

---

# Human Motion Generation From Text

---

**Xinran Li**

ShanghaiTech University

lixr@shanghaitech.edu.cn

**Yuyao Chen**

ShanghaiTech University

chenyy6@shanghaitech.edu.cn

## Abstract

1 Expressive and natural human motion generation is a rewarding area of computer  
2 vision, which is because the generation is really a challenging task on account of  
3 the complex diversity of human motion, human perception of given text, and the  
4 difficulty of accurately describing human motion. However, we have many current  
5 methods to generate human motion such as CLIP, diffusion model, they have some  
6 limitations on generation based on more specific description, and are maybe not  
7 consistent for out-of-domain reference examples. We introduce a new generation  
8 model combining CLIP and DDPM training methods to train a more accurate and  
9 diverse model from multiple text input. Our goal is to generate a dynamic 3D  
10 motion that shows a continuous movement in a short period of time based on the  
11 user's text input, and to generate a random diversity of motions based on the same  
12 text input, showing randomness and diversity to meet different user needs. We  
13 evaluate our model on a dataset of human motion descriptions and compare it with  
14 a baseline approach. Our decoder conditional on the motion representation can  
15 also produce variants of the motion while preserving its semantics and style, while  
16 changing non-essential details that are not present in the motion representation. In  
17 addition, CLIP's joint embedding space supports language-guided motion manip-  
18 ulation in a zero-sample fashion. We used a diffusion model for the decoder and  
19 experimented with an autoregressive model and a diffusion model for the a priori  
20 model, finding that the latter was computationally more efficient and produced  
21 higher quality samples. As we demonstrate, our new generation model is a generic  
22 approach, enabling different modes of conditioning, and different generation tasks.  
23 We show that our model is trained with lightweight resources and yet achieves  
24 state-of-the-art results on leading benchmarks for text-to-motion.

## 25 1 Introduction

26 Motion generation from text is a recently emerging field of research. It involves the task of generating  
27 a sequence of frames representing a human motion from a text description. This can be used for  
28 applications such as virtual reality, animation, and video games. The task is challenging as it requires  
29 understanding of natural language, as well as the ability to generate realistic motion.

30 Recent progress in computer vision has been driven by scaling models on large datasets of captioned  
31 images collected from the internet. Within this framework, CLIP[26] has emerged as a successful  
32 image representation learner. CLIP embedding has many desirable properties: they are robust to  
33 image distribution shifts, have impressive zero-sample capabilities, and are fine-tuned to achieve  
34 state-of-the-art results on a variety of visual and linguistic tasks. Meanwhile, diffusion models have  
35 emerged as a promising framework for generative modelling, driving the latest developments in  
36 image and video generation tasks. To obtain optimal results, diffusion models utilise a bootstrapping  
37 technique that improves sample fidelity (for image, photo-level realism) at the expense of sample  
38 diversity. In this work, we combine these two approaches to solve the problem of text conditional  
39 motion generation. We first train a diffusion decoder to invert the CLIP encoder. Our inverter is  
40 non-deterministic and can generate multiple motions corresponding to a given motion embedding.

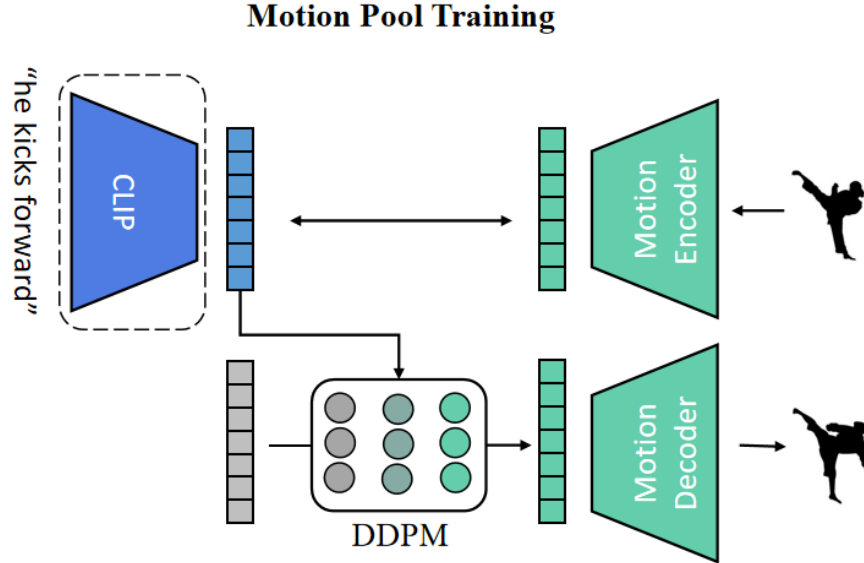


Figure 1: A overview of our human text-to-motion model

41 The presence of the encoder and its approximate inverse (decoder) allows for the ability to go beyond  
 42 text-to-motion translation.

43 In this work, we combine these two approaches to solve the problem of text conditional motion  
 44 generation. We first train a diffusion decoder to invert the CLIP encoder. Our inverter is non-  
 45 deterministic and can generate multiple motions corresponding to a given motion embedding. The  
 46 presence of the encoder and its approximate inverse (decoder) allows for the ability to go beyond  
 47 text-to-motion translation.

48 To obtain a complete model for image generation, we combine a CLIP motion embedding decoder  
 49 with a prior model that generates possible CLIP motion embedding from a given text caption. We  
 50 also develop methods for training diffusion priors in the latent space and show that they achieve  
 51 comparable performance to auto regressive priors, while being more computationally efficient.

## 52 2 Related Work

### 53 2.1 Human Motion Generation

54 Neural motion generation learned from motion-capture data can be conditioned by any signal  
 55 describing motion. Any signal that describes motion. So many methods use parts of the movement  
 56 itself for specific guidance. Some of the related previous work predict human motion from its prefix  
 57 poses[1, 2, 3, 4]. Others[5, 6, 7, 8] make use of bi-directional GRU[9] and Transformer[10] to solve  
 58 super-resolution and in-betweening tasks. Other work[11] uses the automatic encoder to learn the  
 59 latent representation of motion, and then use it to edit and control the motion with spatial constraints,  
 60 such as root locus and bone length.

61 Recently, another generative model that has recently attracted a lot of attention is NeRF([33, 34],  
 62 which has had considerable success in rendering realistic images. An implicit neural representation  
 63 (INR) is a series of neural networks that optimise their parameters to fit a sample rather than an entire  
 64 distribution. A major advantage is the technique’s ability to generalise extremely well in the spatial  
 65 or temporal dimension. For example, Cervantes[35] proposed an implicit scheme which models both  
 66 action categories and timestamps. Similar to the original NeRF, the timestamps are represented by  
 67 sine values. After supervised training, the proposed method can generate a variable-length motion  
 68 sequence for each action category.

69 Apart from input text, human motion can also be controlled with a high-level guidance from natural  
 70 language[12, 13], action class[14, 15, 16], audio[17]. For instance, recent works[18] generated dance

71 moves conditioned on music and the motion prefix, Edwards[19] generate facial expressions to fit  
72 spoken audio sequences. In most cases, the authors recommend a specialized approach to map each  
73 regulatory domain to human motion.

## 74 2.2 Text to Motion

75 In recent years, the dominant approach for text-to-motion tasks is to learn a shared latent space of  
76 language and motion. A motion-language dataset named KIT[20] offers about 11 hours of motion-  
77 capture sequences, each paired with a sentence that clearly describes the action being made. KIT  
78 sentences describe movement type, direction, and sometimes speed, but lack details about movement  
79 style and do not include abstract descriptions of movement. So Currently most of the researches are  
80 based on KIT. Yamada[21] learns both mappings by simultaneously training the text and motion  
81 auto-encoders to bind their latent spaces using both text and motion pairs. JL2P[22] learns the KIT  
82 motion-language dataset with an auto-encoder, constrained to a one-to-one mapping from text to  
83 motion, which has been improved in terms of the subtle concepts of text (i.e. speed, trajectory and  
84 type of action). They also learned to joint motion-text latent space and apply training curriculum to  
85 ease optimization. Lin[22] has further improved trajectory prediction by adding dedicated layers.

86 Another dataset BABEL[23] provides per-frame text labels sorted by 260 classes for the larger  
87 AMASS dataset [24], including approximately 40 hours of motion capture. Although a clear  
88 description of the action is provided, any detail beyond the type of action is usually missing, but this  
89 data covers a wider variety of human motions. MotionCLIP[25] extends text-to-motion data limits  
90 and enables latent space editing using shared text image latent spaces learned by CLIP[26].

## 91 2.3 Diffusion Generation Model

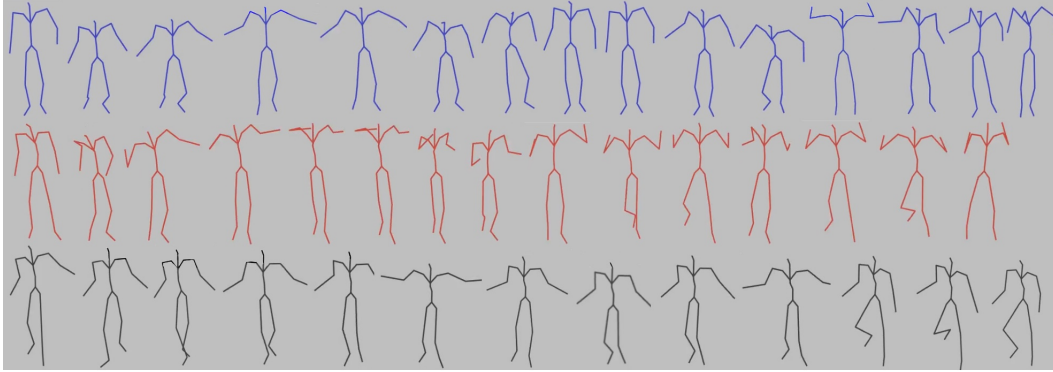
92 Diffusion models[27] are a class of neural generative models based on the stochastic diffusion process  
93 modelled in thermodynamics. In this setup, samples from the data distribution are gradually noised  
94 by the diffusion process. The neural model then learns the reverse process of gradually denoising the  
95 samples. For conditioned generation, Dhariwal & Nichol[29] to enable conditioning on CLIP textual  
96 representations. More recently, Zhang[30] and Kim[31] have suggested diffusion models for motion  
97 generation. For the MotionDiffuse model proposed by Zhang et al., unlike autoregressive inference  
98 schemes that typically require many long motion sequences for training, MotionDiffuse is able to  
99 model correlations between successive movements without introducing additional training costs.  
100 MDM[32] is a diffusion-based generative model carefully tuned for the human motion domain. Being  
101 diffusion-based, the MDM benefits from the local many-to-many domain representation described  
102 above, as evidenced by the quality and diversity of the resulting movements. Furthermore, MDM  
103 incorporates insights already established in the field of motion generation, helping it to become lighter  
104 and more controllable.

## 105 3 Criticism

106 During recent period of time, the MotionCLIP[25] method generated in 2022, it proposes a motion  
107 generation network that makes use of the knowledge encapsulated in CLIP to allow intuitive operations  
108 such as text conditional motion generation and editing. But it still have limitations in understanding  
109 directions like left, right and counter-clockwise. It also has difficulty in capture some styles(e.g.  
110 heavy and proud), and is of course not consistent of some cultural reference examples. For examples,  
111 this model fails to produce *Cristiano Ronaldo's* goal celebration and *batman's* signature pose.

112 MDM[32], a method suitable for a variety of human motion generation tasks. MDM is an atypical  
113 classifier-free diffusion model with a transformer encoder backbone and predicts the signal, rather  
114 than the noise. A significant limitation of the diffusion approach is the long inference time, with a  
115 single result requiring about 1000 forward passes. Since our motion model is small in any case, using  
116 a dimension an order of magnitude smaller than the image reduces its inference time from less than a  
117 second to about a minute.

118 Another obvious shortcoming is that although MDM used CLIP in the generation phase, it only  
119 masked CLIP randomly for classifier-free learning and did not joint embedding space of CLIP with  
120 text, in fact MDM did not link semantics and motion together well, their semantic generation was  
121 poor.



A person moves his hands together and hops in the air.

Figure 2: A complex example of motion generation of same input text

122 MotionDiffuse[30], the first text-driven motion generation method based on a diffusion model.  
 123 MotionDiffuse demonstrates three main advantages: probabilistic mapping to enhance diversity,  
 124 realistic synthesis to ensure rationalisation of motion sequences, and multi-level manipulation to  
 125 allow manipulation of each part and long sequence generation. Although MotionDiffuse pushes the  
 126 performance boundaries of motion generation tasks forward, a number of issues remain. Firstly,  
 127 diffusion models require a large number of diffusion steps during inference and generating motion  
 128 sequences in real time is challenging. Secondly, the current pipeline only accepts a single form of  
 129 motion representation. A more general pipeline that also adapts to all datasets would be more suitable  
 130 for a variety of scenarios.

#### 131 4 Numerical results

132 Our model successfully generate many diverse and interesting human motions given same input text,  
 133 I will show some image of this motion in this report(2, 3, 4) and specific demonstration video in our  
 134 code zip.

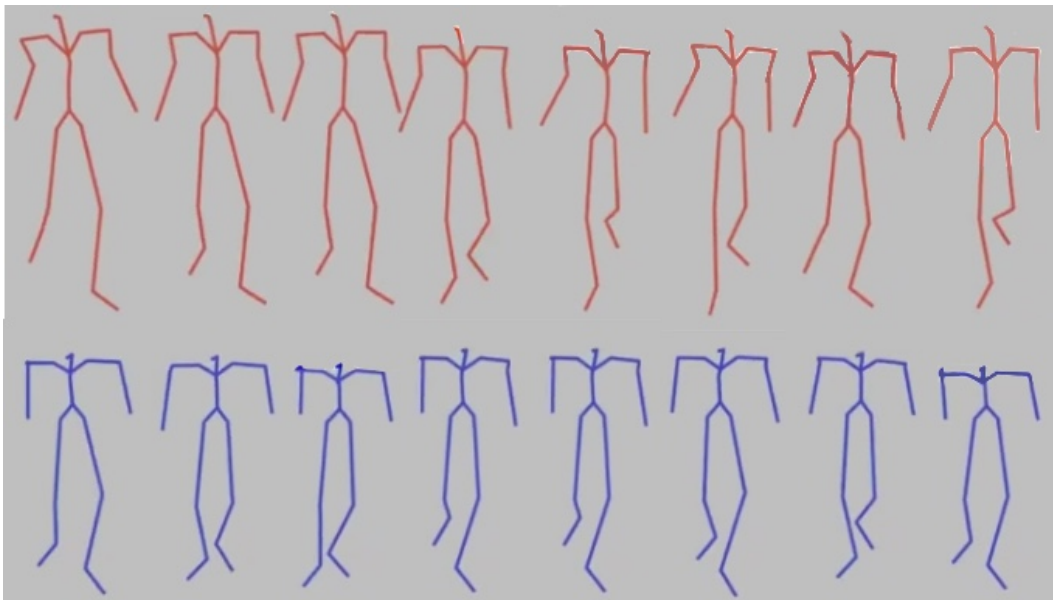
135 We use the HumanML3D[36] dataset for text-to-motion generation, examine MotionCLIP’s ability to  
 136 convert text into animation. Since the latent space of motion is aligned with CLIP, we use CLIP’s  
 137 pre-trained text encoder to process the input text and use MotionCLIP’s decoder to convert the  
 138 resulting latent embedding into motion. We then put the trained latent space into the DDPM model  
 139 to train the motion decoder, which is better than the original MDM model because we learn the  
 140 semantics and style of the given text better by using CLIP.

#### 141 5 Conclusion

142 In this paper, we proposed a novel approach for human motion generation from text. Our model takes  
 143 as input a text description of a motion, and outputs a sequence of motion frames that represent the  
 144 motion described in the text. To obtain a complete motion generation model, we combine the CLIP  
 145 motion embedding decoder with an a prior model that generates possible CLIP motion embedding  
 146 from a given text heading, and then this embedding is used to condition a diffusion decoder which  
 147 produces a final motion. Nevertheless, we see that the same autoencoder with the same data can  
 148 understand motion manifolds and their semantics significantly better, simply by aligning them with  
 149 well-behaved, knowledge-rich latent spaces.

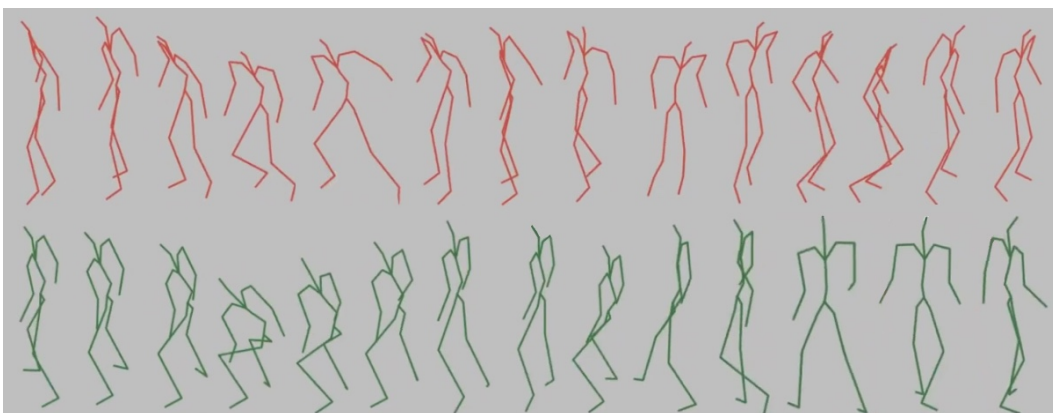
#### 150 References

151 [1] Alejandro Hernandez, Jurgen Gall, and Francesc Moreno-Noguer. Human motion prediction via spatio-  
 152 temporal inpainting. *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp.  
 153 7134–7143, 2019.



A person is walking backwards.

Figure 3: A simple example of motion generation of same input text



A person jumps 360 degree and turns back.

Figure 4: A simple example of motion generation of same input text

- 154 [2] Wen Guo, Yuming Du, Xi Shen, Vincent Lepetit, Xavier Alameda-Pineda, and Francesc Moreno Noguera.  
155 Back to mlp: A simple baseline for human motion prediction. *arXiv preprint arXiv:2207.01567*, 2022b.
- 156 [3] Mathis Petrovich, Michael J. Black and, Gül Varol. Action-Conditioned 3D Human Motion Synthesis with  
157 Transformer VAE. *International Conference on Computer Vision (ICCV)*. 10985–10995, 2021.
- 158 [4] Katerina Fragkiadaki, Sergey Levine, Panna Felsen, and Jitendra Malik. Recurrent network models for  
159 human dynamics. *Proceedings of the IEEE international conference on computer vision*, pp. 4346–4354, 2015.
- 160 [5] Manuel Kaufmann, Emre Aksan, Jie Song, Fabrizio Pece, Remo Ziegler, and Otmar Hilliges. Con volutional  
161 autoencoders for human motion infilling. *International Conference on 3D Vision (3DV)*, pp. 918–927. IEEE,  
162 2020.
- 163 [6] Felix G Harvey and Christopher Pal. Recurrent transition networks for character locomotion. *SIGGRAPH*  
164 *Asia 2018 Technical Briefs*, pp. 1–4, 2018.
- 165 [7] Felix G Harvey, Mike Yurick, Derek Nowrouzezahrai, and Christopher Pal. Robust motion in between-  
166 ing. *ACM Transactions on Graphics (TOG)*, 39(4):60–1, 2020.
- 167 [8] Yinglin Duan, Tianyang Shi, Zhengxia Zou, Yenan Lin, Zhehui Qian, Bohan Zhang, and Yi Yuan. Single-shot  
168 motion completion with transformer. *arXiv preprint arXiv:2103.00776*, 2021.
- 169 [9] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger  
170 Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical  
171 machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- 172 [10] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz  
173 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30,  
174 2017.
- 175 [11] Daniel Holden, Jun Saito, and Taku Komura. A deep learning framework for character motion synthesis  
176 and editing. *ACM Transactions on Graphics (TOG)*, 35(4):1–11, 2016.
- 177 [12] Chaitanya Ahuja and Louis-Philippe Morency. Language2pose: Natural language grounded pose forecasting.  
178 *International Conference on 3D Vision (3DV)*, pp. 719–728. IEEE, 2019.
- 179 [13] Mathis Petrovich, Michael J. Black, and Gül Varol. TEMOS: Generating diverse human motions from  
180 textual descriptions. *European Conference on Computer Vision (ECCV)*, 2022.
- 181 [14] Mathis Petrovich, Michael J. Black, and Gül Varol. Action-conditioned 3D human motion synthesis with  
182 transformer VAE. *International Conference on Computer Vision (ICCV)*, pp. 10985–10995, 2021.
- 183 [15] Chuan Guo, Xinxin Zuo, Sen Wang, Shihao Zou, Qingyao Sun, Annan Deng, Minglun Gong, and Li Cheng.  
184 Action2motion: Conditioned generation of 3d human motions. *Proceedings of the 28th ACM International*  
185 *Conference on Multimedia*, pp. 2021–2029, 2020.
- 186 [16] Pablo Cervantes, Yusuke Sekikawa, Ikuro Sato, and Koichi Shinoda. Implicit neural representations for  
187 variable length human motion generation. *arXiv preprint arXiv:2203.13694*, 2022.
- 188 [17] Ruilong Li, Shan Yang, David A. Ross, and Angjoo Kanazawa. Ai choreographer: Music conditioned 3d  
189 dance generation with aist++. *The IEEE International Conference on Computer Vision (ICCV)*, 2021.
- 190 [18] Andreas Aristidou, Anastasios Yiannakidis, Kfir Aberman, Daniel Cohen-Or, Ariel Shamir, and Yiorgos  
191 Chrysanthou. Rhythm is a Dancer: Music-Driven Motion Synthesis with Global Structure. *arXiv preprint*  
192 *arXiv:2111.12159*, 2021.
- 193 [19] Pif Edwards, Chris Landreth, Eugene Fiume and Karan Singh. JALI: an animator centric viseme model for  
194 expressive lip synchronization. *ACM Transactions on graphics (TOG)* 35, 4, 1–11, 2016.
- 195 [20] Matthias Plappert, Christian Mandery, and Tamim Asfour. The kit motion-language dataset. *Big data*,  
196 4(4):236–252, 2016.
- 197 [21] Tatsuro Yamada, Hiroyuki Matsunaga, and Tetsuya Ogata. Paired recurrent autoencoders for bidirectional  
198 translation between robot actions and linguistic descriptions. *IEEE Robotics and Automation Letters* 3, 4 (2018),  
199 3441–3448, 2018
- 200 [22] Chaitanya Ahuja and Louis-Philippe Morency. Language2pose: Natural language grounded pose forecasting.  
201 *2019 International Conference on 3D Vision (3DV)*, pp. 719–728. IEEE, 2019.
- 202 [22] Angela S Lin, Lemeng Wu, Rodolfo Corona, Kevin Tai, Qixing Huang, and Raymond J Mooney. Generating  
203 animated videos of human activities from natural language descriptions. *Learning2018*, 1, 2018.

- 204 [23] Abhinanda R. Punnakkal, Arjun Chandrasekaran, Nikos Athanasiou, Alejandra QuirosRamirez, and Michael  
205 J. Black. BABEL: Bodies, Action and Behavior with English Labels. *Proceedings IEEE/CVF Conf. on Computer  
206 Vision and Pattern Recognition (CVPR)*. 722–731. 2021.
- 207 [24] Naureen Mahmood, Nima Ghorbani, Nikolaus F. Troje, Gerard Pons-Moll, and Michael J. Black. AMASS:  
208 Archive of Motion Capture as Surface Shapes. *International Conference on Computer Vision*. 5442–5451, 2019.
- 209 [25] Guy Tevet, Brian Gordon, Amir Hertz, Amit H Bermano, and Daniel Cohen-Or. Motionclip: Exposing  
210 human motion generation to clip space. *it arXiv preprint arXiv:2203.08063*, 2022.
- 211 [26] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish  
212 Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural  
213 language supervision. *International Conference on Machine Learning*, pp. 8748–8763. PMLR, 2021.
- 214 [27] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning  
215 using nonequilibrium thermodynamics. *International Conference on Machine Learning*, pp. 2256–2265. PMLR,  
216 2015.
- 217 [28] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in  
218 Neural Information Processing Systems*, 34:8780–8794, 2021.
- 219 [29] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya  
220 Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided  
221 diffusion models. *arXiv preprint arXiv:2112.10741*, 2021.
- 222 [30] Mingyuan Zhang, Zhongang Cai, Liang Pan, Fangzhou Hong, Xinying Guo, Lei Yang, and Ziwei Liu.  
223 Motiondiffuse: Text-driven human motion generation with diffusion model. *arXiv preprint arXiv:2208.15001*,  
224 2022.
- 225 [31] Jihoon Kim, Jiseob Kim, and Sungjoon Choi. Flame: Free-form language-based motion synthesis & editing.  
226 *arXiv preprint arXiv:2209.00349*, 2022.
- 227 [32] Tevet G, Raab S, Gordon B, et al. Human motion diffusion model[J]. *arXiv preprint arXiv:2209.14916*,  
228 2022.
- 229 [33] Mildenhall B, Srinivasan PP, Tancik M, Barron JT, Ramamoorthi R, Ng R. Nerf: Representing scenes as  
230 neural radiance fields for view synthesis. *European conference on computer vision*, Springer, pp 405–421,2020.
- 231 [34] Jain A, Tancik M, Abbeel P. Putting nerf on a diet: Semantically consistent few-shot view synthesis.  
232 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp 5885–5894, 2021.
- 233 [35] Cervantes P, Sekikawa Y, Sato I, Shinoda K. Implicit neural representations for variable length human  
234 motion generation. *arXiv preprint arXiv:220313694*, 2022.
- 235 [36] Chuan Guo, Shihao Zou, Xinxin Zuo, Sen Wang, Wei Ji, Xingyu Li, and Li Cheng. Generating diverse and  
236 natural 3d human motions from text. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
237 Recognition*, pages 5152–5161, 2022.