

# Self-Supervised 3D Human Pose Estimation in Static Video Via Neural Rendering

Luca Schmidtke<sup>1,2</sup>, Benjamin Hou<sup>1</sup>, Athanasios Vlontzos<sup>1</sup>, and Bernhard Kainz<sup>1,2</sup>

<sup>1</sup>Imperial College London, UK

<sup>2</sup>Friedrich-Alexander-Universität Erlangen-Nürnberg, DE

**Abstract.** Inferring 3D human pose from 2D images is a challenging and long-standing problem in the field of computer vision with many applications including motion capture, virtual reality, surveillance or gait analysis for sports and medicine. We present preliminary results for a method to estimate 3D pose from 2D video containing a single person and a static background without the need for any manual landmark annotations. We achieve this by formulating a simple yet effective self-supervision task: our model is required to reconstruct a random frame of a video given a frame from another timepoint and a rendered image of a transformed human shape template. Crucially for optimisation, our ray casting based rendering pipeline is fully differentiable, enabling end to end training solely based on the reconstruction task.

**Keywords:** self-supervised learning, 3D human pose estimation, 3D pose tracking, motion capture

## 1 Introduction

Inferring 3D properties of our world from 2D images is an intriguing open problem in computer vision, even more so when no direct supervision is provided in the form of labels. Although this problem is inherently ill-posed, humans are able to derive accurate depth estimates, even when their vision is impaired, from motion cues and semantic prior knowledge about the perceived world around them. This is especially true for human pose estimation. Self-supervised learning has proven to be an effective technique to utilise large amounts of unlabelled video and image sources. On a more fundamental note, self-supervised learning is hypothesised to be an essential component in the emergence of intelligence and cognition. Moreover, self-supervised approaches allow for more flexibility in domains such as the medical sector where labels are often hard to come by. In this paper we focus on self-supervised 3D pose estimation from monocular video, a key element of a wide range of applications including motion capture, visual surveillance or gait analysis.

Inspired by previous work, we model pose as a factor of variation throughout different frames of a video of a single person and a static background. More

formally, self-supervision is provided by formulating a conditional image reconstruction task: given a pose input different from the current image, what would that image look like if we condition it on the given pose? Differently from previous work, we choose to represent pose as a 3D template consisting of connected parts which we transform and project to two-dimensional image space, thereby inferring 3D pose from monocular images without explicit supervision.

More specifically, our method builds upon the recent emergence and success of combining deep neural networks with an explicit 3D to 2D image formation process through fully differentiable rendering pipelines. This inverse-graphics approach follows the analysis by synthesis principle of generative models in a broader context: We hope to extract information about the 3D properties of objects in our world by trying to recreate their perceived appearance on 2D images. Popular rendering techniques rely on different representations including meshes and polygons, point clouds or implicit surfaces. In our work we make use of volume rendering with a simple occupancy function or density combined with a texture field that assign an occupancy between  $[0, 1]$  and RGB colour value  $c \in \mathbb{R}^3$  for every point defined on a regular 3D grid.

## 2 Related Work

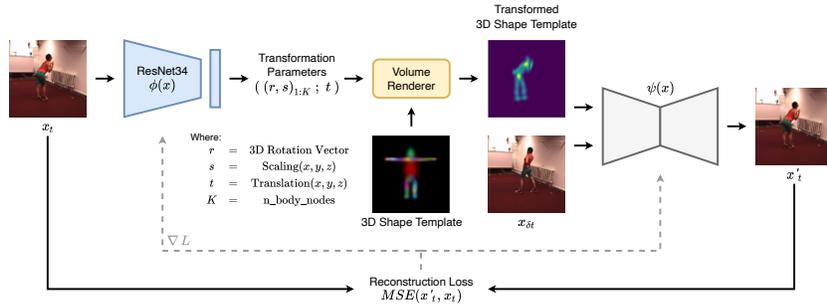
**Monocular 3D Human Pose Estimation** Human pose estimation in general is a long standing problem in computer vision with an associated large body of work and substantial improvements since the advent of deep-learning based approaches. Inferring 3D pose from monocular images however remains a challenging problem tackled by making use of additional cues in the image or video such as motion or multiple views from synchronised cameras or introducing prior knowledge about the hierarchical part based structure of the human body.

**Lifting from 2D to 3D** Many works break down the problem into first estimating 2D pose and subsequently estimate 3D pose either directly [19], by leveraging self-supervision through transformation and reprojection [15] or a kd-tree to find corresponding pairs of detected 2D pose and stored 3D pose [4].

**Motion Cues From Video** Videos provide a rich source of additional temporal information that can be exploited to limit the solution space. [16], [8], [2] and [10] use recurrent architectures in the form of LSTMs or GRUs to incorporate temporal context while [23] employ temporal convolutions and a reprojection objective.

**Multiple Views** Other approaches incorporate images from multiple, synchronised cameras to alleviate the ill-posedness of the problem. [22], [31] and [24] fuse multiple 2D heatmaps while [26] and [27] utilize multi-view consistency as a form of additional supervision in the objective function.

**Human Body Prior** Using non-parametric belief propagation, [29] estimate the 2D pose of loosely-linked human body parts from image features and use a mixture of experts to estimate a conditional distribution of 3D poses. Many more recent approaches rely on features extracted from convolutional neural



**Fig. 1.** Our method – left to right– An input frame  $x_t$  passes through the pose extractor encoder  $\phi$  and produces the transformation parameters for each skeletal node of the shape template  $\mathbf{T}$ . The transformed template is then rendered and concatenated with a random frame of the same sequence,  $x_{\delta t}$ , and is passed into an auto-encoder that’s tasked to reconstruct the original frame.

networks [14]. Many works such as [7], [11] and [10] make use of SMPL [18], a differentiable generative model that produces a 3D human mesh based on disentangled shape and pose parameters. [32] leverage kinematic constraints to improve their predictions while [12] leverage a forward kinematics formulation in combination with the transformation of a 2D part-based template to formulate self-supervision in form of image reconstruction similar in some ways to our approach.

**Human Neural Rendering** Recently, neural rendering approaches, *ie.* fully differentiable rendering pipelines, have gained a lot of attention. Volume rendering techniques [[17], [20]] have been demonstrated to be powerful tools to infer 3D properties of objects from 2D images when used in combination with neural networks. The end-to-end differentiability offers the intriguing opportunity to directly leverage pixel-wise reconstruction losses as a strong self-supervision signal. This has sparked a number of very recent works estimating human 3D shape and pose via neural radiance fields [20] [[13], [30]] or signed-distance function based rendering [6].

### 3 Method

Our approach relies on self-supervision through image reconstruction conditioned on a transformed and rendered shape template. The images are sampled from a video containing a single person moving in front of a static background. More formally, the goal is to reconstruct a number of frames  $(\mathbf{x}_{t_1}, \mathbf{x}_{t_n})$  from random time points  $t_1, \dots, t_n$  in a video with access to *one* frame  $\mathbf{x}_{t_k}$ , again sampled randomly, and rendered images of transformed templates  $\mathbf{T}_1, \dots, \mathbf{T}_n$ .

Our method can be viewed in two distinct steps; regression of template transformation parameters and image reconstruction, where both steps are parameter-

ized using deep convolutional neural networks. An encoder network  $\phi$  regresses rotation, translation and scale parameters from frame  $\mathbf{x}_t$  in order to transform each skeletal node of a 3D shape template,  $\mathbf{T}$ . The generator network  $\psi$  takes as input; (a) a frame  $\mathbf{x}_{\delta t}$ , from a different time instance of the same sequence where the same person assumes a different pose, and (b) a rendered image of the transformed 3D template while being tasked to reconstruct frame  $\mathbf{x}_t$ .

The encoder  $\phi$  consists of a convolutional neural network for feature extraction followed by a number of linear layers and a reshape operation. The generative network  $\psi$  resembles a typical convolutional encoder-decoder structure utilised for image translation, where feature maps are subsequently down-sampled via strided convolutions and the number of features increases. For the decoder we utilise bilinear upsampling and spatially adaptive instance normalisation (SPADE) [21] to facilitate semantic inpainting of the rendered template image.

### 3.1 Template and Volume Rendering

**Shape Template:** A shape template,  $\mathbf{T}$ , consists of  $K$  Gaussian ellipsoids that are arranged in the shape of a human. Each skeletal node, denoted as  $\mathbf{T}_k$ , is defined on a regular volumetric grid and represents a single body part. All ellipsoids are parameterized by their mean  $\mu_k$  and co-variance  $\Sigma_k$ . On the volumetric grid we define two functions: a scalar field  $f : \mathbb{R}^3 \rightarrow [0, 1]$  that assigns a value to each point  $(x, y, z)$  on the grid — in the volume rendering literature it is commonly referred to as the occupancy function, and a vector field  $c : \mathbb{R}^3 \rightarrow \mathcal{C} \subset \mathbb{R}^3$  specifying the RGB-colour for each point, commonly referred to as the colouring function.

**Raycasting and Emission Absorption Function:** We make use of an existing implementation of the raycasting algorithm shipped with the PyTorch3D package [25] to render the template image. Given a camera location  $\mathbf{r}_0 \in \mathbb{R}^3$ , rays are “emitted” from  $\mathbf{r}_0$  that pass through each pixel  $\mathbf{u}_i \in \mathbb{R}^3$  lying on a 2D view plane  $\mathcal{S}$  by sampling uniformly spaced points along each ray starting from the intersecting pixel:

$$\mathbf{p}_j = \mathbf{u}_i + j\delta s, \quad (1)$$

where  $j$  is the step and  $\delta s$  the step size that depends on the maximum depth and number of points along each ray.

The colour value at each pixel location  $\mathbf{u}_i$  is then determined by a weighted sum of all colour values of the points sampled along the ray:

$$\mathbf{c}_i = \sum_{j=0}^J w_j \mathbf{c}_j \quad (2)$$

The weights  $w_j$  are computed by multiplying the occupancy function  $f(x)$  with the transmission function  $T(x)$  evaluated at each point  $\mathbf{p}$  along the ray:

$$w_j = f(\mathbf{p}_j) \cdot T(\mathbf{p}_j), \quad (3)$$

where  $T(\mathbf{x})$  can be interpreted as the probability that a given ray is not terminated, *i.e.* fully absorbed, at a given point  $\mathbf{x}$  and is computed as the cumulative product of the complement of the occupancy function of all  $k$  points up until  $\mathbf{p}_j$ :

$$T(\mathbf{x}) = \prod_k (1 - f(\mathbf{x}_k)) \quad (4)$$

Repeating this for all pixels in the view plane results in a 2D projection of our 3D object representing the rendered image  $\mathbf{f}_r \in \mathbb{R}^{3 \times h \times w}$ .

### 3.2 Pose regression and shape transformation

In order to estimate the skeletal pose of a given frame, we use the encoder network,  $\phi : \mathbb{R}^{3 \times h \times w} \rightarrow \mathbb{R}^{3K+3}$  based on the ResNet-34 architecture [1].

The encoder maps a color input image of size  $h \times w$  to  $K$  rotation and scale vectors,  $(\mathbf{r}, \mathbf{s})_{1:K} \in \mathbb{R}^3$ , and a single global translation vector,  $\mathbf{t} \in \mathbb{R}^3$  for the camera.  $K$  denotes the number of transformable parts in the template. Here, rotation is parameterised via axis-angle representation, and is subsequently converted to 3D transformation matrices using the Rodrigues' rotation formula. Combined with the scaling parameter for each axis, the resulting matrix defines the affine mapping, excluding the sheer component, for spatial transformation of each skeletal node.

After construction of the 3D transformation matrix, each Gaussian ellipsoid of the template  $\mathbf{T}$ , with occupancy  $f_k(\mathbf{x})$  and colour field  $c_k(\mathbf{x})$ , gets transformed according to the regressed parameters. Finally, utilising the aforementioned ray-tracing method we render an image based upon our transformed template by summing together all transformed occupancy and colour fields and clipping to a maximum value of 1:

$$\tilde{\mathbf{T}}_k = \Omega_k(\mathbf{T}_k) \quad (5)$$

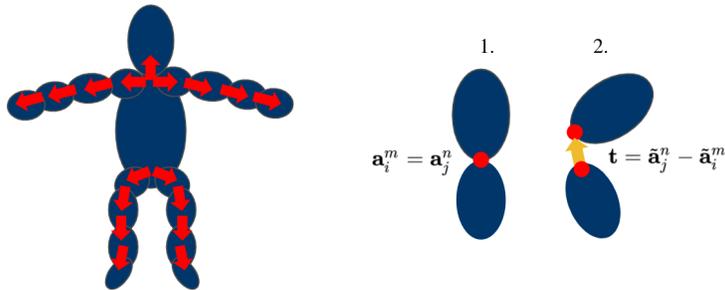
$$\mathbf{f}_r = \mathcal{R}\left(\sum_k f_k(\mathbf{x}), \sum_k c_k(\mathbf{x})\right), \quad (6)$$

,where  $\mathcal{R}$  denotes the rendering operation.

### 3.3 Kinematic chain

Instead of relying on an additional loss to enforce connectivity between body parts as in [28], we define a kinematic chain along which each body part is reconnected to its parent via a translation after rotation and scale have been applied.

Given a parent and child body part with indices  $n$  and  $m$  respectively, we define anchor-points  $\mathbf{a}_i^n$  and  $\mathbf{a}_j^m \in \mathbb{R}^3$  on each part representing the area of overlap in the non-transformed template (see Figure 2 right). If body part  $m$  is being transformed, the position of the anchor point changes:  $\tilde{\mathbf{a}}_j^m = H\mathbf{a}_j^m$ , where  $H \in \mathbb{R}^{3 \times 3}$  specifies a transformation matrix. To ensure continuous connectivity, we



**Fig. 2.** Illustration of the kinematic chain. Red circles denote anchor-points. Following a transformation of the upper part, the translation  $\mathbf{t} = \tilde{\mathbf{a}}_j^n - \tilde{\mathbf{a}}_i^m$  is applied to enforce continuity

apply the transformation for the child body part in an analogous way, it is reconnected with the parent node by applying translation  $\mathbf{t} = \tilde{\mathbf{a}}_j^n - \tilde{\mathbf{a}}_i^m$ . We first transform the core, and then proceed with all other parts in an iterative fashion along the kinematic chain as depicted in Figure 2 left.

### 3.4 Loss function

The loss function is a sum of individual components: the reconstruction loss as the pixel-wise  $l^2$ -norm between the decoder output and the original image and a boundary loss of the form

$$\mathcal{L}_{bx}^i = \begin{cases} |\hat{a}_{x,i}|, & \text{if } |\hat{a}_{x,i}| > 1 \\ 0, & \text{otherwise} \end{cases} \quad \mathcal{L}_{bx} = \sum_i \mathcal{L}_{bx}^i \quad \mathcal{L}_b = \mathcal{L}_{bx} + \mathcal{L}_{by} \quad (7)$$

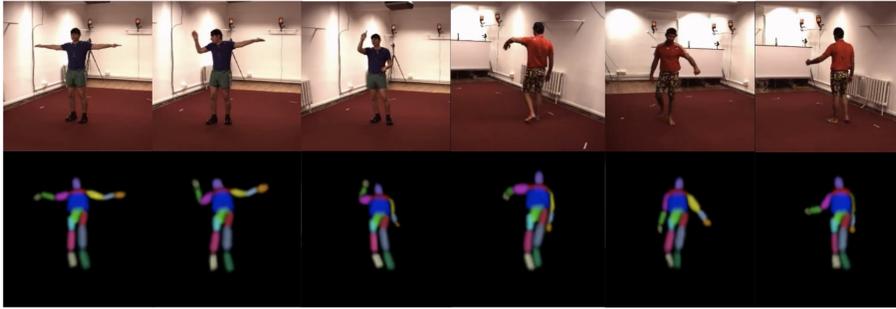
,where  $\hat{a}_{x,i}$  is the x-component of a projected and transformed anchor point. Note that we normalise image coordinates to  $(-1, 1)$ .

We also regularise the pose regression via the  $l^2$ -norm of the rotation vector  $\mathbf{r}$  and decay this term linearly to 0 after 500 iterations. Overall, our objective function is:

$$\mathcal{L} = \mathcal{L}_{\text{recon}} + \mathcal{L}_b + \alpha * \sum_K \|\mathbf{r}_k\|_2, \quad \alpha = \min(1 - 0.02 * \text{iter}) \quad (8)$$

## 4 Experiments

We train and evaluate our model on Human 3.6M [3], a motion-capture dataset including 11 actors performing various activities while being filmed by four different cameras in a studio setting. Following [5] and [28], we train on subjects 1, 5, 6, 7, test on 9 and 11 and restrict activities to mostly upright poses, resulting in roughly 700,000 images for training. We sample video frames in pairs containing the same person in different poses, but with the same background utilising bounding boxes derived from the masks and utilise the Adam optimizer [9].



**Fig. 3.** Results on the two evaluation subjects. Top row: input image. Bottom row: predicted pose in the form of a transformed and rendered 3d shape template.

## 5 Results

We restrict our evaluation to qualitative results in figure 3. These demonstrate that the concept of self-supervision through conditional image translation can be extended to 3D pose estimation. However, there are several issues that still need to be solved: The model is currently not able to distinguish left and right, as can be observed in figure 3 (fourth image from the right), where the subject is facing away from the camera, but the template remains in the front-facing configuration. The model also mostly generates limbs facing away from the camera (third image from the right). We hypothesise that due to depth ambiguity in 2D and limitations in pose variety due to the restricted sampling the decoder can perfectly reconstruct the image despite the wrong orientation of the limb.

## 6 Conclusion

We presented preliminary results for a method to estimate human pose in 3d from monocular images without relying on any landmark labels. Despite issues with depth ambiguity the qualitative results are encouraging and demonstrate the feasibility of combining differentiable rendering techniques and self-supervision. A straightforward improvement would be weak supervision in the form a small labelled dataset. Replacing the image translation task with a purely generative approach with separate fore- and background similarly to [33] might prove to be very successful in extending the approach to non-static backgrounds as well.

**Acknowledgements:** supported by EPSRC EP/S013687/1.

## References

1. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR) pp. 770–778 (2016)
2. Hossain, M.R.I., Little, J.J.: Exploiting temporal information for 3d human pose estimation. In: Ferrari, V., Hebert, M., Sminchisescu, C., Weiss, Y. (eds.) *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part X. Lecture Notes in Computer Science*, vol. 11214, pp. 69–86. Springer (2018)
3. Ionescu, C., Papava, D., Olaru, V., Sminchisescu, C.: Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **36**(7), 1325–1339 (jul 2014)
4. Iqbal, U., Doering, A., Yasin, H., Krüger, B., Weber, A., Gall, J.: A dual-source approach for 3d human pose estimation from a single image. *Comput. Vis. Image Underst.* **172**, 37–49 (2018)
5. Jakab, T., Gupta, A., Bilen, H., Vedaldi, A.: Self-supervised learning of interpretable keypoints from unlabelled videos. 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) pp. 8784–8794 (2020)
6. Jiang, B., Hong, Y., Bao, H., Zhang, J.: Selfrecon: Self reconstruction your digital avatar from monocular video. *ArXiv abs/2201.12792* (2022)
7. Kanazawa, A., Black, M.J., Jacobs, D.W., Malik, J.: End-to-end recovery of human shape and pose. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. pp. 7122–7131. IEEE Computer Society (2018)
8. Katircioglu, I., Tekin, B., Salzmann, M., Lepetit, V., Fua, P.: Learning latent representations of 3d human pose with deep neural networks. *Int. J. Comput. Vis.* **126**(12), 1326–1341 (2018)
9. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. *CoRR abs/1412.6980* (2015)
10. Kocabas, M., Athanasiou, N., Black, M.J.: Vibe: Video inference for human body pose and shape estimation. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (June 2020)
11. Kolotouros, N., Pavlakos, G., Black, M.J., Daniilidis, K.: Learning to reconstruct 3d human pose and shape via model-fitting in the loop. In: *ICCV* (2019)
12. Kundu, J.N., Seth, S., Jampani, V., Rakesh, M., Babu, R.V., Chakraborty, A.: Self-supervised 3d human pose estimation via part guided novel image synthesis. 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) pp. 6151–6161 (2020)
13. Kwon, Y., Kim, D., Ceylan, D., Fuchs, H.: Neural human performer: Learning generalizable radiance fields for human performance rendering. *ArXiv abs/2109.07448* (2021)
14. LeCun, Y., Bengio, Y.: *Convolutional Networks for Images, Speech, and Time Series*, p. 255–258 (1998)
15. Li, Y., Li, K., Jiang, S., Zhang, Z., Huang, C., Xu, R.Y.D.: Geometry-driven self-supervised method for 3d human pose estimation. vol. 34, pp. 11442–11449 (Apr 2020), <https://ojs.aaai.org/index.php/AAAI/article/view/6808>
16. Lin, M., Lin, L., Liang, X., Wang, K., Cheng, H.: Recurrent 3d pose sequence machines. *CoRR abs/1707.09695* (2017), <http://arxiv.org/abs/1707.09695>

17. Lombardi, S., Simon, T., Saragih, J.M., Schwartz, G., Lehrmann, A.M., Sheikh, Y.: Neural volumes. *ACM Transactions on Graphics (TOG)* **38**, 1 – 14 (2019)
18. Loper, M., Mahmood, N., Romero, J., Pons-Moll, G., Black, M.J.: SMPL: A skinned multi-person linear model. *ACM Trans. Graphics (Proc. SIGGRAPH Asia)* **34**(6), 248:1–248:16 (Oct 2015)
19. Martinez, J., Hossain, R., Romero, J., Little, J.J.: A simple yet effective baseline for 3d human pose estimation. In: *Proceedings IEEE International Conference on Computer Vision (ICCV)*. IEEE, Piscataway, NJ, USA (Oct 2017)
20. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. In: *ECCV (2020)*
21. Park, T., Liu, M.Y., Wang, T.C., Zhu, J.Y.: Semantic image synthesis with spatially-adaptive normalization. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* pp. 2332–2341 (2019)
22. Pavlakos, G., Zhou, X., Derpanis, K.G., Daniilidis, K.: Coarse-to-fine volumetric prediction for single-image 3d human pose. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* pp. 1263–1272 (2017)
23. Pavlo, D., Feichtenhofer, C., Grangier, D., Auli, M.: 3d human pose estimation in video with temporal convolutions and semi-supervised training. In: *Conference on Computer Vision and Pattern Recognition (CVPR)* (2019)
24. Qiu, H., Wang, C., Wang, J., Wang, N., Zeng, W.: Cross view fusion for 3d human pose estimation. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)* pp. 4341–4350 (2019)
25. Ravi, N., Reizenstein, J., Novotny, D., Gordon, T., Lo, W.Y., Johnson, J., Gkioxari, G.: Accelerating 3d deep learning with pytorch3d. *arXiv:2007.08501* (2020)
26. Rhodin, H., Salzmann, M., Fua, P.V.: Unsupervised geometry-aware representation for 3d human pose estimation. *ArXiv abs/1804.01110* (2018)
27. Rhodin, H., Spörri, J., Katircioglu, I., Constantin, V., Meyer, F., Müller, E., Salzmann, M., Fua, P.V.: Learning monocular 3d human pose estimation from multi-view images. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition* pp. 8437–8446 (2018)
28. Schmidtke, L., Vlontzos, A., Ellershaw, S., Lukens, A., Arichi, T., Kainz, B.: Unsupervised human pose estimation through transforming shape templates. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)* (2021)
29. Sigal, L., Black, M.J.: Predicting 3d people from 2d pictures. In: Perales, F.J., Fisher, R.B. (eds.) *Articulated Motion and Deformable Objects*. pp. 185–195 (2006)
30. Su, S.Y., Yu, F., Zollhoefer, M., Rhodin, H.: A-nerf: Articulated neural radiance fields for learning human shape, appearance, and pose (2021)
31. Tomè, D., Toso, M., de Agapito, L., Russell, C.: Rethinking pose in 3d: Multi-stage refinement and recovery for markerless motion capture. *2018 International Conference on 3D Vision (3DV)* pp. 474–483 (2018)
32. Xu, J., Yu, Z., Ni, B., Yang, J., Yang, X., Zhang, W.: Deep kinematics analysis for monocular 3d human pose estimation. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (June 2020)
33. Yang, Y., Bilen, H., Zou, Q., Cheung, W.Y., Ji, X.W.: Learning foreground-background segmentation from improved layered gans. *2022 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)* pp. 366–375 (2022)