

Skeleton-aided Articulated Motion Generation

Yichao Yan, Jingwei Xu, Bingbing Ni, Xiaokang Yang

Abstract—This work makes the first attempt to generate articulated human motion sequence from a single image. On one hand, we utilize paired inputs including human skeleton information as motion embedding and a single human image as appearance reference, to generate novel motion frames based on the conditional GAN infrastructure. On the other hand, a triplet loss is employed to pursue appearance smoothness between consecutive frames. As the proposed framework is capable of jointly exploiting the image appearance space and articulated/kinematic motion space, it generates realistic articulated motion sequence, in contrast to most previous video generation methods which yield blurred motion effects. We test our model on two human action datasets including KTH and Human3.6M, and the proposed framework generates very promising results on both datasets.

Index Terms—Motion Generation, Skeleton Aid, Video Analysis.

I. INTRODUCTION

Object motion prediction and generation in the videos is a key factor in video analysis, and it has potential application to smart surveillance, human-computer interaction and other applications. Generative models such as GAN [1] have achieved great success on image generation, but how to generate videos with motion dynamics is rarely explored. Although recent developments of convolutional neural network (CNN) and recurrent neural network (RNN/LSTM [2]) have made great success on action classification task [3], [4], [5], motion generation is still challenging because it often involves high-dimensional data with complex temporal dynamics. In particular, previous video generation methods [6], [7], [8], [9] are only good at simulating rigid movement of objects. In the case of articulated movement (e.g., human motion), these methods mostly yield blurred effects for various body parts.

Existing video generation methods mainly focus on two tasks. The first one is video prediction [9], [6], [10], [11], [12], [13], i.e., the models need to learn the motion patterns from a sequence of observed frames and to predict/generate the next frames. These methods are usually based on a recurrent structure (RNN or LSTM), despite of the good ability of the RNN/LSTM to model sequential data, they usually achieve good results only for short-term predictions where the videos are simple and quiet predictable. While the long-term prediction results usually suffer from low image quality, such as blur and object deformation. The second type of methods aim to directly generate a sequence of frames based on a single input [14] or only the scene types [8]. This task is more challenging as the motion patterns can no longer be observed during test phase. These methods employ the GAN model to generate the spatio-temporal cuboids or employ the Variational Autoencoders [15] to forecast the dense trajectory of pixels in the scene. However, the objects in the scene can move

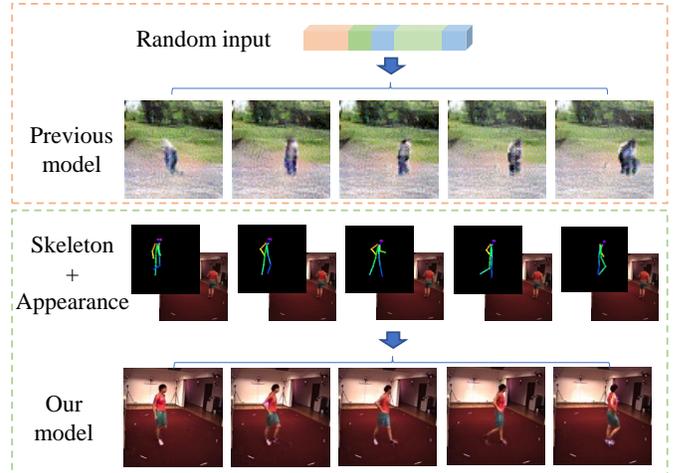


Fig. 1. Overview of the proposed framework. Previous video generation methods do not explicitly model the articulated structure of the foreground motion objects, the generated videos usually suffer from great deformation (see the top part). The bottom part gives an illustration of our approach, and we utilize the appearance information as well as the skeleton information to generate motion sequence.

arbitrarily if no geometric constrains are given to foreground objects, which will result in great deformation on the generated objects. Furthermore, we observe a shared limitation over these two type of methods, i.e., the articulated structures of the foreground motion objects (i.e., human) are not well modeled in the generation model. As previous generative methods only take the whole appearance as input, it will be difficult for the model to learn the structural relationship among the articulations/parts if given no supervision, thus resulting in great deformation during motion. Limited by this constraint, the quality of the generated videos are far from satisfaction.

In this work, we propose to use skeleton information to help generate articulated motion, which is motivated by the following observations. On the one hand, the articulated motion is usually under strong structure/geometric constrain, which can be well represented by the skeletons. On the other hand, compared to images with high dimension, the skeleton (coordinates of body parts) serves as a very good low dimensional embedding for human motion. Therefore it could be used as underlying status parameters to generate flexible poses. Also, skeletons can be mapped to image one-by-one, this avoids the long-term prediction problem shared by previous methods. Moreover, recent development on human pose estimation techniques has made skeleton data easy to access, thus avoiding heavy human annotations.

Our problem is defined as follows. Suppose we have a sequence of skeletons and a single image of human appearance, the task is to generate a sequence of articulated motion images.

This task can be further decomposed into two sub-problems: 1) how to generate realistic articulated motion (i.e., instead of blobs); and 2) how to generate appearance, which is adapted to every generated image frame. On one hand, motion generation (i.e., generating different human pose) is addressed by a GAN-like network. Recently, GAN has achieved great success in image generation, domain adaption and most importantly, image-to-image translation. The motion generation process can be naturally transformed into skeleton-to-image translation problem, which can be naturally handled by a conditional GAN model (i.e., in this work, we employ a GAN loss and an L1 loss to ensure smooth image-to-image translation). On the other hand, if given no appearance information (cloth color, bodily form), appearance of the generated image sequences cannot be controlled, i.e., the generated appearances might differ from image to image. This violates the rule that the appearance of an object should be consistent during the entire motion sequence. To address this issue, we choose to generate the motion sequence based on both skeleton sequence and an appearance image, which is realized with a specially designed generator. Furthermore, in order to ensure inter frame continuity, and we also employ a triplet loss which aims to penalize the generation loss if the adjacent frames have larger distance than the non-adjacent frames.

This work is among the very few works for complete video generation. To the best of our knowledge, this is the first work to employ skeleton information to help generate videos. Our key contribution is that the proposed method can generate images/videos with large scale geometry change, where previous methods achieve little success. We test our model on two human action datasets including KTH and Human3.6M, and the proposed framework generates very promising results on both datasets.

II. RELATED WORK

Human Motion Analysis. In computer vision, human motion analysis is a broad concept which focuses on the understanding and applications of human motion patterns. It has been receiving increasing attention for a long time and has been applied to enormous applications, such as content-based video retrieval, visual surveillance, and man-computer interfaces [16], [17], [18]. Among the researches of human motion analysis, variant human body models (e.g., stick figure model [19], [20], cardboard model [21], 3D volumetric model [22]) play an important role. These models involve low-level processes on human body structures and cover the kinematic properties of the body, which build the foundation in solving different problems, including human motion tracking [23], [24], [25], action recognition [26], [27], [28], [29], and pose estimation [30], [31], [32]. Motivated by these successful applications based on well-defined human body models, in this paper, we pay attention to another more challenging task that attempts to generate consistent human motion sequences based on the correspondent skeleton and appearance information.

Image Generation. Early works for image generation usually make efforts on simple texture synthesis with hand-crafted features [33]. During the past few years, two generation

models have been attracting more and more attention, i.e., the variational autoencoder (VAE) [15] and the generative adversarial network (GAN) [1]. VAE is a classical method which aims to model complicated distribution and it has been widely applied in various generative tasks. Gregor et al. [34] propose a sequential generative model which extends the original VAE with recurrent neural networks and attention mechanism. Another interesting model is proposed by Yan et al. [35], they develop a layered generative model based on conditional VAE. GAN is also a popular generative model and many recent works are built on it. Some works improve the architecture of original GAN for better performances [36], [37], [38], [39]. Conditional generative adversarial network (CGAN) [40] gives extra information to the input as condition, and the output is constrained by the input conditions. CGAN has been further extended by [41], [42], [43] to solve the image-to-image translation problem, which gives inspiration of our model proposed in this paper. Other applications such as image super-resolution [44], image edition [45] and unsupervised representation learning [46] also show impressive results.

Video Generation. Our problem is closely related to video generation or prediction. Video texture based methods [47], [48], [49] can generate periodic motion sequences if an input reference video is given. Lotter et al. [10] propose a predictive neural network motivated by the concepts from neuroscience. Finn et al. [7] develop an action-conditioned video prediction model which concentrates on pixel motion and Mathieu et al. [11] introduce multi-scale architecture to reduce the deformation in the predictions. Instead of focusing on pixel level prediction, Van et al. [50] attempt to predict the transformations between frames. Some works also utilize GAN or VAE in video generation. Vondrick et al. [8] propose a GAN model which generates static background and dynamic foreground sequences separately. Xue et al. [51] introduce conditional VAE and build a cross convolutional network which encodes image and motion information for generation. Although variant models are proposed, the results of these methods are usually limited by two issues: 1) the deformation of the foreground object is serious, 2) and the inter-frame consistency cannot be well maintained. These problems inspire us that in order to generate more realistic videos, strong motion constrains are needed during the generation process. Therefore, we employ skeleton information to guide our model for motion generation.

III. METHOD

Our problem is defined as follows. Given an image x containing the foreground person (the appearance reference image), we would like to generate a sequence of images $Y = \{y_1, \dots, y_n\}$ that share the same appearance, and the foreground objects should keep a specific motion pattern as well (e.g., walking, running). In other words, we would like to generate articulated motion from a single static image. This is challenging because a person could have infinite move patterns. The first step is to choose a specific move pattern for the sequence. As skeleton can well represent a person's

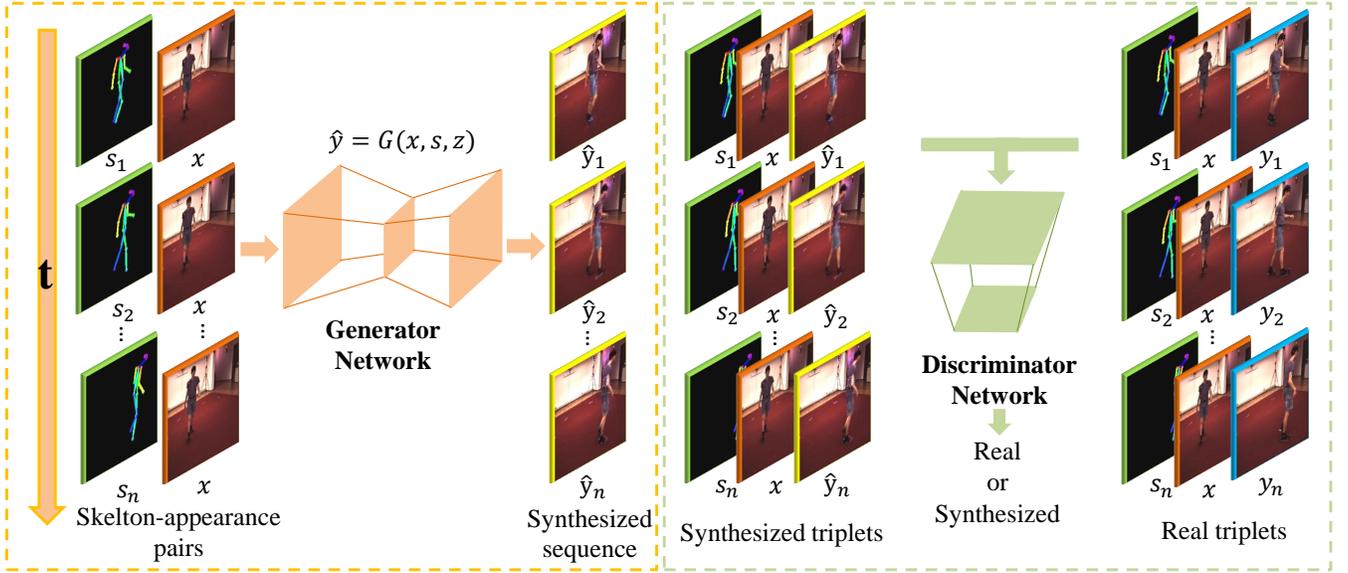


Fig. 2. Architecture of the generation and discrimination network. The inputs for the generator are the skeleton-appearance pairs (s, x) and generate the synthesized sequence $\hat{Y} = \{\hat{y}_1, \dots, \hat{y}_n\}$. The discriminator D tries to distinguish real triplets (x, s, y) and synthesized triplets (x, s, \hat{y}) .

motion, we employ a sequence of skeletons $S = \{s_1, \dots, s_n\}$ as prior knowledge for the motion. The remaining problem is to generate the output sequence based on the appearance image and skeleton sequence: $\{x, S\} \rightarrow Y$. Here, we propose a skeleton conditioned GAN to model this mapping. The second problem is to maintain appearance-smoothness between consecutive frames, we employ a triplet loss on the generator for this purpose. The details of the proposed method are given in the following of this section.

A. Skeleton Conditioned GAN

Different from the previous image-to-image translation model, where only a single input image is mapped to the output, our generative model is conditioned on two inputs, i.e., the appearance reference image x and the skeleton image s . The value function of the Conditional GAN (CGAN) model is expressed as follows:

$$\mathcal{L}_c(G, D) = \mathbb{E}_{x, s, y \sim p_{data}(x, s, y)} [\log D(x, s, y)] + \mathbb{E}_{x, s \sim p_{data}(x, s), z \sim p_z(z)} [\log(1 - D(x, s, G(x, s, z)))], \quad (1)$$

where the generator G tries to produce a new frame, and a discriminator D tries to distinguish real triplets (x, s, y) and synthesized triplets $(x, s, G(x, s, z))$. The architectures of the generator and discriminator model are illustrated in Figure 2.

Previous methods [41], [11] observe that using the CGAN loss alone will give sharp results with artifacts, and it's beneficial adding a contractive loss such as L1 loss. Although this may cause blur effect, mixing these two losses will generate overall better results. The L1 loss can be expressed as:

$$\mathcal{L}_{L_1}(G) = \mathbb{E}_{x, s \sim p_{data}(x, s), z \sim p_z(z)} [\|y - G(x, s, z)\|_1]. \quad (2)$$

Notice that our goal is to generate continuous motions rather than individual images. Thus it is important to consider the

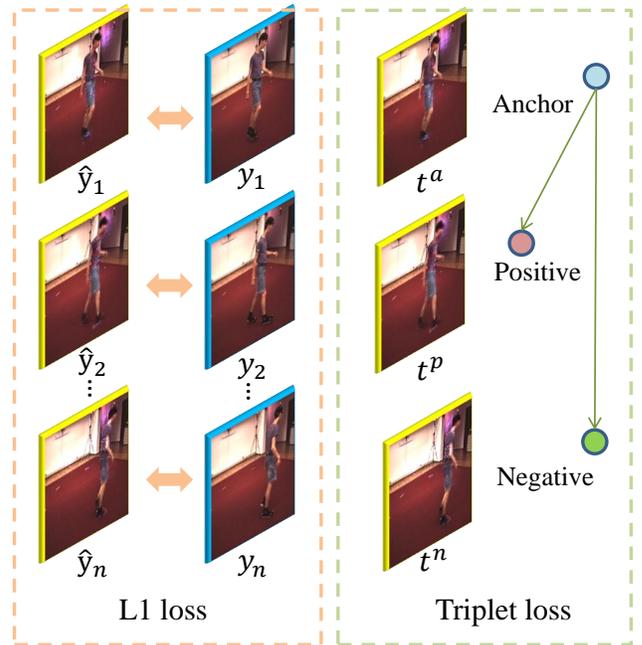


Fig. 3. Loss terms of the model. Despite of the GAN loss, we take two additional loss terms, i.e., the L1 loss to enhance the image-to-image translation quality, and the triplet loss to guarantee the continuity of the generated motion sequence.

motion continuity and appearance smoothness of the adjacent frames. Although this can be achieved by training a perfect generator that precisely maps the input appearance image into a new pose specified by the input skeleton, the perfect generator cannot be achieved in practice because a single appearance image does not contain the complete information of the moving object. e.g., for an image of a person, some parts are inevitably occluded, which can't be generated perfectly through training. To address this issue, we propose to utilize a

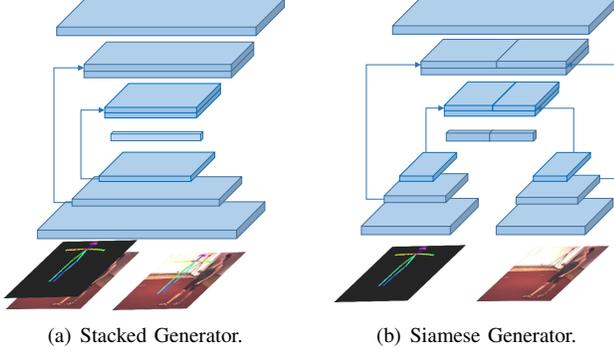


Fig. 4. Structures of the generator. For both the stacked generator and the Siamese generator, we take the U-Net structure for both generators.

triplet loss that motivates adjacent frames to be more similar than the far-away frames. First, we need to construct the triplet set $\mathcal{T} = \{t_i^a, t_i^p, t_i^n\}_{i=1}^m$ from the generated samples: $\hat{Y} = \{\hat{y}_1, \dots, \hat{y}_n\}$, where $\hat{y}_j = G(x, s_j, z)$, $1 \leq j \leq n$, and m is the number of triplets selected for training. The anchor of the triplet can be randomly chosen from the generated samples, say $t_i^a = \hat{y}_j$, the positive sample can select a sample adjacent to the anchor (e.g., $t_i^p = \hat{y}_{j+1}$), and the negative sample can choose a far-away sample (e.g., $t_i^n = \hat{y}_{j+5}$). We would like that the distance between anchor and positive is smaller than that of anchor and negative, thus the loss function can be expressed as:

$$\mathcal{L}_{tri}(G) = \sum_{i=1}^m [\|t_i^a - t_i^p\|_2^2 - \|t_i^a - t_i^n\|_2^2 + \alpha]_+. \quad (3)$$

In our experiments, we also tried to replace the L2 norm with L1 norm in the triplet loss, but we don't observe performance gain. The loss terms are illustrated in Figure 3.

The overall objective function is:

$$\mathcal{L}(G, D) = \mathcal{L}_c(G, D) + \lambda \mathcal{L}_{L_1}(G) + \beta \mathcal{L}_{tri}(G), \quad (4)$$

where λ and β are the weights for different loss terms. We aim to solve

$$G^* = \arg \min_G \max_D \mathcal{L}(G, D). \quad (5)$$

B. Generator architecture

The basic structure of our network is built upon [41], in which the generator usually follows the encoder-decoder structure. Specifically, the U-net structure [52] has proven to be more effective for image-to-image translation tasks, because it adds skip connections between encoder and decoder. Such a structure enables the information to share between the inputs and outputs, and thus achieving success for tasks such as image colorization, image segmentation, etc. We also adopt this structure for the generator. In our case, the generator needs to translate two inputs (i.e., the appearance reference image x and the skeleton image s) into a single output y . Therefore, the encoder takes two inputs simultaneously and shuttles the information to the decoder.

Stacked Generator. There are multiple options to design the structure of the encoder. The most straightforward way is to

TABLE I
DETAILED STRUCTURE OF THE GENERATOR.

Encoder	
Layer	Input size: $256 \times 256 \times 6(3)$
1	Conv-(64, $K4 \times 4, S2$), lReLU-(0.2)
2	Conv-(128, $K4 \times 4, S2$), BN, lReLU-(0.2)
3	Conv-(256, $K4 \times 4, S2$), BN, lReLU-(0.2)
4-8	Conv-(512, $K4 \times 4, S2$), BN, lReLU-(0.2)
Decoder	
Layer	Input size: $60 \times 60 \times 1(2)$
1	FConv-(512, $K4 \times 4, S2$), BN, Dropout-(0.5), ReLU
2-5	FConv-(1024, $K4 \times 4, S2$), BN, Dropout-(0.5), ReLU
6	FConv-(512, $K4 \times 4, S2$), BN, Dropout-(0.5), ReLU
7	FConv-(256, $K4 \times 4, S2$), BN, Dropout-(0.5), ReLU
8	FConv-(128, $K4 \times 4, S2$), BN, Dropout-(0.5), ReLU

stack the two input images as a single input (i.e., $i = [x, s]$), thus the standard encoding network can be applied directly. We denote this structure as **stacked generator 1**. We further observe that the skeleton images have completely black ground which is not informative at all, and all the motion information is contained in the pose of the skeleton. Therefore, we can directly draw the skeleton on the appearance image instead. The resulted image sequence can be viewed as a skeleton moving on the appearance image. In this case, the inputs for the encoder are standard images, and have no distinction with the traditional image-to-image task. We denote this structure as **stacked generator 2**. The two stacked structures are illustrated in Figure 4(a).

The stacked encoders provide a simple solution for simultaneously encoding two input images by stacking the two input as an ensemble. However, the two inputs usually contain different information, which needs to be modeled separately. To this end, we design a second structure.

Siamese Generator. In particular, each input image can be modeled by an encoding network, and the features are concatenated at the bottleneck layer for the decoder. In this way, the different information of both the appearance image and the skeleton image can be well modeled respectively. We denote this structure as the **Siamese generator**, because it has Siamese structure. See Figure 4(b) for an illustration.

The detailed structure of the generator is illustrated in Table I, where the encoder is composed of convolution (Conv), batch normalization (BN) layers and leaky rectified linear unit (lReLU) layers, and the decoder is composed of fractional length convolutional (FConv) layers, BN layer, Dropout layer and the ReLU layers. For gray-scale inputs, we replicate their channel 3 times so that the network doesn't distinguish RGB inputs and gray-scale inputs. We resize all the input images into a fix size, i.e., $256 \times 256 \times 3$. Therefore, the input size for the stacked generators is $256 \times 256 \times 6$ because they stack two inputs. And for the Siamese generator, the number of channel is 3 all the inputs. Also notice that each of the Siamese structure has exactly the same structure as the encoder, the output of each encoder is concatenated as inputs ($60 \times 60 \times 2$) for the decoder.

C. Discriminator architecture

The discriminator needs to be able to classify the realistic triplets (x, s, y) from the synthesized triplets $(x, s, G(x, s, z))$,

TABLE II
DETAILED STRUCTURE OF THE DISCRIMINATOR.

Discriminator	
Layer	Input size: $256 \times 256 \times 9$
1	Conv-(64, $K4 \times 4, S2$), IReLU-(0.2)
2	Conv-(128, $K4 \times 4, S2$), BN, IReLU-(0.2)
3	Conv-(256, $K4 \times 4, S2$), BN, IReLU-(0.2)
4-6	Conv-(512, $K4 \times 4, S2$), BN, IReLU-(0.2)

thus the inputs for the discriminator network are three images. Similar to the encoding network, we can also stack all the three inputs as a 9-dimensional input, or we can extract features from each inputs, and then combine them to make a decision. However, we observe little difference using these two structures during our experiments, therefore we only report the results of the stacked structure. The discriminator structure is illustrated in Table II.

D. Learning and Implementation

Our implementation is based on a modified version of image-to-image translation network [41] on Tensorflow [53]. We report results for several architectures. For all the models, we alternatively train the discriminator and then the generator. We train the generator and discriminator with stochastic gradient descent, with a fixed learning rate as 0.0002 and the Adam optimizer. All the models are trained for 30 epochs. We find that small batch size leads to more appealing results. So the batch size is set as 10 for all the experiments to balance the generation quality and training time. All the videos are re-scaled into range $[-1, 1]$ as normalization. The random noise z is not explicitly sampled from Gaussian distribution, it appears only in the form of dropout, which is in consistent with [41]. The skeletons are extracted using the real-time human pose estimator [54]. And we also use the same estimator to detect the generated human motion sequence to compare with the ground truth pose estimation. Because we use three kinds of loss function to train the network, we try different weights of each term in training phase to get optimal results. More detailed analysis on loss function is in section 4.3.

IV. EXPERIMENTS

In this section, we present extensive experimental evaluations and in-depth analysis of the proposed method. The evaluations are performed on the following two human action datasets:

KTH dataset. This dataset contains several types of human actions in the outdoor environment, and all the videos are gray scale. We experiment on three types of actions, i.e., *walking*, *running* and *hand waving*. And we choose these actions because their motion parts are different, *walking* mainly involve the movement of legs, *hand waving* involve the movement of arms, and *running* involves both the movement of arms and legs. Each type of action contains 100 videos, and we divide the dataset into training set containing 80 videos and test set containing 20 videos.

Human3.6M dataset. This dataset was collected in an indoor environment, there are 4 cameras working simultaneously, i.e., we have access to 4 views of each action. We use

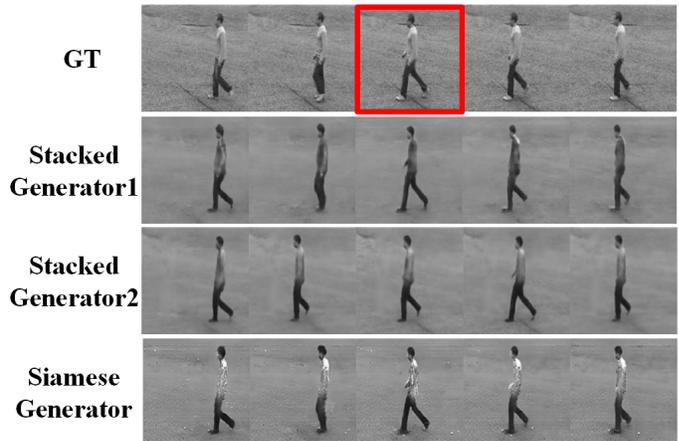


Fig. 5. Examples generated by different generator structure. The first row contains the ground truth motion sequence, the image with red bounding box is the appearance reference image.

the *walking* scenario in this dataset. The dataset also contains foreground segmentations, therefore we experiment on both videos with and without background.

As there exists no standard evaluation protocol for image and video generation tasks, in order to evaluate the quality of the generated videos, we report both qualitative and quantitative results.

A. Structure Analysis

We propose the *stacked structure* and the *Siamese structure* for the generator. We evaluate the performance of the two structure in this part, and the results of different generator structures are compared in Figure 5. The image with red bounding box is the appearance reference image. Benefitting from the U-Net structure, all the generators achieve to generate realistic motion patterns. The two *stacked generators* generate similar results. However, they suffer from a major issue, i.e., the generated sequences do not preserve the appearance of the reference image any more. For example, the subject in the reference image is in white T-shirt and black pants. The results generated by the stacked generators are in dark cloth. In contrast, the results generated by the *Siamese generator* are more similar to the ground truth. This is mainly because that the stacked structure does not distinct the pose and appearance information, which is encoded by the same network. However, the appearance images and the skeleton images are under different distributions, encoding them with two separate networks will better fit the two different distributions. Therefore, the *Siamese generator* achieves better results. It not only encodes the motion pattern, but also successfully encodes the appearance information, thus the generated sequence shares high appearance similarity with the reference image. In the following of this paper, we report the results of the *Siamese generator*.

B. Motion Generation

Some generated sequences on KTH dataset are visualized in Figure 6. For each example, the first row is the ground-truth

sequence, and the second row represents the detected skeleton sequence. The generated results are given in the third row. We re-run the skeleton detector on the generated sequence, and the resulted skeleton sequence is shown in the last row. We have three observations: 1) the foreground subjects are naturally generated in the scene. Different from [8], which employs a two stream model to separately generate the foreground and background, our model is a unified model and does not distinct foreground from backgrounds. Moreover, for each of the generated image, the boundary of the foreground subject is sharp and looks natural in the scene. Qualitatively, this is better than the results in [8], where the people in the scene are often blobs. 2) The model successfully generates motions patterns with high quality. We find that the generated motion sequences are highly recognizable for humans. Moreover, the skeletons extracted from the generated sequence are very similar from ground-truth skeletons. We make two remarks here. On one hand, this demonstrates that the generated pose is close to the objective. On the other hand, it in turn demonstrates that our model has successfully generated humans that can be recognized by the pose detector. 3) The identity of the appearance reference image is well preserved. For example, in Figure 6(a) and Figure 6(c), the appearance reference image of the subjects are with white and black clothes respectively. We can observe that the generated sequence share similar appearances with their reference images, and the appearance of the person is consistent in the generated sequence.

The generation results on Human3.6M dataset are visualized in Figure 7. Note that different from KTH, Human3.6M dataset contains color videos which are much more difficult to encode the appearance into latent space. We can observe that the appearance of generated sequence is close to the ground truth. The results demonstrate that the generator effectively transforms the color space into latent space, and finds good representation for the relationship between different body parts and their corresponding color. And we also present the results on Human3.6M without background in Figure 7(b). We can observe that the network indeed transforms the appearance from the target image to generated sequence, without referring to the background information.

Overall, the proposed method shows promising results in generating human motions. For more examples, please refer to our supplementary materials.

C. Component Analysis

In this part, we study how the loss terms in Equation 4 influence the generation results. The GAN loss and the L1 loss have been analyzed in [41] for image-to-image translation task, which demonstrates that using GAN loss alone will generate artifacts, and using L1 loss alone tends to generate the color averaged over the training set. The results in our experiments are illustrated in Figure 8. We can observe that the L1 loss and GAN loss have different effects for the generated videos. Specially, only using the GAN loss will result in severe artifacts in the video, some parts of the person are missing and color of the person keeps changing in the video. While only using L1 loss will degrade the generalization

TABLE III
RECOGNITION ACCURACIES (%) ON THE GENERATED SEQUENCES. GT DENOTES THE GROUND-TRUTH RESULTS. L1 DENOTES L1 LOSS, G DENOTES GAN LOSS, T DENOTES TRIPLET LOSS.

Action	GT	L1	G	L1+G	L1+G+T
Walking	100	83.1	43.1	93.8	93.8
Running	98.5	76.9	41.5	92.3	92.3
Hand waving	100	95.4	47.7	100	100

ability of the network. The third row of Figure 8 shows these results. No matter how we change the reference image, all the generated sequences have almost the same appearance, i.e., the average color of all the training samples. These adversary effects in these two kinds of losses are critical for training the generator. Combing L1 loss and GAN loss will generate better results, but the color is not consistent along the frames. Furthermore, introducing the triplet loss described in Section 3.1 will stabilize the performance of generator. As shown in bottom two rows of Figure 8, the generated sequence in bottom has consistent clothing color compared to the results which lack the help of triplet loss.

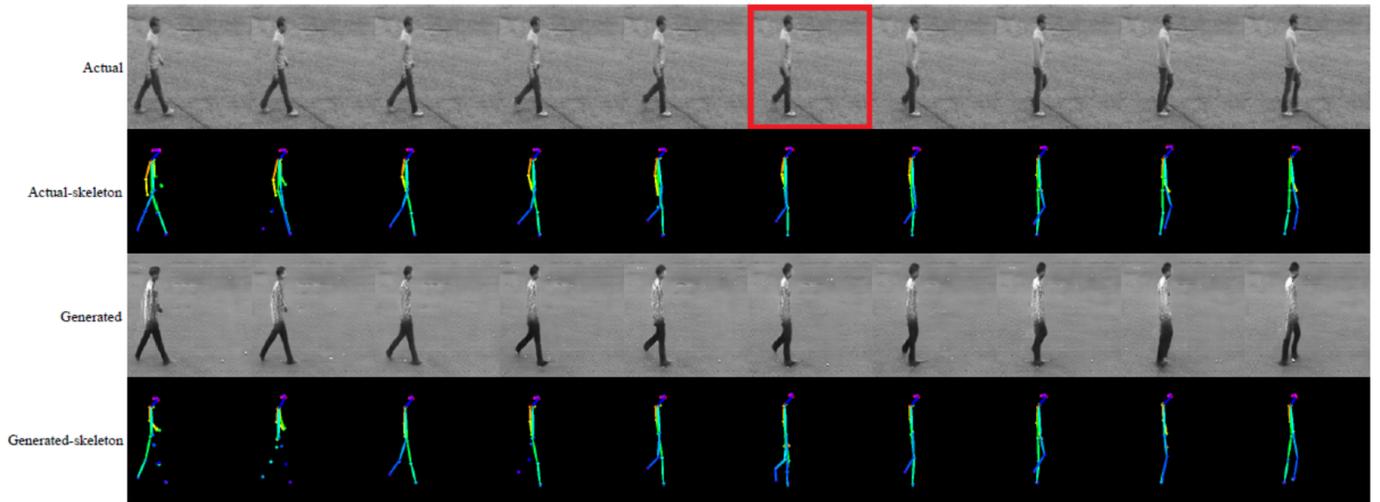
We also perform quantitative evaluations over the loss terms. We train an action classification framework [5] on the KTH dataset, and then use it to classify the generated sequences. If the generated sequences are classified into the right action type, it demonstrates that the generated sequences can be recognized by the off-the-shelf classifier. The results are shown in Table III. As KTH dataset is a relatively easy dataset for action recognition, the trained classifiers achieve very good performance ($\sim 100\%$) on the ground-truth test sequences. We observe that only using L1 or GAN loss leads to great performance drop, and combining the two loss terms significantly improves the recognition accuracy. Further adding triplet loss does not bring in further improvements, because the motion patterns are already clear enough for the classifier.

V. CONCLUSION

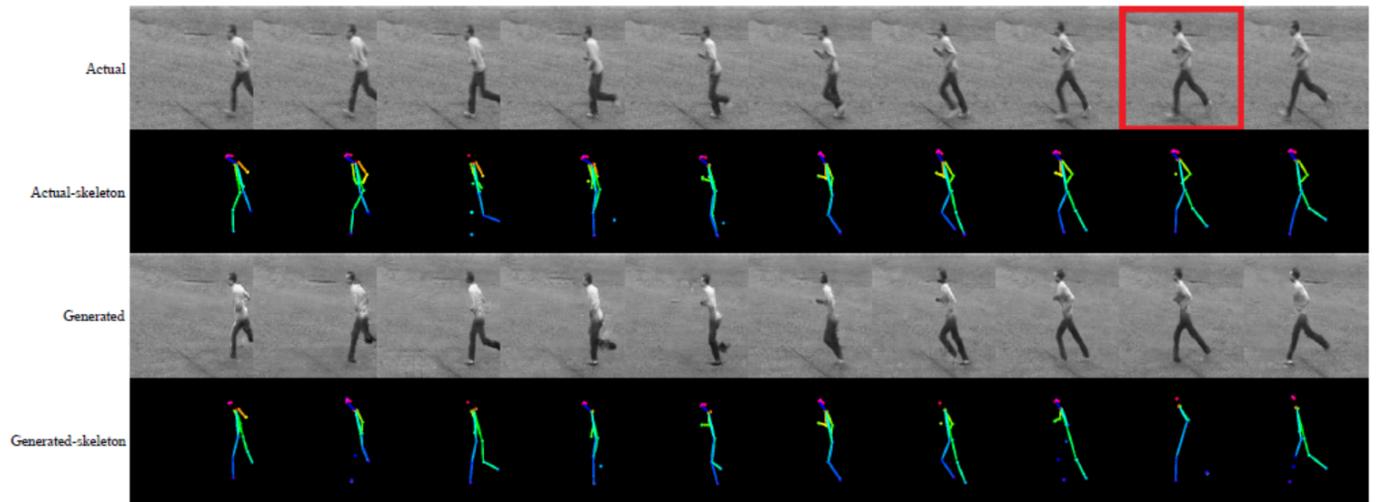
In this work, we propose a skeleton-aided articulated motion generation method. In contrast to existing video generation methods, which usually lack geometric constraints for the foreground object, our method utilizes skeleton information as a guidance for the geometric change during the motion. Experimental results show that by giving an appearance reference image and the skeleton sequence, our model can produce high-quality video sequences that not only preserve the appearance of the reference image, but also have clear motion pattern as the skeletons. The generated motion sequences are also recognizable by the off-the-shelf action recognition framework.

REFERENCES

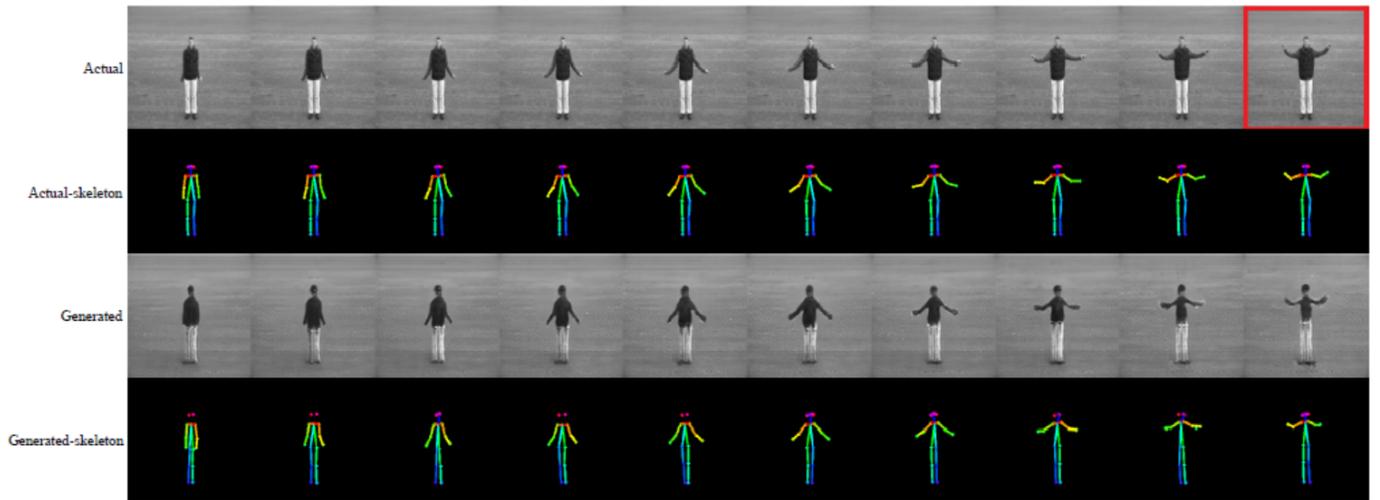
- [1] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. "Generative adversarial nets," in *Advances in neural information processing systems*, 2014, pp. 2672–2680.
- [2] Sepp Hochreiter and Jürgen Schmidhuber. "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [3] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Fei-Fei Li. "Large-scale video classification with convolutional neural networks," in *CVPR*, 2014, pp. 1725–1732.



(a) Walking.



(b) Running.



(c) Hand waving.

Fig. 6. Generated samples on KTH dataset. The first rows contain the ground truth motion sequences, and the images with red bounding boxes are the appearance reference images. The second rows are the ground truth reference skeletons. The third rows are the samples generated by our model, and the last rows are the detected skeletons of the generated sequences.

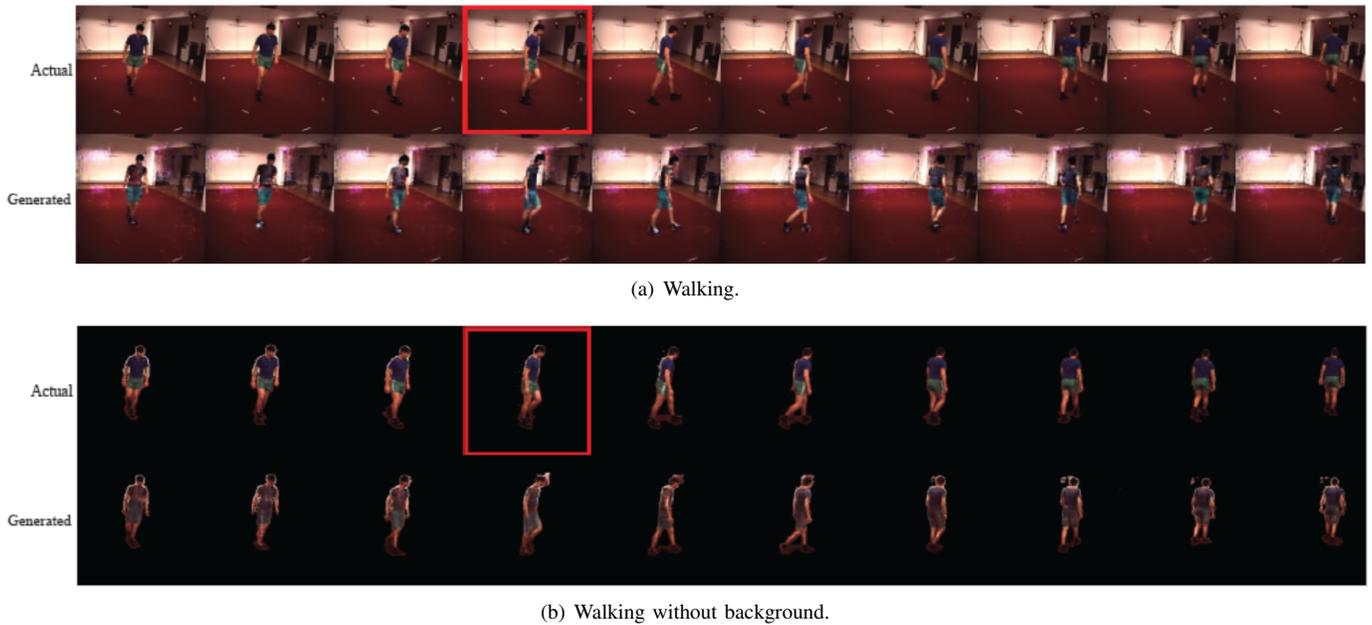


Fig. 7. Generation results on Human3.6M dataset. (a) Walking sequence with background. (b) The same sequence without background.

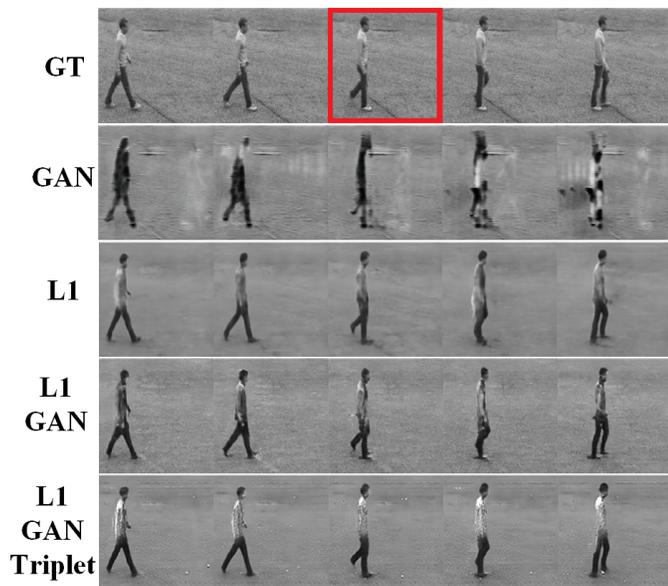


Fig. 8. Examples generated with different loss terms.

[4] Zuxuan Wu, Xi Wang, Yu-Gang Jiang, Hao Ye, and Xiangyang Xue, “Modeling spatial-temporal clues in a hybrid deep learning framework for video classification,” in *ACMMM*, 2015, pp. 461–470.

[5] Jeff Donahue, Lisa Anne Hendricks, Marcus Rohrbach, Subhashini Venugopalan, Sergio Guadarrama, Kate Saenko, and Trevor Darrell, “Long-term recurrent convolutional networks for visual recognition and description,” *TPAMI*, vol. 39, no. 4, pp. 677–691, 2017.

[6] Nitish Srivastava, Elman Mansimov, and Ruslan Salakhutdinov, “Unsupervised learning of video representations using lstms,” in *ICML*, 2015, pp. 843–852.

[7] Chelsea Finn, Ian Goodfellow, and Sergey Levine, “Unsupervised learning for physical interaction through video prediction,” in *NIPS*, 2016, pp. 64–72.

[8] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba, “Generating videos with scene dynamics,” in *NIPS*, 2016, pp. 613–621.

[9] Junhyuk Oh, Xiaoxiao Guo, Honglak Lee, Richard L. Lewis, and

Satinder P. Singh, “Action-conditional video prediction using deep networks in atari games,” in *NIPS*, 2015, pp. 2863–2871.

[10] William Lotter, Gabriel Kreiman, and David Cox, “Deep predictive coding networks for video prediction and unsupervised learning,” *CoRR*, vol. abs/1605.08104, 2016.

[11] Michaël Mathieu, Camille Couprie, and Yann LeCun, “Deep multi-scale video prediction beyond mean square error,” *CoRR*, vol. abs/1511.05440, 2015.

[12] Joost R. van Amersfoort, Anitha Kannan, Marc’Aurelio Ranzato, Arthur Szlam, Du Tran, and Soumith Chintala, “Transformation-based models of video sequences,” *CoRR*, vol. abs/1701.08435, 2017.

[13] Marc’Aurelio Ranzato, Arthur Szlam, Joan Bruna, Michaël Mathieu, Ronan Collobert, and Sumit Chopra, “Video (language) modeling: a baseline for generative models of natural videos,” *CoRR*, vol. abs/1412.6604, 2014.

[14] Jacob Walker, Carl Doersch, Abhinav Gupta, and Martial Hebert, “An uncertain future: Forecasting from static images using variational autoencoders,” in *ECCV*, 2016, pp. 835–851.

[15] Diederik P. Kingma and Max Welling, “Auto-encoding variational bayes,” *CoRR*, vol. abs/1312.6114, 2013.

[16] Ronald Poppe, “Vision-based human motion analysis: An overview,” *CVIU*, vol. 108, no. 1, pp. 4–18, 2007.

[17] Jake K Aggarwal and Quin Cai, “Human motion analysis: A review,” in *Nonrigid and Articulated Motion Workshop, 1997. Proceedings., IEEE*, 1997, pp. 90–102.

[18] Xiaofei Ji and Honghai Liu, “Advances in view-invariant human motion analysis: A review,” *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, vol. 40, no. 1, pp. 13–24, 2010.

[19] Shanon X Ju, Michael J Black, and Yaser Yacoob, “Cardboard people: A parameterized model of articulated image motion,” in *Automatic Face and Gesture Recognition, 1996., Proceedings of the Second International Conference on*, 1996, pp. 38–44.

[20] Nicholas R Howe, Michael E Leventon, and William T Freeman, “Bayesian reconstruction of 3d human motion from single-camera video,” in *NIPS*, 1999, pp. 820–6.

[21] Richard F Rashid, “Towards a system for the interpretation of moving light displays,” *TPAMI*, , no. 6, pp. 574–581, 1980.

[22] Roland Kehl and Luc Van Gool, “Markerless tracking of complex human motions from multiple views,” *CVIU*, vol. 104, no. 2, pp. 190–209, 2006.

[23] Yu Huang and Thomas S Huang, “Model-based human body tracking,” in *ICPR*, 2002, vol. 1, pp. 552–555.

[24] Eng-Jon Ong, Antonio S Micilotta, Richard Bowden, and Adrian Hilton, “Viewpoint invariant exemplar-based 3d human tracking,” *CVIU*, vol. 104, no. 2, pp. 178–189, 2006.

- [25] Ivana Mikić, Mohan Trivedi, Edward Hunter, and Pamela Cosman, “Human body model acquisition and tracking using voxel data,” *IJCV*, vol. 53, no. 3, pp. 199–223, 2003.
- [26] Yong Du, Wei Wang, and Liang Wang, “Hierarchical recurrent neural network for skeleton based action recognition,” in *CVPR*, 2015, pp. 1110–1118.
- [27] Georgios Th Papadopoulos, Apostolos Axenopoulos, and Petros Daras, “Real-time skeleton-tracking-based human action recognition using kinect data,” in *International Conference on Multimedia Modeling*, 2014, pp. 473–483.
- [28] Alexander Grushin, Derek D Monner, James A Reggia, and Ajay Mishra, “Robust human action recognition via long short-term memory,” in *IJCNN*, 2013, pp. 1–8.
- [29] Dian Gong, Gerard Medioni, and Xuemei Zhao, “Structured time series analysis for human action segmentation and recognition,” *TPAMI*, vol. 36, no. 7, pp. 1414–1427, 2014.
- [30] Mykhaylo Andriluka, Stefan Roth, and Bernt Schiele, “Pictorial structures revisited: People detection and articulated pose estimation,” in *CVPR*, 2009, pp. 1014–1021.
- [31] Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele, “2d human pose estimation: New benchmark and state of the art analysis,” in *CVPR*, 2014, pp. 3686–3693.
- [32] Alexander Toshev and Christian Szegedy, “Deeppose: Human pose estimation via deep neural networks,” in *CVPR*, 2014, pp. 1653–1660.
- [33] Javier Portilla and Eero P Simoncelli, “A parametric texture model based on joint statistics of complex wavelet coefficients,” *IJCV*, vol. 40, no. 1, pp. 49–70, 2000.
- [34] Karol Gregor, Ivo Danihelka, Alex Graves, Danilo Jimenez Rezende, and Daan Wierstra, “DRAW: A recurrent neural network for image generation,” in *ICML*, 2015, pp. 1462–1471.
- [35] Xinchen Yan, Jimei Yang, Kihyuk Sohn, and Honglak Lee, “Attribute2image: Conditional image generation from visual attributes,” in *ECCV*, 2016, pp. 776–791.
- [36] Junbo Jake Zhao, Michaël Mathieu, and Yann LeCun, “Energy-based generative adversarial network,” *CoRR*, vol. abs/1609.03126, 2016.
- [37] Xi Chen, Yan Duan, Rein Houthoofd, John Schulman, Ilya Sutskever, and Pieter Abbeel, “Infogan: Interpretable representation learning by information maximizing generative adversarial nets,” in *NIPS*, 2016, pp. 2172–2180.
- [38] Guo-Jun Qi, “Loss-sensitive generative adversarial networks on lipschitz densities,” *CoRR*, vol. abs/1701.06264, 2017.
- [39] Martin Arjovsky, Soumith Chintala, and Léon Bottou, “Wasserstein GAN,” *CoRR*, vol. abs/1701.07875, 2017.
- [40] Mehdi Mirza and Simon Osindero, “Conditional generative adversarial nets,” *CoRR*, vol. abs/1411.1784, 2014.
- [41] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A. Efros, “Image-to-image translation with conditional adversarial networks,” *CoRR*, vol. abs/1611.07004, 2016.
- [42] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A. Efros, “Unpaired image-to-image translation using cycle-consistent adversarial networks,” *CoRR*, vol. abs/1703.10593, 2017.
- [43] Taeksoo Kim, Moonsu Cha, Hyunsoo Kim, Jung Kwon Lee, and Jiwon Kim, “Learning to discover cross-domain relations with generative adversarial networks,” in *ICML*, 2017, pp. 1857–1865.
- [44] Christian Ledig, Lucas Theis, Ferenc Huszar, Jose Caballero, Andrew P. Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, and Wenzhe Shi, “Photo-realistic single image super-resolution using a generative adversarial network,” *CoRR*, vol. abs/1609.04802, 2016.
- [45] Andrew Brock, Theodore Lim, James M. Ritchie, and Nick Weston, “Neural photo editing with introspective adversarial networks,” *CoRR*, vol. abs/1609.07093, 2016.
- [46] Alec Radford, Luke Metz, and Soumith Chintala, “Unsupervised representation learning with deep convolutional generative adversarial networks,” *CoRR*, vol. abs/1511.06434, 2015.
- [47] Arno Schödl, Richard Szeliski, David H Salesin, and Irfan Essa, “Video textures,” in *SIGGRAPH*, 2000, pp. 489–498.
- [48] Aseem Agarwala, Ke Colin Zheng, Chris Pal, Maneesh Agrawala, Michael Cohen, Brian Curless, David Salesin, and Richard Szeliski, “Panoramic video textures,” in *ACM Transactions on Graphics (TOG)*, 2005, vol. 24, pp. 821–827.
- [49] Zicheng Liao, Neel Joshi, and Hugues Hoppe, “Automated video looping with progressive dynamism,” *ACM Transactions on Graphics (TOG)*, vol. 32, no. 4, pp. 77, 2013.
- [50] Joost R. van Amersfoort, Anitha Kannan, Marc’Aurelio Ranzato, Arthur Szlam, Du Tran, and Soumith Chintala, “Transformation-based models of video sequences,” *CoRR*, vol. abs/1701.08435, 2017.
- [51] Tianfan Xue, Jiajun Wu, Katherine Bouman, and Bill Freeman, “Visual dynamics: Probabilistic future frame synthesis via cross convolutional networks,” in *NIPS*, 2016, pp. 91–99.
- [52] Olaf Ronneberger, Philipp Fischer, and Thomas Brox, “U-net: Convolutional networks for biomedical image segmentation,” in *MICCAI*, 2015, pp. 234–241.
- [53] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Gregory S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian J. Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Józefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mané, Rajat Monga, Sherry Moore, Derek Gordon Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul A. Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda B. Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng, “Tensorflow: Large-scale machine learning on heterogeneous distributed systems,” *CoRR*, vol. abs/1603.04467, 2016.
- [54] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh, “Realtime multi-person 2d pose estimation using part affinity fields,” *CoRR*, vol. abs/1611.08050, 2016.