

# Motion Representations for Articulated Animation

Aliaksandr Siarohin<sup>1\*</sup>, Oliver J. Woodford\*, Jian Ren<sup>2</sup>, Menglei Chai<sup>2</sup> and Sergey Tulyakov<sup>2</sup>

<sup>1</sup>DISI, University of Trento, Italy, <sup>2</sup>Snap Inc., Santa Monica, CA

aliaksandr.siarohin@unitn.it, {jren, mchai, stulyakov}@snap.com, \*Work done while at Snap Inc.

## Abstract

We propose novel motion representations for animating articulated objects consisting of distinct parts. In a completely unsupervised manner, our method identifies object parts, tracks them in a driving video, and infers their motions by considering their principal axes. In contrast to the previous keypoint-based works, our method extracts meaningful and consistent regions, describing locations, shape, and pose. The regions correspond to semantically relevant and distinct object parts, that are more easily detected in frames of the driving video. To force decoupling of foreground from background, we model non-object related global motion with an additional affine transformation. To facilitate animation and prevent the leakage of the shape of the driving object, we disentangle shape and pose of objects in the region space. Our model<sup>1</sup> can animate a variety of objects, surpassing previous methods by a large margin on existing benchmarks. We present a challenging new benchmark with high-resolution videos and show that the improvement is particularly pronounced when articulated objects are considered, reaching 96.6% user preference vs. the state of the art.

## 1. Introduction

Animation—bringing static objects to life—has broad applications across education and entertainment. Animated characters and objects, such as those in Fig. 1, increase the creativity and appeal of content, improve the clarity of material through storytelling, and enhance user experiences<sup>2</sup>.

Until very recently, animation techniques necessary for achieving such results required a trained professional, specialized hardware, software, and a great deal of effort. Quality results generally still do, but vision and graphics communities have attempted to address some of these limitations by training data-driven methods [39, 6, 27, 11, 10] on object classes for which prior knowledge of object shape and pose can be learned. This, however, requires ground truth pose and shape data to be available during training.

<sup>1</sup>Our source code is publicly available at <https://github.com/snap-research/articulated-animation>.

<sup>2</sup>Visit our project [website](#) for more qualitative samples.

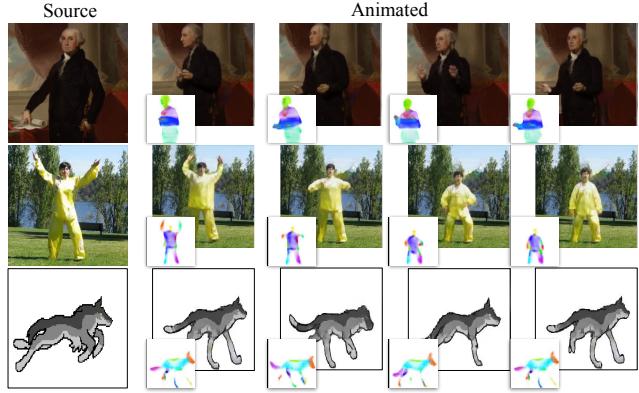


Figure 1: Our method animates still source images via unsupervised region detection (inset).

Recent works have sought to avoid the need for ground truth data through *unsupervised* motion transfer [42, 30, 31]. Significant progress has been made on several key challenges, including training using image reconstruction as a loss [42, 30, 31], and disentangling motion from appearance [20]. This has created the potential to animate a broader range of object categories, without any domain knowledge or labelled data, requiring only videos of objects in motion during training [30]. However, two key problems remain open. The first is how to represent the parts of an articulated or non-rigid moving object, including their shapes and poses. The second is given the object parts, how to animate them using the sequence of motions in a driving video.

Initial attempts used end-to-end frameworks [42, 31] to first extract unsupervised keypoints [20, 17], then warp a feature embedding of a source image to align its keypoints with those of a driving video. Follow on work [30] further modelled the motion around each keypoint with local, affine transformations, and introduced a generation module that both composites warped source image regions and inpaints occluded regions, to render the final image. This enabled a variety of creative applications<sup>3</sup>, for example needing only one source face image to generate a near photo-realistic animation, driven by a video of a different face.

<sup>3</sup>E.g. a [music video](#) in which images are animated using prior work [30].

However, the resulting unsupervised keypoints are detected on the boundary of the objects. While points on edges are easier to identify, tracking such keypoints between frames is problematic, as any point on the boundary is a valid candidate, making it hard to establish correspondences between frames. A further problem is that the unsupervised keypoints do not correspond to semantically meaningful object parts, representing location and direction, but not shape. Due to this limitation, animating articulated objects, such as bodies, remains challenging. Furthermore, these methods assume static backgrounds, i.e. no camera motion, leading to leakage of background motion information into one or several of the detected keypoints. Finally, absolute motion transfer, as in [30], transfers the shape of the driving object into the generated sequence, decreasing the fidelity of the source identity. These remaining deficiencies limit the scope of previous works [30, 31] to more trivial object categories and motions, especially when objects are articulated.

This work introduces three contributions to address these challenges. First, we redefine the underlying motion representation, using *regions* from which first-order motion is *measured*, rather than regressed. This enables improved convergence, more stable, robust object and motion representations, and also empirically captures the shape of the underpinning object parts, leading to better motion segmentation. This motion representation is inspired by Hu moments [8]. Fig. 3(a) contains several examples of region vs. keypoint-based motion representation.

Secondly, we explicitly model background or camera motion between training frames by predicting the parameters of a global, affine transformation explaining non-object related motions. This enables the model to focus solely on the foreground object, making the identified points more stable, and further improves convergence. Finally, to prevent shape transfer and improve animation, we disentangle the shape and pose of objects in the space of unsupervised regions. Our framework is self-supervised, does not require any labels, and is optimized using reconstruction losses.

These contributions further improve unsupervised motion transfer methods, resulting in higher fidelity animation of articulated objects in particular. To create a more challenging benchmark for such objects, we present a newly collected dataset of TED talk speakers. Our framework scales better in the number of unsupervised regions, resulting in more detailed motion. Our method outperforms previous unsupervised animation methods on a variety of datasets, including talking faces, taichi videos and animated pixel art being preferred by 96.6% of independent raters when compared with the state of the art [30] on our most challenging benchmark.

## 2. Related work

Image animation methods can be separated into supervised, which require knowledge about the animated object during training, and unsupervised, which do not. Such

knowledge typically includes landmarks [4, 44, 26, 12], semantic segmentations [24], and parametric 3D models [11, 35, 9, 22, 19]. As a result, supervised methods are limited to a small number of object categories for which a lot of labelled data is available, such as faces and human bodies. Early face reenactment work [35] fitted a 3D morphable model to an image, animating and rendering it back using graphical techniques. Further works used neural networks to get higher quality rendering [16, 41], sometimes requiring multiple images per identity [11, 25]. A body of works treats animation as an image-to-image [33] or video-to-video [39, 6, 27] translation problem. Apart from some exceptions [38], these works further constrain the problem to animating a single instance of an object, such as a single face [16, 2] or a single human body [6, 27, 39], requiring retraining [2, 6, 27] or fine-tuning [44] for each new instance. Despite promising results, generalizing these methods beyond a limited range of object categories remains challenging. Additionally, they tend to transfer not only the motion but also the shape of the driving object [16, 44].

Unsupervised methods address some of these limitations. They do not require any labelled data regarding the shape or landmarks of the animated object. Video-generation-based animation methods predict future frames of a video, given the first frame and an animation class label, such as “make a happy face”, “do jumping jack”, or “play golf” [36, 29, 7]. Recently Menapace *et al.* [21] introduce playable video generation, where action could be selected at each timestamp. A further group of works re-target animation from a driving video to a source frame. X2Face [42] builds a canonical representation of an input face, and generates a warp field conditioned on the driving video. Monkey-Net [31] learns a set of unsupervised keypoints to generate animations. Follow-up work substantially improves the quality of animation by considering a first order motion model (FOMM) [30] for each keypoint, represented by regressing a local, affine transformation. Both of these works apply to a wider range of objects including faces, bodies, robots, and pixel art animations. Empirically, these methods extract keypoints on the boundary of the animated objects. Articulated objects such as human bodies are therefore challenging, as internal motion, for example, an arm moving across the body, is not well modeled, producing unconvincing animations.

This work presents an unsupervised method. We argue that the limitations of previous such methods in animating articulated objects is due to an inability of their internal representations to capture complete object parts, their shape and pose. X2Face [42] assumes an object can be represented with a single RGB texture, while other methods find keypoints on edges [31, 30]. Our new region motion representation resembles the construction of a motion history image whose shape is analyzed using principal components (namely Hu moments [8]). In [8] the authors construct motion descrip-

tors by computing temporal image differences, aggregate them into motion history images, and use Hu moments to build a motion recognition system. In this manner, blob statistics are used to discriminate between different actions.

### 3. Method

We propose three contributions over FOMM [30], namely our PCA-based motion estimation (Sec. 3.2), background motion representation (Sec. 3.3) and animation via disentanglement (Sec. 3.6). To make the manuscript self-contained, we first give the necessary technical background on the original first order motion representation [30] (Sec. 3.1).

#### 3.1. First Order Motion Model

FOMM [30] consists of the two main parts: motion estimation and image generation, where motion estimation further contains coarse motion estimation and dense motion prediction. Coarse motion is modelled as sparse motions between separate object parts, while dense motion produces an optical flow along with the confidence map for the entire image. We denote by  $\mathbf{S}$  and  $\mathbf{D}$  the source and the driving frames extracted from the same video respectively.

The first step is to estimate coarse motions from  $\mathbf{S}$  and  $\mathbf{D}$ . Motions for each object part are represented by affine transformations,  $A_{\mathbf{X} \leftarrow \mathbf{R}}^k \in \mathcal{R}^{2 \times 3}$ , to an abstract, common reference frame,  $\mathbf{R}$ ;  $\mathbf{X}$  is either  $\mathbf{S}$  or  $\mathbf{D}$ . Motions are estimated for  $K$  distinct parts. An encoder-decoder *keypoint* predictor network outputs  $K$  heatmaps,  $\mathbf{M}^1, \dots, \mathbf{M}^K$  for the input image, followed by softmax, s.t.  $\mathbf{M}^k \in [0, 1]^{H \times W}$ , where  $H$  and  $W$  are the height and width of the image respectively, and  $\sum_{z \in \mathcal{Z}} \mathbf{M}^k(z) = 1$ , where  $z$  is a pixel location ( $x, y$  coordinates) in the image, the set of all pixel locations being  $\mathcal{Z}$ , and  $\mathbf{M}^k(z)$  is the  $k$ -th heatmap weight at pixel  $z$ . Thus, the translation component of the affine transformation (which is the last column of  $A_{\mathbf{X} \leftarrow \mathbf{R}}^k$ ) can be estimated using softargmax:

$$\mu^k = \sum_{z \in \mathcal{Z}} \mathbf{M}^k(z) z. \quad (1)$$

In FOMM [30] the remaining affine parameters are regressed per pixel and form 4 additional channels  $P_{ij}^k \in \mathcal{R}^{H \times W}$ , where  $i \in \{0, 1\}, j \in \{0, 1\}$  indexes of the affine matrix  $A_{\mathbf{X} \leftarrow \mathbf{R}}^k$ . The latter is estimated using weighted pooling:

$$A_{\mathbf{X} \leftarrow \mathbf{R}}^k[i, j] = \sum_{z \in \mathcal{Z}} \mathbf{M}^k(z) P_{ij}^k(z). \quad (2)$$

We refer to this way of computing motion as *regression-based*, where affine parameters are *predicted* by a network and pooled to compute  $A_{\mathbf{X} \leftarrow \mathbf{R}}^k$ . Motion between  $\mathbf{D}$  and  $\mathbf{S}$  for part  $k$  is then computed via the common reference frame:

$$A_{\mathbf{S} \leftarrow \mathbf{D}}^k = A_{\mathbf{S} \leftarrow \mathbf{R}}^k \begin{bmatrix} A_{\mathbf{D} \leftarrow \mathbf{R}}^k \\ 0 & 0 & 1 \end{bmatrix}^{-1}. \quad (3)$$

Given the coarse motion, the next step is to predict the optical flow and the confidence map. Since the transformations between  $\mathbf{D}$  and  $\mathbf{S}$  are known for each part, the task is to combine them to obtain a single optical flow field. To this end, the flow predictor selects the appropriate coarse transformation for each pixel, via a weighted sum of coarse motions. Formally,  $K+1$  assignment maps are output, one per region, plus background, which FOMM [30] assumes is *motionless*, and softargmax is applied pixelwise across them, s.t.  $\mathbf{W}^k \in [0, 1]^{H \times W}$ ,  $\mathbf{W}^0$  corresponds to the background, and  $\sum_{k=0}^K \mathbf{W}^k(z) = 1, \forall z$ . Optical flow per pixel,  $\mathbf{O}(z) \in \mathbb{R}^2$ , is then computed as:

$$\mathbf{O}(z) = \mathbf{W}^0(z)z + \sum_{k=1}^K \mathbf{W}^k(z) A_{\mathbf{S} \leftarrow \mathbf{D}}^k \begin{bmatrix} z \\ 1 \end{bmatrix}. \quad (4)$$

With such a model, animation becomes challenging when there is even slight background motion. The model automatically adapts by assigning several of the available keypoints to model background as shown in the first row of Fig. 3(a).

A confidence map,  $\mathbf{C} \in [0, 1]^{H \times W}$ , is also predicted using the same network, to handle parts missing in the source image. Finally  $\mathbf{S}$  is passed through an encoder, followed by warping the resulting feature map using the optical flow (Eq. 4) and multiplied by the confidence map. A decoder then reconstructs the driving image  $\mathbf{D}$ .

At test time FOMM [30] has two modes of animating  $\mathbf{S}$ : *standard* and *relative*. In both cases the input is the source image  $\mathbf{S}$  and the driving video  $\mathbf{D}_1, \mathbf{D}_2, \dots, \mathbf{D}_t, \dots, \mathbf{D}_T$ . In the *standard* animation the motion between source and driving is computed frame-by-frame using Eq. (3). For *relative* animation, in order to generate a frame  $t$  the motion between  $\mathbf{D}_1$  and  $\mathbf{D}_t$  is computed first and then applied to  $\mathbf{S}$ . Both of these modes are problematic when the object in question is articulated, as we show in Sec. 3.6.

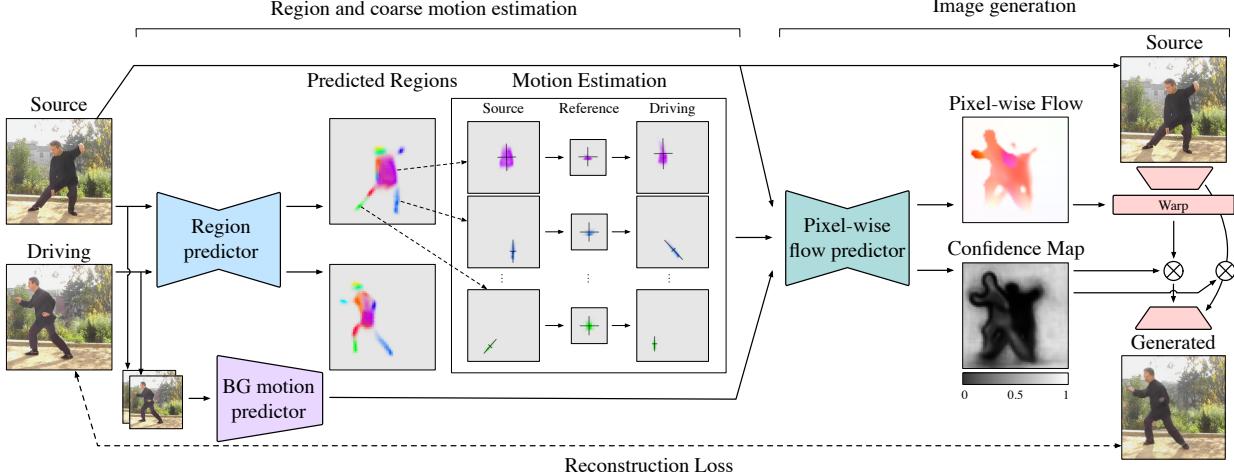
#### 3.2. PCA-based motion estimation

Accurate motion estimation is the main requirement for high-quality image animation. As mentioned previously, FOMM regresses the affine parameters. This requires higher capacity networks, and generalizes poorly (see Sec. 4.1). We propose a different motion representation: *all* motions are *measured* directly from the heatmap  $\mathbf{M}^k$ . We compute the translation as before, while in-plane rotation and scaling in x- and y-directions are computed via a principal component analysis (PCA) of the heatmap  $\mathbf{M}^k$ . Formally, the transformation  $A_{\mathbf{X} \leftarrow \mathbf{R}}^k \in \mathbb{R}^{2 \times 3}$ , of the  $k$ <sup>th</sup> region from the reference frame to the image is computed as:

$$\mu^k = \sum_{z \in \mathcal{Z}} \mathbf{M}^k(z) z, \quad (5)$$

$$U^k S^k V^k = \sum_{z \in \mathcal{Z}} \mathbf{M}^k(z) (z - \mu^k) (z - \mu^k)^T \text{ (SVD)}, \quad (6)$$

$$A_{\mathbf{X} \leftarrow \mathbf{R}}^k = \left[ U^k S^{k \frac{1}{2}}, \mu^k \right]. \quad (7)$$



**Figure 2: Overview of our model.** The region predictor returns heatmaps for each part in the source and the driving frame. We then compute principal axes of each heatmap, to transform each region from the source to the driving frame through a whitened reference frame. Region and background transformations are combined by the pixel-wise flow prediction network. The target image is generated by warping the source image in a feature space using the pixel-wise flow, and inpainting newly introduced regions, as indicated by the confidence map.

Here the singular value decomposition (SVD) approach to computing PCA [37] is used, Eq. (6) decomposing the covariance of the heatmap into unitary matrices  $U^k$  and  $V^k$ , and  $S^k$ , the diagonal matrix of singular values. We call this approach *PCA-based*, in contrast to *regression-based* for Eq. (1). Despite using the same region representation and encoder here, the encoded regions differ significantly (see Fig. 3(a)), ours mapping to meaningful object parts such as the limbs of an articulated body, due to our novel foreground motion representation, described above. Note that in our *PCA-based* approach, shear is not captured, therefore our transform is not fully affine, with only five degrees of freedom instead of six. Nevertheless, as we later show empirically, it captures sufficient motion, with shear being a less significant component of the affine transform for this task. The reference frame in both is used only as an intermediate coordinate frame between the source and driving image coordinate frames. However, here (in contrast to FOMM) it is not in fact abstract, corresponding to the coordinate frame where the heatmap is whitened (i.e. has zero mean and identity covariance); see Fig. 2.  $A_{S \leftarrow D}^k$  is computed per Eq. (3).

### 3.3. Background motion estimation

Background occupies a large portion of image. Hence even small background motion between frames, e.g. due to camera motion, negatively affects the animation quality. FOMM [30] does not treat background motion separately, therefore must model it using keypoints. This has two negative consequences: (i) additional network capacity is required, since several keypoints are used to model the background instead of the foreground; (ii) overfitting to the training set, since these keypoints focus on specific

parts of the background, which may not appear in the test set. Hence, we additionally predict an affine background transformation,  $A_{S \leftarrow D}^0$ , using an encoder network assuming  $S$  and  $D$  as input and predicting six real values,  $a_1, \dots, a_6$ , such that  $A_{S \leftarrow D}^0 = [a_1, a_2, a_3, a_4, a_5, a_6]$ . Since our framework is unsupervised, the background network can include parts of the foreground into the background motion. In practice this does not happen, since it is easier for the network to use a more appropriate *PCA-based* motion representation for the foreground. It is also simpler for the network to use  $S$  and  $D$  to predict background movement, instead of encoding it in the heatmaps modelling the foreground. We verify this empirically, demonstrating that the proposed motion representations can separate background and foreground in a completely unsupervised manner (see Fig. 6 and Sec. 4.5 for comparisons).

### 3.4. Image generation

Similarly to FOMM [30], we render the target image in two stages: a pixel-wise flow generator converts coarse motions to dense optical flow, then the encoded features of the source are warped according to the flow, followed by inpainting the missing regions. The input of the dense flow predictor is a  $H \times W \times (4K + 3)$  tensor, with four channels per region, three for the source image warped according to the region's affine transformation, and one for a heatmap of the region, which is a gaussian approximation of  $M^k$ , and a further three channels for the source image warped according to the background's affine transformation. In contrast to FOMM, which uses constant variances, we estimate covariances from heatmaps. Our dense optical flow is given

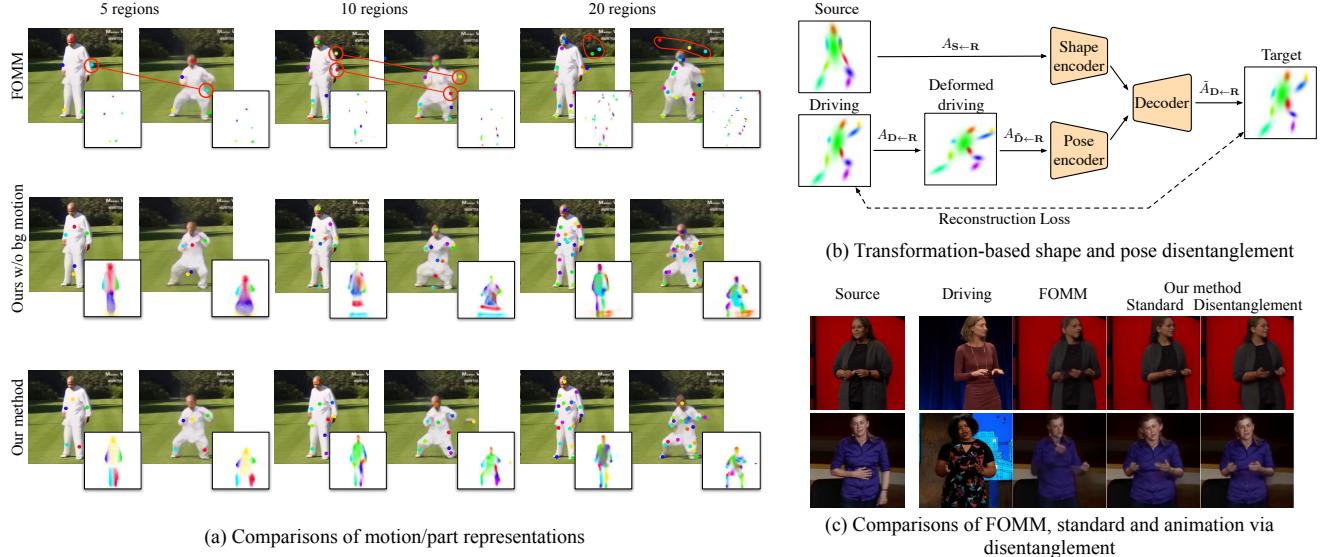


Figure 3: **Motion representation and disentanglement.** (a) A comparison of part regions estimated by regression-based (FOMM [30]) and PCA-based (our method) motion representations. Each row shows a generated sequence along with the detected keypoints and region heatmaps (inset). (b) Framework for animating motion driving frame whilst retaining shape from the source frame, by disentangling shape and pose from motion representations. (c) Qualitative animation results of articulated objects, using FOMM [30], and our method using standard and disentangled motion transfer (sec. 3.6).

	5 regions			10 regions			20 regions		
	$\mathcal{L}_1$	(AKD, MKR)	AED	$\mathcal{L}_1$	(AKD, MKR)	AED	$\mathcal{L}_1$	(AKD, MKR)	AED
FOMM [30]	0.062	(7.34, 0.036)	0.181	0.056	(6.53, 0.033)	0.172	0.062	(8.29, 0.049)	0.196
Ours w/o bg	0.061	(6.67, 0.030)	0.175	0.059	<b>(5.55, 0.026)</b>	0.165	0.057	<b>(5.47, 0.026)</b>	0.155
Ours	<b>0.049</b>	<b>(6.04, 0.029)</b>	<b>0.162</b>	<b>0.047</b>	(5.59, 0.027)	<b>0.152</b>	<b>0.046</b>	<b>(5.17, 0.026)</b>	<b>0.141</b>

Table 1: Comparing our model with FOMM [30] on TaiChiHD (256), for  $K = 5, 10$  and  $20$ . (Best result in bold.)

by (cf. Eq. (4)):

$$\mathbf{O}(z) = \sum_{k=0}^K \mathbf{W}^k(z) A_{\mathbf{S} \leftarrow \mathbf{D}}^k \begin{bmatrix} z \\ 1 \end{bmatrix}. \quad (8)$$

The predicted optical flow and confidence map are used as per FOMM [30]. However in contrast to FOMM [30], but similar to Monkey-Net [31], here deformable skip connections [33] are used between the encoder and decoder.

### 3.5. Training

The proposed model is trained end-to-end using a reconstruction loss in the feature space of the pretrained VGG-19 network [15, 40]. We use a multi-resolution reconstruction loss from previous work [30, 34]:

$$\mathcal{L}_{\text{rec}}(\hat{\mathbf{D}}, \mathbf{D}) = \sum_l \sum_i \left| V_i(\mathbf{F}_l \odot \hat{\mathbf{D}}) - V_i(\mathbf{F}_l \odot \mathbf{D}) \right|, \quad (9)$$

where  $\hat{\mathbf{D}}$  is the generated image,  $V_i$  is the  $i^{\text{th}}$ -layer of the VGG-19 pretrained network, and  $\mathbf{F}_l$  is a downsampling op-

erator. Per FOMM [30], we also use an equivariance loss,

$$\mathcal{L}_{\text{eq}} = \left| A_{\mathbf{X} \leftarrow \mathbf{R}}^k - \tilde{A} A_{\tilde{\mathbf{X}} \leftarrow \mathbf{R}}^k \right|, \quad (10)$$

where  $\tilde{\mathbf{X}}$  is image  $\mathbf{X}$  transformed by  $\tilde{A}$ , and  $\tilde{A}$  is some random geometric transformation. The final loss is the sum of terms,  $\mathcal{L} = \mathcal{L}_{\text{rec}} + \mathcal{L}_{\text{eq}}$ .

### 3.6. Animation via disentanglement

Image animation using both *standard* and *relative* methods has limitations. The *standard* method directly transfers object shape from the driving frame into the generated video, while *relative* animation is only applicable to a limited set of inputs, e.g. it requires that objects be in the same pose in the source  $\mathbf{S}$  and initial driving  $\mathbf{D}_1$  frames. To address this, we learn disentangled shape and pose encoders, as shown in Fig. 3(b). The pose encoder takes in the set of *driving* motions,  $\{A_{\mathbf{D} \leftarrow \mathbf{R}}^k\}_{k=1}^K$ , while the shape encoder takes in the set of *source* motions,  $\{A_{\mathbf{S} \leftarrow \mathbf{R}}^k\}_{k=1}^K$ . A decoder then uses the concatenated latent representations (each in  $\mathbb{R}^{64}$ ) of these two encoders, to produce a set of modified driving motions,

$\{\tilde{A}_{D \leftarrow R}^k\}_{k=1}^K$  encoding the motion of the former and the shape of the latter. These are then used to render the output.

The encoders and the decoder are implemented using fully connected layers, and trained separately from (and after) other blocks, using an  $\mathcal{L}_1$  reconstruction loss on the motion parameters. As with earlier training, source and driving frames come from the same video (i.e. the object has the same shape), therefore to ensure that shape comes from the correct branch, random horizontal and vertical scaling deformations are applied to the driving motions during training, as shown in Fig. 3(b). This forces shape information to come from the other (shape) branch. However, since the shape branch has a different pose, the pose must still come from the pose branch. Thus shape and pose are disentangled. The deformations are not applied at test time.

## 4. Evaluation

We now discuss the datasets, metrics and experiments used to evaluate the proposed method. Later we compare with prior work, as well as ablate our contributions.

### 4.1. Toy Motion Representation Experiment

To demonstrate the benefit of the proposed PCA-based motion representation, we devise an experiment on rotated rectangles (see Sup. Mat.): the task is to predict the rotation angle of a rectangle in an image. To fully isolate our contribution, we consider a supervised task, where three different architectures learn to predict angles under the  $\mathcal{L}_1$  loss. The first, a Naive architecture, directly regresses the angle using an encoder-like architecture. The second is Regression-based, as in FOMM [30]. The third uses our PCA-based approach (see Sup. Mat.). Test results are presented in Fig. 5, against training set size. The Naive baseline struggles to produce meaningful results for any size of training set, while Regression-based performance improves with more data. However, the PCA-based approach significantly improves accuracy over the Regression-based one, being over an order of magnitude better with a large number of samples. This shows that it is significantly easier for the network to infer geometric parameters of the image, such as angle, using our proposed PCA-based representation.

### 4.2. Benchmarks

We evaluate our method on several benchmark datasets for animating human faces, bodies and animated cartoons. Each dataset has separate training and test videos. The datasets are as follows:

- *VoxCeleb* [23] consists of interview videos of different celebrities. We extract square, face regions and down-scale them to  $256 \times 256$ , following FOMM [30]. The number of frames per video ranges from 64 to 1024.
- *TaiChiHD* [30] consists of cropped videos of full human bodies performing Tai Chi actions. We evaluate

on two resolutions of the dataset:  $256 \times 256$  (from FOMM [30]), and a new,  $512 \times 512$  subset, removing videos lacking sufficient resolution to support that size.

- *MGif* [31] is a dataset of .gif files, that depicts 2D cartoon animals. The dataset was collected using google searches.
- *TED-talks* is a new dataset, collected for this paper in order to demonstrate the generalization properties of our model. We cropped the upper part of the human body from the videos, downscaling to  $384 \times 384$ . The number of frames per video ranges from 64 to 1024.

Since video animation is a relatively new problem, there are not currently many effective ways of evaluating it. For quantitative metrics, prior works [42, 31, 30] use video reconstruction accuracy as a proxy for image animation quality. We adopt the same metrics here:

- $\mathcal{L}_1$  error is the mean absolute difference between reconstructed and ground-truth video pixel values.
- *Average keypoint distance* (AKD) and *missing keypoint rate* (MKR) evaluate the difference between poses of reconstructed and ground truth video. Landmarks are extracted from both videos using public, body [5] (for TaiChiHD and TED-talks) and face [3] (for VoxCeleb) detectors. AKD is then the average distance between corresponding landmarks, while MKR is the proportion of landmarks present in the ground-truth that are missing in the reconstructed video.
- *Average Euclidean distance* (AED) evaluates how well identity is preserved in reconstructed video. Public re-identification networks for bodies [13] (for TaiChiHD and TED-talks) and faces [1] extract identity from reconstructed and ground truth frame pairs, then we compute the mean  $\mathcal{L}_2$  norm of their difference across all pairs.

### 4.3. Comparison with the state of the art

We compare our method with the current state of the art for unsupervised animation, FOMM [30], on both reconstruction (the training task) and animation (the test-time task). We used an extended training schedule compared to [30], with 50% more iterations. To compare fairly with FOMM [30], we also re-trained it with the same training schedule. Also, for reference, we include comparisons with X2Face [42] and Monkey-Net [31] on video reconstruction.

**Reconstruction quality** Quantitative reconstruction results are reported in Table 2. We first show that our method reaches state-of-the-art results on a dataset with non-articulated objects such as faces. Indeed, when compared with FOMM [30] on VoxCeleb, our method shows on-par results. The situation changes, however, when articulated objects are considered, such as human bodies in



Figure 4: **Qualitative comparisons.** We show representative examples of articulated animation using our method and FOMM [30], on two datasets of articulated objects: TED-talks (left) and TaiChiHD (right). Zoom in for greater detail.

	TaiChiHD (256) (AKD, MKR)			TaiChiHD (512) (AKD, MKR)			TED-talks (AKD, MKR)			VoxCeleb			MGif $\mathcal{L}_1$
	$\mathcal{L}_1$	AED		$\mathcal{L}_1$	AED		$\mathcal{L}_1$	AED	$\mathcal{L}_1$	AKD	AED		
X2Face	0.080	(17.65, 0.109)	0.27	-	-	-	-	-	-	0.078	7.69	0.405	-
Monkey-Net	0.077	(10.80, 0.059)	0.228	-	-	-	-	-	-	0.049	1.89	0.199	-
FOMM	0.056	(6.53, 0.033)	0.172	0.075	(17.12, 0.066)	0.203	0.033	(7.07, 0.014)	0.163	0.041	<b>1.27</b>	0.134	0.0223
Ours	<b>0.047</b>	<b>(5.58, 0.027)</b>	<b>0.152</b>	<b>0.064</b>	<b>(13.86, 0.043)</b>	<b>0.172</b>	<b>0.026</b>	<b>(3.75, 0.007)</b>	<b>0.114</b>	<b>0.040</b>	1.28	<b>0.133</b>	<b>0.0206</b>

Table 2: Video reconstruction: comparison with the state of the art on five different datasets. For all methods we use  $K = 10$  regions. (Best result in bold.)

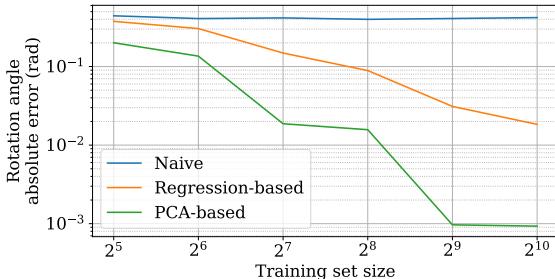


Figure 5: Mean test-time absolute rotation error, as a function of training set size.

TaiChiHD and TED-talks datasets, on which our improved motion representations boost all the metrics. The advantage over the state of the art holds at different resolutions, for TaiChiHD (256), TaiChiHD (512) and TED-talks, as well as for different numbers of selected regions (discussed later).

**Animation quality** Fig. 3(a,c) & 4 show selected and representative animations respectively, using our method and FOMM [30], on articulated bodies. Here for FOMM [30] we use the *standard* method, while ours uses animation via disentanglement. The results show clear improvements, in most cases, in animation quality, especially of limbs.

Animation quality was evaluated quantitatively through a user preference study similar to that of [30]. AMT users were presented with the source image, driving video, and the output from our method and FOMM [30], and asked which of the two videos they preferred. 50 such videos were

Dataset	Our vs FOMM (%)	Our vs Standard (%)
TaiChiHD (256)	83.7%	53.6%
TED-talks	96.6%	65.6%

Table 3: User study: second column - the proportion (%) of users that prefer our method over FOMM [30]; third column - the proportion (%) of users that prefer animation via disentanglement over *standard* animation for our model.

evaluated, by 50 users each, for a total of 2500 preferences per study. The results further support the reconstruction scores in Tab. 2. When the animated object is not articulated (VoxCeleb), the method delivers results comparable to the previous work, e.g. 52% preference in favour of our method. However, when bodies are animated (TaiChiHD & TED-talks), FOMM [30] fails to correctly detect and animate the articulated body parts such as hands. Our method renders them in the driving pose even for extreme cases, leading to a high preference in favor of it (see Tab. 3, middle column). Moreover since it is not possible to demonstrate the benefits of the animation via disentanglement using reconstruction metrics, we run an additional user study to compare our method with *standard* animation and animation via disentanglement. Since animation via disentanglement preserves shape of the object much better (see Fig. 3(c)), users prefer it more often. It especially pronoun in case of the TED-talks dataset, since shape of the objects differs more in that case (see Tab. 3, last column).

	$\mathcal{L}_1$	(AKD, MKR)	AED
No pca or bg model	0.060	(6.14, 0.033)	0.163
No pca	0.056	(9.58, 0.034)	0.206
No bg model	0.059	<b>(5.55, 0.026)</b>	0.165
Full method	<b>0.048</b>	(5.59, 0.027)	<b>0.152</b>

Table 4: Ablation study on TaiChiHD (256) dataset with  $K = 10$ . (Best result in bold.)

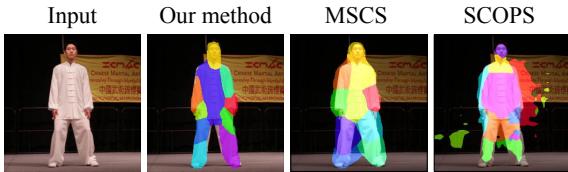


Figure 6: Qualitative co-part segmentation comparison with recent methods.

Finally, we applied animation from a TED-talks video to a photograph of George Washington, shown in Fig. 1, demonstrating animation of out-of-domain data.

#### 4.4. Ablations

In order to understand how much benefit each of our contributions bring, we ran a number of ablation experiments, detailed in Tab. 4.

**PCA-based vs. regression-based representations** First we compare the PCA-based motion model with the previous, regression-based one [30]. From the qualitative heatmap depictions in Fig. 3(a), we observe that the regression-based method localizes one edge of each corresponding part, while our method predicts regions that roughly correspond to the segmentation of the object into its constituent, articulated parts. This meaningful segmentation arises completely unsupervised.

From Tab. 1 we note that adding the PCA-based representation alone (second row) had marginal impact on the  $\mathcal{L}_1$  score (dominated by the much larger background region), but it had a much larger impact on other metrics, which are more sensitive to object-part-related errors on articulated objects. This is corroborated by Tab. 4. Moreover, we observed that when our PCA-based formulation is not used, the network encodes some of the movement into the background branch, leading to significant degradation of keypoint quality, which in turn leads to degradation of the AKD and AED scores (No-pca, Tab. 4).

We intuit that PCA-based estimation both captures regions and improves performance because it is much easier for the convolutional network to assign pixels of an object part to the corresponding heatmap than to directly regress motion parameters to an abstract reference frame. This is borne out by our toy experiment (sec. 4.1). In order to esti-

mate the heatmap, it need only learn all appearances of the corresponding object part, while regression-based networks must learn the joint space of all appearances of a part in all possible geometric configurations (e.g. rotated, scaled etc.).

One of the most important hyper-parameters of our model is the number of regions,  $K$ . The qualitative and quantitative ablations of this parameter are shown in Fig. 3(a) and Tab. 1 respectively. We can observe that, while the regression-based representation fails when the number of keypoints grows to 20, our PCA-based representation scales well with the number of regions.

**Modeling background motion** Background motion modeling significantly lowers  $\mathcal{L}_1$  error (see Tab. 4, Full method vs. No bg Model). Since background constitutes a large portion of the image, and  $\mathcal{L}_1$  treats all pixels equally, this is to be expected. AED was also impacted, suggesting that the identity representation captures some background appearance. Indeed, we observe (Fig. 3(a), second row) that having no background model causes a reduction in region segmentation quality. However, since AKD & MKR metrics evaluate object pose only, they are not improved by background modelling.

#### 4.5. Co-part Segmentation

While designed for articulated animation, our method produces meaningful object parts. To evaluate this capability of our method, we compare it against two recent unsupervised co-part segmentation works: MSCS [32] and SCOPS [14]. Following MSCS, we compute *foreground segmentation IoU* scores on TaiChiHD. Despite not being optimized for this task, our method achieves superior performance reaching 0.81 *IoU* vs. 0.77 for MSCS [32] and 0.55 for SCOPS [14]. See Fig. 6 for qualitative results.

### 5. Conclusion

We have argued that previous unsupervised animation frameworks' poor results on articulated objects are due to their representations. We propose a new, PCA-based, region motion representation, that we show both makes it easier for the network to learn region motion, and encourages it to learn semantically meaningful object parts. In addition, we propose a background motion estimation module to decouple foreground and background motion. Qualitative and quantitative results across a range of datasets and tasks demonstrate several key benefits: improved region distribution and stability, improved reconstruction accuracy and user perceived quality, and an ability to scale to more regions. We also introduce a new, more challenging dataset, TED-talks, for benchmarking future improvements on this task.

While we show some results on out of domain data (Fig. 1), generalization remains a significant challenge to making this method broadly practical in articulated animation of inanimate objects.

## References

- [1] Brandon Amos, Bartosz Ludwiczuk, and Mahadev Satyanarayanan. Openface: A general-purpose face recognition. 2016. [6](#)
- [2] Aayush Bansal, Shugao Ma, Deva Ramanan, and Yaser Sheikh. Recycle-gan: Unsupervised video retargeting. In *Proceedings of the European Conference on Computer Vision*, 2018. [2](#)
- [3] Adrian Bulat and Georgios Tzimiropoulos. How far are we from solving the 2d & 3d face alignment problem? (and a dataset of 230,000 3d facial landmarks). In *ICCV*, 2017. [6](#)
- [4] Chen Cao, Qiming Hou, and Kun Zhou. Displaced dynamic expression regression for real-time facial tracking and animation. *ACM Transactions on Graphics*, 2014. [2](#)
- [5] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *CVPR*, 2017. [6](#)
- [6] Caroline Chan, Shiry Ginosar, Tinghui Zhou, and Alexei A Efros. Everybody dance now. In *Proceedings of the IEEE International Conference on Computer Vision*, 2019. [1, 2](#)
- [7] A Clark, J Donahue, and K Simonyan. Adversarial video generation on complex datasets. *arXiv preprint arXiv:1907.06571*, 2019. [2](#)
- [8] J. W. Davis and A. F. Bobick. The representation and recognition of human movement using temporal templates. In *CVPR*, 1997. [2](#)
- [9] Yu Deng, Jiaolong Yang, Dong Chen, Fang Wen, and Xin Tong. Disentangled and controllable face image generation via 3d imitative-contrastive learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020. [2](#)
- [10] Oran Gafni, Lior Wolf, and Yaniv Taigman. Vid2game: Controllable characters extracted from real-world videos. *arXiv preprint arXiv:1904.08379*, 2019. [1](#)
- [11] Zhenglin Geng, Chen Cao, and Sergey Tulyakov. 3d guided fine-grained face manipulation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. [1, 2](#)
- [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*, 2020. [2](#)
- [13] Alexander Hermans, Lucas Beyer, and Bastian Leibe. In defense of the triplet loss for person re-identification. *arXiv:1703.07737*, 2017. [6](#)
- [14] Wei-Chih Hung, Varun Jampani, Sifei Liu, Pavlo Molchanov, Ming-Hsuan Yang, and Jan Kautz. Scops: Self-supervised co-part segmentation. In *CVPR*, 2019. [8, 11, 12](#)
- [15] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *Proceedings of the European Conference on Computer Vision*, 2016. [5, 13](#)
- [16] Hyeongwoo 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*, 2018. [2](#)
- [17] Yunji Kim, Seonghyeon Nam, In Cho, and Seon Joo Kim. Unsupervised keypoint learning for guiding class-conditional video prediction. In *Proceedings of the Neural Information Processing Systems Conference*, 2019. [1](#)
- [18] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2014. [13](#)
- [19] Wen Liu, Zhixin Piao, Jie Min, Wenhan Luo, Lin Ma, and Shenghua Gao. Liquid warping gan: A unified framework for human motion imitation, appearance transfer and novel view synthesis. In *Proceedings of the IEEE International Conference on Computer Vision*, 2019. [2](#)
- [20] Dominik Lorenz, Leonard Bereska, Timo Milbich, and Bjorn Ommer. Unsupervised part-based disentangling of object shape and appearance. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. [1](#)
- [21] Willi Menapace, Stéphane Lathuilière, Sergey Tulyakov, Aliaksandr Siarohin, and Elisa Ricci. Playable video generation. In *CVPR*, 2021. [2](#)
- [22] Koki Nagano, Jaewoo Seo, Jun Xing, Lingyu Wei, Zimo Li, Shunsuke Saito, Aviral Agarwal, Jens Furund, and Hao Li. pagan: real-time avatars using dynamic textures. *ACM Transactions on Graphics*, 2018. [2](#)
- [23] A. Nagrani, J. S. Chung, and A. Zisserman. Voxceleb: a large-scale speaker identification dataset. In *INTERSPEECH*, 2017. [6](#)
- [24] Yuval Nirkin, Yosi Keller, and Tal Hassner. Fsgan: Subject agnostic face swapping and reenactment. In *Proceedings of the IEEE International Conference on Computer Vision*, 2019. [2](#)
- [25] Albert Pumarola, Antonio Agudo, Aleix M Martinez, Alberto Sanfeliu, and Francesc Moreno-Noguer. Ganimation: Anatomically-aware facial animation from a single image. In *Proceedings of the European Conference on Computer Vision*, 2018. [2](#)

- [26] Shengju Qian, Kwan-Yee Lin, Wayne Wu, Yangxi-aokang Liu, Quan Wang, Fumin Shen, Chen Qian, and Ran He. Make a face: Towards arbitrary high fidelity face manipulation. In *Proceedings of the IEEE International Conference on Computer Vision*, 2019. 2
- [27] Jian Ren, Menglei Chai, Sergey Tulyakov, Chen Fang, Xiaohui Shen, and Jianchao Yang. Human motion transfer from poses in the wild. *arXiv preprint arXiv:2004.03142*, 2020. 1, 2
- [28] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015. 13
- [29] Masaki Saito, Eiichi Matsumoto, and Shunta Saito. Temporal generative adversarial nets with singular value clipping. In *Proceedings of the IEEE International Conference on Computer Vision*, 2017. 2
- [30] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. First order motion model for image animation. In *Proceedings of the Neural Information Processing Systems Conference*, 2019. 1, 2, 3, 4, 5, 6, 7, 8, 11, 12, 13
- [31] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. Animating arbitrary objects via deep motion transfer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 1, 2, 5, 6
- [32] Aliaksandr Siarohin, Subhankar Roy, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. Motion-supervised co-part segmentation. In *Proceedings of the International Conference on Pattern Recognition*, 2020. 8, 11, 12
- [33] 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*, 2018. 2, 5
- [34] Hao Tang, Dan Xu, Wei Wang, Yan Yan, and Nicu Sebe. Dual generator generative adversarial networks for multi-domain image-to-image translation. In *ACCV*, 2018. 5
- [35] Justus Thies, Michael Zollhofer, Marc Stamminger, Christian Theobalt, and Matthias Nießner. Face2face: Real-time face capture and reenactment of rgb videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016. 2
- [36] Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, and Jan Kautz. Mocogan: Decomposing motion and content for video generation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 2
- [37] Michael E Wall, Andreas Rechtsteiner, and Luis M Rocha. Singular value decomposition and principal component analysis. In *A practical approach to microarray data analysis*, pages 91–109. Springer, 2003. 4
- [38] Ting-Chun Wang, Ming-Yu Liu, Andrew Tao, Guilin Liu, Jan Kautz, and Bryan Catanzaro. Few-shot video-to-video synthesis. In *Proceedings of the Neural Information Processing Systems Conference*, 2019. 2
- [39] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Guilin Liu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. Video-to-video synthesis. In *Proceedings of the Neural Information Processing Systems Conference*, 2018. 1, 2
- [40] 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*, 2017. 5
- [41] Wei Wang, Xavier Alameda-Pineda, Dan Xu, Pascal Fua, Elisa Ricci, and Nicu Sebe. Every smile is unique: Landmark-guided diverse smile generation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 2
- [42] Olivia Wiles, A Sophia 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*, 2018. 1, 2, 6
- [43] Yuxin Wu, Alexander Kirillov, Francisco Massa, Wan-Yen Lo, and Ross Girshick. Detectron2. <https://github.com/facebookresearch/detectron2>, 2019. 13
- [44] Egor Zakharov, Aliaksandra Shysheya, Egor Burkov, and Victor Lempitsky. Few-shot adversarial learning of realistic neural talking head models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 2

In this supplementary material we report additional details on the toy experiment in Sec. A. In Sec. B we provide additional details for the co-part segmentation experiment. We provide additional implementation details in Sec. C. Additionally in Sec. D we visually demonstrate the ability of the model to control the background. Finally in Sec. E we describe the TED-talks data collection procedure.

## A. Toy Experiment Details

The rotated rectangles dataset consists of images of rectangles randomly rotated from  $0^\circ$  to  $90^\circ$ , along with labels that indicate the angle of rotation. The rectangles have different, random colors. Visual samples are shown in Fig. 7.

We tested three different networks: Naive, Regression-based and PCA-based. The Naive network directly predicts an angle from an image using an encoder and a fully-connected layer. Regression-based is similar to FOMM [30]; the angle is regressed per pixel an using hourglass network, and pooled according to heatmap weights predicted using the same hourglass network. PCA-based is our method described in Sec. 3.2 we predict the heatmap using an hourglass network, PCA is performed according to Eq. (6), and the angle is computed from matrix  $U$  as  $\arctan(U_{10}/U_{00})$ .

Each of the networks was trained, on subsets of the dataset of varying sizes, to minimize the  $\mathcal{L}_1$  loss between predicted and ground truth rotation angle. All models were trained for 100 epochs, with batch size 8. We used the Adam optimizer, with a learning rate of  $10^{-4}$ . We varied the size of the training set from 32 to 1024. Results, on a separate, fixed test set of size 128, were then computed, shown in Fig. 5.

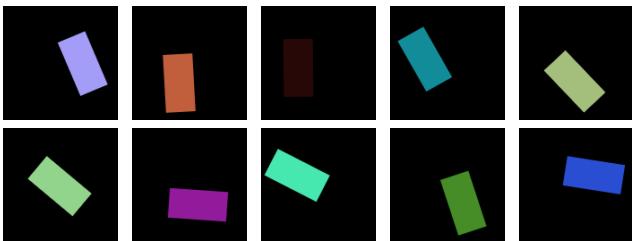


Figure 7: Examples of synthetic rectangle dataset.

## B. Co-part segmentation details

To perform co-part segmentation we use  $M^k$ . A pixel  $z$  is assigned to a part that has the maximum response of  $M^k$  in that pixel, e.g.  $\text{argmax}_k M^k(z)$ . Moreover since region predictor did not explicitly predict background region, we assign pixel  $z$  to the background iff  $\sum_k M^k(z) < 0.001$ . We demonstrate additional qualitative comparisons with MSCS [32] and SCOPS [14] in Fig. 9. Fig. 9 shows that our method produces more meaningful co-part segmentations compared to SCOPS [14]. Furthermore, our method

separates the foreground object from the background more accurately than MSCS [32].

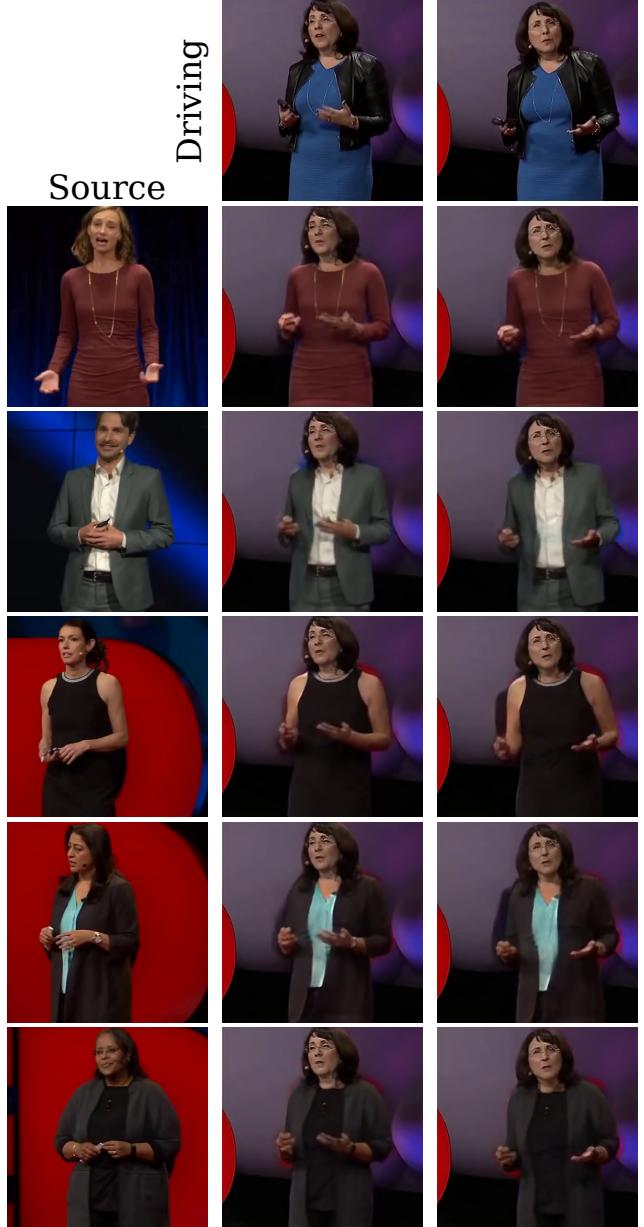


Figure 8: Examples of cloth swap performed using our model. First column depicts sources from which cloth is taken, while the first row shows a driving video to which we put the cloth. Rest demonstrates images generated with our model.

Similarly to MSCS [32] we can exploit produced segmentations in order to perform a part swap. In Fig. 8 we copy the cloth from the person in the source image on to the person in the driving video.

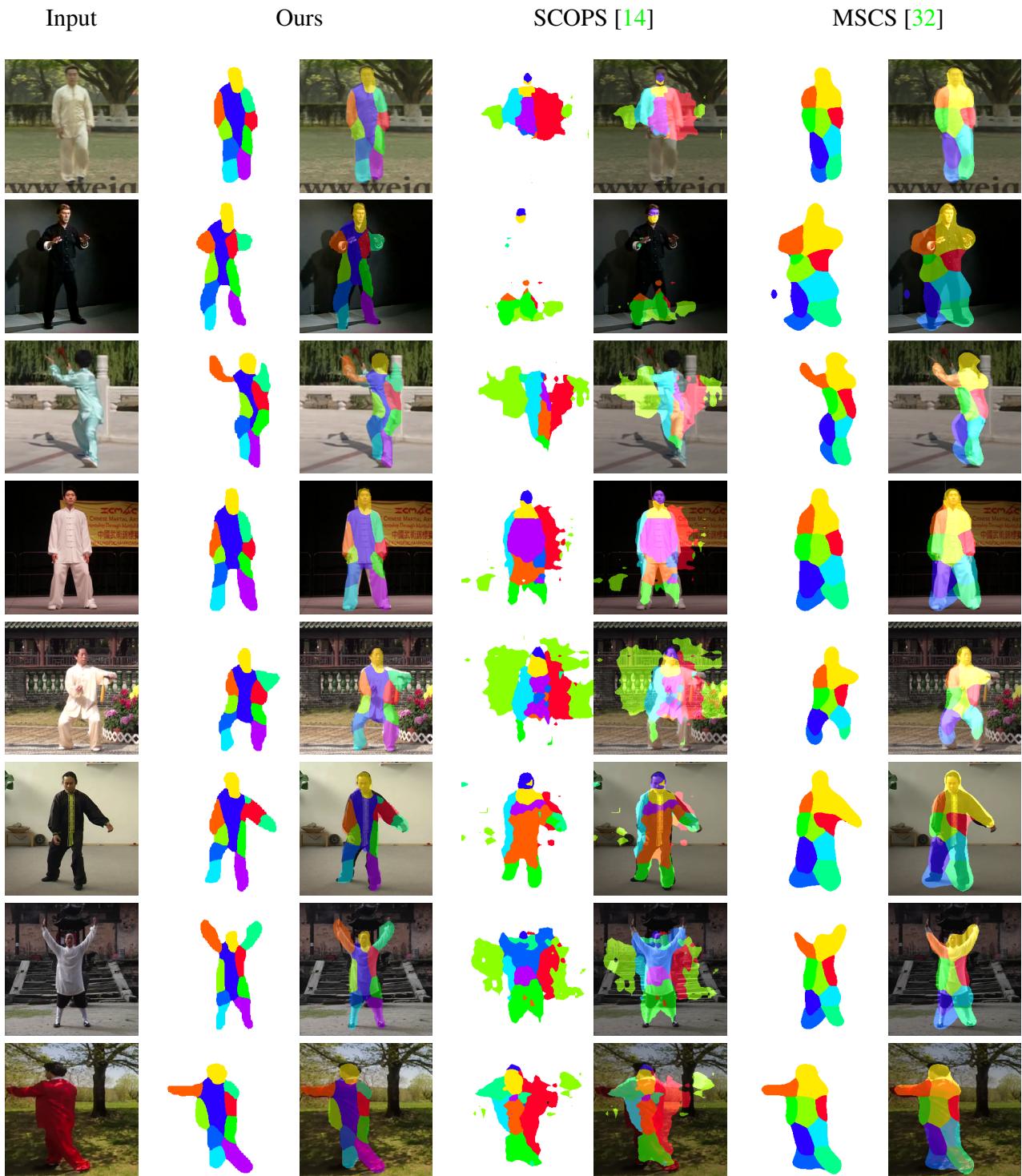


Figure 9: Additional qualitative co-part segmentation comparisons with recent methods. First column is an input. In next columns, for every method segmentation mask and image with overlayed segmentation are shown.

## C. Implementation details

For a fair comparison, in order to highlight our contributions, we mostly follow the architecture design of

FOMM [30]. Similar to FOMM, our region predictor, background motion predictor and pixel-wise flow predictor operate on a quarter of the original resolution, e.g.  $64 \times 64$  for

$256 \times 256$  images,  $96 \times 96$  for  $384 \times 384$  and  $128 \times 128$  for  $512 \times 512$ . We use the U-Net [28] architecture with five "convolution - batch norm - ReLU - pooling" blocks in the encoder and five "upsample - convolution - batch norm - ReLU" blocks in the decoder for both the region predictor and the pixel-wise flow predictor. For the background motion predictor, we use only five block encoder part. Similarly to FOMM [30], we use the Johnson architecture [15] for image generation, with two down-sampling blocks, six residual-blocks, and two up-sampling blocks. However, we add skip connections that are warped and weighted by the confidence map. Our method is trained using Adam [18] optimizer with learning rate  $2e - 4$  and batch size 48, 20, 12 for  $256 \times 256$ ,  $384 \times 384$  and  $512 \times 512$  resolutions respectively. During the training process, the networks observe 3M source-driving pairs, each pair selected at random from a random video chunk, and we drop the learning rate by a factor of 10 after 1.8M and 2.7M pairs. We use 4 Nvidia P100 GPUs for training.

The shape-pose disentanglement network consists of 2 identical encoders and 1 decoder. Each encoder consists of 3 "linear - batch norm - ReLU" blocks, with a number of hidden units equal to 256, 512, 1024, and another linear layer with number of units equal to 64. Decoder takes a concatenated input from encoders and applies 3 "linear - batch norm - ReLU" blocks, with sizes 1024, 512, 256. The network is trained on 1M source-driving pairs, organized in batches of 256 images. We use Adam optimizer with learning rate  $1e - 3$  and drop the learning rate at 660K and 880K pairs.

## D. Background movement

While the primary purpose of background modelling is to free up the capacity of the network to better handle the object. For animating articulated objects, background motion is usually unnecessary. Thus, though we estimate background motion, we set it to zero during animation. However, nothing in our framework prevents us from controlling camera motion. Below we show a still background, then move it left, right, and rotate counterclockwise.

## E. TED-talks dataset creation

In order to create the TED-talks dataset, we downloaded 3,035 YouTube videos, shared under the "CC BY – NC – ND 4.0 International" license<sup>4</sup> using the "TED talks" query. From these initial candidates, we selected the videos in which the upper part of the person is visible for at least 64 frames, and the height of the person bounding box was at least 384 pixels. After that, we manually filtered out static videos and videos in which a person is doing something other than presenting. We ended up with 411 videos, and split these videos in 369



Figure 10: Visualizations of background movement. From top to bottom we show driving frame, still background, background that moves left, moves right and rotates counterclockwise.

training and 42 testing videos. We then split each video into chunks without significant camera changes (i.e. with no cuts to another camera), and for which the presenter did not move too far from their starting position in the chunk. We cropped the a square region around the presenter, such that they had a consistent scale, and downscaled this region to  $384 \times 384$  pixels. Chunks that lacked sufficient resolution to be downscaled, or had a length shorter than 64 frames, were removed. Both the distance moved and the region cropping were achieved using a bounding box estimator for humans [43]. We obtained 1,177 training video chunks and 145 test videos chunks.

<sup>4</sup>This license allows for non-commercial use.