

# Animating Arbitrary Objects via Deep Motion Transfer

Aliaksandr Siarohin<sup>1</sup>, Stéphane Lathuilière<sup>1</sup>, Sergey Tulyakov<sup>2</sup>, Elisa Ricci<sup>1,3</sup> and Nicu Sebe<sup>1</sup>

<sup>1</sup>DISI, University of Trento, Italy, <sup>2</sup> Snap Inc., <sup>3</sup> Fondazione Bruno Kessler (FBK), Trento, Italy

{aliaksandr.siarohin, stephane.lathuiliere, e.ricci, niculae.sebe}@unitn.it, stulyakov@snap.com

## Abstract

This paper introduces a novel deep learning framework for image animation. Given an input image with a target object and a driving video sequence depicting a moving object, our framework generates a video in which the target object is animated according to the driving sequence. This is achieved through a deep architecture that decouples appearance and motion information. Our framework consists of three main modules: (i) a Keypoint Detector unsupervisedly trained to extract object keypoints, (ii) a Dense Motion prediction network for generating dense heatmaps from sparse keypoints, in order to better encode motion information and (iii) a Motion Transfer Network, which uses the motion heatmaps and appearance information extracted from the input image to synthesize the output frames. We demonstrate the effectiveness of our method on several benchmark datasets, spanning a wide variety of object appearances, and show that our approach outperforms state-of-the-art image animation and video generation methods. Our source code is publicly available <sup>1</sup>.

## 1. Introduction

This paper introduces a framework for motion-driven image animation to automatically generate videos by combining the appearance information derived from a *source image* (e.g. depicting the face or the body silhouette of a certain person) with motion patterns extracted from a *driving video* (e.g. encoding the facial expressions or the body movements of another person). Several examples are given in Fig. 1. Generating high-quality videos from static images is challenging, as it requires learning an appropriate representation of an object, such as a 3D model of a face or a human body. This task also requires accurately extracting the motion patterns from the driving video and mapping them on the object representation. Most approaches are object-specific, using techniques from computer graphics [7, 34]. These methods also use an explicit object representation, such as a 3D morphable model [5], to facilitate animation, and therefore only consider faces.

<sup>1</sup> <https://github.com/AliaksandrSiarohin/monkey-net>

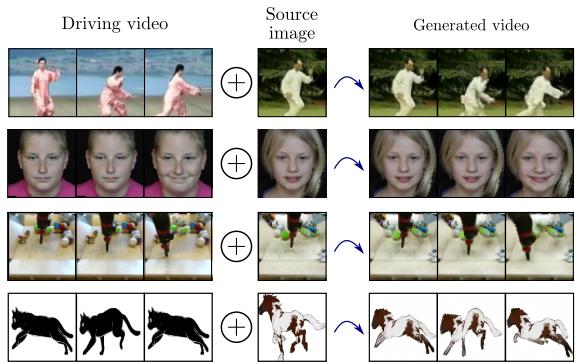


Figure 1: Our deep motion transfer approach can animate arbitrary objects following the motion of the driving video.

Over the past few years, researchers have developed approaches for automatic synthesis and enhancement of visual data. Several methods derived from Generative Adversarial Networks (GAN) [16] and Variational Autoencoders (VAE) [27] have been proposed to generate images and videos of different objects [38, 31, 35]. These approaches use additional information such as conditioning labels (e.g. indicating a facial expression, a body pose) [41, 32]. More specifically, they are purely data-driven, leveraging a large collection of training data to learn a latent representation of the visual inputs for synthesis. Noting the significant progress of these techniques, recent research studies have started exploring the use of deep generative models for image animation and video retargeting [42, 9, 4, 39, 3]. These works demonstrate that deep models can effectively transfer motion patterns between human subjects in videos [4], or transfer a facial expression from one person to another [42]. However, these approaches have limitations: for example, they rely on pre-trained models for extracting object representations that require costly ground-truth data annotations [9, 39, 3]. Furthermore, these works do not address the problem of animating arbitrary objects: instead, considering a single object category [42] or learning to translate videos from one specific domain to another [4, 24].

This paper addresses some of these limitations by introducing a novel deep learning framework for animating a static image using a driving video. Inspired by [42], we

propose learning a latent representation of the object in a self-supervised way, leveraging a large collection of video sequences. There are two key distinctions between our work and [42]. Firstly, our approach is not designed for specific object category, but rather is effective in animating arbitrary objects. Secondly, we introduce a novel strategy to model and transfer motion information, using a set of sparse motion-specific keypoints that were learned in an unsupervised way to describe relative pixel movements. Our intuition is that only relevant motion patterns (derived from the driving video) must be transferred for object animation, while other information should not be used. We call the proposed deep framework Monkey-Net, as it enables motion transfer by considering MOViNg KEYpoints.

We demonstrate the effectiveness of our framework by conducting an extensive experimental evaluation on three publicly available datasets, previously used for video generation: the Tai-Chi [35], the BAIR robot pushing [12] and the UvA-NEMO Smile [11] datasets. As shown in our experiments, our image animation method produces high quality videos for a wide range of objects. Furthermore, our quantitative results clearly show that our approach outperforms state-of-the-art methods for image-to-video translation tasks.

## 2. Related work

**Deep Video Generation.** Early deep learning-based approaches for video generation proposed synthesizing videos by using spatio-temporal networks. Vondrick *et al.* [38] introduced VGAN, a 3D convolutional GAN which simultaneously generates all the frames of the target video. Similarly, Saito *et al.* [31] proposed TGAN, a GAN-based model which is able to generate multiple frames at the same time. However, the visual quality of these methods outputs is typically poor.

More recent video generation approaches used recurrent neural networks within an adversarial training framework. For instance, Wang *et al.* [41] introduced a Conditional MultiMode Network (CMM-Net), a deep architecture which adopts a conditional Long-Short Term Memory (LSTM) network and a VAE to generate face videos. Tulyakov *et al.* [35] proposed MoCoGAN, a deep architecture based on a recurrent neural network trained with an adversarial learning scheme. These approaches can take conditional information as input that comprises categorical labels or static images and, as a result, produces high quality video frames of desired actions.

Video generation is closely related to the future frame prediction problem addressed in [33, 29, 14, 36, 44]. Given a video sequence, these methods aim to synthesize a sequence of images which represents a coherent continuation

of the given video. Earlier methods [33, 29, 25] attempted to directly predict the raw pixel values in future frames. Other approaches [14, 36, 2] proposed learning the transformations which map the pixels in the given frames to the future frames. Recently, Villegas *et al.* [37] introduced a hierarchical video prediction model consisting of two stages: it first predicts the motion of a set of landmarks using an LSTM, then generates images from the landmarks.

Our approach is closely related to these previous works since we also aim to generate video sequences by using a deep learning architecture. However, we tackle a more challenging task: image animation requires decoupling and modeling motion and content information, as well as a re-combining them.

**Object Animation.** Over the years, the problems of image animation and video re-targeting have attracted attention from many researchers in the fields of computer vision, computer graphics and multimedia. Traditional approaches [7, 34] are designed for specific domains, as they operate only on faces, human silhouettes, *etc.* In this case, an explicit representation of the object of interest is required to generate an animated face corresponding to a certain person’s appearance, but with the facial expressions of another. For instance, 3D morphable models [5] have been traditionally used for face animation [45]. While especially accurate, these methods are highly domain-specific and their performance drastically degrades in challenging situations, such as in the presence of occlusions.

Image animation from a driving video can be interpreted as the problem of transferring motion information from one domain to another. Bansal *et al.* [4] proposed RecycleGAN, an approach which extends conditional GAN by incorporating spatio-temporal cues in order to generate a video in one domain given a video in another domain. However, their approach only learns the association between two specific domains, while we want to animate an image depicting one object without knowing at training time which object will be used in the driving video. Similarly, Chan *et al.* [9] addressed the problem of motion transfer, casting it within a per-frame image-to-image translation framework. They also proposed incorporating spatio-temporal constraints. The importance of considering temporal dynamics for video synthesis was also demonstrated in [39]. Wiles *et al.* [42] introduced X2Face, a deep architecture which, given an input image of a face, modifies it according to the motion patterns derived from another face or another modality, such as audio. They demonstrated that a purely data-driven deep learning-based approach is effective in animating still images of faces without demanding explicit 3D representation. In this work, we design a self-supervised deep network for animating static images, which is effective for generating arbitrary objects.

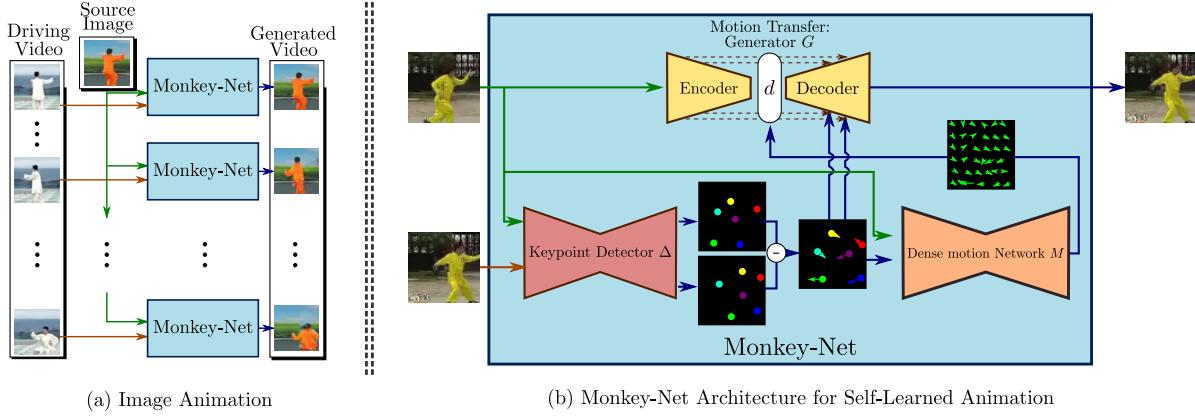


Figure 2: A schematic representation of the proposed motion transfer framework for image animation. At testing time (Fig. (a)), the model generates a video with the object appearance of the source image but with the motion from the driving video. Monkey-Net (Fig. (b)) is composed of three networks: a motion-specific keypoint detector  $\Delta$ , a motion prediction network  $M$  and an image generator  $G$ .  $G$  reconstructs the image  $x'$  from the keypoint positions  $\Delta(x)$  and  $\Delta(x')$ . The optical flow computed by  $M$  is used by  $G$  to handle misalignments between  $x$  and  $x'$ . The model is learned with a self-supervised learning scheme.

### 3. Monkey-Net

The architecture of the Monkey-Net is given in Fig. 2. We now describe it in detail.

#### 3.1. Overview and Motivation

The objective of this work is to animate an object based on the motion of a similar object in a driving video. Our framework is articulated into three main modules (Fig. 2). The first network, named Keypoint Detector, takes as input the source image and a frame from the driving video and automatically extracts sparse keypoints. The output of this module is then fed to a Dense Motion prediction network, which translates the sparse keypoints into motion heatmaps. The third module, the Motion Transfer network, receives as input the source image and the dense motion heatmap and recombines them producing a target frame.

The output video is generated frame-by-frame as illustrated in Fig. 2.a. At time  $t$ , the Monkey-Net uses the source image and the  $t^{th}$  frame from the driving video. In order to train a Monkey-Net one just needs a dataset consisting of videos of objects of interest. No specific labels, such as keypoint annotations, are required. The learning process is fully self-supervised. Therefore, at test time, in order to generate a video sequence, the generator requires only a static input image and a motion descriptor from the driving sequence. Inspired by recent studies on unsupervised landmark discovery for learning image representations [22, 43], we formulate the problem of learning a motion representation as an unsupervised motion-specific keypoint detection task. Indeed, the keypoints locations differences between two frames can be seen as a compact motion representation.

In this way, our model generates a video by modifying the input image according to the landmarks extracted from the driving frames. Using a Monkey-Net at inference time is detailed in Sec. 3.6.

The Monkey-Net architecture is illustrated in Fig. 2.b. Let  $x$  and  $x' \in \mathcal{X}$  be two frames of size  $H \times W$  extracted from the same video. The  $H \times W$  lattice is denoted by  $\mathcal{U}$ . Inspired by [22], we jointly learn a keypoint detector  $\Delta$  together with a generator network  $G$  according to the following objective:  $G$  should be able to reconstruct  $x'$  from the keypoint locations  $\Delta(x) \in \mathcal{U}$ ,  $\Delta(x') \in \mathcal{U}$ , and  $x$ . In this formulation, the motion between  $x$  and  $x'$  is implicitly modeled. To deal with large motions, we aim to learn keypoints that describe motion as well as the object geometry. To this end, we add a third network  $M$  that estimates the optical flow  $F \in \mathbb{R}^{H \times W \times 2}$  between  $x'$  and  $x$  from  $\Delta(x)$ ,  $\Delta(x')$  and  $x$ . The motivation for this is twofold. First, this forces the keypoint detector  $\Delta$  to predict keypoint locations that capture not only the object structure but also its motion. To do so, the learned keypoints must be located especially on the object parts with high probability of motion. For instance, considering the human body, it is important to obtain keypoints on the extremities (as in feet or hands) in order to describe the body movements correctly, since these body-parts tend to move the most. Second, following common practises in conditional image generation, the generator  $G$  is implemented as an encoder-decoder composed of convolutional blocks [21]. However, standard convolutional encoder-decoders are not designed to handle large pixel-to-pixel misalignment between the input and output images [32, 3, 15]. To this aim, we introduce a deformation module within the generator  $G$  that employs the estimated optical flow  $F$  in order to align the encoder features with  $x'$ .

### 3.2. Unsupervised Keypoint Detection

In this section, we detail the structure employed for unsupervised keypoint detection. First, we employ a standard U-Net architecture that, from the input image, estimates  $K$  heatmaps  $H_k \in [0, 1]^{H \times W}$ , one for each keypoint. We employ softmax activations for the last layer of the decoder in order to obtain heatmaps that can be interpreted as detection confidence map for each keypoint. An encoder-decoder architecture is used here since it has shown good performance for keypoints localization [6].

To model the keypoint location confidence, we fit a Gaussian on each detection confidence map. Modeling the landmark location by a Gaussian instead of using directly the complete heatmap  $H_k$  acts as a bottle-neck layer, and therefore allows the model to learn landmarks in an indirect way. The expected keypoint coordinates  $\mathbf{h}_k \in \mathbb{R}^2$  and its covariance  $\Sigma_k$  are estimated according to:

$$\mathbf{h}_k = \sum_{p \in \mathcal{U}} H_k[p]p; \quad \Sigma_k = \sum_{p \in \mathcal{U}} H_k[p](p - \mathbf{h}_k)(p - \mathbf{h}_k)^\top \quad (1)$$

The intuition behind the use of keypoint covariances is that they can capture not only the location of a keypoint but also its orientation. Again considering the example of the human body: in the case of the legs, the covariance may capture their orientation. Finally, we encode the keypoint distributions as heatmaps  $H_k^i \in [0, 1]^{H \times W}$ , such that they can be used as inputs to the generator and to the motion networks. Indeed, the advantage of using a heatmap representation, rather than considering directly the 2D coordinates  $\mathbf{h}_k$ , is that heatmaps are compatible with the use of convolutional neural networks. Formally, we employ the following Gaussian-like function:

$$\forall p \in \mathcal{U}, H_k(\mathbf{p}) = \frac{1}{\alpha} \exp(-(\mathbf{p} - \mathbf{h}_k)\Sigma_k^{-1}(\mathbf{p} - \mathbf{h}_k)) \quad (2)$$

where  $\alpha$  is normalization constant. This process is applied independently on  $\mathbf{x}$  and  $\mathbf{x}'$  leading to two sets of  $K$  keypoints heatmaps  $H = \{H_k\}_{k=1..K}$  and  $H' = \{H'_k\}_{k=1..K}$ .

### 3.3. Generator Network with Deformation Module

In this section, we detail how we reconstruct the target frame  $\mathbf{x}'$  from  $\mathbf{x}$ ,  $\Delta(\mathbf{x}) = H$  and  $\Delta(\mathbf{x}') = H'$ . First we employ a standard convolutional encoder composed of a sequence of convolutions and average pooling layers in order to encode the object appearance in  $\mathbf{x}$ . Let  $\xi_r \in \mathbb{R}^{H_r \times W_r \times C_r}$  denote the output of the  $r^{th}$  block of the encoder network ( $1 \leq r \leq R$ ). The architecture of this generator network is also based on the U-Net architecture [30] in order to obtain better details in the generated image. Motivated by [32], where it was shown that a standard

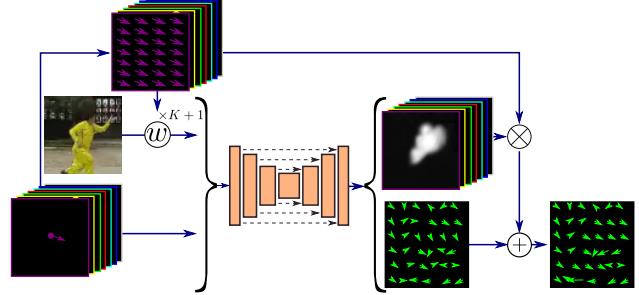


Figure 3: A schematic representation of the adopted part-based model for optical flow estimation from sparse representation. From the appearance of the first frame and the keypoints motion, the network  $M$  predicts a mask for each keypoint and the residual motion (see text for details).

U-Net cannot handle large pixel-to-pixel misalignment between the input and the output images, we propose using a deformation module to align the features of the encoder with the output images. Contrary to [32] that defines an affine transformation for each human body part in order to compute the feature deformation, we propose a deformation module that can be used on any object. In particular, we propose employing the optical flow  $\mathcal{F}$  to align the features  $\xi_r$  with  $\mathbf{x}'$ . The deformation employs a warping function  $f_w(\cdot, \cdot)$  that warps the feature maps according to  $\mathcal{F}$ :

$$\xi'_r = f_w(\xi_r, \mathcal{F}) \quad (3)$$

This warping operation is implemented using a bilinear sampler, resulting in a fully differentiable model. Note that  $\mathcal{F}$  is down-sampled to  $H_r \times W_r$  via nearest neighbour interpolation when computing Eq. (3). Nevertheless, because of the small receptive field of the bilinear sampling layer, encoding the motion only via the deformation module leads to optimization problems. In order to facilitate network training, we propose inputting the decoder the difference of the keypoint locations encoded as heatmaps  $\dot{H} = H' - H$ . Indeed, by providing  $\dot{H}$  to the decoder, the reconstruction loss applied on the  $G$  outputs (see Sec. 3.5) is directly propagated to the keypoint detector  $\Delta$  without going through  $M$ . In addition, the advantage of the heatmap difference representation is that it encodes both the locations and the motions of the keypoints. Similarly to  $\mathcal{F}$ , we compute  $R$  tensors  $\dot{H}_r$  by down-sampling  $\dot{H}$  to  $H_r \times W_r$ . The two tensors  $\dot{H}_r$  and  $\xi'_r$  are concatenated along the channel axis and are then treated as skip-connection tensors by the decoder.

### 3.4. From Sparse Keypoints to Dense Optical Flow

In this section, we detail how we estimate the optical flow  $\mathcal{F}$ . The task of predicting a dense optical flow only from the displacement of a few keypoints and the appearance of the first frame is challenging. In order to facilitate

the task of the network, we adopt a part base formulation. We make the assumption that each keypoint is located on an object part that is locally rigid. Thus, the task of computing the optical flow becomes simpler since, now, the problem consists in estimating masks  $M_k \in \mathcal{R}^{H \times W}$  that segment the object in rigid parts corresponding to each keypoint. A first coarse estimation of the optical flow can be given by:

$$\mathcal{F}_{\text{coarse}} = \sum_{k=1}^{K+1} M_k \otimes \rho(h_k) \quad (4)$$

where  $\otimes$  denotes the element-wise product and  $\rho(\cdot) \in \mathcal{R}^{H \times W \times 2}$  is the operator that returns a tensor by repeating the input vector  $H \times W$  times. Additionally, we employ one specific mask  $M_{K+1}$  without deformation (which corresponds to  $\rho([0, 0])$ ) to capture the static background. In addition to the masks  $M_k$ , the motion network  $M$  also predicts the residual motion  $\mathcal{F}_{\text{residual}}$ . The purpose of this residual motion field is to refine the coarse estimation by predicting non-rigid motion that cannot be modeled by the part-based approach. The final estimated optical flow is:  $\mathcal{F} = \mathcal{F}_{\text{coarse}} + \mathcal{F}_{\text{residual}}$ .

Concerning the inputs of the motion network,  $M$  takes two tensors,  $\dot{H}$  and  $\mathbf{x}$  corresponding respectively to the sparse motion and the appearance. However, we can observe that, similarly to the generator network,  $M$  may suffer from the misalignment between the input  $\mathbf{x}$  and the output  $\mathcal{F}$ . Indeed,  $\mathcal{F}$  is aligned with  $\mathbf{x}'$ . To handle this problem, we use the warping operator  $f_w$  according to the motion field of each keypoint  $\rho(h_k)$ , e.g.  $\mathbf{x}_k = f_w(\mathbf{x}, \rho(h_k))$ . This solution provides images  $\mathbf{x}_k$  that are locally aligned with  $\mathcal{F}$  in the neighborhood of  $h'_k$ . Finally, we concatenate  $H' - H$ ,  $\{\mathbf{x}_k\}_{k=1..K}$  and  $\mathbf{x}$  along the channel axis and feed them into a standard U-Net network. Similarly to the keypoint and the generator network, the use of U-Net architecture is motivated by the need of fine-grained details.

### 3.5. Network Training

We propose training the whole network in an end-to-end fashion. As formulated in Sec. 3.1, our loss ensures that  $\mathbf{x}'$  is correctly reconstructed from  $\Delta(\mathbf{x}) \in \mathcal{U}$ ,  $\Delta(\mathbf{x}') \in \mathcal{U}$  and  $\mathbf{x}$ . Following the recent advances in image generation, we combine an adversarial and the feature matching loss proposed in [40] in order to learn to reconstruct  $\mathbf{x}'$ . More precisely, we use a discriminator network  $D$  that takes as input  $H'$  concatenated with either the real image  $\mathbf{x}'$  or the generated image  $\hat{\mathbf{x}}'$ . We employ the least-square GAN formulation [28] leading to the two following losses used to train the discriminator and the generator:

$$\begin{aligned} \mathcal{L}_{\text{gan}}^D(D) = & \mathbb{E}_{\mathbf{x}' \in \mathcal{X}} [(D(\mathbf{x}' \oplus H') - 1)^2] \\ & + \mathbb{E}_{(\mathbf{x}, \hat{\mathbf{x}}') \in \mathcal{X}^2} [D(\hat{\mathbf{x}}' \oplus H')]^2 \end{aligned}$$

$$\mathcal{L}_{\text{gan}}^G(G) = \mathbb{E}_{(\mathbf{x}, \hat{\mathbf{x}}') \in \mathcal{X}^2} [(D(\hat{\mathbf{x}}' \oplus H') - 1)^2] \quad (5)$$

where  $\oplus$  denotes the concatenation along the channel axis. Note that in Eq (5), the dependence on the trained parameters of  $G$ ,  $M$ , and  $\Delta$  appears implicitly via  $\hat{\mathbf{x}}'$ . Note that we provide the keypoint locations  $H'$  to the discriminator to help it to focus on moving parts and not on the background. However, when updating the generator, we do not propagate the discriminator loss gradient through  $H'$  to avoid that the generator tends to fool the discriminator by generating meaningless keypoints.

The GAN loss is combined with a feature matching loss that encourages the output image  $\hat{\mathbf{x}}'$  and  $\mathbf{x}'$  to have similar feature representations. The feature representations employed to compute this loss are the intermediate layers of the discriminator  $D$ . The feature matching loss is given by:

$$\mathcal{L}_{\text{rec}} = \mathbb{E}_{(\mathbf{x}, \mathbf{x}')} [\|D_i(\hat{\mathbf{x}}' \oplus H') - D_i(\mathbf{x}' \oplus H')\|_1] \quad (6)$$

where  $D_i$  denotes the  $i^{\text{th}}$ -layer feature extractor of the discriminator  $D$ .  $D_0$  denotes the discriminator input. The main advantage of the feature matching loss is that, differently from other perceptual losses, [32, 23], it does not require the use of an external pre-trained network. Finally the overall loss is obtained by combining Eqs. (6) and (5),  $\mathcal{L}_{\text{tot}} = \lambda_{\text{rec}} \mathcal{L}_{\text{rec}} + \mathcal{L}_{\text{gan}}^G$ . In all our experiments, we chose  $\lambda_{\text{rec}} = 10$  following [40]. Additional details of our implementation are shown in the Supplementary Material.

### 3.6. Generation Procedure

At test time, our network receives a driving video and a source image. In order to generate the  $t^{\text{th}}$  frame,  $\Delta$  estimates the keypoint locations  $h_k^s$  in the source image. Similarly, we estimate the keypoint locations  $h_k^1$  and  $h_k^t$  from first and the  $t^{\text{th}}$  frames of the driving video. Rather than generating a video from the absolute positions of the keypoints, the source image keypoints are transferred according to the relative difference between keypoints in the video. The keypoints in the generated frame are given by:

$$h_k^{s'} = h_k^s + (h_k^t - h_k^1) \quad (7)$$

The keypoints  $h_k^{s'}$  and  $h_k^s$  are then encoded as heatmaps using the covariance matrices estimated from the driving video, as described in Sec. 3.2. Finally, the heatmaps are given to the dense motion and the generator networks together with the source image (see Secs. 3.3 and 3.4). Importantly, one limitation of transferring relative motion is that it cannot be applied to arbitrary source images. Indeed, if the driving video object is not roughly aligned with the source image object, Eq. (7) may lead to absolute keypoint positions that are physically impossible for the considered object.

	<i>Tai-Chi</i> (AKD, MKR)			Nemo			Bair
	$\mathcal{L}_1$	AED	$\mathcal{L}_1$	AKD	AED	$\mathcal{L}_1$	
X2Face	0.068 (4.50, 35.7%)	0.27	0.022	0.47	0.140	0.069	
Ours	<b>0.050</b> ( <b>2.53, 17.4%</b> )	<b>0.21</b>	<b>0.017</b>	<b>0.37</b>	<b>0.072</b>	<b>0.025</b>	

Table 1: Video reconstruction comparisons

## 4. Experiments

In this section, we present a in-depth evaluation on three problems, tested on three very different datasets and employing a large variety of metrics.

**Datasets.** The Uva-*Nemo* dataset [11] is a facial dynamics analysis dataset composed of 1240 videos. We follow the same pre-processing as in [41]. Specifically, faces are aligned using the OpenFace library [1] before re-sizing each frame to  $64 \times 64$  pixels. Each video starts from a neutral expression and lasts 32 frames. As in [41], we use 1110 videos for training and 124 for evaluation.

The *Tai-Chi* dataset [35] is composed of 4500 tai-chi video clips downloaded from YouTube. We use the data as pre-processed in [35]. In particular, the frames are resized to  $64 \times 64$  pixels. The videos are split into 3288 and 822 videos for training and testing respectively. The video length varies from 32 to 100 frames.

The *BAIR* robot pushing dataset [12] contains videos collected by a Sawyer robotic arm pushing a variety of objects over a table. It contains 40960 training and 256 test videos. Each video is  $64 \times 64$  pixels and has 30 frames.

**Evaluation Protocol.** Evaluating the results of image animation methods is a difficult task, since ground truth animations are not available. In addition, to the best of our knowledge, X2Face [42] is the only previous approach for data-driven model-free image animation. For these two reasons, we evaluate our method also on two closely related tasks. As proposed in [42], we first evaluate Monkey-Net on the task of video reconstruction. This consists in reconstructing the input video from a representation in which motion and content are decoupled. This task is a “proxy” task to image animation and it is only introduced for the purpose of quantitative comparison. In our case, we combine the extracted keypoints  $\Delta(x)$  of each frame and the first frame of the video to re-generate the input video. Second, we evaluate our approach on the problem of Image-to-Video translation. Introduced in [38], this problem consists of generating a video from its first frame. Since our model is not directly designed for this task, we train a small recurrent neural network that predicts, from the keypoint coordinates in the first frame, the sequence of keypoint coordinates for the other 32 frames. Additional details can be found in the Supplementary Material. Finally, we evaluate our model on image animation. In all experiments we use K=10.

**Metrics.** In our experiments, we adopt several metrics in

order to provide an in-depth comparison with other methods. We employ the following metrics.

- $\mathcal{L}_1$ . In the case of the video reconstruction task where the ground truth video is available, we compare the average  $\mathcal{L}_1$  distance between pixel values of the ground truth and the generated video frames.
- AKD. For the *Tai-Chi* and *Nemo* datasets, we employ external keypoint detectors in order to evaluate whether the motion of the generated video matches the ground truth video motion. For the *Tai-Chi* dataset, we employ the human-pose estimator in [8]. For the *Nemo* dataset we use the facial landmark detector of [6]. We compute these keypoints for each frame of the ground truth and the generated videos. From these externally computed keypoints, we deduce the *Average Keypoint Distance* (AKD), *i.e.* the average distance between the detected keypoints of the ground truth and the generated video.
- MKR. In the case of the *Tai-Chi* dataset, the human-pose estimator returns also a binary label for each keypoint indicating whether the keypoints were successfully detected. Therefore, we also report the *Missing Keypoint Rate* (MKR) that is the percentage of keypoints that are detected in the ground truth frame but not in the generated one. This metric evaluates the appearance quality of each video frame.
- AED. We compute the feature-based metric employed in [13] that consists in computing the *Average Euclidean Distance* (AED) between a feature representation of the ground truth and the generated video frames. The feature embedding is chosen such that the metric evaluates how well the identity is preserved. More precisely, we use a network trained for facial identification [1] for *Nemo* and a network trained for person re-id [18] for *Tai-Chi*.
- FID. When dealing with Image-to-video translation, we complete our evaluation with the *Frechet Inception Distance* [19] (FID) in order to evaluate the quality of individual frames.

Furthermore, we conduct a user study for both the Image-to-Video translation and the image animation tasks (see Sec. 4.3).

### 4.1. Ablation Study

In this section, we present an ablation study to empirically measure the impact of each part of our proposal on the performance. First, we describe the methods obtained by “amputating” key parts of the model described in Sec. 3.1: (i) *No  $\mathcal{F}$*  - the dense optical flow network  $M$  is not used; (ii) *No  $\mathcal{F}_{coarse}$*  - in the optical flow network  $M$ , we do not use the part based-approach; (iii) *No  $\mathcal{F}_{residual}$*  - in the Optical Flow network  $M$ , we do not use  $\mathcal{F}_{residual}$ ; (iv) *No  $\Sigma_k$*  - we do not estimate the covariance matrices  $\Sigma_K$  in the keypoint detector  $\Delta$  and the variance is set to  $\Sigma_k = 0.01$  as in [22]; (v) the source image is not given to the motion network  $M$ ,

Tai-Chi			
	$\mathcal{L}_1$	(AKD, MKR)	AED
No $\mathcal{F}$	0.057	(3.11, 23.8%)	0.24
No $\mathcal{F}_{\text{residual}}$	0.051	(2.81, 18.0%)	0.22
No $\mathcal{F}_{\text{coarse}}$	0.052	(2.75, 19.7%)	0.22
No $\Sigma_k$	0.054	(2.86, 20.6%)	0.23
No $\mathbf{x}$	0.051	(2.71, 19.3%)	<b>0.21</b>
Full	<b>0.050</b>	( <b>2.53, 17.4%</b> )	<b>0.21</b>

Table 2: Video reconstruction ablation study *TaiChi*.



Figure 4: Qualitative ablation evaluation of video reconstruction on *Tai-Chi*.

$M$  estimates the dense optical flow only from the keypoint location differences; (vi) *Full* denotes the full model as described in Sec. 3.

In Tab. 2, we report the quantitative evaluation. We first observe that our full model outperforms the baseline method without deformation. This trend is observed according to all the metrics. This illustrates the benefit of deforming the features maps according to the estimated motion. Moreover, we note that *No  $\mathcal{F}_{\text{coarse}}$*  and *No  $\mathcal{F}_{\text{residual}}$*  both perform worse than when using the full optical flow network. This illustrates that  $\mathcal{F}_{\text{coarse}}$  and  $\mathcal{F}_{\text{residual}}$  alone are not able to estimate dense motion accurately. A possible explanation is that  $\mathcal{F}_{\text{coarse}}$  cannot estimate non rigid motions and that  $\mathcal{F}_{\text{residual}}$ , on the other hand, fails in predicting the optical flow in the presence of large motion. The qualitative results shown in Fig. 4 confirm this analysis. Furthermore, we observe a drop in performance when covariance matrices are replaced with static diagonal matrices. This shows the benefit of encoding more information when dealing with videos with complex and large motion, as in the case of the *TaiChi* dataset. Finally, we observe that if the appearance is not provided to the deformation network  $M$ , the video reconstruction performance is slightly lower.

Tai-Chi				
	FID	AED	MKR	
MoCoGAN [35]	54.83	0.27	46.2%	
Ours	<b>19.75</b>	<b>0.17</b>	<b>30.3%</b>	
Nemo		Bair		
	FID	AED	FID	
MoCoGAN [35]	51.50	0.33	MoCoGAN [35]	244.00
CMM-Net [41]	27.27	0.13	SV2P [2]	57.90
Ours	<b>11.97</b>	<b>0.12</b>	Ours	<b>23.20</b>

Table 3: Image-to-video translation comparisons.

## 4.2. Comparison with Previous Works

**Video Reconstruction.** First, we compare our results with the X2Face model [42] that is closely related to our proposal. Note that this comparison can be done since we employ image and motion representation of similar dimension. In our case, each video frame is reconstructed from the source image and 10 landmarks, each one represented by 5 numbers (two for the location and three for the symmetric covariance matrix), leading to a motion representation of dimension 50. For X2face, motion is encoded into a driving vector of dimension 128. The quantitative comparison is reported in Tab. 1. Our approach outperforms X2face, according to all the metrics and on all the evaluated datasets. This confirms that encoding motion via motion-specific keypoints leads to a compact but rich representation.

**Image-to-Video Translation:** In Tab. 3 we compare with the state of the art Image-to-Video translation methods: two unsupervised methods MoCoGAN [35] and SV2P [2], and CMM-Net which is based on keypoints [41]. CMM-Net is evaluated only on *Nemo* since it requires facial landmarks. We report results SV2P on the *Bair* dataset as in [2]. We can observe that our method clearly outperforms the three methods for all the metrics. This quantitative evaluation is confirmed by the qualitative evaluation presented in the Supplementary material. In the case of MoCoGAN, we observe that the AED score is much higher than the two other methods. Since AED measures how well the identity is preserved, these results confirm that, despite the realism of the video generated by MoCoGAN, the identity and the person-specific details are not well preserved. A possible explanation is that MoCoGAN is based on a feature embedding in a vector, which does not capture spatial information as well as the keypoints. The method in [41] initially produces a realistic video and preserves the identity, but the lower performance can be explained by the apparition of visual artifacts in the presence of large motion (see the Supplementary material for visual examples). Conversely, our method both preserves the person identity and performs

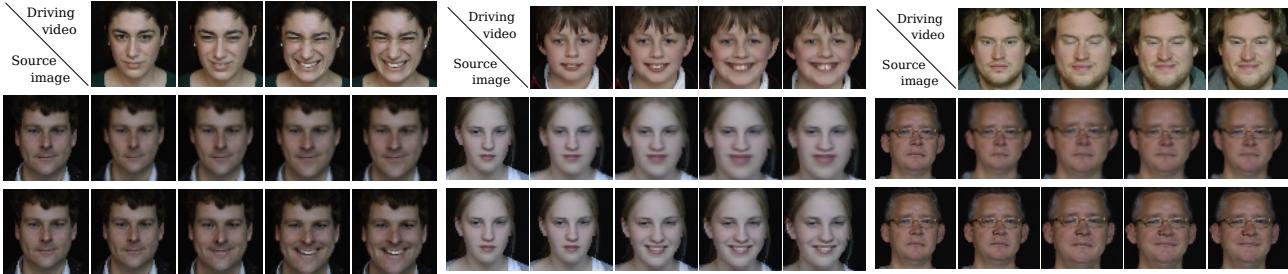


Figure 5: Qualitative results for image animation on the *Nemo* dataset: X2face (2-nd row) against our method (3-rd row).



Figure 6: Qualitative results for image animation on the *Tai-Chi* dataset: X2face (2-nd row) against our method (3-rd row)

<i>Tai-Chi</i>	<i>Nemo</i>	<i>Bair</i>
85.0%	79.2%	90.8%

Table 4: User study results on image animation. Proportion of times our approach is preferred over X2face [42].

well even under large spatial deformations.

**Image Animation.** In Fig. 5, we compare our method with X2Face [42] on the *Nemo* dataset. We note that our method generates more realistic smiles on the three randomly selected samples despite the fact that the XFace model is specifically designed for faces. Moreover, the benefit of transferring the relative motion over absolute locations can be clearly observed in Fig. 5 (column 2). When absolute locations are transferred, the source image inherits the face proportion from the driving video, resulting in a face with larger cheeks. In Fig. 6, we compare our method with X2Face on the *Tai-Chi* dataset. X2Face [42] fails to consider each body-part independently and, consequently, warps the body in such a way that its center of mass matches the center of mass in the driving video. Conversely, our method successfully generates plausible motion sequences that match the driving videos. Concerning the *Bair* dataset, exemplar videos are shown in the Supplementary material. The results are well in line with those obtained on the two other datasets.

### 4.3. User Evaluation

In order to further consolidate the quantitative and qualitative evaluations, we performed user studies for both the Image-to-Video translation (see the Supplementary Material) and the image animation problems using Amazon Me-

chanical Turk.

For the image animation problem, our model is again compared with X2face [42] according to the following protocol: we randomly select 50 pairs of videos where objects in the first frame have a similar pose. Three videos are shown to the user: one is the driving video (reference) and 2 videos from our method and X2Face. The users are given the following instructions: *Select the video that better corresponds to the animation in the reference video.* We collected annotations for each video from 10 different users. The results are presented in Tab. 4. Our generated videos are preferred over X2Face videos in almost more than 80% of the times for all the datasets. Again, we observe that the preference toward our approach is higher on the two datasets which correspond to large motion patterns.

## 5. Conclusion

We introduced a novel deep learning approach for image animation. Via the use of motion-specific keypoints, previously learned following a self-supervised approach, our model can animate images of arbitrary objects according to the motion given by a driving video. Our experiments, considering both automatically computed metrics and human judgments, demonstrate that the proposed method outperforms previous work on unsupervised image animation. Moreover, we show that with little adaptation our method can perform Image-to-Video translation. In future work, we plan to extend our framework to handle multiple objects and investigate other strategies for motion embedding.

## References

- [1] B. Amos, B. Ludwiczuk, and M. Satyanarayanan. Openface: A general-purpose face recognition. 2016. 6
- [2] M. Babaeizadeh, C. Finn, D. Erhan, R. H. Campbell, and S. Levine. Stochastic variational video prediction. In *ICLR*, 2017. 2, 7, 11, 12
- [3] G. Balakrishnan, A. Zhao, A. V. Dalca, F. Durand, and J. Guttag. Synthesizing images of humans in unseen poses. In *CVPR*, 2018. 1, 3
- [4] A. Bansal, S. Ma, D. Ramanan, and Y. Sheikh. Recycle-gan: Unsupervised video retargeting. In *ECCV*, 2018. 1, 2
- [5] V. Blanz and T. Vetter. A morphable model for the synthesis of 3d faces. In *SIGGRAPH*, 1999. 1, 2
- [6] A. Bulat and G. 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. 4, 6
- [7] C. Cao, Q. Hou, and K. Zhou. Displaced dynamic expression regression for real-time facial tracking and animation. *TOG*, 2014. 1, 2
- [8] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *CVPR*, 2017. 6
- [9] C. Chan, S. Ginosar, T. Zhou, and A. A. Efros. Everybody dance now. In *ECCV*, 2018. 1, 2
- [10] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. In *NIPS*, 2014. 11
- [11] H. Dibeklioğlu, A. A. Salah, and T. Gevers. Are you really smiling at me? spontaneous versus posed enjoyment smiles. In *ECCV*, 2012. 2, 6
- [12] F. Ebert, C. Finn, A. X. Lee, and S. Levine. Self-supervised visual planning with temporal skip connections. In *CoRL*, 2017. 2, 6
- [13] P. Esser, E. Sutter, and B. Ommer. A variational u-net for conditional appearance and shape generation. In *CVPR*, 2018. 6
- [14] C. Finn, I. Goodfellow, and S. Levine. Unsupervised learning for physical interaction through video prediction. In *NIPS*, 2016. 2
- [15] Y. Ganin, D. Kononenko, D. Sungatullina, and V. Lempitsky. Deepwarp: Photorealistic image resynthesis for gaze manipulation. In *ECCV*, 2016. 3
- [16] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *NIPS*, 2014. 1
- [17] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*. 11
- [18] A. Hermans, L. Beyer, and B. Leibe. In defense of the triplet loss for person re-identification. *arXiv:1703.07737*, 2017. 6
- [19] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *NIPS*. 6
- [20] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. 11
- [21] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. In *CVPR*, 2017. 3
- [22] T. Jakab, A. Gupta, H. Bilen, and A. Vedaldi. Unsupervised learning of object landmarks through conditional image generation. In *NIPS*, 2018. 3, 6
- [23] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *ECCV*, 2016. 5
- [24] D. Joo, D. Kim, and J. Kim. Generating a fusion image: One's identity and another's shape. In *CVPR*, 2018. 1
- [25] N. Kalchbrenner, A. v. d. Oord, K. Simonyan, I. Danihelka, O. Vinyals, A. Graves, and K. Kavukcuoglu. Video pixel networks. In *ICML*, 2016. 2
- [26] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *ICLR*, 2014. 11
- [27] D. P. Kingma and M. Welling. Auto-encoding variational bayes. In *ICLR*, 2014. 1
- [28] X. Mao, Q. Li, H. Xie, R. Y. Lau, Z. Wang, and S. P. Smolley. Least squares generative adversarial networks. In *ICCV*. 5
- [29] J. Oh, X. Guo, H. Lee, R. L. Lewis, and S. Singh. Action-conditional video prediction using deep networks in atari games. In *NIPS*, 2015. 2
- [30] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*. Springer, 2015. 4
- [31] M. Saito, E. Matsumoto, and S. Saito. Temporal generative adversarial nets with singular value clipping. In *ICCV*, 2017. 1, 2
- [32] A. Siarohin, E. Sangineto, S. Lathuilière, and N. Sebe. Deformable gans for pose-based human image generation. In *CVPR*, 2018. 1, 3, 4, 5
- [33] N. Srivastava, E. Mansimov, and R. Salakhudinov. Unsupervised learning of video representations using lstms. In *ICML*, 2015. 2
- [34] J. Thies, M. Zollhofer, M. Stamminger, C. Theobalt, and M. Nießner. Face2face: Real-time face capture and reenactment of rgb videos. In *CVPR*, 2016. 1, 2
- [35] S. Tulyakov, M.-Y. Liu, X. Yang, and J. Kautz. Mocogan: Decomposing motion and content for video generation. In *CVPR*, 2018. 1, 2, 6, 7, 11, 12
- [36] J. Van Amersfoort, A. Kannan, M. Ranzato, A. Szlam, D. Tran, and S. Chintala. Transformation-based models of video sequences. *arXiv preprint arXiv:1701.08435*, 2017. 2
- [37] R. Villegas, J. Yang, Y. Zou, S. Sohn, X. Lin, and H. Lee. Learning to generate long-term future via hierarchical prediction. In *ICML*, 2017. 2
- [38] C. Vondrick, H. Pirsiavash, and A. Torralba. Generating videos with scene dynamics. In *NIPS*, 2016. 1, 2, 6
- [39] T.-C. Wang, M.-Y. Liu, J.-Y. Zhu, G. Liu, A. Tao, J. Kautz, and B. Catanzaro. Video-to-video synthesis. In *NIPS*, 2018. 1, 2
- [40] T.-C. Wang, M.-Y. Liu, J.-Y. Zhu, A. Tao, J. Kautz, and B. Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In *CVPR*, 2017. 5, 11
- [41] W. Wang, X. Alameda-Pineda, D. Xu, P. Fua, E. Ricci, and N. Sebe. Every smile is unique: Landmark-guided diverse smile generation. In *CVPR*, 2018. 1, 2, 6, 7, 11, 12

- [42] O. Wiles, A. S. Koepke, and A. Zisserman. X2face: A network for controlling face generation using images, audio, and pose codes. In *ECCV*, 2018. [1](#), [2](#), [6](#), [7](#), [8](#), [12](#)
- [43] Y. Zhang, Y. Guo, Y. Jin, Y. Luo, Z. He, and H. Lee. Unsupervised discovery of object landmarks as structural representations. In *CVPR*, 2018. [3](#)
- [44] L. Zhao, X. Peng, Y. Tian, M. Kapadia, and D. Metaxas. Learning to forecast and refine residual motion for image-to-video generation. In *ECCV*. [2](#)
- [45] M. Zollhöfer, J. Thies, P. Garrido, D. Bradley, T. Beeler, P. Pérez, M. Stamminger, M. Nießner, and C. Theobalt. State of the art on monocular 3d face reconstruction, tracking, and applications. In *Computer Graphics Forum*. [2](#)

## A. Implementation details

As described in Sec. 3, each module employs a U-Net architecture. We use the exact same architecture for all the networks. More specifically each block of each of the encoder consists of a  $3 \times 3$  convolution, batch normalization [20], ReLU and average pooling. The first convolution layers have 32 filters and each subsequent convolution doubles the number of filters. Each encoder is composed of a total of 5 blocks. The decoder blocks have similar structure:  $3 \times 3$  convolution, batch normalization and ReLU followed by nearest neighbour up-sampling. The first block of the decoder has 512 filters. Each consequent block has the reduced number of filters by a factor of 2.

As described in Sec 3, the keypoint detector  $\Delta$  produces  $K$  heatmaps followed by softmax. In particular, we employ softmax activations with 0.1 temperature. Indeed, thanks to the use of a low temperature for softmax, we obtain sharper heatmaps and avoid uniform heatmaps that would lead to keypoints constantly located in the image center.

For  $G$ , we employ 4 additional Residual Blocks [17] in order to remove possible warping artifacts produced by  $M$ . The output of  $G$  is a 3 channel feature map followed by the sigmoid. We use the discriminator architecture described in [40].

The framework is trained for  $T$  epochs where  $T$  equals 250, 500 and 10 for *Tai-Chi*, *Nemo* and *Bair* respectively. Epoch involves training the network on 2 randomly sampled frames from each training video. We use the Adam optimizer [26] with learning rate  $2e-4$  and then with learning rate  $2e-5$  for another  $\frac{T}{2}$  epochs.

As explained in Sec. 4.2, for *Image-to-Video* translation, we employ a single-layer GRU network in order to predict the keypoint sequence used to generate the video. This recurrent network [10] has 1024 hidden units and is trained via  $\mathcal{L}_1$  minimization.

## B. MGif dataset

We collected an additional dataset of videos containing movements of different cartoon animals. Each video is a moving *gif* file. Therefore, we called this new dataset *MGif*. The dataset consists of 1000 videos, we used 900 videos for training and 100 for evaluation. Each video has size  $128 \times 128$  and contains from 5 to 100 frames. The dataset is particularly challenging because of the high appearance variation and motion diversity. Note that in the experiments on this dataset, we use absolute keypoint locations from the driving video instead of the relative keypoint motion detailed in Sec. 3.6.

## C. Additional experimental results

In this section, we report additional results. In Sec C.1, we complete the ablation study and, in Secs. C.2 and C.3, we report qualitative results for both the image-to-video and image animation problems. Finally, in Sec. C.4, we visualize the keypoint predicted by our self-supervised approach.

### C.1. Additional ablation study

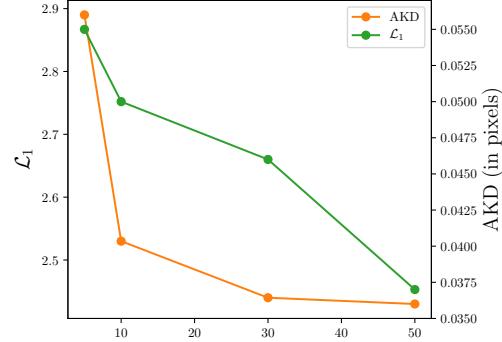


Figure 7: Reconstruction Errors as functions of the number of keypoints. Results obtained on *Tai-Chi*.

We perform experiments to measure the impact of the number of keypoints on video reconstruction quality. We report results on *Tai-Chi* dataset in Fig. 7. We computed  $\mathcal{L}_1$  and  $AKD$  metrics as described in the paper. As expected, increasing the number of keypoints leads to a lower reconstruction error, but additional keypoints introduce memory and computational overhead. We use 10 keypoints in all our experiments, since we consider this to be a good trade-off.

### C.2. Image-to-Video translation

As explained in Sec. 4.2 of the main paper, we compare with the three state of the art methods for Image-to-Video translation: MoCoGAN [35] and SV2P [2], and CMM-Net [41]. CMM-Net is evaluated only on *Nemo* and SV2P only on the *Bair* dataset. We report a user study and qualitative results.

**User Study.** We perform a user study for the image-to-video translation problem. As explained in Sec. 4.3, we perform pairwise comparisons between our method and the competing methods. We employ the following protocol: we randomly select 50 videos and use the first frame of each of video as the reference frames to generate new videos. For each of the 50 videos the initial frame, and two videos generated by our and one of the competing methods are shown to the user. We provide the following instructions: "Select a more realistic animation of the reference image". As in

Sec. 4.2 of the main paper, our method is compared with MoCoGAN [35], Sv2p [2], and CMM-Net [41]. The results of the user study are presented in Table 5. On average, users preferred the videos generated by our approach over those generated by other methods. The preference gap is especially evident for the *Tai-Chi* and *Bair* datasets that contain a higher amount of large motion. This supports the ability of our approach to handle driving videos with large motion.

	<i>Tai-Chi</i>	<i>Nemo</i>	<i>Bair</i>
MoCoGAN [35]	88.2%	68.2%	90.6%
CMM-Net [41]	-	63.6%	-
SV2P [2]	-	-	98.8%

Table 5: User study results on image-to-video translation. Proportion of times our approach is preferred over the competitors methods

**Qualitative results.** We report additional qualitative results in Figs. 8, 9 and 10. These qualitative results further support the ability of our method to generate realistic videos from source images and driving sequences.

In particular, for the *Nemo* dataset (Fig. 8), MoCoGAN and CMM-Net suffer from more artifacts. In addition, the videos generated by MoCoGAN do not preserve the identity of the person. This issue is particularly visible when comparing the first and the last frames of the generated video. CMM-Net preserves better the identity but fails in generating realistic eyes and teeth. In contrast to these works, our method generates realistic smiles while preserving the person identity.

For *Tai-Chi* (Fig. 10), MoCoGAN [35] produces videos where some parts of the human body are not clearly visible (see rows 3,4 and 6). This is again due to the fact that visual information is embedded in a vector. Conversely, our method generates realistic human body with richer details.

For *Bair* (Fig. 9), [2] completely fails to produce videos where the robotic is sharp. The generated videos are blurred. MoCoGAN [35] generates videos with more details but containing many artifacts. In addition, the backgrounds generated by MoCoGAN are not temporally coherent. Our method generates realistic robotic arm moving in front of detailed and temporally coherent backgrounds.

### C.3. Image animation

As explained in the main paper, we compare our method with X2Face [42]. Results are reported in Figs. 12, 13 and 14 on the *Nemo*, *Tai-Chi* and *Bair* datasets respectively.

When tested using the *Nemo* dataset (Fig. 12), our method generates more realistic smiles on most of the randomly selected samples despite the fact that the XFace

model is specifically designed for faces. Similarly to the main paper, the benefit of transferring the relative motion over absolute locations can be clearly observed in the bottom right example where the video generated by X2face inherits the large cheeks of the young boy in the driving video.

For *Tai-Chi* (Fig. 13), X2face is not able to handle the motion of the driving video and simply warps the human body in the source image as a single blob.

For *Bair* (Fig. 14), we observe a similar behavior. X2face generates unrealistic videos where the robotic arm is generally not distinguishable. On the contrary, our model is able to generate a realistic robotic arm moving according to the driving video motion.

Finally in Fig 11, we report results on the *MGif* dataset. First, these examples illustrate high diversity of *MGif* dataset. Second, we observe that our model is able to transfer the motion of the driving video even if the appearance of the source frame is very different from the driving video. In particular, in all the generated sequences, we observe that the legs are correctly generated and follow the motion of the driving video. The model preserves the rigid parts of the animals as, for instance, the abdomen. In the last row, we see that the model is also able to animate the fox tail according to the motion of the cheetah tail.

### C.4. Keypoint visualization

Finally, we report visual examples of keypoints learned by our model in Figs. 15, 16, 17 and 18. On the *Nemo* dataset, we observe that the obtained keypoints are semantically consistent. For instance, the cyan and light green keypoints constantly correspond to the nose and the chin respectively. For *Tai-Chi*, the keypoints are also semantically consistent: light green for the chin and yellow for the left-side arm (right arm in frontal views and left arm in back views), for instance. For the *Bair* dataset, we observe that two keypoints (light green and dark blue) correspond to the robotic arm. The other keypoints are static and can correspond to the background. Finally, concerning the *MGif* dataset, we observe that each keypoint corresponds to two different animal parts depending if the animal is going towards left or right. In the case of animals going right (last three rows), the keypoints are semantically consistent (red for the tail, dark blue for the head etc.). Similarly, the keypoints are semantically consistent among images of animal going left (red for the head, dark blue for the tail etc.). Importantly, we observe that a keypoint is associated to each highly moving part, as legs and tails.

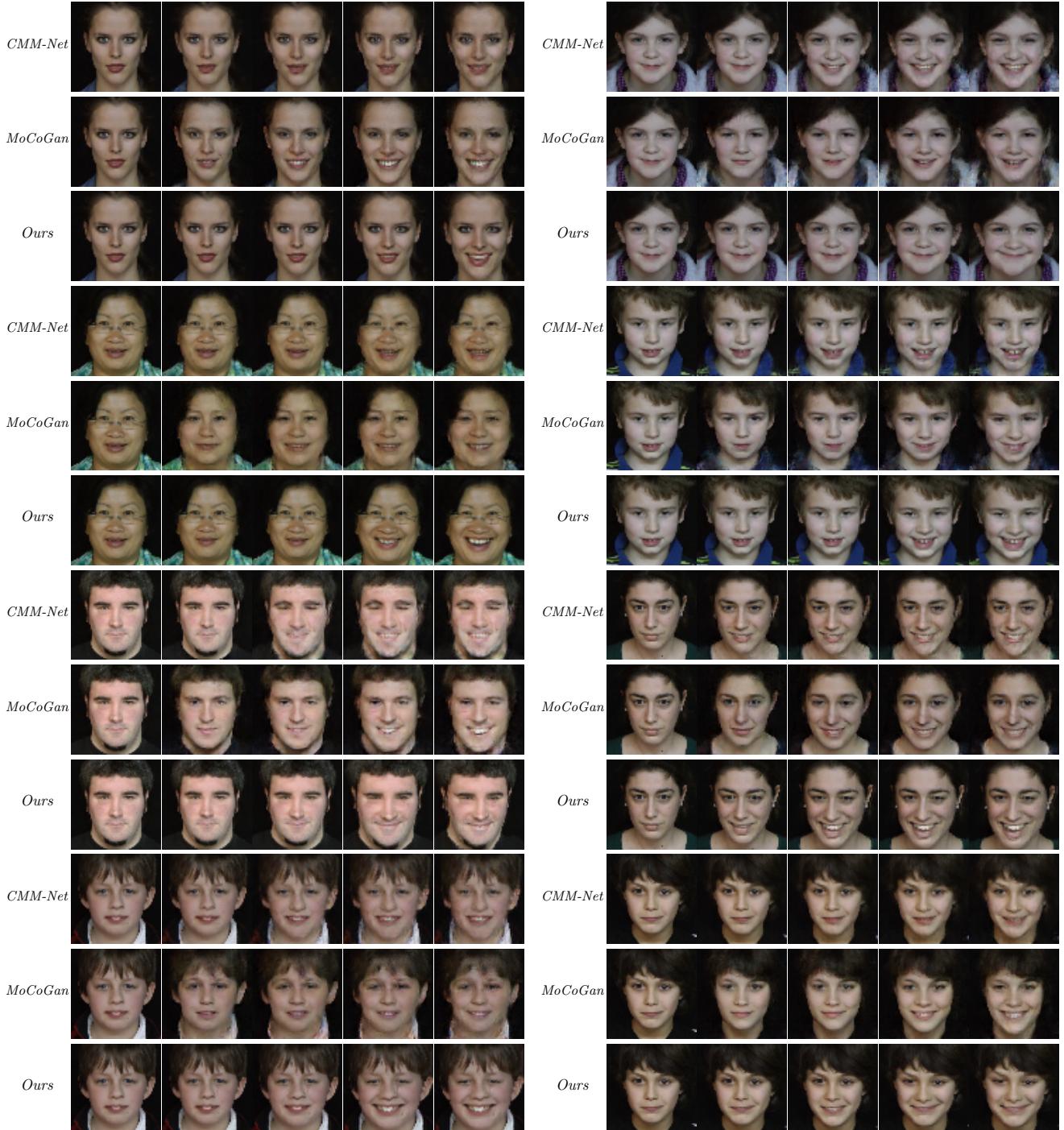


Figure 8: Qualitative results for *Image-to-video* translation on the *Nemo* dataset.

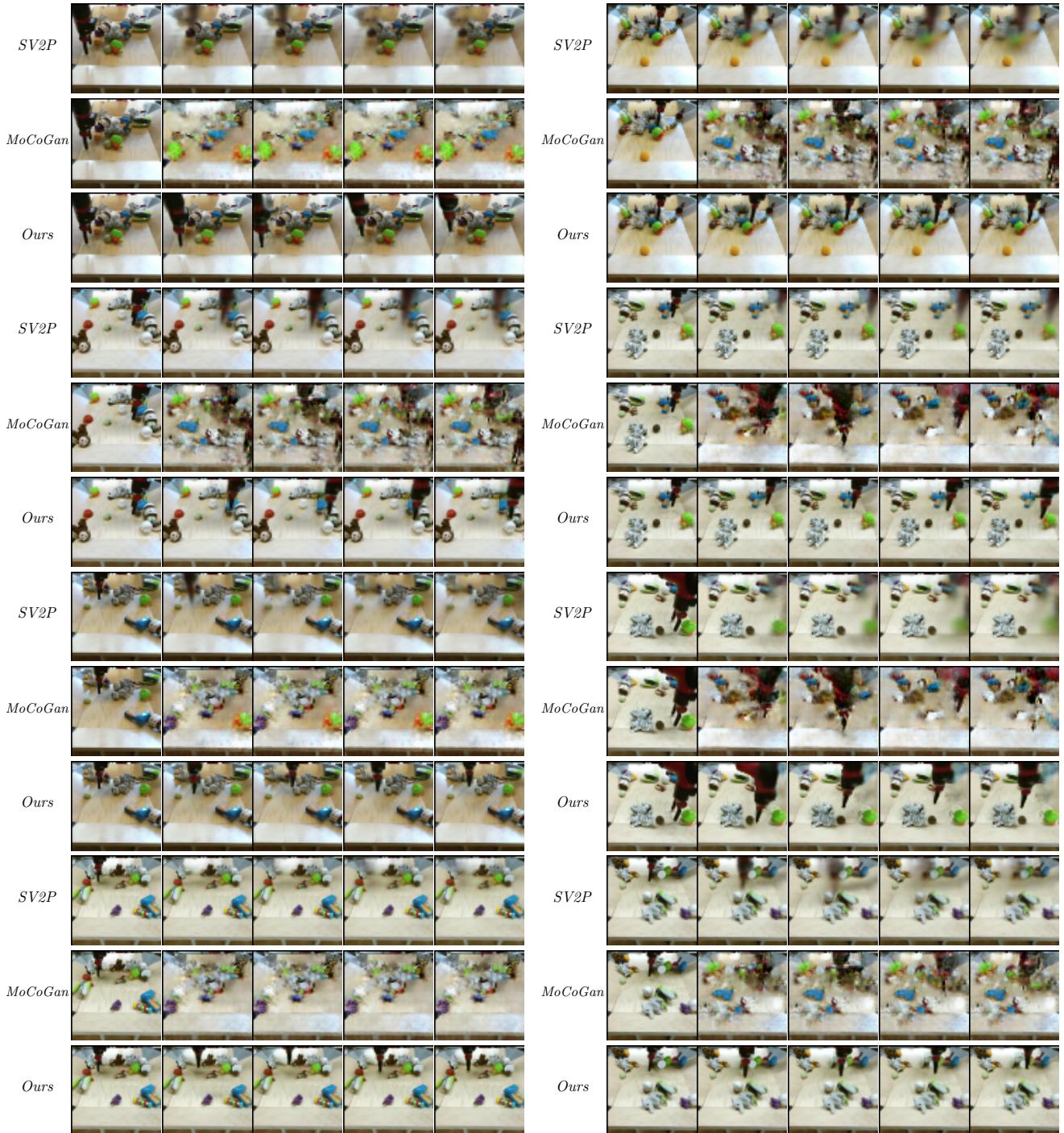


Figure 9: Qualitative results for *Image-to-video* translation on the *Bair* dataset.

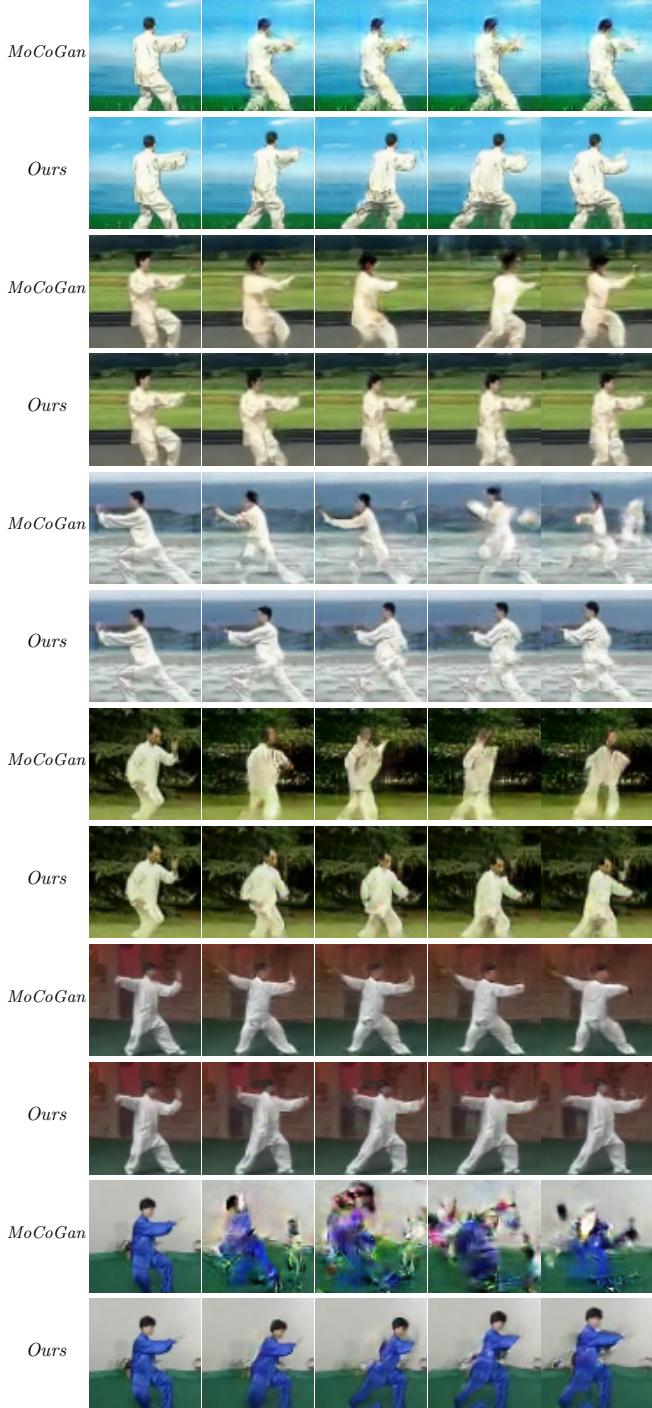


Figure 10: Qualitative results for *Image-to-video* translation on the *Tai-Chi* dataset.

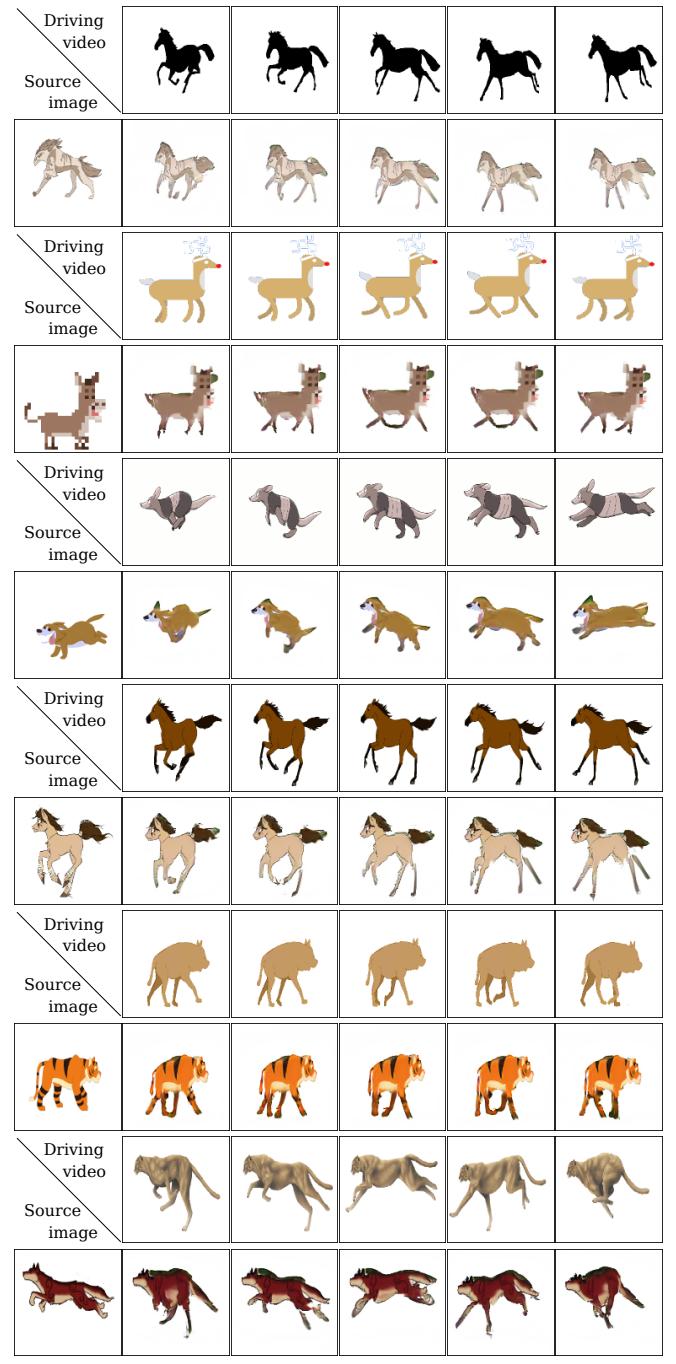


Figure 11: Qualitative results for image animation on the MGif.

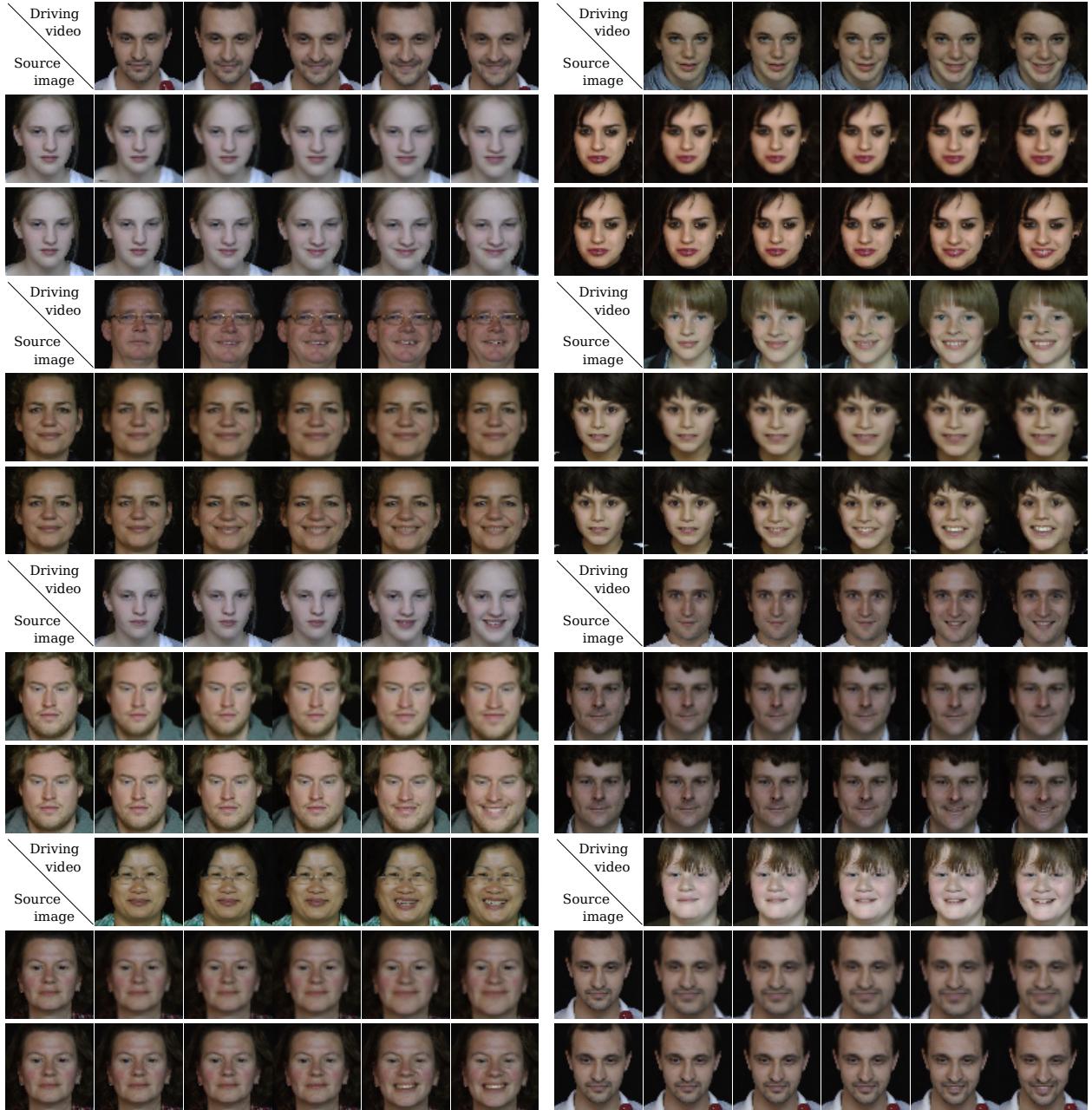


Figure 12: Additional qualitative results for image animation on the *Nemo* dataset: X2face (first) against our method (second).

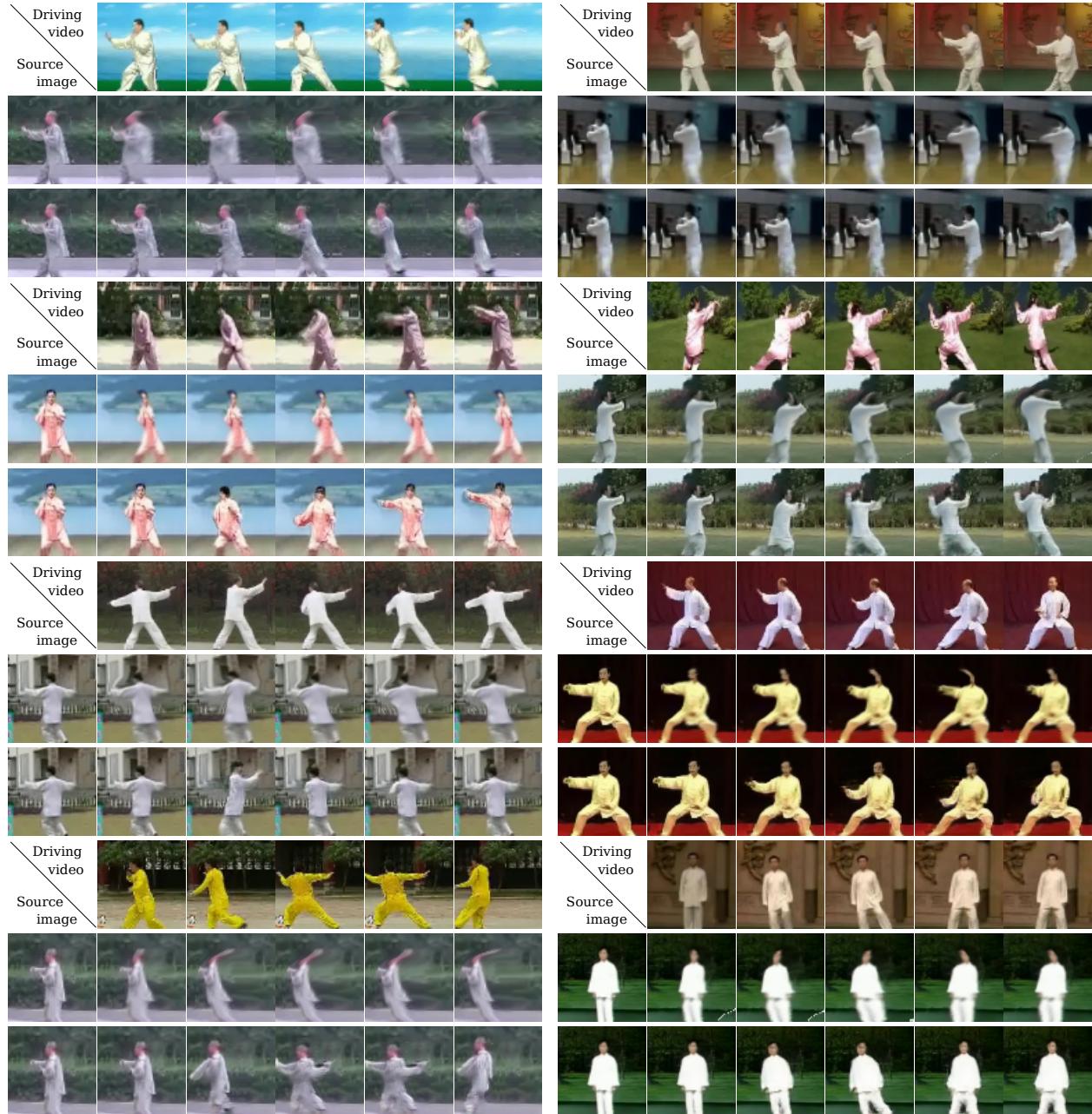


Figure 13: Additional qualitative results for image animation on the *Tai-Chi* dataset: X2face (first) against our method (second).



Figure 14: Additional qualitative results for image animation on the *Bair* dataset: X2face (first) against our method (second).



Figure 15: Keypoints predicted on the *Nemo* dataset.

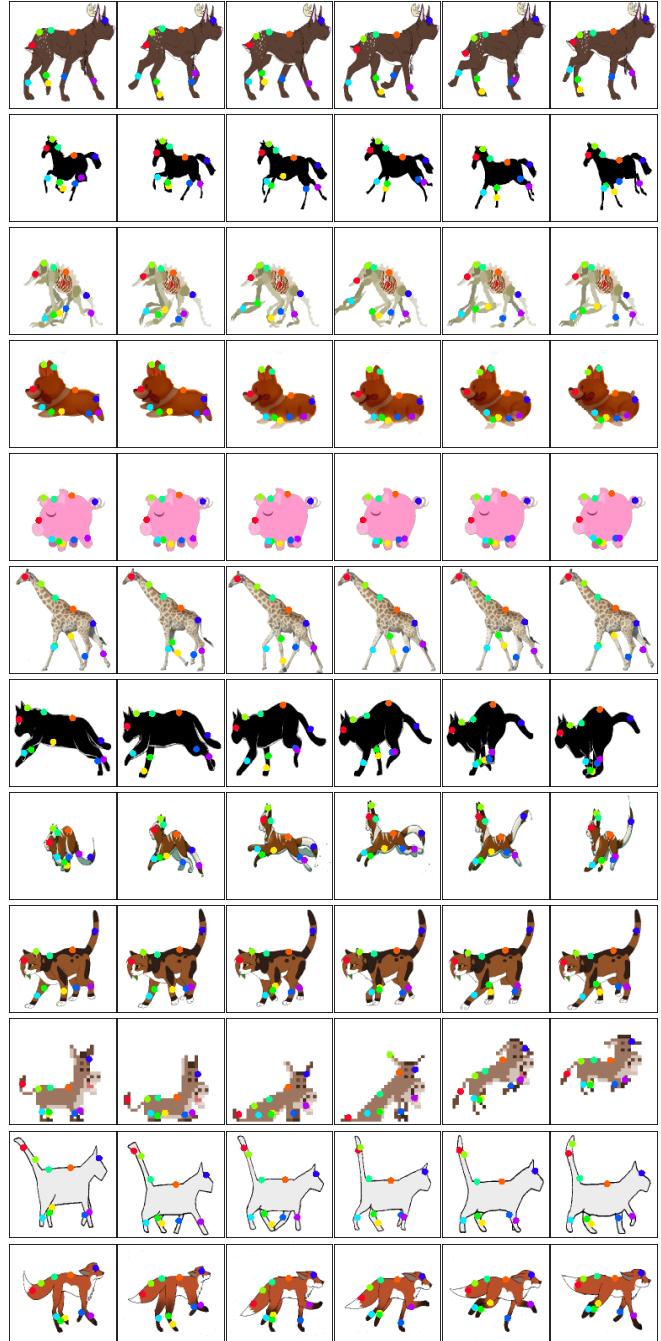


Figure 16: Keypoints predicted on the *MGif* dataset.

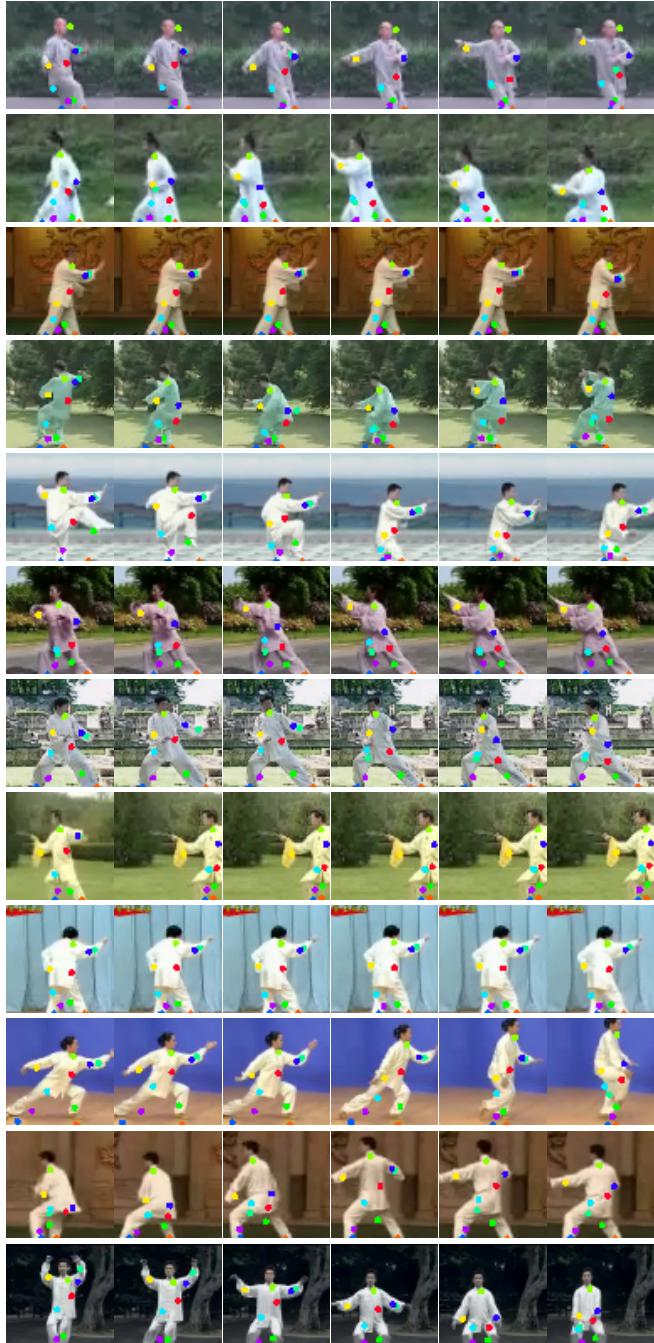


Figure 17: Keypoints predicted on the *Tai-Chi* dataset.

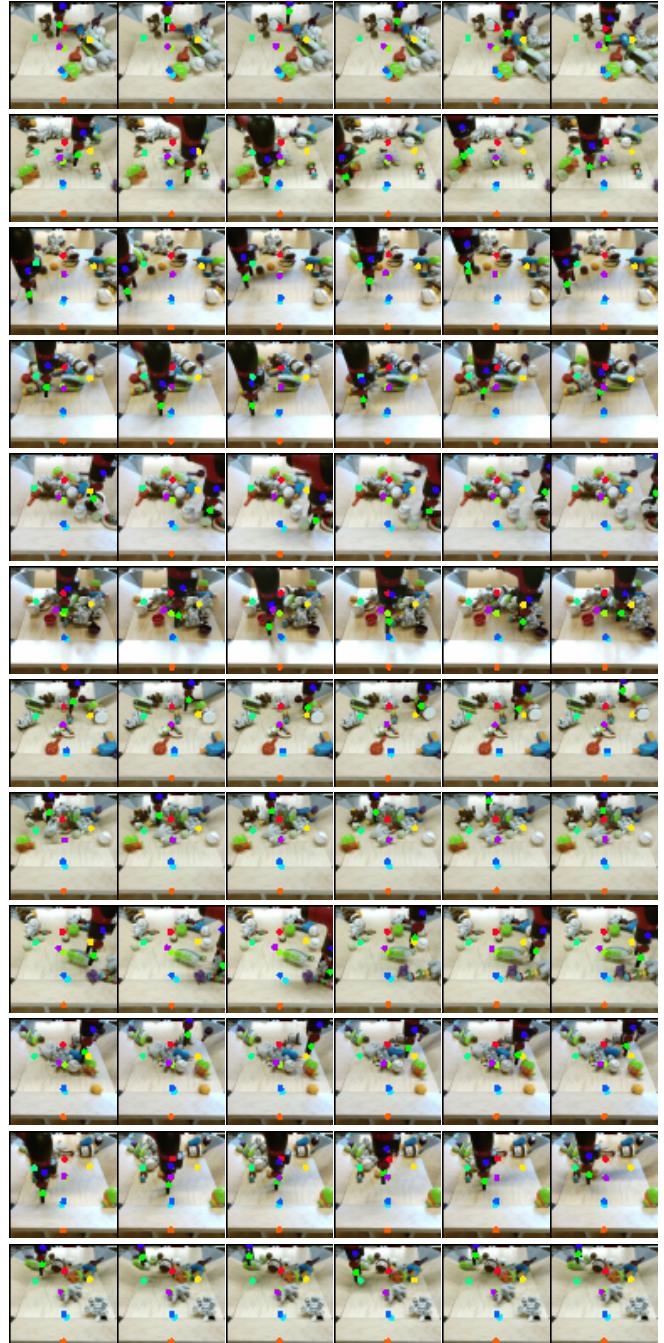


Figure 18: Keypoints predicted on the *Bair* dataset.