

Leveraging Temporal Joint Depths for Improving 3D Human Pose Estimation in Video

Naoki Kato

Mobility Technologies Co., Ltd.
naoki.kato@mo-t.com

Hiroto Honda

Mobility Technologies Co., Ltd.
hiroto.honda@mo-t.com

Yusuke Uchida

Mobility Technologies Co., Ltd.
yusuke.uchida@mo-t.com

Abstract—The effectiveness of the approaches to predict 3D poses from 2D poses estimated in each frame of a video has been demonstrated for 3D human pose estimation. However, 2D poses without appearance information of persons have much ambiguity with respect to the joint depths. In this paper, we propose to estimate a 3D pose in each frame of a video and refine it considering temporal information. The proposed approach reduces the ambiguity of the joint depths and improves the 3D pose estimation accuracy.

Index Terms—video analysis, 3D human pose estimation

I. INTRODUCTION

3D human pose estimation aims to localize human joints in a 3D coordinate system. It is important for many applications including AR/VR, human action analysis, and in-vehicle camera video analysis. One of the difficulties of this task is that there is ambiguity in the estimation of the depth of the joints from an image. In recent years, there has been an increasing number of studies on the use of temporal information in video. The most common approach in those studies is to utilize 2D poses estimated in each frame to predict final 3D pose [1]–[3].

However, since 2D pose does not include the appearance information, there is large ambiguity with respect to the depths of the joints. The experiments by Pavllo et al. [2] suggest that estimating the 3D pose from the 2D pose is likely to have limited accuracy, even if the ground-truth 2D poses are used.

In this paper, we propose a novel 3D pose estimation approach that first estimates a 3D pose for each frame (first stage) and then aggregates the multi-frame predictions for estimating the final 3D poses (second stage). This method is based on the assumption that the joint depths can be estimated more correctly from an image, which has more information than 2D pose. Unlike using 2D poses intermediately, it is possible to perform inferences that do not lose the joint depth information. We demonstrate the effectiveness of the proposed method through experiments on a public dataset.

II. PROPOSED METHOD

We show a schematic diagram of the proposed method in Figure 1. Given a video $\{I_t\}$ ($t \in \{1, \dots, T\}$) of length T containing a human image $I_t \in \mathbb{R}^{H \times W \times 3}$ of size $H \times W$ at frame t , our objective is to estimate 3D coordinates of the human joints $y_t \in \mathbb{R}^{3J}$ for each frame where J is the number of the joints. The target 3D coordinates of the joints are defined as relative coordinates from a root joint (usually hip is employed [2], [4]).

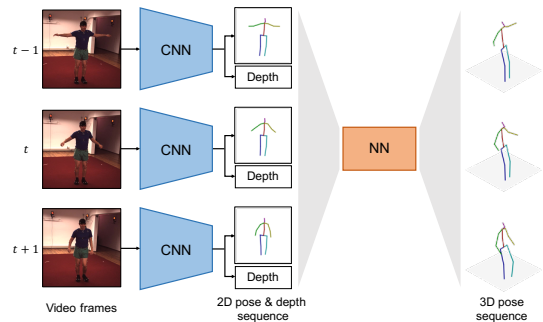


Fig. 1. Overview of our approach.

One of the existing approaches for estimating 3D poses in video is to estimate the 3D poses from the 2D pose sequences estimated at each frame [1]–[3]. This simplifies the task and the effectiveness of exploiting temporal information has been demonstrated. However, since 2D poses used as intermediate representations of human poses in this approach lack the depth information of the joints, 3D pose estimation from 2D pose has the large ambiguity of joint depths. To resolve the issue, it is a natural choice to estimate the joint depths from images which have rich person appearance information.

Based on the above considerations, we propose to use the depths of joints as an intermediate representation in the framework of the above two-stage method. The first stage predicts a 3D pose $x_t \in \mathbb{R}^{3J}$ comprised of image coordinates and depths of joints in each frame I_t , then the second stage predicts the final 3D pose y_t from the 3D pose sequences $\{x_t\}$ ($t \in \{1, \dots, T\}$) estimated by the first stage. This approach can be regarded as refining 3D poses using temporal information. The advantage of the method is that the joint depth information is not lost in the middle of the inference, unlike the approach which uses only 2D poses as an intermediate representation.

First stage. We employ Integral Regression [5] for the first stage, which predicts 3D pose in each frame. Given a human image, this model outputs a 3D heatmap of size $w \times h \times d$ for each joint. Here, the x - y grid uniformly discretize the image coordinates and the z grid uniformly discretize $[-D/2, D/2]$ mm centered at the root joint. The final estimates of the joint coordinates are computed from the centroids of the heatmaps. The model can be trained using both 2D and 3D datasets. By utilizing a 2D dataset which is diverse in appearance, we can mitigate overfitting which is problematic when training

TABLE I
IMPACT ON MPJPE BY INPUT AND RECEPTIVE FIELD.

Input	Receptive field			
	1	27	81	243
2D pose	49.5	48.8	47.7	47.6
2D pose + depth	50.4	48.5	48.4	47.8
2D pose + depth ($\sigma = 0.1$)	48.4	46.4	45.9	45.6
GT 2D pose	38.7	37.3	-	36.3
GT 2D pose + GT depth	20.2	13.4	-	15.9

on 3D datasets. We use ResNet-50 [6] as the backbone and apply three deconvolution layers of kernel size 4 to output 3D heatmaps where the number of channels for the z axis is $d = 72$. We set $D = 1500$ mm.

Second stage. We employ 1D ConvNet as the second stage, which is effective for exploiting 2D pose sequences [2]. The input layer is a temporal convolution with kernel size W and output channels C . This is followed by B residual blocks, each of which is composed of a 1D convolution with kernel size W , dilation $D = W^B$ and another convolution with kernel size 1. The output layer is a convolution with kernel size 1, which outputs an estimated pose with $3J$ channels for each frame. All the convolutional layers except the last layer are followed by batch normalization, ReLU and a dropout layer with a dropout rate p . The temporal receptive fields of this model is W^{B+1} . Unless otherwise noted, we set $W = 3$, $B = 4$, $C = 1024$, and $p = 0.25$.

We use the predictions of the first stage to train the second stage. In this case, the input of the second stage at training time is the predictions from the model which has been supervised by the training data for the first stage. However at test time, the input is the predictions on unseen data, so there is a concern that the distribution of the input data at training and test time may deviate and the performance of the model may be degraded. Therefore, we apply data augmentation that adds Gaussian noise to the input 2D poses and depths during training of the second stage, to reproduce the test-time predictions of the first stage to enhance model performance.

III. EXPERIMENTS

Datasets and Evaluation Metrics. We conduct experiments using Human3.6M [7]. Following the standard evaluation protocol [2], [4], [5], we use 5 subjects (S1, S5, S6, S7, S8) for training and 2 subjects (S9, S11) for evaluation. Along with Human3.6M, we use 2D datasets of COCO [8] and MPII [9] to train Integral Regression.

We consider two evaluation protocols commonly used for the evaluation on Human3.6M. Protocol 1 is Mean Per Joint Position Error (MPJPE) calculated by averaging the distances between the predicted and ground-truth coordinates. Protocol 2 is also MPJPE but after alignment with the ground-truth in translation, rotation, and scale.

Results. The evaluation results for different input types and temporal receptive fields of the second-stage model are shown in Table I. The models with receptive fields of 1, 27, 81, and 243 in the table have (W, B) of $(1, 2)$, $(3, 2)$, $(3, 3)$ and $(3, 4)$ respectively. When joint depths are directly added to the inputs

TABLE II
THE EVALUATION RESULTS ON HUMAN3.6M DATASET.

Method	Protocol 1	Protocol 2
Martinez et al. [4]	62.9	47.7
Sun et al. [5]	49.6	40.6
Cai et al. [3]	48.8	39.0
Pavlo et al. [2]	46.8	36.5
Ours	45.6	34.8

without data augmentation, the error is comparable to the result when only 2D poses are used as input. This may be due to overfitting caused by a discrepancy between the distribution of depths during training and test time. On the other hand, the model trained with Gaussian noise on depths resulted in lower errors in all receptive fields than 2D poses. This result demonstrates the effectiveness of using depths as inputs of the second stage with an appropriate data augmentation.

Our implementation of Integral Regression has MPJPE of 48.9 mm at the first stage, and the errors are reduced by up to 3.3 mm at the second stage, showing that the prediction results of a single-frame 3D pose estimator can be properly refined by utilizing temporal information.

The lowest error when using the ground-truth 2D pose as input is 36.3 mm, showing that the error of the system is greatly reduced by the accurate 2D pose. Adding ground-truth depths further reduce errors by 63%, indicating that exploiting the depths significantly reduces the lower bound of errors.

The evaluation results of our method and existing monocular methods are shown in Table II. For both protocols, our approach outperforms all the comparative methods.

IV. CONCLUSION

In this paper we proposed a two-stage 3D pose estimation pipeline in video that uses a joint depth sequence as an intermediate representation for the human pose in addition to a 2D pose sequence. In the evaluation experiments, we observe that adding depth to the input of the second stage reduces the 3D joint localization error, indicating that our pipeline appropriately refine 3D poses leveraging temporal information.

REFERENCES

- [1] M. Rayat Imtiaz Hossain and J. J. Little, "Exploiting temporal information for 3d human pose estimation," in *ECCV*, 2018.
- [2] D. Pavlo, C. Feichtenhofer, D. Grangier, and M. Auli, "3d human pose estimation in video with temporal convolutions and semi-supervised training," in *CVPR*, 2019.
- [3] Y. Cai, L. Ge, J. Liu, J. Cai, T.-J. Cham, J. Yuan, and N. M. Thalmann, "Exploiting spatial-temporal relationships for 3d pose estimation via graph convolutional networks," in *ICCV*, 2019.
- [4] J. Martinez, R. Hossain, J. Romero, and J. J. Little, "A simple yet effective baseline for 3d human pose estimation," in *CVPR*, 2017.
- [5] X. Sun, B. Xiao, F. Wei, S. Liang, and Y. Wei, "Integral human pose regression," in *ECCV*, 2018.
- [6] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016.
- [7] C. Ionescu, D. Papava, V. Olaru, and C. Sminchisescu, "Human3. 6m: Large scale datasets and predictive methods for 3d human sensing in natural environments," *TPAMI*, 2013.
- [8] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *ECCV*, 2014.
- [9] M. Andriluka, L. Pishchulin, P. Gehler, and B. Schiele, "2d human pose estimation: New benchmark and state of the art analysis," in *CVPR*, 2014.