cases}$$

$\mathbf{M}_i$  and  $\mathbf{Y}_j$  represent the  $i$ -th music data and the  $j$ -th motion data respectively, where  $\oplus$  denotes the feature concatenation operation. Both  $E_{music}$  and  $E_{motion}$  denote the music and motion encoders respectively and  $f$  represents a dense layer.

At the second stage, we follow (Ramesh et al. 2022; Tevet et al. 2022) to directly predict the motion  $x_0$  rather than predicting the noise  $\epsilon$  as formulated by (Ho, Jain, and Abbeel 2020), for improved generation performance. The loss can be expressed as follows:

$$\mathcal{L}_{ddim} = \|x_0 - G(\mathbf{x}_t, t, E_{music}(M))\|_2^2 \quad (7)$$

where  $x_0$  is the original motion sequence and  $G(\mathbf{x}_t, t, E_{music}(M))$  denotes the final step of motion sequence generated by the diffusion model.

In addition, a geometric loss is employed to regularize the generative model, enforcing physical properties and preventing artifacts in order to generate natural and coherent motion. This consists of a velocity loss and an elbow loss; the former ensures that the velocity of the generated motion coincides with the ground-truth motion and the latter encourages more intensive arm swing for more vivid motion.

$$\mathcal{L}_{vel} = \frac{1}{N-1} \sum_{i=1}^{N-1} \|(x_0^{i+1} - x_0^i) - (\hat{x}_0^{i+1} - \hat{x}_0^i)\|_2^2 \quad (8)$$

$$\mathcal{L}_{elbow} = -\frac{1}{N-1} \sum_{i=1}^{N-1} \|\hat{x}_{0_{elbow}}^{i+1} - \hat{x}_{0_{elbow}}^i\|_2^2 \quad (9)$$

Overall, our training objective is:

$$\mathcal{L} = \mathcal{L}_{ddim} + \lambda_{vel} \mathcal{L}_{vel} + \lambda_{elbow} \mathcal{L}_{elbow} \quad (10)$$

## Experiment

In this section, we will first present the training datasets and evaluation metrics. Subsequently, we will conduct quantitative and qualitative experiments that are compared to our baseline method, followed by providing some ablation studies intended to demonstrate the efficacy of our frameworks.

### Datasets

We leverage the ConductorMotion100 dataset (Chen et al. 2021) for training purposes. It consists of a training set, validation set and test set, with respective durations of 90, 5 and 5 hours. Since the motion of the conductor’s lower body contains very little useful information and is often occluded or outside of the camera’s view, ConductorMotion only preserves 13 2D keypoints of the upper body in the MS COCO format. All motion data is re-sampled to 30 fps, with corresponding music motion encoding at 90 Hz.

### Evaluation Metrics

We use four metrics that are commonly utilized in motion generation and relative fields to evaluate our method.

**Mean Squared Error (MSE).** Mean squared error (MSE) is the most direct way to measure how closely the generated motion corresponds to the ground-truth motion and has been widely used as an evaluation metric in music-to-motion tasks. The representation of MSE is defined as follows:

$$MSE(Y, \hat{Y}) = \|Y - \hat{Y}\|_2^2$$

where  $Y$  denotes the ground-truth motion and  $\hat{Y}$  denotes the generated motion.

**Fréchet Gesture Distance (FGD).** FGD is frequently used to measure the distance between synthesized gesture distribution and real data distribution. Since gesture motion and conducting motion are closely related, both being represented as keypoints, we employ FGD to evaluate the distance of the generated conducting motion distribution and the ground-truth conducting motion distribution. FGD is expressed as follows:

$$FGD(Y, \hat{Y}) = \|\mu_{gt} - \mu_{gen}\|_2^2 + \text{Tr}(\Sigma_{gt} + \Sigma_{gen} - 2(\Sigma_{gt} \Sigma_{gen})^{\frac{1}{2}})$$

where  $\mu_{gt}$  and  $\Sigma_{gt}$  stand for the mean and variance of the latent feature distribution of the ground-truth motion  $Y$ , while  $\mu_{gen}$  and  $\Sigma_{gen}$  are the mean and variance of the latent feature distribution of the generated motion  $\hat{Y}$ .

**Beat Consistency Score (BC).** Beat Consistency Score is a metric to evaluate motion-music correlation in terms of the similarity between the motion beats and music beats. We follow (Li et al. 2021) to define motion beats as the local minima of kinetic velocity and use librosa (McFee et al. 2015) to extract music beats. Beat Consistency Score computes the average distance between every music beat and its nearest motion beat:

$$BC = \frac{1}{|\mathcal{B}^x|} \sum_{i=1}^{|\mathcal{B}^x|} \exp \left( - \frac{\min_{t_j^x \in \mathcal{B}^x} \|t_j^x - t_i^y\|_2^2}{2\sigma^2} \right)$$

where  $\mathcal{B}^x = \{t_j^x\}$  represent motion beats and  $\mathcal{B}^y = \{t_i^y\}$  represent music beats, and  $\sigma$  is the parameter to normalize sequences, which is set to 3 empirically.

**Diversity.** Similar to prior works (Zhu et al. 2023; Li et al. 2021), we evaluate our model’s ability to generate diverse conducting motions given various input music. Like (Zhu et al. 2023), we choose 500 generated samples randomly and calculate the mean absolute error between the generated latent motion features and the shuffled features.

### Implementation Details

For the diffusion model, we set the diffusion steps to 1000 and use Adam (Kingma and Ba 2014) for optimization with a learning rate of  $2e-4$  and batch size of 32. We train the diffusion model over 1200 epochs, setting the unconditional rate of random mask to 0.1. Experiments are conducted on two NVIDIA TESLA V100 GPUs.

### Main Results

Table 1: Main results on ConductorMotion100 test set

| Methods                                 | MSE ↓         | FGD ↓         | BC ↑         | Diversity ↑    |
|---|---------------|---------------|--------------|----------------|
| M <sup>2</sup> S-GAN (VirtualConductor) | 0.0054        | 1051.97       | 0.109        | 1012.06        |
| Diffusion-Conductor                     | <b>0.0042</b> | <b>812.01</b> | <b>0.119</b> | <b>1152.06</b> |

As shown in table 1, we report four metrics compared with VirtualConductor (Chen et al. 2021) on ConductorMotion100 test set. It is shown that our method outperforms VirtualConductor on all the four metrics.

We further visualize the beat consistency between the music and generated conducting motion, making a comparison with VirtualConductor. As illustrated in Fig. 2, our generated motion beats are better able to match the given music beats.

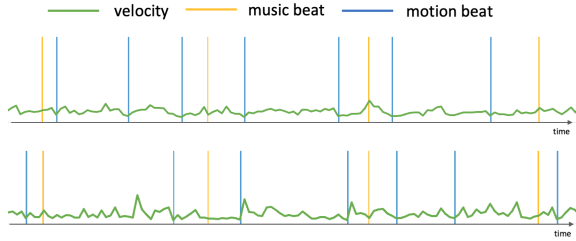


Figure 2: Qualitative comparison of beat consistency between VirtualConductor (top) and ours (bottom).

In addition, we provide visualizations of motion generation conditioned on music which were not included in the training or test sets. We randomly select the following symphonies: Tchaikovsky Piano Concerto No.1, Beethoven’s Symphony No.7, The Marriage of Figaro Overture, and Vivaldi Four Seasons (Spring) (see Fig. 3).

### Ablation Study

**Comparison of predicting  $\epsilon$  and  $x_0$ .** We further investigate the effect of predicting the noise  $\epsilon$  versus the motion  $x_0$  via an additional study. The results as shown in Fig. 4 demonstrate that the model trained by minimizing the loss between the noise  $\epsilon$  performs much worse than one trained by minimizing the loss between motion  $x_0$ , which fails to generate plausible motion sequences in longer frames, whereas predicting  $x_0$  successfully produces stable and plausible motion sequences. These results demonstrate

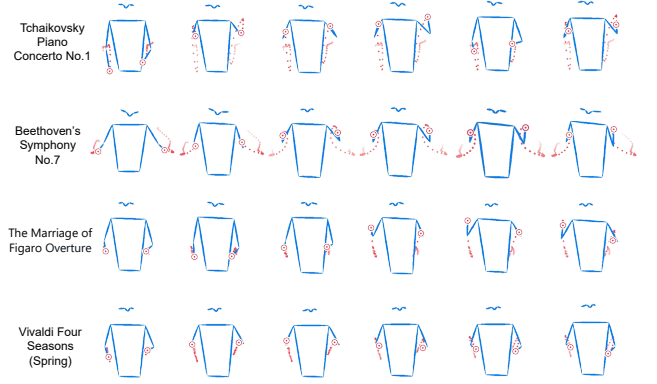


Figure 3: Visualization of the four symphonies.

the effectiveness of our design-choice to predict the motion rather than noise for each diffusion step.

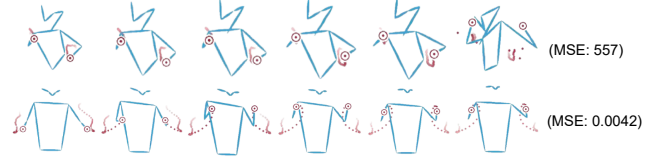


Figure 4: Qualitative and MSE comparison of generated motion of predicting  $\epsilon$  (top) and  $x_0$  (bottom) on ConductorMotion100 test set.

**Effect of geometric loss.** We examine the effect of incorporating a geometric loss in the training objective and compare it with one trained without its use. The results indicated in Table 2 show that the model trained with geometric loss can achieve better performance than the model trained without it on the test set. Furthermore, as visualized in Fig. 5, the model trained with a geometric loss is able to produce motion with more vivid arm swings and plausible poses, which confirms its effectiveness in yielding high-quality motion.

Table 2: Comparison of five metrics on ConductorMotion100 test set with and without geometric loss

| Method             | MSE ↓         | FGD ↓         | BC ↑         | Diversity ↑    |
|--------------------|---------------|---------------|--------------|----------------|
| w/o geometric loss | 0.0045        | 822.07        | 0.116        | 1127.90        |
| w geometric loss   | <b>0.0042</b> | <b>812.01</b> | <b>0.119</b> | <b>1152.06</b> |



Figure 5: Qualitative comparison of generated motion of w/o (left) and w (right) geometric loss.

### Conclusion

In this paper, we presents Diffusion-Conductor, a novel DDIM-based approach for music-driven conducting motion generation, which integrates the diffusion model to a two-stage learning framework. And extensive experiments on several metrics, including Fréchet Gesture Distance (FGD) and Beat Consistency Score (BC) demonstrated the superiority of our approach.

## References

- Chen, D.; Liu, F.; Li, Z.; and Xu, F. 2021. VirtualConductor: Music-driven Conducting Video Generation System. *CoRR*, abs/2108.04350.
- Goodfellow, I.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A.; and Bengio, Y. 2020. Generative adversarial networks. *Communications of the ACM*, 63(11): 139–144.
- Ho, J.; Jain, A.; and Abbeel, P. 2020. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33: 6840–6851.
- Ho, J.; and Salimans, T. 2022. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*.
- Kingma, D. P.; and Ba, J. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Li, R.; Yang, S.; Ross, D. A.; and Kanazawa, A. 2021. Ai choreographer: Music conditioned 3d dance generation with aist++. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 13401–13412.
- Liu, F.; Chen, D.-L.; Zhou, R.-Z.; Yang, S.; and Xu, F. 2022. Self-supervised music motion synchronization learning for music-driven conducting motion generation. *Journal of Computer Science and Technology*, 37(3): 539–558.
- McFee, B.; Raffel, C.; Liang, D.; Ellis, D. P.; McVicar, M.; Battenberg, E.; and Nieto, O. 2015. librosa: Audio and music signal analysis in python. In *Proceedings of the 14th python in science conference*, volume 8, 18–25.
- Mourot, L.; Hoyet, L.; Le Clerc, F.; Schnitzler, F.; and Hellier, P. 2022. A Survey on Deep Learning for Skeleton-Based Human Animation. In *Computer Graphics Forum*, volume 41, 122–157. Wiley Online Library.
- Nichol, A.; Dhariwal, P.; Ramesh, A.; Shyam, P.; Mishkin, P.; McGrew, B.; Sutskever, I.; and Chen, M. 2021. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*.
- Ramesh, A.; Dhariwal, P.; Nichol, A.; Chu, C.; and Chen, M. 2022. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*.
- Rombach, R.; Blattmann, A.; Lorenz, D.; Esser, P.; and Ommer, B. 2021. High-Resolution Image Synthesis with Latent Diffusion Models. *arXiv:2112.10752*.
- Rombach, R.; Blattmann, A.; Lorenz, D.; Esser, P.; and Ommer, B. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 10684–10695.
- Saharia, C.; Chan, W.; Saxena, S.; Li, L.; Whang, J.; Denton, E.; Ghasemipour, S. K. S.; Ayan, B. K.; Mahdavi, S. S.; Lopes, R. G.; et al. 2022. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*.
- Song, J.; Meng, C.; and Ermon, S. 2020. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*.
- Tevet, G.; Raab, S.; Gordon, B.; Shafir, Y.; Cohen-Or, D.; and Bermano, A. H. 2022. Human motion diffusion model. *arXiv preprint arXiv:2209.14916*.
- Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, L. u.; and Polosukhin, I. 2017. Attention is All you Need. In Guyon, I.; Luxburg, U. V.; Bengio, S.; Wallach, H.; Fergus, R.; Vishwanathan, S.; and Garnett, R., eds., *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Yan, S.; Xiong, Y.; and Lin, D. 2018. Spatial temporal graph convolutional networks for skeleton-based action recognition. In *Thirty-second AAAI conference on artificial intelligence*.
- Zhong, Z.; Zheng, L.; Kang, G.; Li, S.; and Yang, Y. 2020. Random erasing data augmentation. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, 13001–13008.
- Zhu, L.; Liu, X.; Liu, X.; Qian, R.; Liu, Z.; and Yu, L. 2023. Taming Diffusion Models for Audio-Driven Co-Speech Gesture Generation. *arXiv preprint arXiv:2303.09119*.