

Pose-guided Generative Adversarial Net for Novel View Action Synthesis

Xianhang Li¹, Junhao Zhang², Kunchang Li³, Shruti Vyas¹, and Yogesh S Rawat¹

¹ CRCV, University of Central Florida, ² National University of Singapore
³ Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences

Abstract

We focus on the problem of novel-view human action synthesis. Given an action video, the goal is to generate the same action from an unseen viewpoint. Naturally, novel view video synthesis is more challenging than image synthesis. It requires the synthesis of a sequence of realistic frames with temporal coherency. Besides, transferring different actions to a novel target view requires awareness of action category and viewpoint change simultaneously. To address these challenges we propose a novel framework named Pose-guided Action Separable Generative Adversarial Net (PAS-GAN), which utilizes pose to alleviate the difficulty of this task. First, we propose a **recurrent pose-transformation module** which transforms actions from the source view to the target view and generates novel view pose sequence in 2D coordinate space. Second, a well-transformed pose sequence enables us to separate the action and background in the target view. We employ a novel **local-global spatial transformation module** to effectively generate sequential video features in the target view using these action and background features. Finally, the generated video features are used to synthesize human action with the help of a 3D decoder. Moreover, to focus on dynamic action in the video, we propose a novel **multi-scale action-separable loss** which further improves the video quality. We conduct extensive experiments on two large-scale multi-view human action datasets, NTU-RGBD and PKU-MMD, demonstrating the effectiveness of PAS-GAN which outperforms existing approaches. The codes and models will be available on <https://github.com/xhl-video/PAS-GAN>.

1. Introduction

Video generation is an interesting problem with a wide range of applications in robotics [32, 6], data augmentation

[45], and augmented reality [12]. We have seen some recent efforts in video generation, which include stochastic video synthesis [33, 24, 36] and conditional video generation [35, 29, 39]. The unconditional stochastic approaches attempt to learn the distribution of data directly, whereas the conditional generation approaches utilize the guidance of some priors, such as semantic segmentation labels [39], pose of target video [1] to simplify the problem.

Motivated by the success of conditional video generation [35, 29, 39], we have recently seen some efforts [15, 38, 25, 16, 27] which can generate the novel view video by utilizing priors in feature level. These methods can obtain the novel view of the motion and the appearance features through some differential operations. We first conclude that previous works generally transfer the action into target view with three steps: Firstly, using an encoder to represent the action feature from source-view video in latent feature space, which can enhance the learning capacity of the network. Secondly, transforming the viewpoint of this action features by the cooperation power of the priors and differential operations in the target view, significantly reducing the complexity of the problem. Finally, feeding the transformed features into a decoder, aiming at generating a set of sequential frames. However, although many elaborate designs proved to be effective, we still clearly found significant motion blur and complete stillness in the generated video. Thus, there are three questions that are worth investigating: First, what kind of **action features** do we prefer to translate into the target view more easily? Second, what type of **priors and corresponding operations** are most beneficial for generating the motion and appearance details of the target view? Third, what **constraints** can improve the balance of motion and static detail generation for the novel view video synthesis model?

Regarding the **action features**, in [15, 38, 25, 16, 27], the authors use 3D convolutional neural network to encode source video. However, the action features from RGB space may not be the best choice. Because the RGB modality contains redundant background information in addition to the rich foreground action, it is challenging to extract ex-

¹The project page is on https://xhl-video.github.io/xianhangli/pas_gan.html.

²The authors email are: xli421@ucsc.edu, junhao.zhang@u.nus.edu, kc.li@siat.ac.cn, (shruti.yogesh)@crcv.ucf.edu

clusive action dynamics without explicit supervision. More importantly, view transformation in such redundant feature space leads to inaccurate pixel rendering. It has been shown that pose information can model the action dynamics better than other modalities [1, 39]. Motivated by this, we explore the use of pose information from the source viewpoint to represent the action features. Specifically, we use the pose sequence from the source video to ignore redundant background features. We propose a novel *recurrent pose-transformation module* to map this source pose to the target viewpoint.

Second, there are various types of *priors* used, such as RGB frame [15, 25], depth sequence [15], and pose sequence [15] from the target viewpoint. We argue that these priors provide different types of information, e.g., RGB can provide rich appearance information and pose can provide rich action dynamics. Therefore, it is worth investigating their appropriate utilization instead of treating them in the same manner [15]. The existing approaches [15, 25] perform a global transformation in latent space, which make it challenging to recover the fine action details for action prediction. Thus in our work, we only take an image prior and use it to decouple pose and appearance features. The extracted pose is later used for transforming action dynamics and the transformed action dynamics is used along with the extracted appearance features to generate latent video features. Specifically, we propose a novel *recurrent local-global spatial transformation module* which transforms latent features independently for each joints in the target pose and jointly learned with a *global transformation*. The co-operation power of local and global transformation modules can generate comprehensive latent features for video decoder.

Furthermore, we also investigate the *constraints* used for video synthesis. We argue that the commonly-used reconstruction loss (mean squared error) is the main reason that causes the motion blur [7, 10] because motion can not belong to the same Gaussian distribution. Although using the mean squared error (MSE) loss can lift the score of peak signal-to-noise ratio (PSNR), it can result in heavy motion blur for the visual quality. More discussion can be found in Section 3.4. To address this issue, we propose to separate the action from the background in the generated video with the help of a transformed pose sequence. We utilize a multi-scale perceptual loss as the main objective to force the network to focus on the action and background details separately. In addition, we also use the adversarial loss to make the network learn the distribution of actions, so it can focus on action dynamics in addition to the static details.

In conclusion, we propose PAS-GAN, a novel approach for novel-view action synthesis which is trained end-to-end jointly optimizing multiple objectives. We validate our approach on two large-scale multi-view human ac-

tion datasets, NTU-RGBD [26] and PKU-MMD [4]. We demonstrate: (1) the effectiveness of the proposed framework in novel view action synthesis; (2) the benefits of local spatial transformation module; (3) the capability of multi-scale action-separable perpetual and adversarial loss formulation to generate better static details and coherent motion.

2. Related Works

Conditional Video Generation Recent video generation approaches mainly focus on generative adversarial networks (GANs) [8, 33, 24, 36, 5]. video prediction belongs to the conditional category, which [35, 29, 39] can benefit from the priors. In [14], the author proposed a deep video prediction model conditioned on a single image and an action class. In [28, 44, 18], the authors use a representation consisting of a set of learned key-points along with their local affine transformations to represent complex motion. These methods utilize the key-points information to generate videos. However, the key-points and first frame are in the same view. More recently, methods in [39, 1, 19, 43] use pose as auxiliary information to help generate a more realistic video. However, our approach focuses on video generation from an unseen viewpoint. Unlike these methods that pose auxiliary information, we use pose to first separate appearance features and then recurrently generate action features. Second, we use the position information provided by pose to make the perceptual loss and adversarial loss more focused on the generation of action. Moreover, we need to estimate the key-points/pose for an unseen viewpoint, which is more challenging.

Novel-View Synthesis. Most of the current work in novel view synthesis is focused on images where the focus is either on geometry-based [23, 21] or learning-based methods [42]. with the help of a generative query network in [6], the authors utilize multiple views to render an image from unseen views. In [42], the authors proposed an encoder to extract view-independent features and a decoder hallucinate the image of a novel view. However, it is challenging to adopt the novel view image-synthesis methods into video synthesis due to several challenges. The image synthesis approaches focus on the transformation of appearance and the lack of temporal coherence inevitably leads to a failure in modeling the motion transformation. In [25, 15, 16], we have seen some efforts on novel view video generation, however, due to multiple unknown parameters of the target-view, it remains unclear how much information on target view is required for video synthesis. Moreover, the issue of motion blur and static-frames is more evident in this task than conditional video generation.

3. Method

Given a source action video and an image prior from a target viewpoint, our goal is to generate a video of the same

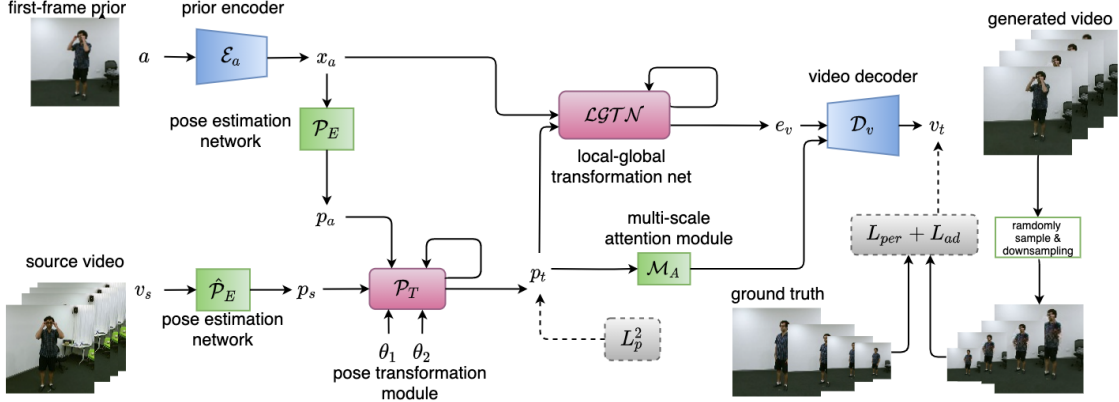


Figure 1. Overview of the proposed framework. Given an image prior a , which is the first frame of the target view video and an action video v_s from the source viewpoint, the proposed network generates the video v_t from the target viewpoint. Pose-transformation module \mathcal{P}_T transforms the extracted pose information p_s using p_a to generate the target view pose sequence p_t . The global and local transformation module \mathcal{LGTN} focus on coarse and fine action features, which uses appearance prior x_a and p_t to generate the latent features e_v . Different scales heat-map of p_t which generated by multi-scale attention from \mathcal{M}_A and e_v are used to generate the target view video v_t . Finally, our model relies on the supervision of the multi-scale action-separable perceptual loss L_{per} and adversarial loss L_{ad} .

action from the target viewpoint. Formally, given an action video v_s from a source viewpoint and an actor image a from a novel target view (the first frame of the target view video), we aim to generate an action video v_t from the target viewpoint. One of the key challenges is to transform relevant action dynamics from the source viewpoint to the target viewpoint. We propose a recurrent pose-transformation module \mathcal{P}_T to overcome this challenge, which uses pose information p_s and p_a from the source video and target prior along with the corresponding viewpoint angles θ_1 and θ_2 . The pose-transformation module \mathcal{P}_T generates pose-sequence p_t from the target view capturing action information. The next challenge is to focus on fine action details which are hard to capture in latent representation. We propose a recurrent local and global transformation module \mathcal{LGTN} which utilize the transformed pose-sequence p_t to transform the appearance prior x_a and generate latent features e_v for target view action. The generated latent features e_v are used to synthesize the target action video v_t with a video decoder \mathcal{D}_v . The last challenge is how to design an effective loss to supervise the model to learn a more detailed motion and reduce the generation of blurry motion. To this end, we designed multi-scale action-separable loss. It contains perceptual loss L_{per} and adversarial loss L_{ad} where the input of both loss is one randomly sampled frame from the generated video and ground truth frame at different scales. Moreover, the transformed pose is used to separate the action region which helps in generating details of motion and also preserve the appearance. An overview of the proposed approach is shown in Figure 1.

3.1. Recurrent Pose Transformation

Given a pose-sequence p_s from the source view and the prior-pose p_a from the target viewpoint, the goal of pose-

transformation module \mathcal{P}_T is to estimate the pose-sequence p_t from the target viewpoint. The \mathcal{P}_T module has a recurrent structure that also utilizes the viewpoint information θ_1 and θ_2 corresponding to the source and the target viewpoints to perform this transformation.

$$p_t = \mathcal{P}_T(p_s, p_a, \theta_1, \theta_2), \quad (1)$$

where p_t is the transformed pose-sequence, p_a is the prior pose, p_s is the source pose sequence, θ_1 is the source viewpoint, and θ_2 is the target viewpoint. The prior pose p_a is extracted from the prior image a using a 2D convolution based model \mathcal{E}_a [30] and a pose estimation head \mathcal{P}_E . The appearance encoder \mathcal{E}_a extracts latent appearance features x_a and these latent features are passed to a fully connected network \mathcal{P}_E for pose estimation. The source pose sequence p_s corresponds to the source video v_s and can be extracted using any existing pose estimation method $\hat{\mathcal{P}}_E$ [2]. In our experiments we utilize the pose information available with the training data.

A detailed overview of the pose-transformation module \mathcal{P}_T is shown in Figure 2. The recurrent structure of \mathcal{P}_T consists of three main components; a motion estimator \mathcal{M}_E , change in viewpoint estimator θ_E , and a motion transformer Δ_T . The motion estimator \mathcal{M}_E aims at predicting pose changes between subsequent frames. It takes two subsequent poses p_s^i and p_s^{i+1} and predicts ΔP , which represents the change in pose.

$$\Delta P = N(< N(p_s^i), N(p_s^{i+1}) >) \quad (2)$$

Here, N represents a two-layered neural network and $<>$ indicates concatenation operation. Similarly, the θ_E module estimates the change in viewpoint between the source and the target view. It estimates $\Delta\theta = \theta_E(\theta_1, \theta_2)$ with the help of a neural network same as described in equation 2.

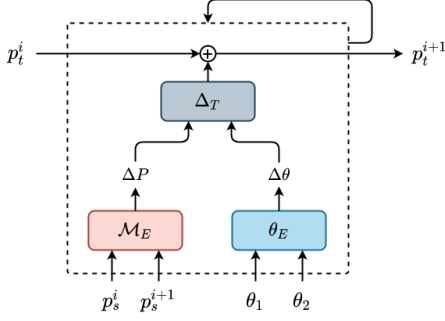


Figure 2. Overview of pose-transformation module.

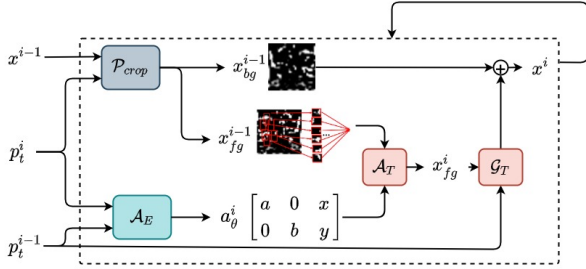


Figure 3. Overview of the local-global transformation network.

The third component, motion transformer Δ_T , takes the change in pose ΔP and change in viewpoint $\Delta\theta$ as input and transforms this motion to the target viewpoint. Finally, the transformed motion estimate is added to the hidden pose variable p_t^i to generate the target pose for the next time-step. Here p_t^i represents the pose in the previous time-step. The Δ_T module also uses the same architecture defined in equation 2 and the hidden pose variable is initialized with the prior pose p_a .

3.2. Pose-guided the Motion Integration

The estimated transformed pose p_t is passed to a two-stream motion integration module which generates action features e_v in latent space. This motion integration module consists of a local-global transformation network \mathcal{LGTN} .

Local-Global Transformation Network. A global transformation to generate action features in latent space can be effective, but it can be challenging to recover fine action details with such global transformation. To address this issue, we propose a local-global transformation network \mathcal{LGTN} which allows us to transform different key-regions in action independently. We utilize the joints present in the pose to define these key-regions. \mathcal{LGTN} is based on Spatial Transformer Network (STN) [11] with three key differences. First, STN utilizes a common source to determine the transformation parameters and perform the transformations. \mathcal{LGTN} , on the other hand, utilizes the pose information to estimate the transformation parameters and the transformation is performed on the appearance features. Second, in STN a global transformation is performed on the source,

whereas in \mathcal{LGTN} we propose to perform a different transformation on different key-regions. Finally, STN does not have a notion of recurrent transformation, and \mathcal{LGTN} has a recurrent structure that performs these transformations for a sequence of action features. A detailed overview of \mathcal{LGTN} is shown in Figure 3. Moreover, after transforming the feature locally, we feed the foreground transformed feature into a global transformation based [3].

The local-global transformation network \mathcal{LGTN} consists of four main components; a parameter estimator \mathcal{A}_E , key-region separator \mathcal{P}_{crop} , a transformation module \mathcal{A}_T and a global transformation module \mathcal{G}_T . The parameter estimator \mathcal{A}_E predicts the transformation parameters A_θ with the help of a small convolution neural network which takes two subsequent pose heat-maps as input. It aims at learning the set of transformation parameters that are required to move from the initial pose p_t^{i-1} to the next pose p_t^i . We utilize a 2D affine transformation where the estimated parameters can be represented as,

$$a_\theta^i = \begin{bmatrix} a^i & 0 & x^i \\ 0 & b^i & y^i \end{bmatrix} \quad (3)$$

where varying a^i, b^i, x^i , and y^i allow for translation, rotation, and scaling.

The key-region separator \mathcal{P}_{crop} takes the pose from current time-step and extracts key-regions from the latent appearance features x^{i-1} based on the joints. All the extracted key-regions represent foreground features x_{fg}^{i-1} and the remaining features x_{bg}^{i-1} are marked as background. The transformation module \mathcal{A}_T applies the learned parameters A_θ on the foreground features x_{fg}^{i-1} to generate the foreground appearance features x_{fg}^i for the next time-step. This pointwise transformation is defined as,

$$\begin{bmatrix} \alpha^i \\ \beta^i \end{bmatrix} = a_\theta^i \begin{bmatrix} \alpha^{i-1} \\ \beta^{i-1} \\ 1 \end{bmatrix} = \begin{bmatrix} a^i & 0 & x^i \\ 0 & b^i & y^i \end{bmatrix} \begin{bmatrix} \alpha^{i-1} \\ \beta^{i-1} \\ 1 \end{bmatrix} \quad (4)$$

where (α^i, β^i) are the target coordinates in the transformed foreground appearance feature map, $(\alpha^{i-1}, \beta^{i-1})$ are the source coordinates in the appearance feature map x_{fg}^{i-1} from the previous time-step, and a_θ^i is the estimated affine transformation matrix. Finally the transformed foreground features x_{fg}^i are aggregated and integrated back with the background appearance features x_{bg}^{i-1} to compute video features x_i for next time-step. Then we feed the foreground feature x_{fg}^i into the \mathcal{G}_T , based on the convolutional gated recurrent unit (Conv-GRU) [3], which takes current subsequent pose heat-maps as input to estimate motion and transform the x_{fg}^i . This transformation is performed recurrently for each time-step. Finally, we can obtain a sequence of features. The final video features can be obtained $e_v = \mathcal{LGTN}(p_t, x_i)$, are for target view video generation.

3.3. Video Decoder

We utilize a 3D Conv based video decoder \mathcal{D}_v , which takes the generated video features e_v and synthesize corresponding action video v_t . The pose information provides crucial details regarding the foreground region where the action occurs. Therefore we explore the use of estimated pose for attention mechanism in the video decoder. The transformed pose sequence p_t is used to create a sequence of heatmaps h_a with the help of a Gaussian kernel over joints. These pose heatmaps are integrated with the generated video features e_v in the video decoder \mathcal{D}_v which encourage the decoder to focus on foreground action regions.

Multi-scale Learning. It can be challenging to recover fine action dynamics from a compressed latent representation; therefore we take a multi-scale approach while decoding the action video v_t . We extract the prior appearance features x_a at different resolutions from the earlier layers of the prior encoder \mathcal{E}_a . Also, the motion integration using the local-global transformation network \mathcal{LGTN} is performed at multiple scales using these multi-scale appearance features. The \mathcal{LGTN} modules share weights for feature transformation across all scales. The pose heatmaps h_a for soft attention are also computed at multiple scales using linear interpolation using \mathcal{M}_A . The transformed video features e_v and pose heatmaps h_a at multiple scales are then utilized by the video decoder at different stages to generate the video.

3.4. Training Objective

Overall, the proposed framework is jointly optimized for two different training objectives; transformed pose prediction L_p^2 , and action video generation L_a . The transformed pose prediction L_p^2 is optimized using a mean squared loss defined as,

$$L_p^2 = \frac{1}{NT} \sum_{j=1}^T \sum_{i=1}^N ((x_i^j - \hat{x}_i^j)^2 + (y_i^j - \hat{y}_i^j)^2) \quad (5)$$

where, T is the total number of time-steps, N is total number of skeleton joints with coordinates x_i^j and y_i^j each. We experimented with multiple loss formulations for L_a to improve the quality of video prediction.

Multi-scale Action Separable Loss. The multi-scale action separable loss contains one perceptual loss and adversarial loss of generated frames and generated actions at different scales. First, we use a perceptual loss L_{per} to improve the visual quality of the generated video frames. [13]. We use a pre-trained vgg16-net [30] to extract the frame-level features at different layers. We employed the ground

truth frame $\hat{\mathbf{F}}$ and generated frames \mathbf{F} as input.

$$L_{per} = \sum_{c=1}^C |V_c(\hat{\mathbf{F}}) - V_c(\mathbf{F})| + \sum_{c=1}^C |V_c(\mathbf{F}_{\text{crop}}) - V_c(\hat{\mathbf{F}}_{\text{crop}})| \quad (6)$$

Where V represents the pre-trained vgg16-net and c is the current channel. We adopt the key-region separator \mathcal{P}_{crop} to obtain \mathbf{F}_{crop} . To reduce the computational cost, the perceptual loss L_{per} is calculated by only randomly sampling one of the generated video frames each step. Finally, we adopt a downsampling layer for the input F to employ the multi-scale F at 4 scales [1, 0.5, 0.25, 0.125]. Thus we have 16 items in total in L_{per} . In addition to this, we also explore the use of adversarial loss [8] to make the generated videos realistic. We use vgg16-net based discriminator with four scale inputs and adopt the \mathcal{P}_{crop} to obtain actions as input similarly. The final training objective for our framework is,

$$L = \lambda_1 L_p^2 + \lambda_2 L_{per} + \lambda_3 L_{ad} \quad (7)$$

where λ_1, λ_2 and λ_3 are weights which are estimated experimentally. Motivated by some excellent super-resolution works [17, 40], we use a MSE loss pre-trained network to initialize our model parameters.

4. Experiments

We aim to demonstrate the effectiveness of the proposed approach in novel view video synthesis and highlight the benefit of its components including, recurrent pose-transformation module, local-global transformation module, and multi-scale action separable loss.

4.1. Experiment Setup

Dataset. We conduct our experiments on NTU-RGBD [26] and PKU-MMD [4], which are both large scale multi-view action datasets. The NTU-RGBD contains over 56,000 videos and 60 action classes in total. The videos are recorded using three different cameras and there are a total of 80 different viewpoints. The PKU-MMD contains over 3500 action clips and 51 action classes which are also recorded by three cameras. We use the 2D RGB pose, which indicates key points in the original RGB resolution and the viewpoint angles. The pose information can also be extracted using some advanced pose estimation methods [31, 20] for each frame.

Implementation Details. We train our network in an end-to-end manner optimizing multiple objectives. We use $\lambda_1 = \lambda_2 = 1, \lambda_3 = 1e - 4$ in our training objective. The network is trained using Adam optimizer with a learning rate of $1e-4$ with a batch-size 12. For NTU-RGBD, we use

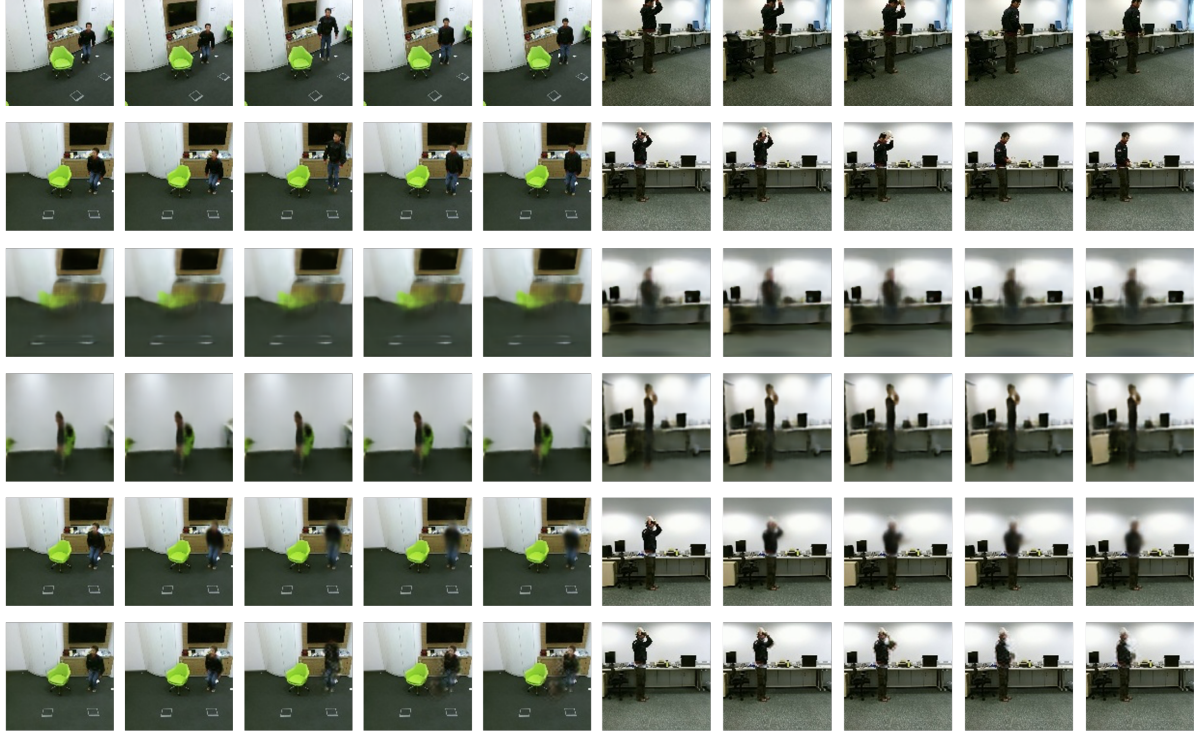


Figure 4. Comparison of the generated video frames using PAS-GAN with existing methods on NTU-RGBD [26] dataset. Row 1: source; Row 2: target; Row 3: VRNet [37]; Row 4: VDNet [15]; Row 5: RTNet [25]; Row 6: PAS-GAN (ours).

Model	Pair-view SSIM Score						Average SSIM	Average PSNR
	$v1 \rightarrow v2$	$v1 \rightarrow v3$	$v2 \rightarrow v1$	$v2 \rightarrow v3$	$v3 \rightarrow v1$	$v3 \rightarrow v2$		
VDG[9]	.502 ± .058	.543 ± .068	.584 ± .060	.563 ± .062	.611 ± .077	.522 ± .063	.554 ± .075	-
PG ² [19]	.499 ± .071	.561 ± .060	.600 ± .064	.557 ± .071	.598 ± .075	.543 ± .055	.560 ± .076	-
VRNet[37]	-	-	-	-	-	-	.68	19.8
ResNet[15]	.705 ± .115	.735 ± .095	.717 ± .130	.690 ± .122	.734 ± .127	.669 ± .150	.708 ± .127	-
VDNet[15]	.789 ± .076	.791 ± .069	.800 ± .076	.765 ± .079	.797 ± .067	.756 ± .089	.783 ± .078	-
RTNet[25]	.974 ± .021	.975 ± .021	.975 ± .019	.971 ± .021	.974 ± .017	.971 ± .022	.973 ± .020	27.5 ± 2.70
PAS-GAN(ours)	.977 ± .007	.975 ± .006	.974 ± .008	.975 ± .006	.978 ± .007	.977 ± .006	.976 ± .006	28.07 ± 1.39

Table 1. Comparison of SSIM and average PSNR scores with existing methods. We report scores from all the combinations of three testing views. The scores for VDG [9] and PG² [19] are shown as reported by the authors of VDNet[15].

Model	MSE ↓	PSNR ↑	SSIM ↑	FVD ↓	Pose MSE ↓
VRNet[37]	.00281	25.54 ± 1.45	.923 ± .010	26.30 ± .020	-
RTNet[25]	.00082	31.56 ± 2.46	.974 ± .021	3.90 ± .12	-
BasicNet	.00406	24.05 ± 1.03	.938 ± .016	16.64 ± .022	0.032
PAS-GAN(ours)	.00058	32.58 ± 1.27	.984 ± .005	3.61 ± .10	0.009

Table 2. Key-region based MSE, PSNR, SSIM, & FVD scores along with MSE score for pose estimation.

a resolution of 112×112 in our experiments with 8 frames. For PKU-MMD, we use 224×224 . We extract four-scales feature with spatial resolution 112×112 & (224×224) , 56×56 & (112×112) , 28×28 & (56×56) , 14×14 & (28×28) respectively from a modified vgg16-net [30]. We implement our framework on PyTorch [22] and use a Tesla V100.

Evaluation Metrics. To evaluate the performance of the proposed methods, we follow [15, 25] and use Structural Similarity (SSIM) [41] and Peak Signal to Noise Ratio (PSNR) for per-frame quantitative results. In addition, we adopt the Fréchet Video Distance [34] for evaluating the

quality of generated videos. Moreover, we also compute MSE, SSIM and PSNR scores only on the foreground region to focus more on the generated action. We utilize the pose information to crop out the foreground region in the generated and ground-truth videos.

Baseline To obtain our baseline model (BasicNet), we use a modified vgg16-net [30] to encode the image prior and a 3D convolution-based video decoder to generate the video. The input is the image prior and we expand the image encodings to add temporal dimension before passing it to the video decoder. Based on this BasicNet, we use this as a basic model in our approach and evaluate the effectiveness of various components on top of this. Another important baseline method is VD-Net [15]. Instead of using a single image as prior, view-lstm utilizes the sequence of depth heatmap and skeleton from the target view, which already contains sufficient motion information from the target view.

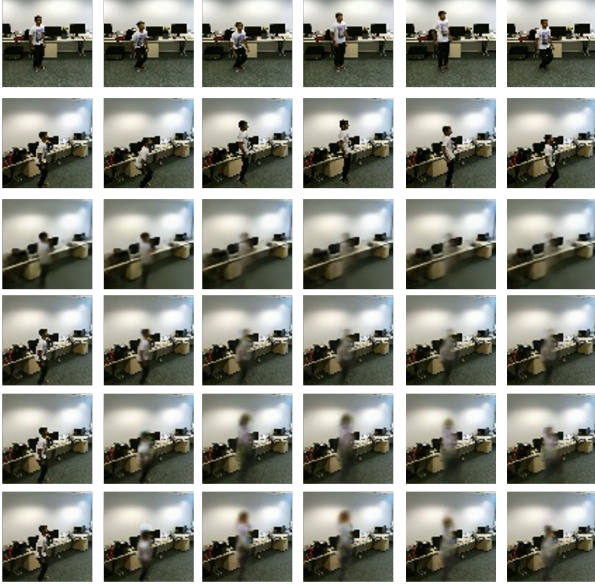


Figure 5. Ablation to evaluate different components of PAS-GAN. Row 1: source; Row 2: target; Row 3: BasicNet; Row 4: w/ multi-scale learning; Row 5: w/ recurrent pose transformation module; Row 6: w/ local-global transformation module. **It is important to note that row 6 is not generated using the full model.**

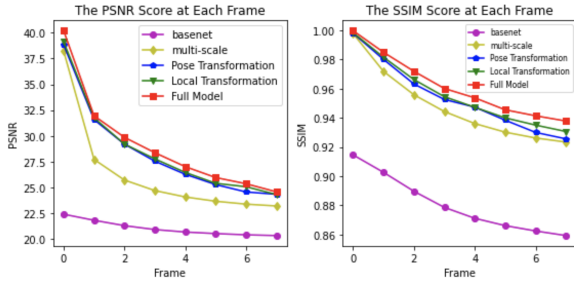


Figure 6. Ablation to compare different components using frame-level PSNR and SSIM scores.

Model	MSE ↓	PSNR ↑	SSIM ↑	FVD ↓
BasicNet	.00394	24.38 ± 1.63	.924 ± .030	17.87 ± .67
w/ Multi-Scale	.00272	26.32 ± 2.39	.950 ± .029	5.94 ± .084
w/ Recurrent Pose Transformation	.00226	27.23 ± 2.58	.958 ± .028	4.83 ± .127
w/ local-global Transformation Network	.00182	27.85 ± 1.89	.968 ± .016	4.11 ± .067

Table 3. Key-region based evaluation of components in PAS-GAN.

Model	MSE ↓	PSNR ↑	SSIM ↑	FVD ↓
w/ MSE Only	.00182	27.85 ± 1.89	.968 ± .016	4.11 ± .067
w/ Perceptual Loss [13] Only	.00197	27.43 ± 1.75	.967 ± .013	3.94 ± .090
w/ Multi-Scale Action Separable Loss	.00177	28.09 ± 1.87	.970 ± .015	3.83 ± .067

Table 4. Key-region based evaluation of different loss functions.

On the other hand, our method focuses on transforming the motion from the source view to the target view. We compare the proposed method with these baselines and other state-of-the-art methods, including [25] both qualitatively as well as quantitatively.

4.2. Evaluation

We perform an extensive evaluation of PAS-GAN on NTU-RGBD and PKU-MMD datasets. The quantitative evaluation on NTU-RGBD dataset is shown in Table 1 and 2. In Table 1, we have shown pair-wise SSIM scores for

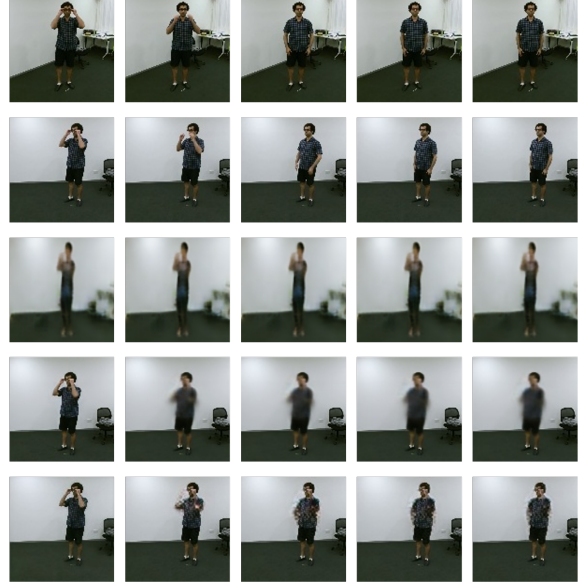


Figure 7. Generated novel view videos using existing methods highlighting motion blur and static video frames. Row 1: Source Video; Row 2: Target Video; Row 3: VNet [15]; Row 4: RTNet [25]; Row 5: PAS-GAN (Proposed method). Compared with previous state-of-the-art methods, our model can alleviate the motion blur and maintain the temporal coherence to some extent.

each pair of viewpoints. We observe that PAS-GAN consistently performs well on all pairs of viewpoints independent of the change in viewpoint. In Table 2, we have shown scores for the action region in the video and we can observe that scores are comparable (slightly better) compared with the evaluation on the full video. We also show the generated video frames for qualitative evaluation on NTU-RGBD and PKU-MMD datasets. The generated frames are shown in Figure 4 and 8. We observe that the generated frames have visible action dynamics capturing the target action.

Quantitative Comparison. We first perform a quantitative comparison of PAS-GAN with existing methods using SSIM, PSNR, and FVD metrics on the NTU-RGBD dataset. The comparison is shown in Table 1 and 2. In Table 1, we compare the SSIM and PSNR scores on the full generated video. We observe that PAS-GAN outperforms all the existing methods in both the evaluation metrics. Table 2 shows the evaluation only on the activity region. We can observe that the proposed method provides a significant improvement over the baseline model BasicNet and also outperforming the other existing methods.

Qualitative Comparison. We also perform qualitative comparison on NTU-RGBD and PKU-MMD. The generated video frames are shown in Figure 4, in Figure 7, and in Figure 8. Although the quantitative scores of PAS-GAN and previous SOTA RTNet [25] is close, we can see a significant improvement our PAS-GAN made in terms of the visual quality. The action dynamics can be visible in the



Figure 8. Comparison of the generated video frames using PAS-GAN with existing methods on PKU-MMD [4]. Row 1: source; Row 2: target; Row 3: RTNet [25]; Row 4: PAS-GAN (proposed method).

generated frames and it is much better when compared with the other models on the same dataset.

4.3. Ablations

We conduct several ablation experiments from two different perspectives. First is to validate the effectiveness of each component in PAS-GAN and second is to explore the impact of different loss functions.

Effectiveness of Components. To study the impact of different components in PAS-GAN, we sequentially add proposed components in the BasicNet model. The comparison is shown in Table 3 and in Figure 5 where we observe that the largest improvement results from the multi-scale learning. We also observe a consistent performance gain as each module is added to the network. However, the performance gain diminishes eventually and we argue that this is due the use of a smaller resolution (56×56) in our ablations, which limits the performance to some extent. Moreover, we plot the per-frame PSNR and SSIM respectively in Figure 5. It clearly shows that different modules contributes distinctively. Beside, due to the use of an image prior, the results are at the peak for the first frame.

Influence of Loss Functions. PAS-GAN is trained using multiple objectives. To study the impact of these loss functions, we perform some ablations which are shown in Table 4. First surprising finding is that, although using only perceptual results in worse scores than MSE, there is a huge difference in visual quality between them as shown in Figure 4. Next, we observe that the proposed multi-scale action-separable loss improves the performance further. It is due to the use of multiple scales and separation of action and background. They significantly help the model to focus more on fine motion and appearance details.

4.4. Discussion

We would like to explore in more depth the similarities and differences between our approach and RTNet [25]. Our problem definition is the same as theirs: given the first frame of the target view, the goal is to transform the source video to the target viewpoint. The first difference is that we do not use the source video to represent the action features. Instead, we use the 2D pose obtained from each frame of the source video. Transforming the viewpoint in 2D pose coordinate space greatly reduces our transformation difficulty as non-action features are automatically ignored. Secondly, RTNet only focuses on the global information when generating the target view features, which leads to a lack of local information which leads to motion blur. Our local-global transformation module uses the 2D pose coordinate information to automatically focus on the local information and use a global transformation module to integrate further. Finally, we use the multi-scale perceptual loss as the primary loss instead of MSE, and separate the action with the transformed pose coordinates for its supervision. This enables our approach to address the issue of motion blur and generate frames with coherent motion.

5. Conclusion

We present PAS-GAN for novel view video synthesis where we explore the use of pose to solve this problem. Our method utilizes a pose transformation module to capture action dynamics from a source viewpoint. We show how a local-global feature transformation can be used jointly to learn effective action dynamics. We also propose a novel video specific multi-scale action-separable perceptual and adversarial loss formulation to improve the quality of generated videos. Extensive evaluation on two large-scale action datasets demonstrates the effectiveness of the proposed approach and the benefits of its various components.

References

- [1] C. Chan, S. Ginosar, T. Zhou, and A. Efros. Everybody dance now. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5932–5941, 2019. 1, 2
- [2] Bowen Cheng, Bin Xiao, Jingdong Wang, Honghui Shi, Thomas S. Huang, and Lei Zhang. Higherhrnet: Scale-aware representation learning for bottom-up human pose estimation. In *CVPR*, 2020. 3
- [3] Kyunghyun Cho, Bart van Merriënboer, Çağlar Gülçehre, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder-decoder for statistical machine translation. *CoRR*, abs/1406.1078, 2014. 4
- [4] Liu Chunhui, Hu Yueyu, Li Yanghao, Song Sijie, and Liu Jiaying. Pku-mmd: A large scale benchmark for continuous multi-modal human action understanding. *arXiv preprint arXiv:1703.07475*, 2017. 2, 5, 8
- [5] Aidan Clark, Jeff Donahue, and Karen Simonyan. Efficient video generation on complex datasets. *CoRR*, abs/1907.06571, 2019. 2
- [6] SM Ali Eslami, Danilo Jimenez Rezende, Frederic Besse, Fabio Viola, Ari S Morcos, Marta Garnelo, Avraham Ruderman, Andrei A Rusu, Ivo Danihelka, Karol Gregor, et al. Neural scene representation and rendering. *Science*, 360(6394):1204–1210, 2018. 1, 2
- [7] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*, volume 1. 2016. 2
- [8] 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*, pages 2672–2680, 2014. 2, 5
- [9] Mohamed Ilyes Lakhal, Oswald Lanz, and Andrea Cavallaro. Pose guided human image synthesis by view disentanglement and enhanced weighting loss. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 0–0, 2018. 6
- [10] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1125–1134, 2017. 2
- [11] Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. Spatial transformer networks. In *Advances in neural information processing systems*, pages 2017–2025, 2015. 4
- [12] Varun Jain, Shivam Aggarwal, Suril Mehta, and Ramya Hebbalaguppe. Synthetic video generation for robust hand gesture recognition in augmented reality applications. *CoRR*, abs/1911.01320, 2019. 1
- [13] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*, pages 694–711. Springer, 2016. 5, 7
- [14] Yunji Kim, Seonghyeon Nam, In Cho, and Seon Joo Kim. Unsupervised keypoint learning for guiding class-conditional video prediction. In *Advances in Neural Information Processing Systems*, pages 3814–3824, 2019. 2
- [15] M. Lakhal, O. Lanz, and A. Cavallaro. View- lstm: Novel-view video synthesis through view decomposition. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 7576–7586, 2019. 1, 2, 6, 7
- [16] Mohamed Ilyes Lakhal, Davide Boscaini, Fabio Poiesi, Oswald Lanz, and Andrea Cavallaro. Novel-view human action synthesis. In *Asian Conference on Computer Vision (ACCV)*, December 2020. 1, 2
- [17] Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al. Photorealistic single image super-resolution using a generative adversarial network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4681–4690, 2017. 5
- [18] Tianyu Luan, Yali Wang, Junhao Zhang, Zhe Wang, Zhipeng Zhou, and Yu Qiao. Pc-hmr: Pose calibration for 3d human mesh recovery from 2d images/videos. *arXiv preprint arXiv:2103.09009*, 2021. 2
- [19] Liqian Ma, Xu Jia, Qianru Sun, Bernt Schiele, Tinne Tuytelaars, and Luc Van Gool. Pose guided person image generation. In *Advances in neural information processing systems*, pages 406–416, 2017. 2, 6
- [20] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation, 2016. 5
- [21] David Novotny, Ben Graham, and Jeremy Reizenstein. Perspectivenet: A scene-consistent image generator for new view synthesis in real indoor environments. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32, pages 7601–7612. Curran Associates, Inc., 2019. 2
- [22] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019. 6
- [23] K. Rematas, C. H. Nguyen, T. Ritschel, M. Fritz, and T. Tuytelaars. Novel views of objects from a single image. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(8):1576–1590, 2017. 2
- [24] Masaki Saito, Eiichi Matsumoto, and Shunta Saito. Temporal generative adversarial nets with singular value clipping. In *ICCV*, 2017. 1, 2
- [25] Kara Marie Schatz, Erik Quintanilla, Shruti Vyas, and Y. Rawat. A recurrent transformer network for novel view action synthesis. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020. 1, 2, 6, 7, 8
- [26] Amir Shahrudiy, Jun Liu, Tian-Tsong Ng, and Gang Wang. Ntu rgb+d: A large scale dataset for 3d human activity analysis. In *IEEE Conference on Computer Vision and Pattern Recognition*, June 2016. 2, 5, 6

- [27] Sarah Shiraz, Krishna Regmi, Shruti Vyas, Yogesh S Rawat, and Mubarak Shah. Novel view video prediction using a dual representation. *International Conference on Image Processing*, 2021. 1
- [28] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. First order motion model for image animation. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32, pages 7137–7147. Curran Associates, Inc., 2019. 2
- [29] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. Animating arbitrary objects via deep motion transfer. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 1, 2
- [30] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 3, 5, 6
- [31] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning for human pose estimation. In *CVPR*, 2019. 5
- [32] Joshua Tobin, Wojciech Zaremba, and Pieter Abbeel. Geometry-aware neural rendering. In *Advances in Neural Information Processing Systems*, pages 11559–11569, 2019. 1
- [33] Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, and Jan Kautz. MoCoGAN: Decomposing motion and content for video generation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1526–1535, 2018. 1, 2
- [34] Thomas Unterthiner, Sjoerd van Steenkiste, Karol Kurach, Raphaël Marinier, Marcin Michalski, and Sylvain Gelly. Towards accurate generative models of video: A new metric & challenges. *CoRR*, abs/1812.01717, 2018. 6
- [35] Ruben Villegas, Jimei Yang, Seunghoon Hong, Xunyu Lin, and Honglak Lee. Decomposing motion and content for natural video sequence prediction. *ICLR*, 2017. 1, 2
- [36] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Generating videos with scene dynamics. In *Advances in neural information processing systems*, pages 613–621, 2016. 1, 2
- [37] Shruti Vyas, Yogesh Singh Rawat, and Mubarak Shah. Time-aware and view-aware video rendering for unsupervised representation learning. *CoRR*, abs/1811.10699, 2018. 6
- [38] Shruti Vyas, Yogesh S Rawat, and Mubarak Shah. Multi-view action recognition using cross-view video prediction. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXVII 16*, pages 427–444. Springer, 2020. 1
- [39] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Guilin Liu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. Video-to-video synthesis. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018. 1, 2
- [40] Xintao Wang, Ke Yu, Shixiang Wu, Jinjin Gu, Yihao Liu, Chao Dong, Yu Qiao, and Chen Change Loy. Esrgan: Enhanced super-resolution generative adversarial networks. In *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, pages 0–0, 2018. 5
- [41] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. 6
- [42] Xiaogang Xu, Ying-Cong Chen, and Jiaya Jia. View independent generative adversarial network for novel view synthesis. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 7791–7800, 2019. 2
- [43] Ceyuan Yang, Zhe Wang, Xinge Zhu, Chen Huang, Jianping Shi, and Dahua Lin. Pose guided human video generation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 201–216, 2018. 2
- [44] Junhao Zhang, Yali Wang, Zhipeng Zhou, Tianyu Luan, Zhe Wang, and Yu Qiao. Learning dynamical human-joint affinity for 3d pose estimation in videos. *IEEE Transactions on Image Processing*, 30:7914–7925, 2021. 2
- [45] Yumeng Zhang, Gaoguo Jia, Li Chen, Mingrui Zhang, and Jun-Hai Yong. Self-paced video data augmentation with dynamic images generated by generative adversarial networks. *CoRR*, abs/1909.12929, 2019. 1