

Cross-identity Video Motion Retargeting with Joint Transformation and Synthesis

Haomiao Ni¹

Yihao Liu²

Sharon X. Huang¹

Yuan Xue²

¹The Pennsylvania State University, University Park, PA, USA

²Johns Hopkins University, Baltimore, MD, USA

¹{hfn5052, suh972}@psu.edu ²{yliu236, yuanxue}@jhu.edu



Figure 1: Examples of video motion retargeting, where motion from the driving video (1st column in 2nd&3rd rows) is transferred to the subject in the subject video (1st row). Videos generated by RegionMM [41] for face video and EDN [4] for dance video are shown in the 2nd column of each block, 2nd&3rd rows. Videos generated by our proposed TS-Net are in the 3rd column of each block, 2nd&3rd rows (highlighted with blue boxes).

Abstract

In this paper, we propose a novel dual-branch Transformation-Synthesis network (TS-Net), for video motion retargeting. Given one subject video and one driving video, TS-Net can produce a new plausible video with the subject appearance of the subject video and motion pattern of the driving video. TS-Net consists of a warp-based transformation branch and a warp-free synthesis branch. The novel design of dual branches combines the strengths of deformation-grid-based transformation and warp-free generation for better identity preservation and robustness to occlusion in the synthesized videos. A mask-aware similarity module is further introduced to the transformation branch to reduce computational overhead. Experimental results on face and dance datasets show that TS-Net achieves better performance in video motion retargeting than several state-of-the-art models as well as its single-branch variants. Our code is available at https://github.com/nihaomiao/WACV23_TSNet.

1. Introduction

Motion retargeting aims to transfer motion from a driving video to a target video while maintaining the subject’s identity of the target video. It has become an important topic due to its practical applications in special effects, virtual/augmented reality and video editing, etc. Motion retargeting in the image domain has been explored extensively and compelling results have been shown in many tasks, such as person image generation [1, 27, 32, 37, 41], and facial expression generation [5, 21, 34, 55]. Often formulated as a guided video synthesis task, motion retargeting between videos is known to be more challenging than motion retargeting between images since the temporal dynamics of the motion to be transferred has to be learned [6]. Moreover, synthesizing realistic videos, especially human motion videos, is more challenging than the generation of high-quality images because human perception is sensitive to unnatural temporal changes, and human motion is often highly articulated [52, 54]. In this paper, we mainly focus on video motion retargeting between different human

subjects (Fig. 1). Given one subject video and one driving video, we aim to synthesize a new plausible video with the same identity of the person from the subject video and the same motion as the person in the driving video.

Recent works in video motion retargeting [2, 4, 6, 9, 14, 18, 46, 47, 49, 50, 52, 54] have shown impressive progress. To capture the temporal relationship among video frames, prior works [6, 49, 50] generated frames via warping subject frames by motion flow, which is usually extracted by specifically designed warping field estimators, such as FlowNet [8] or first-order approximation [39]. While warp-based systems can generally preserve subject identity well, traditional flow-based warping may suffer from occlusion and large motion due to its requirement of learning a warp field with point-to-point correspondence between frames [15]. Other methods [2, 4, 12, 20, 52, 54] utilized warp-free (direct) synthesis with a conditional GAN-style structure [16, 30, 48]. To ease the challenging of direct synthesis, they often employed feature disentangle/decomposition [52] or followed the state-of-the-art generator architectures [31, 48] to add various connections among inputs, the encoder, and the decoder network. Unlike warp-based generation, direct synthesis is not limited to only using pixels from reference images, and therefore is easier to synthesize novel pixels for unseen/occluded objects. However, such flexibility can also lead to identity leakage [12], *i.e.*, identity changes in the generated video.

Considering that warp-based synthesis can better preserve identity while warp-free generation helps produce new pixels, in this paper, we propose a novel video motion retargeting framework, termed *Transformation-Synthesis Network*, or *TS-Net* for short, to combine their advantages. TS-Net has a dual branch structure which consists of a transformation branch and a synthesis branch. The network architectures within the two branches are inherently different, thus learning via the two branches can be regarded as a special multi-view learning case [51]. Unlike the popular warp-based methods using specially designed optical flow estimators [38, 39, 41] and inspired by [26], our proposed transformation branch computes deformation flow by weighting the regular grid with a spatial similarity matrix between driving mask features and subject image features. The computation of similarity takes multiple correspondences into consideration; thus it can better alleviate occlusion and handle large motion. We also design a mask-aware similarity to avoid comparing all pairs of points within the feature maps and thus be more efficient than traditional similarity computation methods. In our synthesis branch, we use a fully-convolutional fusion network. Features of two branches are concatenated and fed to the decoder network to generate realistic video frames. Experiments in Sec. 4 shows the effectiveness of this simple concatenation strategy.

Merely based on sparse 2D masks of driving videos, our proposed TS-Net can consistently achieve state-of-the-art results for both face and dance videos, successfully modeling hair and clothes details and their motion. TS-Net also handles large motions and preserves identity better when compared with other state-of-the-art methods, as shown in Fig. 1. Our contributions are summarized as follows:

1. We propose a novel dual branch video motion retargeting network TS-Net to generate identity-preserving and temporally coherent videos via joint learning of transformation and synthesis.
2. We utilize a simple yet effective way to estimate deformation grid based on similarity matrix. Mask-aware similarity is adopted to further reduce computation overhead.
3. Comprehensive experiments on facial motion and body motion retargeting tasks show that TS-Net can achieve state-of-the-art results by only using sparse 2D masks.

2. Related Work

Guided Image Generation. For conditional image generation, many works focused on generation tasks guided by specific conditions such as pose-guided person image synthesis [1, 27, 32, 36, 40, 43] and conditioned facial expression generation [5, 12, 34]. Pose-guided person image generation can produce person images in arbitrary poses, based on a subject image of that person and a novel pose from the driving image. Ma *et al.* [27] proposed a two-staged coarse-to-fine Pose Guided Person Generation Network (PG²), which utilizes pose integration and image refinement to generate high-quality person images. Conditioned facial expression generation aims to generate a reenacted face which shows the same expression as the driving face image while preserving the identity of the subject image. Chen *et al.* [5] proposed a two-stage framework called PuppeteerGAN, which first performs expression retargeting by the sketching network and then executes appearance transformation by the coloring network. Though these works have shown promising results, they are restricted to a specific object category (face or human body). Several recent works [37, 38, 41, 55] have proposed general guided image generation in various domains. Most of works [38, 39, 41, 44] applies motion flow to image animation because it can model the physical dynamics. Siarohin *et al.* [39] proposed a general self-supervised first-order-motion model for estimating dense motion flow to animate arbitrary objects using learned keypoints and local affine transformations. In [41], the authors further improved their network by modeling object movement through unsupervised region detection. Despite of building upon similar motion flow, instead of adopting complicated modeling in [39, 41], the

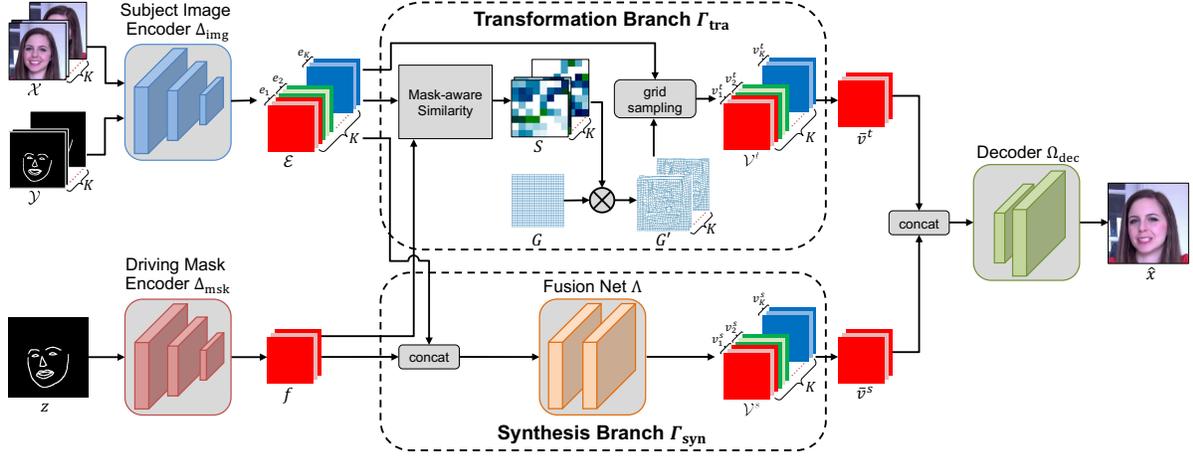


Figure 2: Illustration of the TS-Net generator to generate one frame \hat{x} in the target video.

transformation branch in TS-Net generate deformation flow by weighting regular grid with similarity matrix in feature space, which shows better simplicity and efficiency.

Video Motion Retargeting. Different from image-based generation, video motion retargeting is more challenging due to the additional coherence requirements in the temporal dimension. Most existing literature focused on specific domains such as human pose motion retargeting [4, 52], or facial expression retargeting [9, 12, 20, 49, 50, 54], yet they may lack generality when applied to multiple domains. In contrast, our proposed TS-Net can work well on both face and human body videos. Using off-the-shelf detectors to extract driving motion masks, such as 3D masks [9, 12, 20], 2D dense mask [46, 47], or 2D sparse mask [4, 54], is also popular in current video motion retargeting methods. Due to the simplicity of 2D sparse masks, our proposed TS-Net also utilizes keypoints extracted by Dlib [22] and OpenPose [3] to synthesize videos of 3D human face/body. To learn representation and preserve input information effectively, most recent methods are based on state-of-the-art generators with U-Net structure and AdaIN module [12, 20, 46, 54], feature disentangle/decomposition [49, 52], or specifically designed motion flow estimators [6, 46, 50]. On the contrary, our proposed TS-Net uses a more robust and general GAN generator [19] as backbone to jointly learn transformation and synthesis. Some previous works [46, 47] also performed video motion retargeting by combining warp-based and warp-free generation. However, their warping flows are always applied to previous generated frames, which may lead to the accumulation of synthesis artifacts. Our proposed TS-Net instead computes the warping flow between driving mask and *real* subject images in feature space to avoid this issue.

3. Methodology

3.1. Model Architecture

Given a sequence $\mathcal{X} = \{x_1, x_2, \dots, x_K\}$ with K subject frames, their corresponding mask sequence $\mathcal{Y} = \{y_1, y_2, \dots, y_K\}$, and a mask frame z from driving video, the TS-Net generator can produce a new video frame \hat{x} with the subject from \mathcal{X} and mask from z . Masks are generated by applying off-the-shelf pretrained 2D sparse keypoint detectors, *i.e.*, Dlib [22] for face landmark detection and OpenPose [3] for pose keypoint estimation. As illustrated in Fig. 2, TS-Net generator consists of two branches: a transformation branch Γ_{tra} and a synthesis branch Γ_{syn} for generating the new video frame using warp-based transformation and direct synthesis, respectively.

During training, we concatenate K subject frames \mathcal{X} and their masks \mathcal{Y} and feed them to an image encoder Δ_{img} to extract subject embedding features $\mathcal{E} = \{e_1, e_2, \dots, e_K\}$. A mask encoder Δ_{msk} encodes the input driving mask z into driving embedding feature f . To reduce computational costs of matrix multiplication, TS-Net operates in a low-resolution feature space, where the spatial size of \mathcal{E} and f are only $1/8^2$ of the input frames. We then input \mathcal{E} and f to the transformation branch Γ_{tra} and the synthesis branch Γ_{syn} , as illustrated as follows.

Transformation Branch. Inside Γ_{tra} , we implement warp-based transformation using spatial sampling grids [17]. We first compute the cosine similarity matrix S_k between the driving embedding feature f and the k -th subject feature e_k as

$$S_{k_{pq}} = \frac{e_{k_p} \cdot f_q}{\|e_{k_p}\|_2 \|f_q\|_2}, \quad (1)$$

where $S_{k_{pq}}$ is the affinity value between f_q at position q in map f , and, e_{k_p} at position p in map e_k , and $\|\cdot\|_2$ indicates the L2 norm. Suppose that the size of feature f and e_k are $m \times m$, the size of matrix S_k will be $m^2 \times m^2$, which is

quartic to m . Thus adopting low-resolution feature maps is important to alleviate computational overhead.

We further reduce computational costs by designing a novel mask-aware similarity computation method, as shown in Fig. 3. Given one driving mask z and one subject mask y , we first generate their corresponding bounding box b_z and b_y according to the maximum and minimum keypoint coordinates in masks. Intuitively, most of pixels inside the bounding box b_z will not be warped to the regions outside b_y thus we can skip similarity computation between pixels of these two regions. Based on this observation, we down-sample b_z and b_y to be the same spatial size as feature map f and e and then only compute the affinity values between points of their inside/outside-bounding-box regions.

For the input subject features \mathcal{E} , we now have K similarity matrices $S = \{S_1, S_2, \dots, S_K\}$. We then use similarity matrix S_k to weight the regular grid G and obtain the k -th sampling grid G'_k as

$$G'_{k_p} = \frac{\sum_q (\exp(\tau S_{k_{pq}}) \cdot G_q)}{\sum_q \exp(\tau S_{k_{pq}})}, \quad (2)$$

where G'_{k_p} is the coordinate p of sampling grid G'_k , G_q is the coordinate q of regular grid G , and τ is the coefficient to control the relative difference between affinity values. This results in K sampling grids $G' = \{G'_1, G'_2, \dots, G'_K\}$. By applying sampling grids G' to subject features \mathcal{E} , we acquire K warped features $\mathcal{V}^t = \{v_1^t, v_2^t, \dots, v_K^t\}$. The final warped feature \bar{v}^t is then generated by averaging the K features in \mathcal{V}^t .

Synthesis Branch. Inside Γ_{syn} , we concatenate the k -th subject embedding feature e_k with driving mask feature f and feed them to a fusion network Λ , which consists of a series of fully-convolutional layers, for creating the k -th synthesized warp-free feature map v_k^s . Processing K feature maps in \mathcal{E} will generate K synthesized feature maps $\mathcal{V}^s = \{v_1^s, v_2^s, \dots, v_K^s\}$. We then take average of the K features in \mathcal{V}^s to produce the final synthesized feature \bar{v}^s .

Combination of Branches. We concatenate the feature \bar{v}^t and \bar{v}^s of two branches and adopt a decoder network Ω_{dec} to synthesize the final output \hat{x} . We also tried to combine \bar{v}^t and \bar{v}^s with an attention-based matting function as in [46, 47], yet we found that this strategy fails to generate better results, as later illustrated in Sec. 4.3. More architecture details are in Sec. 4.2.

3.2. Training and Inference

We train our proposed TS-Net generator using a self-supervised way of training. Specifically, input driving mask sequence \mathcal{Z} and subject image sequence \mathcal{X} are from different segments of the same subject video. Thus we have frames in the subject video as ground truth. The overall loss for generating one frame is calculated as

$$l = \mathcal{L}_{\text{GAN}}(\hat{x}, x) + \alpha \mathcal{L}_{\text{VGG}}(\hat{x}, x) + \beta \mathcal{L}_{\text{FM}}(\hat{x}, x) + \lambda \mathcal{L}_{\text{TRA}}(\hat{x}^t, x), \quad (3)$$

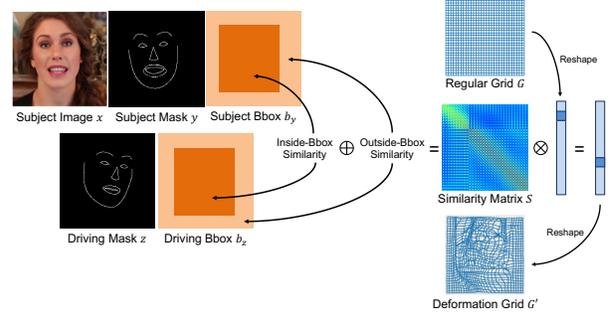


Figure 3: Illustration of our proposed mask-aware similarity computation.

where \mathcal{L}_{GAN} is an adversarial loss [10], \mathcal{L}_{VGG} represents a perceptual loss [19] based on VGG network [42], \mathcal{L}_{FM} is a feature matching loss [48], and \mathcal{L}_{TRA} is extra regularization loss for transformation branch. Here α , β , and λ are balancing factors, \hat{x} is the generated frame, \hat{x}^t is warped subject frame, and x is ground truth real frame.

We now introduce detailed loss terms. The adversarial loss \mathcal{L}_{GAN} is defined by the minimax optimization [10]:

$$\min_G \max_D \mathbb{E}_x [\log D(x)] + \mathbb{E}_{\hat{x}} [\log(1 - D(\hat{x}))]. \quad (4)$$

Discriminator D is designed to distinguish the real video frame x from the synthesized video frame \hat{x} given driving mask frame z . The perceptual loss \mathcal{L}_{VGG} is defined as

$$\sum_{i=1}^N \frac{1}{W_i} [\|F^{(i)}(\hat{x}) - F^{(i)}(x)\|_1], \quad (5)$$

where N is the number of layers in VGG feature extraction network and $F^{(i)}$ denotes the output of i -th layer with W_i elements of the VGG network [42] pretrained on ImageNet [7]. The feature matching loss \mathcal{L}_{FM} is defined as

$$\sum_{i=1}^M \frac{1}{U_i} [\|D^{(i)}(\hat{x}) - D^{(i)}(x)\|_1], \quad (6)$$

where $D^{(i)}$ denotes the i -th layer with U_i elements of our proposed discriminator D . The transformation branch loss \mathcal{L}_{TRA} is calculated as

$$\mathcal{L}_{\text{TRA}} = \|\hat{x}^t - x\|_1, \quad (7)$$

where \hat{x}^t is computed by patch-wise warping subject frame using deformation grid G' . For K input subject frames, we compute \mathcal{L}_{TRA} for each frame and then sum them up. In (3), the first three loss terms (\mathcal{L}_{GAN} , \mathcal{L}_{VGG} , and \mathcal{L}_{FM}) are commonly used in current video generation models [48, 47]. We show the importance of introducing \mathcal{L}_{TRA} to the training of our model in Sec. 4.3.

Inference. Given the subject video \mathcal{X} and the mask sequence of driving video \mathcal{Z} , we randomly select K frames from the subject video to synthesize a new frame \hat{x} .

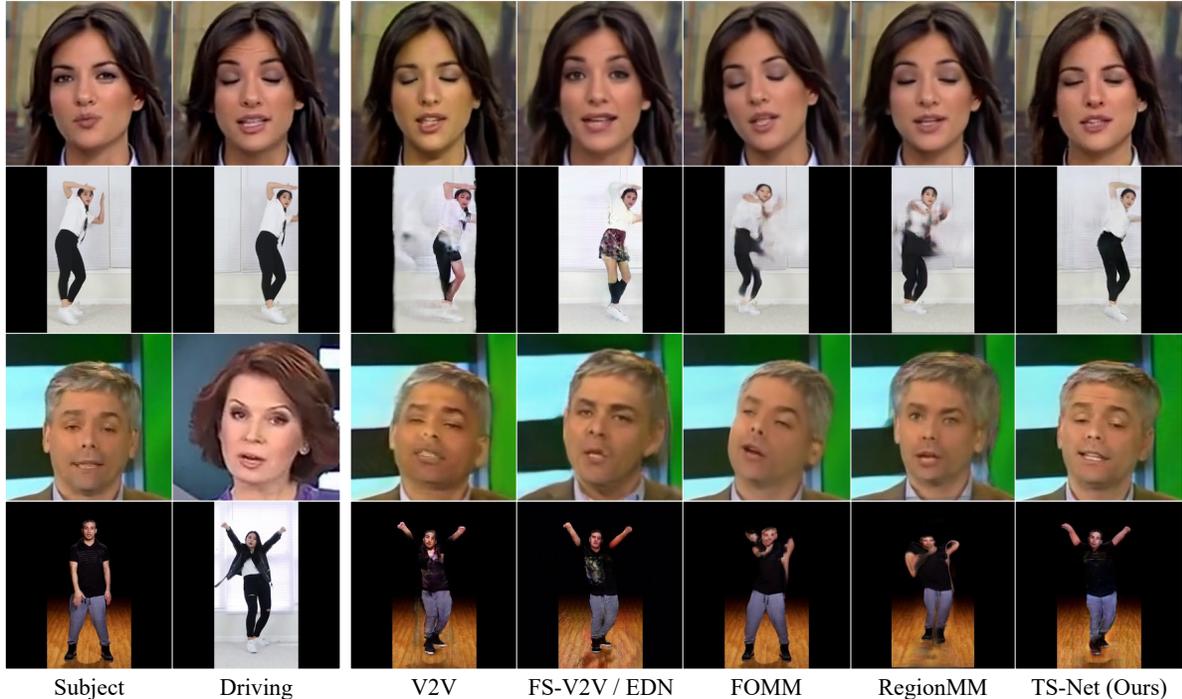


Figure 4: Qualitative comparison with state-of-the-art methods (V2V [47], FS-V2V [46], EDN [4], FOMM [39], and RegionMM [41]) on face and dance video datasets. The top two rows are results for self-reconstruction and the bottom two rows are for cross-identity transfer. Note that FS-V2V is used for face videos and EDN is used for dance videos.

4. Experiments

4.1. Datasets and Metrics

Datasets. We conduct experiments on face videos and dance videos. For face videos, We use the real videos in FaceForensics [35] dataset, which contains 1,004 videos of news briefing from different reporters. We randomly choose 150 videos for training and 150 videos for testing. Since the original videos are long, we randomly selected a short segment of 30 continuous frames from each video, and the selected short videos are used in our experiments. To extract mask sequences from videos, we first apply a face alignment algorithm [22] to localize 68 facial landmarks in each frame. The sparse facial landmarks are then connected to create the face mask. For dance videos, following [46, 47], we downloaded dancing videos from Youtube¹. We randomly chose 100 videos for training and 85 videos for testing and randomly sampled 30 continuous frames containing only one person from each video. We extracted human poses as masks via OpenPose [3]. Face and hand keypoints are kept for better motion retargeting.

Metrics. Following [38, 39], we compute metrics based on two testing settings, *self-reconstruction* and *cross-identity*

transfer. For each setting, we synthesize 100 videos where the size of each frame is 256×256 . For *self-reconstruction*, we segment a video of the same subject to two non-overlapping clips and use one clip as subject video and another one as driving video. In this setting, driving video also serves as ground truth. Similar to [9], we compute the normalized mean L_2 distance and Learned Perceptual Image Patch Similarity (LPIPS) [56] metrics between self-reconstructed results and driving videos. For *cross-identity transfer*, which is more practical in real world applications, subject video and driving video are from different subjects in this setting. Due to the lack of ground truth, we conduct user study to compare our models with state-of-the-art methods. Human evaluators are shown sets of n videos generated by n different models and then are asked to rank videos in each set from 1 (best) to n (worst) based on perceptual similarity and realism. Tied rank scores will be given for videos that are perceived to have comparable quality.

4.2. Implementation

Model Implementation. Our proposed encoder Δ and decoder Θ in TS-Net are general and can have various backbone networks, such as pix2pix [16] and SPADE [31]. We adopt the architecture in [19] due to its simplicity. For

¹The video links are available on the project website of [46]. We obtained permission to use the videos from the video owners.

Dataset	Method	$L_2 \downarrow$	LPIPS \downarrow
Face	V2V [47]	0.0356	0.1123
	FS-V2V [46]	0.0422	0.1064
	FOMM [39]	0.0443	0.1184
	RegionMM [41]	0.0148	0.0532
	TS-Net ($K = 1$)	0.0275	0.0731
	TS-Net ($K = 3$)	0.0271	0.0683
Dance	TS-Net ($K = 5$)	0.0270	0.0673
	V2V [47]	0.0895	0.2622
	EDN [4]	0.0471	0.1718
	FOMM [39]	0.1517	0.3081
	RegionMM [41]	0.1945	0.4081
	TS-Net ($K = 1$)	0.0433	0.1586
TS-Net ($K = 3$)	0.0421	0.1543	
TS-Net ($K = 5$)	0.0423	0.1541	

Table 1: Comparison with state-of-the-art methods under the self-reconstruction setting on face and dance datasets. K is the number of subject frames used in generation.

the encoder Δ_{img} , we use the network with three stride-2 convolutions and 9 residual blocks [13]. For Δ_{msk} , we use three stride-2 convolutions without additional residual blocks since masks contain less information. Thus the spatial size of embedding feature map is only $1/8^2$ size of input image. To encode position-related information for better synthesis, we apply coordinate convolution [25] to inputs. For decoder Θ_{dec} , we employ 4 residual blocks, followed by three up-sampling and convolution layers. For fusion network Λ , we use one residual block and one 1×1 convolution [24] to generate warp-free feature maps \mathcal{V}^s . Instance normalization [45] is adopted in TS-Net. For the discriminator D , we use 70×70 PatchGAN [16, 48, 57], which aims to classify whether the 70×70 overlapping patches are real or fake. To stabilize the training, we use LSGAN [28] for the adversarial loss.

When training TS-Net, we set batch size as 20 videos and train the model for 600 epochs using the Adam optimizer [23] with $(\beta_1, \beta_2) = (0.5, 0.999)$. The learning rate is fixed to 2×10^{-4} in the first 275 epochs and then linearly decayed to zero. The balancing parameters α , β , and λ are all set to be 10 in (3). The coefficient τ in (2) is set to be 100. Data augmentation such as color jitter and flipping are also applied. Hyper-parameters are selected via multiple runs of experiments. When training our models on the face video dataset, we adopt an image gradient difference loss [29] as an extra smoothness constraint to eliminate minor artifacts in the generated videos. When training our models on the dance video dataset, similar to [4, 46, 47], we introduce an extra face discriminator to synthesize better face details. To normalize masks across different subjects, the masks of driving videos are aligned to the masks of subject videos with the similar methods used in [4, 46].

Baselines. For face video dataset, we choose four state-of-the-art video synthesis or image animation models, vid2vid (V2V) [47], few-shot vid2vid (FS-V2V) [46], FOMM [39],

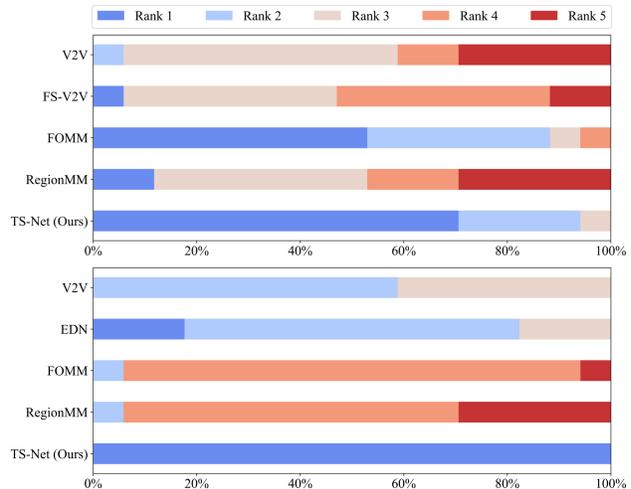


Figure 5: User study of ranking different methods under the cross-identity setting. Ties are allowed. The top chart is for face videos and the bottom one is for dance videos.



Figure 6: Ablation study under the cross-identity transfer setting on the face dataset.

and RegionMM [41] as baselines. For dance video dataset, we compare TS-Net with V2V, FOMM, RegionMM, and Everybody Dance Now (EDN) [4]. FS-V2V is not included for dance videos since it requires DensePose [11] as extra inputs. We follow the default settings in the methods' original implementations wherever possible. The original V2V and EDN train with a single video and test on the same video. For fair comparison, we train V2V and EDN using all available training videos.

4.3. Result Analysis

Comparison with state-of-the-art methods. Table 1 shows a quantitative comparison of our models with state-of-the-art methods under the self-reconstruction setting. TS-Net achieves comparable or better performance when compared with state-of-the-art methods even when using only one subject frame ($K = 1$). Though RegionMM [41]

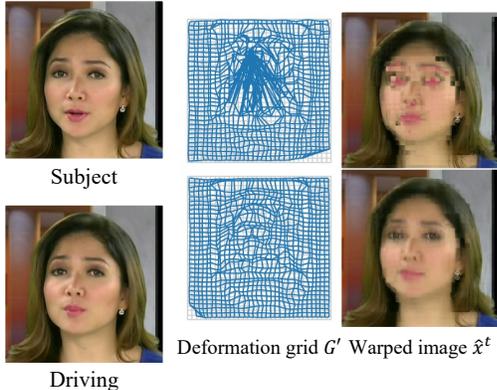


Figure 7: Ablation study for \mathcal{L}_{TRA} . The top row is from T-Net without using \mathcal{L}_{TRA} and the bottom one is with \mathcal{L}_{TRA} . \mathcal{L}_{TRA} clearly learns more reasonable deformation grid and improves warped result.

achieves the best performance in the metrics for face videos under the self-reconstruction setting, it performs worse than TS-Net on all the other tasks (*i.e.*, dance videos and cross-identity transfer settings). The reason may be that RegionMM relies on an unsupervised trained region detection network, which may not be robust enough to handle large motions or fine-grained details in various tasks. In Fig. 4, one can observe that V2V [47] suffers from color/shape distortion, FS-V2V [46] misses some details (*e.g.* opened eyes in first row), EDN fails to preserve some details such as in the face or clothing, FOMM [39] struggles to capture the head/body pose correctly, and RegionMM [41] generates images with some blurry regions and unrealistic appearance details. (Similar results can also be observed in Fig. 1). In contrast, TS-Net can better handle large motion and preserve identity. Table 1 also confirms the effectiveness of using multiple subject frames in TS-Net to collect various appearance information, where most metrics get improved as the number of subject frames increases. For the cross-identity transfer setting, we conduct a user study to compare models with human perception. As shown in Fig. 5, TS-Net gets the most user preference, especially on dance videos.

Ablation study. To analyze the effectiveness of each module in TS-Net, we conduct an ablation study on the face video dataset. The number of input subject frames is fixed to be 3 ($K = 3$) to ease model training and testing. Table 2 shows the quantitative comparison results of the ablation study under the self-reconstruction setting. We first train and test two single branch models, T-Net (Γ_{tra}) and S-Net (Γ_{syn}), which only employ the transformation branch or synthesis branch, respectively. From the results shown in Table 1 and Table 2, one can observe that even a single branch can achieve promising performance when compared to other state-of-the-art methods. However, as shown in Fig. 6, warp-based T-Net fails to generate unseen con-

Method	$L_2 \downarrow$	LPIPS \downarrow
T-Net	0.0276	0.0698
T-Net w/o \mathcal{L}_{TRA}	0.0287	0.0725
S-Net	0.0285	0.0726
TS-Net w/ cross	0.0276	0.0696
TS-Net w/ matting	0.0281	0.0696
TS-Net	0.0271	0.0683

Table 2: Ablation Study under the self-reconstruction setting on face dataset. The number of input subject frames is fixed to be 3 ($K = 3$).

tent (*e.g.*, regions marked by red box) while warp-free S-Net is incompetent to preserve identity. Results demonstrate that a single T-Net or S-Net enables efficient representation learning, and the combination of two branches can complement each other to achieve more satisfactory results. We also compare the transformation branch trained with and without \mathcal{L}_{TRA} , T-Net and [T-Net w/o \mathcal{L}_{TRA}] in Table 2, from which one can observe that removing \mathcal{L}_{TRA} diminishes performance. As Fig. 7 shows, the lack of \mathcal{L}_{TRA} led to a less meaningful deformation grid G' and resulted in a poor warped image \hat{x}^t .

We also evaluate the effectiveness of some common techniques adopted by previous video motion retargeting methods [18, 46, 47], such as adding cross-identity transfer to the training processing [TS-Net w/ cross] or using the matting function to combine different types of features [TS-Net w/ matting]. To enable cross-identity training, we choose input mask sequence \mathcal{Z} and input image sequence \mathcal{X} from different videos. Thus ground truth frames are not available for training. In this case, we only use adversarial loss \mathcal{L}_{GAN} for training, where discriminator D is designed to distinguish the synthesized frame \hat{x} from arbitrary real video frame x . For the matting function, we design an extra attention network with similar architecture as fusion network Λ to generate a matting mask for combining \bar{v}^t and \bar{v}^s . However, both these modules fail to be more effective as Table 2 shows, which demonstrates that the simple design of TS-Net has already achieved sufficient representation learning and synthesis power.

5. Limitations

For most cases, our proposed TS-Net can generate realistic videos by only taking 2D sparse masks (see Fig. 8 and Supp. videos). However, it still suffers from several limitations. First, the input masks of TS-Net are generated from off-the-shelf detectors. Misdetections by the detectors could result in inconsistent motion or incorrect appearances. As shown in Fig. 9, the synthesized face in the top row has an opened mouth, and the generated man in the middle row shows missing hands. Second, TS-Net sometimes struggles to synthesize high-frequency details. One can observe a few



Figure 8: Examples of generated face videos (top block) and dance videos (bottom block) using our proposed TS-Net. For each block, TS-Net synthesizes the new video (3rd row) with the appearance from subject video (1st row) and the motion from driving video (2nd row).

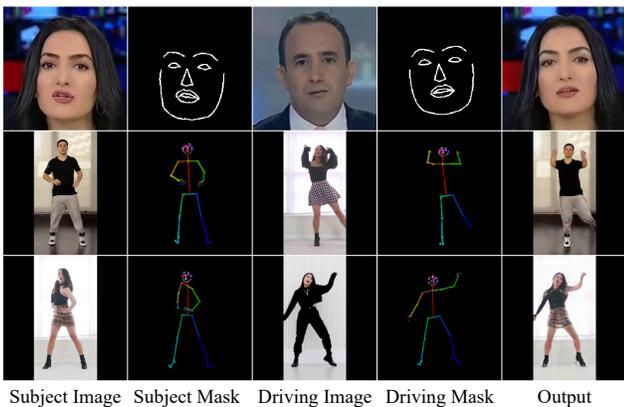


Figure 9: Some failure cases of TS-Net. Driving masks are aligned to match subject masks using the normalization methods in [4, 46].

texture artifacts in the kilt of the last row in Fig. 9. Future work could focus on improving the keypoint detection system and generating more realistic high-frequency textures. **Potential Negative Societal Impact.** Video motion retargeting could be used for unethical purposes [53], *e.g.*, creating videos of celebrities for fake news. We will restrict the usage of our method and model to research purposes only. We also plan to investigate fake video detection techniques [33] that will be effective in detecting fake videos like the ones generated by our proposed method.

6. Conclusion

In this paper, we propose TS-Net to jointly learn transformation and synthesis for video motion transfer. Comprehensive experiments show that TS-Net can achieve state-of-the-art performance on both face and dance videos using only 2D sparse masks. In the future, we plan to investigate TS-Net using different kinds of masks and multi-modal information (*e.g.* audio or text) in motion retargeting.

References

- [1] Guha Balakrishnan, Amy Zhao, Adrian V Dalca, Fredo Durand, and John Gutttag. Synthesizing images of humans in unseen poses. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8340–8348, 2018.
- [2] Aayush Bansal, Shugao Ma, Deva Ramanan, and Yaser Sheikh. Recycle-gan: Unsupervised video retargeting. In *Proceedings of the European conference on computer vision (ECCV)*, pages 119–135, 2018.
- [3] Zhe Cao, Gines Hidalgo, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Openpose: realtime multi-person 2d pose estimation using part affinity fields. *IEEE transactions on pattern analysis and machine intelligence*, 43(1):172–186, 2019.
- [4] Caroline Chan, Shiry Ginosar, Tinghui Zhou, and Alexei A Efros. Everybody dance now. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5933–5942, 2019.
- [5] Zhuo Chen, Chaoyue Wang, Bo Yuan, and Dacheng Tao. Puppeteergan: Arbitrary portrait animation with semantic-aware appearance transformation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13518–13527, 2020.
- [6] Mengyu Chu, You Xie, Jonas Mayer, Laura Leal-Taixé, and Nils Thuerey. Learning temporal coherence via self-supervision for gan-based video generation. *ACM Transactions on Graphics (TOG)*, 39(4):75–1, 2020.
- [7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [8] Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Hausser, Caner Hazirbas, Vladimir Golkov, Patrick Van Der Smagt, Daniel Cremers, and Thomas Brox. FlowNet: Learning optical flow with convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2758–2766, 2015.
- [9] Guy Gafni, Justus Thies, Michael Zollhofer, and Matthias Nießner. Dynamic neural radiance fields for monocular 4d facial avatar reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8649–8658, 2021.
- [10] Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *arXiv preprint arXiv:1406.2661*, 2014.
- [11] Rıza Alp Güler, Natalia Neverova, and Iasonas Kokkinos. Densepose: Dense human pose estimation in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7297–7306, 2018.
- [12] Sungjoo Ha, Martin Kersner, Beomsu Kim, Seokjun Seo, and Dongyoung Kim. Marionette: Few-shot face reenactment preserving identity of unseen targets. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 10893–10900, 2020.
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [14] Fa-Ting Hong, Longhao Zhang, Li Shen, and Dan Xu. Depth-aware generative adversarial network for talking head video generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3397–3406, 2022.
- [15] Yuan-Ting Hu, Jia-Bin Huang, and Alexander G Schwing. Videomatch: Matching based video object segmentation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 54–70, 2018.
- [16] 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.
- [17] Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. Spatial transformer networks. *Advances in neural information processing systems*, 28, 2015.
- [18] Subin Jeon, Seonghyeon Nam, Seoung Wug Oh, and Seon Joo Kim. Cross-identity motion transfer for arbitrary objects through pose-attentive video reassembling. In *European Conference on Computer Vision*, pages 292–308. Springer, 2020.
- [19] 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.
- [20] Hyeonwoo Kim, Pablo Garrido, Ayush Tewari, Weipeng Xu, Justus Thies, Matthias Niessner, Patrick Pérez, Christian Richardt, Michael Zollhöfer, and Christian Theobalt. Deep video portraits. *ACM Transactions on Graphics (TOG)*, 37(4):1–14, 2018.
- [21] Junho Kim, Minjae Kim, Hyeonwoo Kang, and Kwanghee Lee. U-gat-it: unsupervised generative attentional networks with adaptive layer-instance normalization for image-to-image translation. *arXiv preprint arXiv:1907.10830*, 2019.
- [22] Davis E King. Dlib-ml: A machine learning toolkit. *The Journal of Machine Learning Research*, 10:1755–1758, 2009.
- [23] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [24] Min Lin, Qiang Chen, and Shuicheng Yan. Network in network. *arXiv preprint arXiv:1312.4400*, 2013.
- [25] Rosanne Liu, Joel Lehman, Piero Molino, Felipe Petroski Such, Eric Frank, Alex Sergeev, and Jason Yosinski. An intriguing failing of convolutional neural networks and the coordconv solution. *arXiv preprint arXiv:1807.03247*, 2018.
- [26] Yihao Liu, Lianrui Zuo, Shuo Han, Jerry L Prince, and Aaron Carass. Coordinate translator for learning deformable medical image registration. *arXiv preprint arXiv:2203.03626*, 2022.
- [27] Liqian Ma, Xu Jia, Qianru Sun, Bernt Schiele, Tinne Tuytelaars, and Luc Van Gool. Pose guided person image generation. *arXiv preprint arXiv:1705.09368*, 2017.

- [28] Xudong Mao, Qing Li, Haoran Xie, Raymond YK Lau, Zhen Wang, and Stephen Paul Smolley. Least squares generative adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2794–2802, 2017.
- [29] Michael Mathieu, Camille Couprie, and Yann LeCun. Deep multi-scale video prediction beyond mean square error. *arXiv preprint arXiv:1511.05440*, 2015.
- [30] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014.
- [31] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2337–2346, 2019.
- [32] Albert Pumarola, Antonio Agudo, Alberto Sanfeliu, and Francesc Moreno-Noguer. Unsupervised person image synthesis in arbitrary poses. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8620–8628, 2018.
- [33] Ashifur Rahman, Md Mazharul Islam, Mohasina Jannat Moon, Tahera Tasnim, Nipo Siddique, Md Shahiduzzaman, and Samsuddin Ahmed. A qualitative survey on deep learning based deep fake video creation and detection method. *Aust. J. Eng. Innov. Technol*, 4(1):13–26, 2022.
- [34] Yurui Ren, Ge Li, Yuanqi Chen, Thomas H Li, and Shan Liu. Pirenderer: Controllable portrait image generation via semantic neural rendering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13759–13768, 2021.
- [35] Andreas Rössler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, and Matthias Nießner. FaceForensics: A large-scale video dataset for forgery detection in human faces. *arXiv*, 2018.
- [36] Kripasindhu Sarkar, Dushyant Mehta, Weipeng Xu, Vladislav Golyanik, and Christian Theobalt. Neural re-rendering of humans from a single image. In *European Conference on Computer Vision*, pages 596–613. Springer, 2020.
- [37] Aliaksandr Siarohin, Stéphane Lathuilière, Enver Sangineto, and Nicu Sebe. Appearance and pose-conditioned human image generation using deformable gans. *IEEE transactions on pattern analysis and machine intelligence*, 2019.
- [38] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. Animating arbitrary objects via deep motion transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2377–2386, 2019.
- [39] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. First order motion model for image animation. *arXiv preprint arXiv:2003.00196*, 2020.
- [40] Aliaksandr Siarohin, Enver Sangineto, Stéphane Lathuilière, and Nicu Sebe. Deformable gans for pose-based human image generation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3408–3416, 2018.
- [41] Aliaksandr Siarohin, Oliver J Woodford, Jian Ren, Menglei Chai, and Sergey Tulyakov. Motion representations for articulated animation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13653–13662, 2021.
- [42] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [43] Sijie Song, Wei Zhang, Jiaying Liu, and Tao Mei. Unsupervised person image generation with semantic parsing transformation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2357–2366, 2019.
- [44] Jiale Tao, Biao Wang, Borun Xu, Tiezheng Ge, Yuning Jiang, Wen Li, and Lixin Duan. Structure-aware motion transfer with deformable anchor model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3637–3646, 2022.
- [45] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*, 2016.
- [46] Ting-Chun Wang, Ming-Yu Liu, Andrew Tao, Guilin Liu, Jan Kautz, and Bryan Catanzaro. Few-shot video-to-video synthesis. *arXiv preprint arXiv:1910.12713*, 2019.
- [47] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Guilin Liu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. Video-to-video synthesis. *arXiv preprint arXiv:1808.06601*, 2018.
- [48] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8798–8807, 2018.
- [49] Ting-Chun Wang, Arun Mallya, and Ming-Yu Liu. One-shot free-view neural talking-head synthesis for video conferencing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10039–10049, 2021.
- [50] Olivia Wiles, A Koepke, and Andrew Zisserman. X2face: A network for controlling face generation using images, audio, and pose codes. In *Proceedings of the European conference on computer vision (ECCV)*, pages 670–686, 2018.
- [51] Chang Xu, Dacheng Tao, and Chao Xu. A survey on multi-view learning. *arXiv preprint arXiv:1304.5634*, 2013.
- [52] Zhuoqian Yang, Wentao Zhu, Wayne Wu, Chen Qian, Qiang Zhou, Bolei Zhou, and Chen Change Loy. Transmomo: Invariance-driven unsupervised video motion retargeting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5306–5315, 2020.
- [53] Sihyun Yu, Jihoon Tack, Sangwoo Mo, Hyunsu Kim, Junho Kim, Jung-Woo Ha, and Jinwoo Shin. Generating videos with dynamics-aware implicit generative adversarial networks. *arXiv preprint arXiv:2202.10571*, 2022.
- [54] Egor Zakharov, Aliaksandra Shysheya, Egor Burkov, and Victor Lempitsky. Few-shot adversarial learning of realistic neural talking head models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9459–9468, 2019.
- [55] Pan Zhang, Bo Zhang, Dong Chen, Lu Yuan, and Fang Wen. Cross-domain correspondence learning for exemplar-based

- image translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5143–5153, 2020.
- [56] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018.
- [57] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017.