

Epipolar Transformer for Multi-view Human Pose Estimation

Yihui He* Rui Yan* Katerina Fragkiadaki
Carnegie Mellon University
Pittsburgh, PA 15213

{he2@alumni, ruiyan@alumni, katef@cs}.cmu.edu

Shoou-I Yu
Facebook Reality Labs
Pittsburgh, PA 15213

shoou-i.yu@fb.com

Abstract

A common way to localize 3D human joints in a synchronized and calibrated multi-view setup is a two-step process: (1) apply a 2D detector separately on each view to localize joints in 2D, (2) robust triangulation on 2D detections from each view to acquire the 3D joint locations. However, in step 1, the 2D detector is constrained to solve challenging cases which could be better resolved in 3D, such as occlusions and oblique viewing angles, purely in 2D without leveraging any 3D information. Therefore, we propose the differentiable “epipolar transformer”, which empowers the 2D detector to leverage **3D-aware features** to improve 2D pose estimation. The intuition is: given a 2D location p in the reference view, we would like to first find its corresponding point p' in the source view, then combine the features at p' with the features at p , thus leading to a more 3D-aware feature at p . Inspired by stereo matching, the epipolar transformer leverages epipolar constraints and feature matching to approximate the features at p' . The key advantages of the epipolar transformer are: (1) it has minimal learnable parameters, (2) it can be easily plugged into existing networks, moreover (3) it is interpretable, i.e., we can analyze the location p' to understand whether matching over the epipolar line was successful. Experiments on Human3.6M [9] show that our approach has consistent improvements over the baselines. Specifically, in the condition where no external data is used, our Human3.6M model trained with ResNet-50 and image size 256×256 outperforms state-of-the-art by a large margin and achieves MPJPE **26.9 mm**. Code is available¹. This is the workshop version of our CVPR 2020 paper [8]

1. Introduction

The pose estimation task can be divided into two categories: single-view and multi-view 3D pose estimation.

*Equal contribution

¹github.com/yihui-he/epipolar-transformers

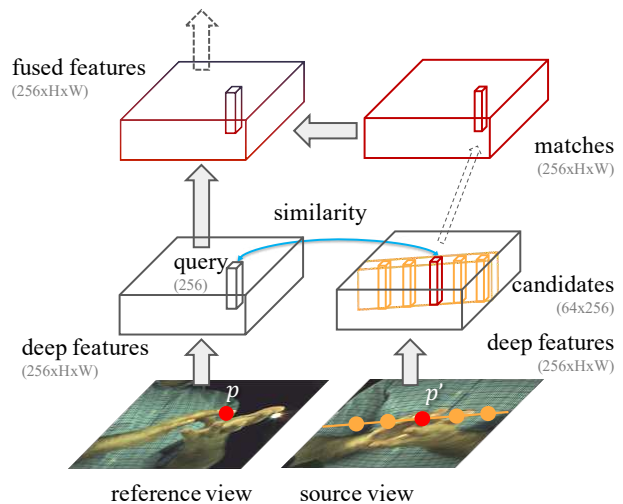


Figure 1: Overview of the proposed epipolar transformer, which enables 2D detectors to leverage 3D-aware features for more accurate keypoint localization. For a query vector (e.g., with length 256) on the deep feature maps ($256 \times H \times W$), we extract K samples along the epipolar line of the source view. Dot-product and softmax are used to compute similarity between the query and sampled vectors to find correspondences. The corresponding features are then fused into the reference view.

Single-view pose estimation [23, 22, 2, 5] directly estimates 3D pose from a monocular image. This is an ill-posed problem due to the ambiguity in depth, which can be alleviated through multi-view pose estimation. This paper focuses on the latter.

Multi-view pose estimation leverages multiple synchronized views with known intrinsic and extrinsic calibration to resolve the depth ambiguity in single-view pose estimation. A common framework [17, 11, 10] to resolve the 3D location of joints follows a two-step process: (1) apply a 2D pose detector on each view separately to localize joints in 2D, and (2) robust triangulation based on camera calibration and 2D detections from each view to acquire the 3D location of joints. Robust triangulation is

required as the prediction of the 2D pose detector could be incorrect or missing due to occlusions. One main disadvantage of this framework is that step 2 only has access to the 2D detections from the 2D detector, which predicts keypoint locations *independently from all other views*. Thus, challenging cases that could potentially be better resolved in 3D, such as occlusions and viewing the scene from oblique angles, are all resolved in 2D by the detector *without leveraging any 3D information*. This could lead to inaccurate detections that are inconsistent in 3D, or the network might require more capacity and training data to resolve these challenging cases.

Therefore, in this paper, we explore the possibility of leveraging 3D information not only on the final 2D detections, but also on the *features* of the 2D detector. The intuition behind our 3D-aware features is shown in Figure 1: given a 2D location p in reference view, we would like to first find its corresponding point p' in source view, then combine the features at p' with the features at p . In this way, the 2D detector can leverage features from neighboring views to compute 3D-aware features, thus empowering the detector to perform 3D reasoning *in the 2D detector itself*, not only during the triangulation phase.

To this end, we propose the “epipolar transformer”, which is a fully differentiable module that takes a feature at p from the reference view, and fuses it with the estimated features at p' from the source view. Inspired by stereo matching, we first leverage the epipolar line generated by p to limit the potential locations of p' . Then, we compute the similarity between the feature at p and features sampled along the epipolar line. Ideally, the feature similarity should be highest when the correct p' is found. However, we do not know where p' is along the epipolar line, so we perform a weighted sum of the features along the line as an approximation of the feature at p' . The weights used for the weighted sum are the feature similarities. Next, given the features at p and p' , we propose multiple methods inspired by [20] to fuse the two features. Finally, we train our network with the epipolar transformer in an end-to-end fashion. The epipolar transformer not only enables features of 2D detectors to be influenced by features from other views, but also potentially promotes features that are coherent across views. Note that the epipolar transformer only operates on the features of the network, so the final output of the 2D detector is still the 2D location of the keypoints.

To evaluate our epipolar transformer, we conduct experiments on Human3.6M [9]. On Human3.6M [9], we achieve **26.9 mm** when using the ResNet-50 backbone on images of resolution 256×256 and trained without external data. This outperforms the state-of-the-art, Qiu *et al.* [16] from ICCV'19 by **4.23 mm**.

The proposed epipolar transformer has multiple

advantages, including (1) can easily be added into existing network architectures, (2) minimal learnable parameters (parameter size is C -by- C , where C is input feature channel size) and (3) is interpretable: one can analyze the feature similarity along the epipolar line to gauge whether matching is successful.

In sum, our contributions are as follows:

1. We propose the epipolar transformer, which is a differentiable module that enables existing 2D pose detectors to gain access to 3D-aware features, thus leading to more accurate predictions.
2. We performed detailed ablation studies to analyze the epipolar transformer, and also understand the effect of different design choices.
3. Experiments show that our proposed model improves upon state-of-the-art on human pose estimation task.

2. The Epipolar Transformer

There are two main components to our epipolar transformer: the epipolar sampler and the feature fusion module. Given a point p in the reference view, the epipolar sampler will, in the source view, compute the locations along the epipolar line from which to sample features. The feature fusion module will then take all the features at the sampled locations in the source view and the feature at p in the reference view to produce a final 3D-aware feature. We now detail each component, and also some implementation details on how to handle image transformations when using the epipolar transformer.

2.1. The Epipolar Sampler

We first define the notations used to describe the epipolar sampler. Given two images captured at the same time but from different views, namely, reference view \mathcal{I} and source view \mathcal{I}' , we denote their projection matrices as $M, M' \in \mathbb{R}^{3 \times 4}$ and camera centers as $C, C' \in \mathbb{R}^4$ in homogeneous coordinates respectively. As illustrated in Figure 1, assuming the camera centers do not overlap, the epipolar line l corresponding to a given query pixel $p = (x, y, 1)$ in \mathcal{I} can be deterministically located on \mathcal{I}' as follows [6].

$$l = [M'C]_{\times} M' M^+ p, \quad (1)$$

where M^+ is the pseudo inverse of M , and $[\cdot]_{\times}$ represents the skew symmetric matrix. p 's correspondence p' should lie on the epipolar line, i.e., $l^T p' = 0$.

The epipolar sampler \mathcal{S} uniformly samples K locations (e.g., 64) on the epipolar line l of the source view, thus forming a set \mathcal{P}' of cardinality K . The function takes as input the query location p on the reference view, and the projection matrices M, M' as shown below.

$$\mathcal{P}' = \mathcal{S}_K(p, M, M') \quad (2)$$

For query points whose epipolar line do not intersect with the source view image plane, we simply skip them.

2.2. Feature Fusion Module

Ideally, if we knew the ground-truth p' in the source view that corresponds to p in the reference view, then all we need to do is sample the feature at p' : $F_{\text{source}}(p')$, and then combine it with the feature at p : $F_{\text{reference}}(p)$. However, we do not know the correct p' . Therefore, inspired by Transformer [19] and non-local network [20], we approximate $F_{\text{source}}(p')$ by a weighted sum of all the features along the epipolar line as follows:

$$\bar{F}_{\text{source}}(p) = \sum_{p' \in \mathcal{P}'} \text{sim}(p, p') F_{\text{source}}(p'), \quad (3)$$

where the pairwise function $\text{sim}(\cdot, \cdot)$ computes the similarity score between two vectors. More specifically, it is the dot-product followed by softmax.

Once we have the feature from the source view: $\bar{F}_{\text{source}}(p)$, we can fuse it with the feature in the reference view: $F_{\text{reference}}(p)$ as follows.

$$F_{\text{fused}}(p) = F_{\text{reference}}(p) + W_z(\bar{F}_{\text{source}}(p)) \quad (4)$$

Note that the output F_{fused} is of the same shape as the input $F_{\text{reference}}$, thus this property enables us to insert the epipolar transformer module into different stages of many existing networks. The weights W_z can be as simple as a 1×1 convolution. In this case, keeping a copy of the original $F_{\text{reference}}$ feature is similar to the design of the residual block.

We explored some other architectures for the feature fusion module, as shown in Figure 2. The one we have just described in Equation 3 and Equation 4 corresponds to Figure 2 (b): the *identity Gaussian* architecture. This architecture is simpler as it only has one learnable convolution layer. Figure 2 (a) is the *bottleneck embedded Gaussian* popularized by non-local network [20]. The reference view and source view are fed into the embedded Gaussian kernel, where the input is down-sampled by two, and the output is up-sampled by two, so that the shape of the fused feature still matches the input's shape.

2.3. Dealing with Image Transformations

As the epipolar transformer relies on camera calibration, any spatial transformations made to the image also need to be reflected in the calibration parameters. More details are as follows.

Data Augmentation: Data augmentation like rotation, scaling and cropping can still be performed with the epipolar transformer. The projection matrix needs to be updated accordingly when the image is transformed with

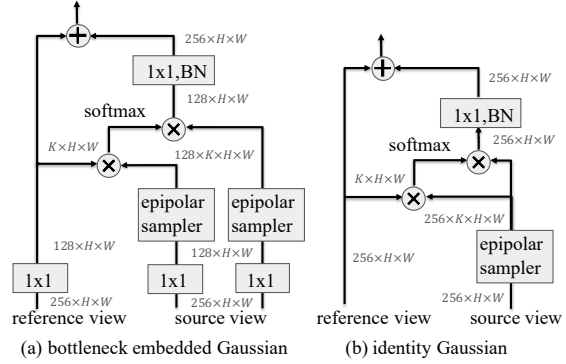


Figure 2: Different feature fusion module architectures. The feature maps are shown as the shape of their tensors, e.g., $256 \times H \times W$ for 256 channels. " \oplus " and " \otimes " denote element-wise sum and matrix multiplication respectively.

an affine transformation parameterized by $A \in \mathbb{R}^{2 \times 2}$ and $b \in \mathbb{R}^2$:

$$M := \begin{bmatrix} A & b \\ \mathbf{0}^T & 1 \end{bmatrix} M \quad (5)$$

Different scaling and cropping parameters can be applied separately to the reference view and source view.

Scaling of projection matrices: Special care is required when we scale the projection matrices due to image resizing or pooling. Suppose the input image is spatially down-sampled s_x and s_y times along the x-axis and y-axis (e.g., $s_x = s_y = 4$ in the hourglass network), the projection matrix is updated as follows:

$$M := \begin{bmatrix} 1/s_x & 0 & (1 - s_x)/2s_x \\ 0 & 1/s_y & (1 - s_y)/2s_y \\ 0 & 0 & 1 \end{bmatrix} M \quad (6)$$

The coordinates are aligned with the center of pixels rather than the top-left corners, which is important for extracting features at precise locations in the epipolar transformer.

3. Experiments

We conducted experiments on a publicly available dataset, Human3.6M [9]. We adopt the same training and testing sets as in [16], where subjects 1, 5, 6, 7, 8 are used for training, and 9, 11 are for testing.

As there are only four views in Human3.6M [9], we choose the closest view as source view. We adopt ResNet-50 with image resolution 256×256 proposed in simple baselines for human pose estimation [21] as our backbone network. We use the ImageNet [4] pre-trained model [14] for initialization. The networks are trained for 20 epochs with batch size 16 and Adam optimizer [13]. Learning rate decays at 10 and 15 epochs. Unless specified, we do

	Net	scale	shlder	elb	wri	hip	knee	ankle	root	belly	neck	nose	head	Avg
-	R152	320	88.50	88.94	85.72	90.37	94.04	90.11	-	-	-	-	-	-
sum over epipolar line[16]	R152	320	91.36	91.23	89.63	96.19	94.14	90.38	-	-	-	-	-	-
max over epipolar line[16]	R152	320	92.67	92.45	91.57	97.69	95.01	91.88	-	-	-	-	-	-
cross-view fusion [16]	R152	320	95.58	95.83	95.01	99.36	97.96	94.75	-	-	-	-	-	-
cross-view fusion [16]*	R50	320	95.6	95.0	93.7	96.6	95.5	92.8	96.7	96.4	96.5	96.4	96.2	95.9
cross-view fusion [16]*	R50	256	86.1	86.5	82.4	96.7	91.5	79.0	100.0	94.1	93.7	95.4	95.5	95.1
epipolar transformer	R50	256	96.44	94.16	92.16	98.95	97.26	96.62	99.89	99.86	99.68	99.78	99.63	97.01
epipolar transformer⁺	R50	256	97.71	97.34	94.85	99.77	98.32	97.55	99.99	99.99	99.76	99.74	99.54	98.25

Table 1: 2D pose estimation accuracy comparison on the Human3.6M [9], where no external training data is used. The metric is joint detection rate, JDR (%). +: indicates using data augmentation mentioned in Section 2.3. "-": We cite numbers from [16] and these entries are absent. *: We trained the models using released code [16]. R50 and R152 are ResNet-50 and ResNet-152 [7] respectively. Scale is the input resolution of the network.

MPJPE (mm)	Dir	Disc	Eat	Greet	Phone	Photo	Pose	Purch	Sit	SitD	Smoke	Wait	WalkD	WalkT	Avg	
Multi-View Martinez [18]	46.5	48.6	54.0	51.5	67.5	70.7	48.5	49.1	69.8	79.4	57.8	53.1	56.7	42.2	45.4	57.0
Pavlakos <i>et al.</i> [15]	41.2	49.2	42.8	43.4	55.6	46.9	40.3	63.7	97.6	119.0	52.1	42.7	51.9	41.8	39.4	56.9
Tome <i>et al.</i> [18]	43.3	49.6	42.0	48.8	51.1	64.3	40.3	43.3	66.0	95.2	50.2	52.2	51.1	43.9	45.3	52.8
Kadkhodamohammadi & Padoy [12]	39.4	46.9	41.0	42.7	53.6	54.8	41.4	50.0	59.9	78.8	49.8	46.2	51.1	40.5	41.0	49.1
R50 256×256+triangulate	38.9	46.1	36.2	59.7	46.4	44.7	44.9	37.7	51.2	72.0	48.2	61.0	46.2	45.7	52.0	48.7
R50 256×256+crossview+triangulate[16]	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	45.5
R50 256×256+ours+triangulate	30.6	33.2	26.7	28.2	32.8	38.4	29.3	28.9	36.6	45.2	34.3	31.7	33.1	34.8	31.2	33.1
R50 256×256+ours+triangulate⁺	29.0	30.6	27.4	26.4	31.0	31.8	26.4	28.7	34.2	42.6	32.4	29.3	27.0	29.3	25.9	30.4
R50 256×256+crossview+RPSM [16]	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	41.2
R50 256×256+ours+RPSM [16]	25.7	27.7	23.7	24.8	26.9	31.4	24.9	26.5	28.8	31.7	28.2	26.4	23.6	28.3	23.5	26.9
R152 320×320+crossview+triangulate[16]	34.8	35.8	32.7	33.5	34.4	38.2	29.7	60.7	53.1	35.2	41.0	41.6	31.9	31.4	34.6	38.3
R152 320×320+crossview+RPSM [16]	28.9	32.5	26.6	28.1	28.3	29.3	28.0	36.8	42.0	30.5	35.6	30.0	29.3	30.0	30.5	31.2

Table 2: Comparison with state-of-the-art multi-view keypoint estimation methods on Human3.6M [9], where no external training data is used. The metric is MPJPE (mm). "+": rotation and scaling augmentation. "-": models trained using released code [16], where the per action MPJPE evaluation were not provided.

not use data augmentation for fair comparisons, following Qiu *et al.* [16]. We follow Qiu *et al.* [16] for other hyper-parameters. Following [16], as there are only four cameras in this dataset, direct linear transformation (DLT) is used for triangulation (Hartley & Zisserman [6], p.312), instead of RANSAC which also needs tuning the inlier/outlier threshold.

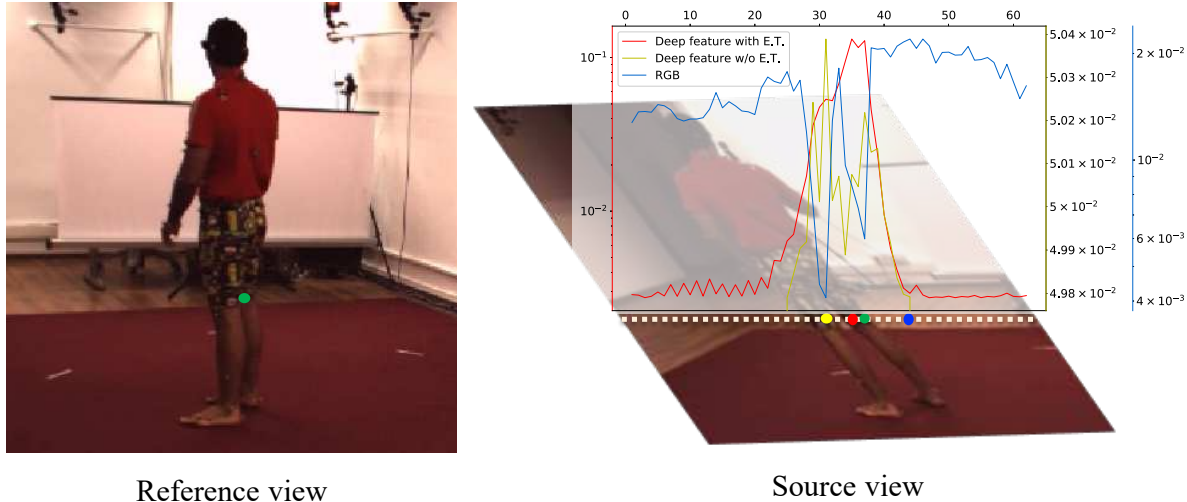
2D Pose Estimation: Following Qiu *et al.* [16], the 2D pose estimation accuracy is measured by Joint Detection Rate (JDR), which measures the percentage of the successfully detected joints. A joint is detected if the distance between the estimated location and the ground truth is smaller than half of the head size [1]. Shown in Table 1, sum or max over the epipolar line does not perform well. The cross-view fusion [16] achieved better performance by fusing with learned global attention, which resembles the non-local network [20]. Epipolar transformer instead attends the features along the epipolar line and fuse them. Using the same backbone ResNet-50 image size 256×256, epipolar transformer achieves

97.01% JDR, which outperforms 95.9% JDR from Qiu *et al.* [16] (ICCV'19) by a large margin. The improvement consolidates the idea that fusing along the epipolar line is better than fusing globally.

We apply data augmentation as is mentioned in Section 2.3, which consists of random scales drawn from a truncated normal distribution $TN(1, 0.25^2, 0.75, 1.25)$ and random rotations from $TN(0^\circ, (30^\circ)^2, -60^\circ, 60^\circ)$ [21]. JDR is further improved to **98.25%** JDR.

Effect of the number of views: Shown in Figure 4, compared with ICCV'19 cross-view [16], epipolar transformer still have better performance when there are fewer views. This shows that epipolar transformer efficiently fuses features from other views.

Compare with state-of-the-art: Table 2 demonstrates state-of-the-art multi-view keypoint estimation methods. Our epipolar transformer outperforms the state-of-the-art by a large margin. Specifically, using triangulation for estimating 3D human poses, epipolar transformer achieves **33.1 mm**, which is \sim **12 mm** better than



(i) Right knee selected, denoted in green.

Figure 3: Visualizations of the matching results along the epipolar line in **more difficult cases** in InterHand. We here use E.T. as a shorthand for Epipolar Transformer. The compared features are (a) deep features learned through the epipolar transformer (deep features with E.T., denoted in red), (b) deep feature learned by ResNet-50 [7] without epipolar transformer (deep features w/o E.T., denoted in yellow), and (c) RGB features (denoted in blue). Green dot on the reference view is the selected joint, and the green dot on the source view is the corresponding point offered by the groundtruth.

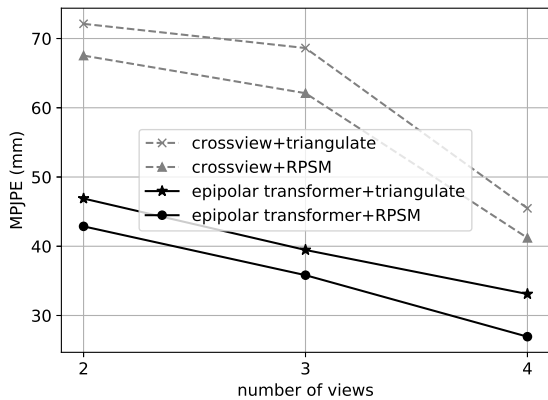


Figure 4: MPJPE by varying the number of views on Human3.6M

the cross-view [16], using the same backbone network (ResNet-50, input size 256×256). Using the recursive pictorial structural model (RPSM [16]) for estimating 3D poses, our epipolar transformer achieves **26.9 mm**, which is ~ 14 mm better than cross-view [16]. More importantly, epipolar transformer on ResNet-50 input size 256×256 even surpasses the state-of-the-art result from cross-view [16] on ResNet-152 input size 320×320 by ~ 4 mm, which is **13%** relative improvement. We argue that epipolar transformer find correspondences and fuse features

based on feature similarity, which is superior than cross-view [16] which use fixed attention for specific cameras settings.

Our model with data augmentation achieves **MPJPE 30.4 mm** with triangulation, which is better than state-of-the-art even without RPSM.

Visualization: As shown in Figure 3, our predictions with epipolar transformer (red dot) are closer to the ground truth points, compared to the features without the awareness of the multi-view information.

References

- [1] Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele. 2d human pose estimation: New benchmark and state of the art analysis. In *Proceedings of the IEEE Conference on computer Vision and Pattern Recognition*, pages 3686–3693, 2014. 4
- [2] Adnane Boukhayma, Rodrigo de Bem, and Philip HS Torr. 3d hand shape and pose from images in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 10843–10852, 2019. 1
- [3] Cristian Sminchisescu Catalin Ionescu, Fuxin Li. Latent structured models for human pose estimation. In *International Conference on Computer Vision*, 2011.
- [4] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Computer Vision and Pattern Recognition*,

2009. *CVPR 2009. IEEE Conference on*, pages 248–255. IEEE, 2009. 3
- [5] Lihao Ge, Zhou Ren, Yuncheng Li, Zehao Xue, Yingying Wang, Jianfei Cai, and Junsong Yuan. 3d hand shape and pose estimation from a single rgb image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 10833–10842, 2019. 1
- [6] Richard Hartley and Andrew Zisserman. *Multiple view geometry in computer vision*. Cambridge university press, 2003. 2, 4
- [7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016. 4, 5
- [8] Yihui He, Rui Yan, Shou-I Yu, and Katerina Fragkiadaki. Epipolar transformer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020. 1
- [9] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. 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, 2013. 1, 2, 3, 4
- [10] Karim Iskakov, Egor Burkov, Victor Lempitsky, and Yuri Malkov. Learnable triangulation of human pose. *arXiv preprint arXiv:1905.05754*, 2019. 1
- [11] Yasamin Jafarian, Yuan Yao, and Hyun Soo Park. Monet: Multiview semi-supervised keypoint via epipolar divergence. *arXiv preprint arXiv:1806.00104*, 2018. 1
- [12] Abdolrahim Kadkhodamohammadi and Nicolas Padoy. A generalizable approach for multi-view 3d human pose regression. *arXiv preprint arXiv:1804.10462*, 2018. 4
- [13] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 3
- [14] Adam Paszke, Sam Gross, Soumith Chintala, and Gregory Chanan. Pytorch, 2017. 3
- [15] Georgios Pavlakos, XiaoWei Zhou, Konstantinos G Derpanis, and Kostas Daniilidis. Harvesting multiple views for marker-less 3d human pose annotations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6988–6997, 2017. 4
- [16] Haibo Qiu, Chunyu Wang, Jingdong Wang, Naiyan Wang, and Wenjun Zeng. Cross view fusion for 3d human pose estimation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4342–4351, 2019. 2, 3, 4, 5
- [17] Tomas Simon, Hanbyul Joo, Iain Matthews, and Yaser Sheikh. Hand keypoint detection in single images using multiview bootstrapping. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 1145–1153, 2017. 1
- [18] Denis Tome, Matteo Toso, Lourdes Agapito, and Chris Russell. Rethinking pose in 3d: Multi-stage refinement and recovery for markerless motion capture. In *2018 International Conference on 3D Vision (3DV)*, pages 474–483. IEEE, 2018. 4
- [19] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017. 3
- [20] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7794–7803, 2018. 2, 3, 4
- [21] Bin Xiao, Haiping Wu, and Yichen Wei. Simple baselines for human pose estimation and tracking. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 466–481, 2018. 3, 4
- [22] Xiong Zhang, Qiang Li, Hong Mo, Wenbo Zhang, and Wen Zheng. End-to-end hand mesh recovery from a monocular rgb image. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2354–2364, 2019. 1
- [23] Christian Zimmermann and Thomas Brox. Learning to estimate 3d hand pose from single rgb images. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4903–4911, 2017. 1