

VoxelPose: Towards Multi-Camera 3D Human Pose Estimation in Wild Environment

Hanyue Tu^{1,2*}, Chunyu Wang¹, and Wenjun Zeng¹

¹ Microsoft Research Asia

{chnuwa,wezeg}@microsoft.com

² University of Science and Technology of China

tuhanyue@mail.ustc.edu.cn

Abstract. We present *VoxelPose* to estimate 3D poses of multiple people from multiple camera views. In contrast to the previous efforts which require to establish cross-view correspondence based on noisy and incomplete 2D pose estimates, *VoxelPose* directly operates in the 3D space therefore avoids making incorrect decisions in each camera view. To achieve this goal, features in all camera views are aggregated in the 3D voxel space and fed into *Cuboid Proposal Network* (CPN) to localize all people. Then we propose *Pose Regression Network* (PRN) to estimate a detailed 3D pose for each proposal. The approach is robust to occlusion which occurs frequently in practice. Without bells and whistles, it outperforms the previous methods on several public datasets. The code is available at <https://github.com/microsoft/voxelpose-pytorch>

Keywords: 3D Human Pose Estimation

1 Introduction

Estimating 3D human pose from multiple cameras separated by wide baselines [1,2,3,4,5,6,7] has been a longstanding problem in computer vision. The goal is to predict 3D positions of the landmark joints for all people in a scene. The successful resolution of the task can benefit many applications such as intelligent sports [5] and retail analysis.

The previous works such as [4,5] propose to address the problem in three steps. They first estimate 2D poses in each camera view independently, for example, by Convolutional Neural Networks (CNN) [8,9]. Then, in the second step, the poses that correspond to the same person in different views are grouped into clusters according to appearance and geometry cues. The final step is to estimate a 3D pose for each person (*i.e.* each cluster) by standard methods such as triangulation [10] or pictorial structure models [11].

In spite of the fact that 2D pose estimation has quickly matured due to the development of CNN models [12,13], the estimation results are still unsatisfactory for challenging cases especially when occlusion occurs which is often

* This work is done when Hanyue Tu is an intern at Microsoft Research Asia.

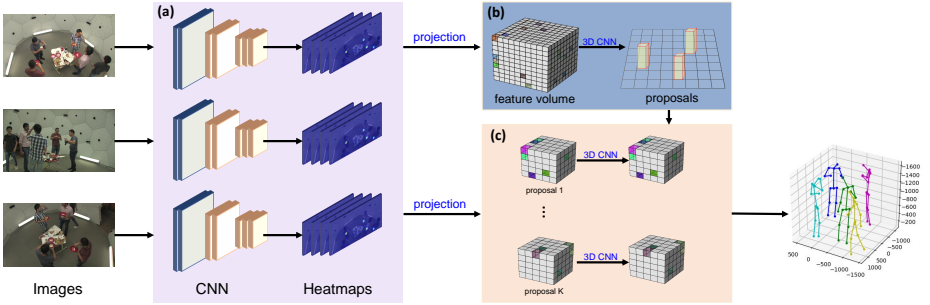


Fig. 1: Overview of our approach. It consists of three modules: (a) we first estimate 2D pose heatmaps for all views; (b) we warp the heatmaps to a common 3D space and construct a feature volume which is fed into a Cuboid Proposal Network to localize all people instances; (c) for each proposal, we construct a finer-grained feature volume and estimate a 3D pose.

the case for natural scenes. See Figure 5 for some example 2D poses estimated by the state-of-the-art method [13]. In addition, it is very difficult to establish cross-view correspondence when 2D poses are inaccurate. All these pose a serious challenge for 3D pose estimation in the wild.

To avoid making incorrect decisions for each camera view, we propose a 3D pose estimator which directly operates in the 3D space by gathering information from all camera views. Figure 1 shows an overview of our approach. It first estimates 2D heatmaps for each view to encode per-pixel likelihood of all joints as shown in Figure 1 (a). Different from the previous works, we do not determine the locations of joints (*e.g.*, by finding the maximum response) nor group the joints into different instances because estimated heatmaps are usually very noisy and incomplete. Instead, we project the heatmaps of all views to a common 3D space as in [6] and obtain a more complete feature volume which allows us to accurately estimate the 3D positions of all joints.

We first present Cuboid Proposal Network (CPN), as shown in Figure 1 (b), to coarsely localize all people in the scene by predicting a number of 3D cuboid proposals from the 3D feature volume. Then for each proposal, we construct a separate *finer-grained* feature volume centered at each proposal, and feed it into a Pose Regression Network (PRN) to estimate a detailed 3D pose. See Figure 1 (c) for illustration. The two networks are composed of basic 3D convolution layers and can be jointly trained.

It is worth noting that our approach implicitly accomplishes two types of association which was previously addressed by post-processing methods. Firstly, the joints of the same person *in a single camera view* are implicitly associated by the cuboid proposal. This was previously addressed in the 2D space either by bottom-up approaches [8,14] or by top-down approaches [13] which would suffer when occlusion occurs. Secondly, the joints that correspond to the same person *in different camera views* are also implicitly associated based on the fact

that the 2D poses which overlap with the projections of a 3D pose belong to the same person. Our approach allows us to avoid the challenging association tasks therefore significantly improves the robustness.

We evaluate our approach on three public datasets including Campus [2], Shelf [2] and CMU Panoptic [15]. It outperforms the state-of-the-arts on the first two datasets. Since no work has reported numerical results on Panoptic, we conduct a series of ablation studies by comparing our approach to several baselines. In addition, we find that CPN and PRN can be accurately trained on automatically generated synthetic heatmaps. They achieve similar results as the models trained on realistic images. This is possible mainly because the heatmap based 3D feature volume representation is a high level abstraction that is disentangled from appearance/lighting, etc. This favorable property dramatically enhances the practical values of the approach.

2 Related Work

In this section, we briefly review the related works on 3D pose estimation for single and multiple people scenarios, respectively. We discuss their main difference from our work and summarize our contributions.

2.1 Single Person 3D Pose Estimation

We briefly classify the existing works into *analytical* and *predictive* approaches based on whether they have learnable parameters. Analytical methods [16,6,11,17] explicitly model the relationship between a 2D and 3D pose according to the camera geometry. **On one hand**, when multiple cameras are available, the 3D pose can be fully determined by simple geometry methods such as triangulation [10] based on the 2D poses in each view. So the bottleneck lies in the inaccuracy of 2D pose estimations. Some works [6,11] propose to model the conditional dependence between the joints and jointly infer their 3D positions to improve their robustness to errors in 2D poses. **On the other hand**, when only one camera is available, the problem is under-determined because multiple 3D poses may correspond to the same 2D pose. The previous works [17,16,18] propose to use low-dimensional pose representations to reduce ambiguities. They optimize the low-dimensional parameters of the representation such that the discrepancy between its projection and the 2D pose is minimized. The improvement in 2D pose estimation has boosted 3D accuracy.

The predictive models [19,20,21,22,23,24,25,26,27] are mainly proposed for the single camera setup aiming to address the ambiguity issue by powerful neural networks. The pioneer works [20,21] propose to regress 3D pose from 2D joint locations by various networks. Some recent works [6,28,19,29,26] also propose to regress a volumetric 3D pose representation from images. In particular, in [6,26], the authors project 2D features or pose heatmaps to 3D space and estimate 3D positions of the body joints. The approach achieves better performance than the triangulation and pictorial structure models on *single* person pose estimation.

However, it requires to address the challenging association problem in order to apply to scenes with multiple people.

2.2 Multiple Person 3D Pose Estimation

There are two challenging association problems in this task. First, it needs to associate the joints of the same person by either top-down [30,9] or bottom-up [8,31,14] strategies. Second, it needs to associate the 2D poses of the same person in different views based on appearance features [4,5] which are unstable when people are occluded. The pictorial structure model is extended to deal with multiple people in [2,1]. The number of people is assumed to be known which is difficult by itself. Besides, the interactions between different people introduce loops into the graph model and complicate the optimization problem. These challenges limit the 3D pose estimation accuracy.

Our work differs from the previous methods [4,5] in that it elegantly avoids the two challenging association problems. This is because a 3D cuboid proposal already naturally associates the joints of the same person in the same and different views by projecting the proposals to image space. Different from the pictorial structure models [2,1], our approach does not suffer from the local optimum and does not need the number of people in each frame to be known as an input. We find in our experiments that our approach outperforms the previous methods on several public datasets.

3 Cuboid Proposal Network

The overview of our approach is shown in Figure 1. It first estimates 2D pose heatmaps for every camera view independently by HRNet [13]. Then we introduce Cuboid Proposal Network (CPN) to localize all people by a number of cuboid proposals. Finally, we present Pose Regression Network (PRN) to regress a detailed 3D pose for each proposal. In this section, we focus on the details of CPN including the input, output and network structures.

3.1 Feature Volume

The input to CPN is a 3D feature volume which contains rich information for detecting people in the 3D space. The feature volume is constructed by projecting the 2D pose heatmaps in all camera views to a common discretized 3D space as will be detailed later. Since the 2D pose heatmaps encode location information of the joints, the resulting 3D feature volume also carries rich information for detecting 3D poses.

We discretize the 3D space, in which people can freely move, by $X \times Y \times Z$ discrete locations $\{\mathbf{G}^{x,y,z}\}$. Each location can be regarded as an anchor of people. In order to reduce the quantization error, we set the distance between the neighboring anchors to be small by adjusting the values of X, Y and Z , respectively. In general, the space is about $8m \times 8m \times 2m$ on the public datasets

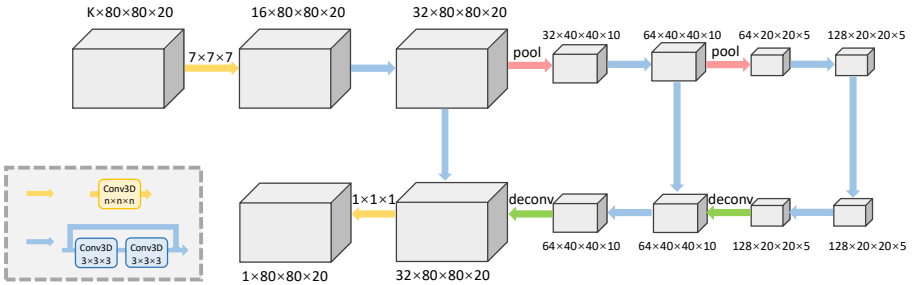


Fig. 2: Network structure of CPN. The input is a feature volume (see section 3.1) and the output is the probability map \mathbf{V} (see section 3.2). The yellow arrow represents a standard 3D convolutional layer and the blue arrow represents a Residual Block of two 3D convolutional layers as shown in the legend.

[15, 2]. So we set X , Y and Z to be 80, 80 and 20, respectively, to strike a good balance between speed and precision. The distance between two neighboring bins is about 100mm which is sufficiently accurate for coarsely localizing people. Note that we will obtain finer-grained 3D poses in PRN.

We compute a feature vector for each anchor by sampling and fusing the 2D heatmap values at its projected locations in all camera views. Denote the 2D heatmap of view v as $\mathbf{M}_v \in \mathcal{R}^{K \times H \times W}$ where K is the number of body joints. For each anchor location $\mathbf{G}^{x,y,z}$, we compute its projected location in view v as $\mathbf{P}_v^{x,y,z}$. The heatmap values at $\mathbf{P}_v^{x,y,z}$ is denoted as $\mathbf{M}_v^{x,y,z} \in \mathcal{R}^K$. Then we compute a feature vector for the anchor as the average heatmap values in all camera views: $\mathbf{F}^{x,y,z} = \frac{1}{V} \sum_{v=1}^V \mathbf{M}_v^{x,y,z}$ where V is the number of cameras. More advanced fusion strategies, for example, assigning a data-dependent weight to reflect the heatmap estimation quality in each camera view, could be explored in the future work. In this work, we stick to the the approach of computing average in order to keep the overall approach as simple as possible. We can see that $\mathbf{F}^{x,y,z}$ actually encodes the likelihood that the K joints are at $\mathbf{G}^{x,y,z}$ which is sufficient to infer people presence. In the following sections, we will describe how we estimate cuboid proposals from the feature volume \mathbf{F} .

3.2 Cuboid Proposals

We represent a cuboid proposal by a 3D bounding box whose orientation and size are fixed in our experiments. This is a reasonable simplification because the sizes of people in 3D space have limited variations which differs from 2D proposals in object detection [32]. So the main task in CPN is to determine the people presence likelihood at each anchor location.

To generate proposals, we slide a small network over the feature volume \mathbf{F} . Each sliding window centered at an anchor is mapped to a lower-dimensional feature which is fed to a fully connected layer to regress a confidence score $\mathbf{V}^{x,y,z}$ representing the likelihood of people presence at the location. The likelihood at

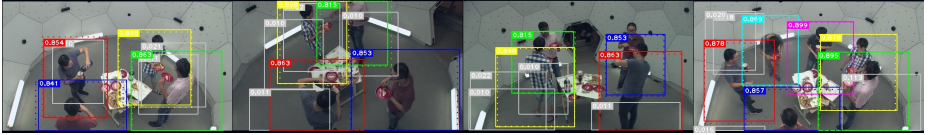


Fig. 3: Estimated cuboid proposals. We project the eight corners of each proposal to the image and compute the minimum and maximum coordinates along the x and y-axis, respectively, which form a bounding box. The numbers represent the estimated confidence scores. The gray boxes denote low confidence proposals. The dashed boxes are the ground-truth.

all anchors form a 3D heatmap $\mathbf{V} \in \mathcal{R}^{X,Y,Z}$. Since we use fixed box orientation and size, we do not estimate them as the 2D object detectors [32,33]. Neither do we estimate center offsets relative to the anchor locations for precise people locations because coarse locations are sufficient.

We compute a Ground-Truth (GT) heatmap score $\mathbf{V}_*^{x,y,z}$ for every anchor according to its distance to the GT poses. Specifically, for each pair of GT pose (root joint) and anchor, we compute a Gaussian score according to their distance. The score decreases exponentially when the distance increases. Note that there could be multiple scores for one anchor if there are multiple people in the scene and we simply keep the largest one. We train CPN by minimizing:

$$\mathcal{L}_{\text{CPN}} = \sum_{x=1}^X \sum_{y=1}^Y \sum_{z=1}^Z \|\mathbf{V}_*^{x,y,z} - \mathbf{V}^{x,y,z}\|_2 \quad (1)$$

The edge length of every proposal is set to be 2000mm which is sufficiently large to cover people in arbitrary poses. The orientations of the proposals are aligned with the world coordinate axes.

3.3 Non-Maximum Suppression

We select the anchors with large regression confidence values as the proposals. On top of the 3D heatmap, we perform Non-Maximum Suppression (NMS) based on the heatmap scores to extract local peaks. Then, we keep the locations of peaks whose heatmap scores are larger than a threshold. Similar to [34], IOU based NMS is not needed for generating proposals because we only have one positive anchor per pose.

3.4 Network Structures of CPN

Inspired by the *voxel-to-voxel* prediction network in [35], we also adopt the 3D convolutions as the basic building blocks for CPN. Since the input feature volume is sparse and has clear semantic meanings, we propose a simpler structure than [35] which is shown in Figure 2. In some scenarios such as football court, the

motion capture space could be larger than that of the public datasets, which would lead to larger feature volume, thus notably reducing the inference speed. We solve the problem by using sparse 3D convolutions [36] because the feature volume only has a small number of non-zero values.

We visualize some estimated proposals in Figure 3. We project the 3D proposals to 2D images using the camera parameters for the sake of simplicity. We can see that most people instances can be accurately retrieved even though some of them are severely occluded in the current view. This is mainly due to the effective fusion of multiview features in a common 3D space. We will numerically evaluate CPN in more details in the experiment section.

4 Pose Regression Network

In this section, we present the details of Pose Regression Network (PRN) which, for each proposal, predicts a complete 3D pose.

4.1 Constructing Feature Volume

Recall that we have already constructed a big feature volume in the previous CPN step, which covers the whole motion space, to coarsely localize people in the environment. However, we do NOT reuse it here in PRN because it is too coarse to accurately estimate the 3D positions of all joints. Instead, we construct a separate finer-grained feature volume centered at each proposal. The size of the feature volume is set to be $2000\text{mm} \times 2000\text{mm} \times 2000\text{mm}$, which is much smaller than that of CPN ($8\text{m} \times 8\text{m} \times 2\text{m}$), but is still large enough to cover people in arbitrary poses. This volume is divided into a discrete grid with $X' \times Y' \times Z'$ bins where $X' = Y' = Z' = 64$. The edge length of a bin is about $\frac{2000}{64} = 31.25\text{mm}$. Note that the precision of our approach is not bounded to 31.25mm because we will use the integration trick [22] to reduce the impact of quantization error as will be described in more detail later. With these definitions, we compute the feature volume following the descriptions in section 3.1.

4.2 Regression of Human Poses

We estimate a 3D heatmap $\mathbf{H}_k \in \mathcal{R}^{X' \times Y' \times Z'}$ for each joint k based on the constructed feature volume. Then the 3D location \mathbf{J}_k of the joint can be obtained by computing the center of mass of \mathbf{H}_k according to the following formula:

$$\mathbf{J}_k = \sum_{x=1}^{X'} \sum_{y=1}^{Y'} \sum_{z=1}^{Z'} (x, y, z) \cdot \mathbf{H}_k(x, y, z) \quad (2)$$

Note that we do not obtain the location \mathbf{J}_k by finding the maximum of \mathbf{H}_k because the quantization error of 31.25mm is still large. Computing the expectation as in the above equation effectively reduces the error. This technique is frequently used in the previous works such as [22].

The estimated joint location is compared to the ground-truth location \mathbf{J}_* to train PRN. Specifically, the L_1 loss is used:

$$\mathcal{L}_{\text{PRN}} = \sum_{k=1}^K \|\mathbf{J}_*^k - \mathbf{J}^k\|_1 \quad (3)$$

The network of PRN is kept the same as CPN as shown in Figure 2 except that the input and output are different. The network weights are shared for different joints. We conducted experiments using different weights but that did not make much difference on current datasets.

4.3 Training Strategies

We first train the 2D pose estimation network for 20 epochs. The initial learning rate is $1e-4$, and decreases to $1e-5$ and $1e-6$ at the 10_{th} and 15_{th} epochs, respectively. Then we jointly train the whole network including CPN and PRN for 10 epochs to convergence. The learning rate is set to be $1e-4$. In some experiments (which will be described clearly), we directly use the backbone network learned on the COCO dataset without finetuning on target datasets.

5 Datasets and Metrics

The Campus Dataset [2] This dataset captures three people interacting with each other in an outdoor environment by three cameras. We follow [2,4] to split the dataset into training and testing subsets. To avoid over-fitting to this small training data, we directly use the 2D pose estimator trained on the COCO dataset and only train CPN and PRN.

The Shelf Dataset [2] It captures four people disassembling a shelf by five cameras. Similar to what we do for Campus, we use the 2D pose estimator trained on COCO and only train CPN and PRN.

The CMU Panoptic Dataset [15] It captures people doing daily activities by dozens of cameras among which five HD cameras (3, 6, 12, 13, 23) are used in our experiments. We also report results for fewer cameras. Following [37], the training set consists of the following sequences: “160422_ultimatum1”, “16022_4_haggling1”, “160226_haggling1”, “161202_haggling1”, “160906_ian1”, “160906_ian2”, “160906_ian3”, “160906_band1”, “160906_band2”, “160906_band3”. The testing set consists of: “160906_pizza1”, “160422_haggling1”, “160906_ian5”, and “160906_band4”.

The Proposal Evaluation Metric We compute recall of proposals at different proposal-groundtruth-distance. It is noteworthy that this metric is only loosely related to the 3D estimation accuracy. We keep ten proposals after NMS for evaluation on the three datasets.

The 3D Pose Estimation Metric Following [4], we use the Percentage of Correct Parts (PCP3D) metric to evaluate the estimated 3D poses. Specifically, for each ground-truth 3D pose, it finds the closest pose estimation and computes

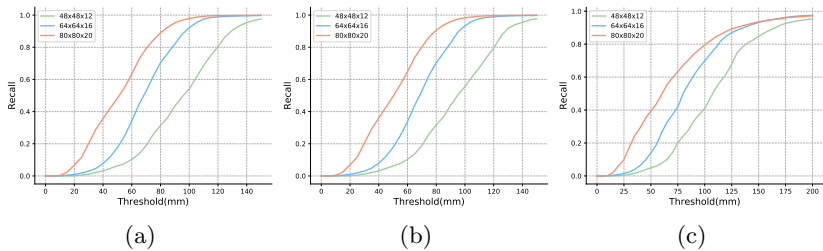


Fig. 4: Recall curves when the motion space is discretized with different numbers of bins on the Panoptic dataset. **(a)** CPN is trained/tested on the real images of five cameras from the Panoptic dataset. **(b)** CPN is trained on the synthetic heatmaps and tested on the real images of five cameras. **(c)** CPN is trained/tested on the real images of one camera from the Panoptic dataset.

the percentage of correct parts. This metric does not penalize false positive pose estimations. To overcome the limitation, we also extend the Average Precision (AP_K) metric [38] to the multi-person 3D pose estimation task which is more comprehensive than PCP3D. If the Mean Per Joint Position Error (MPJPE) of a pose is smaller than K millimeters, we think the pose is accurately estimated.

6 Evaluation of CPN

We first study the impact of the space division granularity to the proposals by setting different values to the X, Y and Z parameters. The results are shown in Figure 4 (a). When we increase the number of bins from $48 \times 48 \times 12$ to $80 \times 80 \times 20$, the recall improves significantly for small thresholds. This is mainly because the quantization error is effectively reduced and the locations of the proposals become more precise. However, the gap becomes smaller for large thresholds. In our experiments, we use $80 \times 80 \times 20$ bins to strike a good balance between the accuracy and speed.

We also consider a practical situation where we do not have sufficient data to train CPN. We propose to address the problem by generating many synthetic heatmaps: we place a number of 3D poses (sampled from the motion capture datasets) at random locations in the space and project them to all views to get the respective 2D locations. Then we generate 2D heatmaps from the locations to train CPN. The experimental results are shown in Figure 2 (b). We can see that the performance is on par with the model trained on the real images. This significantly improves the general applicability of CPN in the wild (we may also need to address the domain adaptation problem in 2D heatmap estimation but it is beyond the scope of this work).

Finally, we study the impact of the number of cameras to the proposals. In general, the recall decreases when fewer cameras are used. In particular, the results of a single camera are shown in Figure 4 (c). We can see that the recall rates at different thresholds are consistently lower than those of the five-camera

Methods		AP	AP ⁵⁰	AP ⁷⁵
HRNet-w48[13]		55.8	67.4	59.0
Ours		98.3	99.5	99.1

Table 1: 2D pose estimation accuracy on the Panoptic dataset. Ours are obtained by projecting the estimated 3D poses to the images.

setup in (a). However, it is still larger than 95% at the threshold of 175mm. It means that it can coarsely retrieve most people using a single camera which demonstrates its practical feasibility. We will report the ultimate 3D pose error using a single camera in the next section.

7 Evaluation of PRN

7.1 2D Pose Estimation Accuracy

We project the 3D poses estimated by our approach to 2D and compare them to the results of HRNet [13]. Since our approach also uses HRNet to estimate heatmaps, they are comparable. The two models are both trained on the Panoptic dataset [15]. The results are shown in Table 1. The AP of HRNet is only 55.8% because there is severe occlusion. Figure 5 shows some 2D poses estimated by HRNet and our approach, respectively. HRNet gets accurate estimates when there is no occlusion which validates its effectiveness. However, it gets inaccurate estimates when people are occluded. In addition, as a top-down method, it may generate false positives if object detector fails. For example, there are two poses mistakenly detected at the dome entrance area in the fourth example.

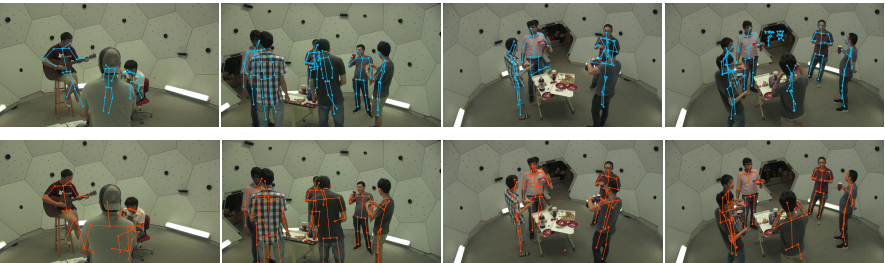


Fig. 5: Comparison of 2D poses estimated by HRNet [13] (top row) and our approach (bottom row). Ours are obtained by projecting the 3D poses to the images. Note that this is only proof-of-concept result rather than rigorous fair comparison as our approach uses multiview images as input.

7.2 Ablation Study on 3D Pose Estimation

We conduct ablation studies to evaluate a variety of factors of our approach. The results on the Panoptic dataset are shown in Table 2.

Space Division Granularity of CPN By comparing the results of (a) and (b), we can see that increasing the number of bins from $64 \times 64 \times 16$ to $80 \times 80 \times 20$ improves accuracy in general. In particular, the AP_{25} metric improves most significantly whereas AP_{50} improves only marginally. The results represent that using finer-grained grids improves the precision but not accuracy which agrees with our expectation. Further increasing the grid size only slightly decreases the error but notably increases the computation time. To strike a good balance, we use $80 \times 80 \times 20$ for the rest of the experiments.

Number of Cameras As shown in (b-d) of Table 2, reducing the number of cameras generally increases the 3D error because the information in the feature volume becomes less complete. In extreme cases, when there is only one camera, the 3D error increases dramatically to 66.95mm as shown in row (d). This is mainly because there is severe ambiguity in monocular 3D pose estimation. If we align the pelvis joints of the estimated poses to the ground-truth (as the previous methods), the 3D pose error decreases to 51.14mm as shown in Table 2 (j). This is comparable to the state-of-the-art monocular 3D pose estimation methods such as [20,39,26]. In addition, we find that AP_{25} drops dramatically but AP_{150} only drops slightly when we reduce the number of cameras from five to one. This means that it can estimate coarse 3D poses using a single camera although they are not as precise as the multiview setup.

Generalization to Different Cameras We train and test our approach on different sets of cameras. Specifically, we randomly select a few cameras from the remaining HD cameras for training and test on the selected five cameras. The 3D error is about 25.51mm (f) which is larger than the situation where training and testing are on the same cameras. But this is still a reasonably good result demonstrating that the approach has strong generalization capability.

Impact of Heatmaps By comparing the results in (b) and (g), we can see that getting accurate 2D heatmaps is critical to the 3D accuracy. When the heatmaps are the ground-truth, the AP s at a variety of thresholds are very high suggesting that the estimated poses are accurate. The MPJPE remarkably decreases to 11.77mm. The main reason for this remaining small error is the quantization error caused by space discretization.

Impact of Proposals By comparing (b) and (h), we can see that replacing CPN by ground-truth proposals does not notably improve the results. The results suggest that the estimated proposals are already very accurate and more attention should be spent on improving the heatmaps and PRN. We do not compute AP s when using ground-truth proposals because the confidence scores of all proposals are all set to be one.

	# Views	Backbone	CPN Size	AP ₂₅	AP ₅₀	AP ₁₀₀	AP ₁₅₀	MPJPE
(a)	5	ResNet-50	64 × 64 × 16	81.54	98.24	99.56	99.85	18.15mm
(b)	5	ResNet-50	80 × 80 × 20	83.59	98.33	99.76	99.91	17.68mm
(c)	3	ResNet-50	80 × 80 × 20	58.94	93.88	98.45	99.32	24.29mm
(d)	1	ResNet-50	80 × 80 × 20	0.860	23.47	80.69	93.32	66.95mm
(e)*	5	ResNet-50	80 × 80 × 20	71.26	96.96	99.12	99.52	20.31mm
(f) ⁺	5	ResNet-50	80 × 80 × 20	50.91	95.25	99.36	99.56	25.51mm
(g)	5	GT Heatmap	80 × 80 × 20	98.61	99.82	99.98	99.99	11.77mm
(h)	5	ResNet-50	GT Proposal	-	-	-	-	16.94mm
(i)	5	GT Heatmap	GT Proposal	-	-	-	-	11.32mm
	# Views	Backbone	CPN Size	AP ₂₅ ^{rel}	AP ₅₀ ^{rel}	AP ₁₀₀ ^{rel}	AP ₁₅₀ ^{rel}	MPJPE ^{rel}
(j)	1	ResNet-50	80 × 80 × 20	1.520	39.86	92.37	96.98	51.14mm

Table 2: Ablation study on the Panoptic dataset. “*” means that CPN and PRN are trained on synthetic heatmaps. “+” means that CPN and PRN are trained and tested with different cameras. “rel” represents that we align the root joints of the estimated poses to the ground-truth.

Qualitative Study We show the estimated 3D poses of three examples in Figure 6. We can see that there are severe occlusions in the images of all camera views. However, by fusing the noisy and incomplete heatmaps from multiple cameras, our approach obtains more comprehensive features which allows us to successfully estimate the 3D poses without bells and whistles. It is noteworthy that we do not need to associate 2D poses in different views based on noisy observations by combining a number of sophisticated techniques. This significantly improves the robustness of the approach. Please see the supplementary video for more examples ³.

Figure 7 (B) shows two examples where our approach did not obtain accurate estimations in the three-camera setup. In the first example, most joints of the lady can be seen from two of the three cameras and our approach accurately estimates the 3D pose. However, the little child is only visible in the first view and, even in that view, many joints are actually occluded by its body. So the resulting 3D pose has large errors. The second example is also interesting. The person is only visible in one view but, fortunately, most joints are visible. We can see that our approach estimates a 3D pose which seems like a translated version of the ground-truth pose plotted in dashed lines. This is reasonable because there is ambiguity for 3D pose estimation from a single image.

Computational Complexity It takes about 300ms on a single Titan X GPU to estimate 3D poses in a five-camera setup. In particular, 93ms is spent on estimating heatmaps and 24ms is spent on generating proposals. The time spent on regressing poses depends on the number of proposals (people). In particular,

³ <https://youtu.be/qZAYHUzdpw>



Fig. 6: Estimated 3D poses and their projections in images. The last column shows the estimated 3D poses.

it takes about 46ms to process one proposal. The inference time has the potential to be further reduced by using sparse 3D convolutions [36].

7.3 Comparison to the State-of-the-arts

Table 3 shows the results of the state-of-the-art methods on the Campus and the Shelf datasets in the top and bottom sections, respectively. On the Campus dataset, we can see that our approach improves PCP3D from 96.3% of [4] to 96.7% which is a decent improvement considering the already very high accuracy. As discussed in Section 5, the PCP3D metric does not penalize false positive estimates. However, it is also meaningless to report AP scores because the GT pose annotations in this dataset are incomplete. So we propose to visualize and publish all of our estimated poses⁴. We find that our approach usually gets accurate estimates as long as joints are visible in at least two views.

Our approach also achieves better results than [4] on the Shelf dataset. In particular, it gets fewer false positives. For example, in Figure 7 (A.2), there is a false positive pose in the pink dashed circle estimated by [4]. In contrast, our approach can suppress most false positives. We find that most errors of our approach are caused by inaccurate GT annotations. For example, as shown in the first column of Figure 7 (A.1), the GT joint locations within the red circle are incorrect. In summary, 66 out of the 301 frames have completely correct annotations and our approach gets accurate estimates on them.

The previous works [2,1,3,4] did not report numerical results on the large scale Panoptic dataset. We encourage future works to do so as in Table 2 (b). We also evaluate our approach on the single person dataset Human3.6M [41]. The MPJPE of our approach is about 19mm which is comparable to [26]. We also visualize and publish our estimated poses⁵.

⁴ <https://youtu.be/AgDQFII5IM>

⁵ <https://youtu.be/S6G3TXaBukw>

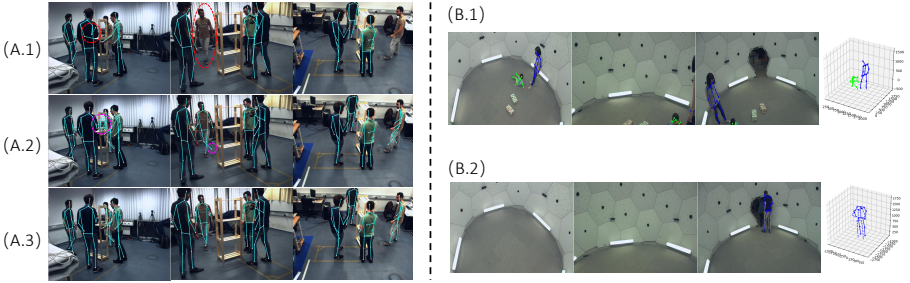


Fig. 7: (A) shows the 3D poses of ground-truth (A.1), estimated by [4] (A.2) and ours (A.3), respectively. The joints in the dashed circles represent the locations are incorrect. (B) shows two typical cases where our approach makes mistakes. The pose plotted by dashed lines in B.2 is the ground-truth.

Campus	Actor 1	Actor 2	Actor 3	Average
Belagiannis <i>et al.</i> [2]	82.0	72.4	73.7	75.8
Belagiannis <i>et al.</i> [3]	83.0	73.0	78.0	78.0
Belagiannis <i>et al.</i> [1]	93.5	75.7	84.4	84.5
Ershadi-Nasab <i>et al.</i> [40]	94.2	92.9	84.6	90.6
Dong <i>et al.</i> [4]	97.6	93.3	98.0	96.3
Ours	97.6	93.8	98.8	96.7
Shelf	Actor 1	Actor 2	Actor 3	Average
Belagiannis <i>et al.</i> [2]	66.1	65.0	83.2	71.4
Belagiannis <i>et al.</i> [3]	75.0	67.0	86.0	76.0
Belagiannis <i>et al.</i> [1]	75.3	69.7	87.6	77.5
Ershadi-Nasab <i>et al.</i> [40]	93.3	75.9	94.8	88.0
Dong <i>et al.</i> [4]	98.8	94.1	97.8	96.9
Ours	99.3	94.1	97.6	97.0

Table 3: Comparison to the state-of-the-art methods on the Campus and the Shelf datasets. The metric is PCP3D.

8 Conclusion

we present a novel approach for multi-person 3D pose estimation. Different from the previous methods, it only makes hard decisions in the 3D space which allows to avoid the challenging association problems in the 2D space. In particular, noisy and incomplete information of all camera views are warped to a common 3D space to form a comprehensive feature volume which is used for 3D estimation. The experimental results on the benchmark datasets validate that the approach is robust to occlusion which has practical values. In addition, the approach has strong generalization capability to different camera setups.

References

1. Belagiannis, V., Amin, S., Andriluka, M., Schiele, B., Navab, N., Ilic, S.: 3d pictorial structures revisited: Multiple human pose estimation. *TPAMI* **38**(10) (2015) 1929–1942
2. Belagiannis, V., Amin, S., Andriluka, M., Schiele, B., Navab, N., Ilic, S.: 3d pictorial structures for multiple human pose estimation. In: *CVPR*. (2014) 1669–1676
3. Belagiannis, V., Wang, X., Schiele, B., Fua, P., Ilic, S., Navab, N.: Multiple human pose estimation with temporally consistent 3d pictorial structures. In: *European Conference on Computer Vision*, Springer (2014) 742–754
4. Dong, J., Jiang, W., Huang, Q., Bao, H., Zhou, X.: Fast and robust multi-person 3d pose estimation from multiple views. In: *CVPR*. (2019) 7792–7801
5. Bridgeman, L., Volino, M., Guillemaut, J.Y., Hilton, A.: Multi-person 3d pose estimation and tracking in sports. In: *CVPRW*. (2019) 0–0
6. Qiu, H., Wang, C., Wang, J., Wang, N., Zeng, W.: Cross view fusion for 3d human pose estimation. In: *ICCV*. (2019) 4342–4351
7. Zhang, Y., An, L., Yu, T., Li, X., Li, K., Liu, Y.: 4d association graph for realtime multi-person motion capture using multiple video cameras. In: *CVPR*. (2020) 1324–1333
8. Cao, Z., Simon, T., Wei, S.E., Sheikh, Y.: Realtime multi-person 2d pose estimation using part affinity fields. In: *CVPR*. (2017) 7291–7299
9. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn. In: *ICCV*. (2017) 2961–2969
10. Hartley, R., Zisserman, A.: *Multiple view geometry in computer vision*. Cambridge university press (2003)
11. Amin, S., Andriluka, M., Rohrbach, M., Schiele, B.: Multi-view pictorial structures for 3d human pose estimation. In: *BMVC*, Citeseer (2013)
12. Newell, A., Yang, K., Deng, J.: Stacked hourglass networks for human pose estimation. In: *ECCV*, Springer (2016) 483–499
13. Sun, K., Xiao, B., Liu, D., Wang, J.: Deep high-resolution representation learning for human pose estimation. In: *CVPR*. (2019) 5693–5703
14. Newell, A., Huang, Z., Deng, J.: Associative embedding: End-to-end learning for joint detection and grouping. In: *NIPS*. (2017) 2277–2287
15. Joo, H., Simon, T., Li, X., Liu, H., Tan, L., Gui, L., Banerjee, S., Godisart, T.S., Nabbe, B., Matthews, I., Kanade, T., Nobuhara, S., Sheikh, Y.: Panoptic studio: A massively multiview system for social interaction capture. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2017)
16. Wang, C., Wang, Y., Lin, Z., Yuille, A.L., Gao, W.: Robust estimation of 3d human poses from a single image. In: *CVPR*. (2014) 2361–2368
17. Ramakrishna, V., Kanade, T., Sheikh, Y.: Reconstructing 3d human pose from 2d image landmarks. In: *ECCV*, Springer (2012) 573–586
18. Zhou, X., Zhu, M., Leonardos, S., Daniilidis, K.: Sparse representation for 3d shape estimation: A convex relaxation approach. *TPAMI* **39**(8) (2016) 1648–1661
19. Pavlakos, G., Zhou, X., Daniilidis, K.: Ordinal depth supervision for 3d human pose estimation. In: *CVPR*. (2018) 7307–7316
20. Martinez, J., Hossain, R., Romero, J., Little, J.J.: A simple yet effective baseline for 3d human pose estimation. In: *ICCV*. (2017)
21. Moreno-Noguer, F.: 3d human pose estimation from a single image via distance matrix regression. In: *CVPR*, IEEE (2017) 1561–1570

22. Sun, X., Xiao, B., Wei, F., Liang, S., Wei, Y.: Integral human pose regression. In: ECCV. (2018) 529–545
23. Fang, H.S., Xu, Y., Wang, W., Liu, X., Zhu, S.C.: Learning pose grammar to encode human body configuration for 3d pose estimation. In: AAAI. (2018)
24. Pavllo, D., Feichtenhofer, C., Grangier, D., Auli, M.: 3d human pose estimation in video with temporal convolutions and semi-supervised training. In: CVPR. (2019)
25. Bogo, F., Kanazawa, A., Lassner, C., Gehler, P., Romero, J., Black, M.J.: Keep it smpl: Automatic estimation of 3d human pose and shape from a single image. In: ECCV, Springer (2016) 561–578
26. Iskakov, K., Burkov, E., Lempitsky, V., Malkov, Y.: Learnable triangulation of human pose. In: ICCV. (2019) 7718–7727
27. Remelli, E., Han, S., Honari, S., Fua, P., Wang, R.: Lightweight multi-view 3d pose estimation through camera-disentangled representation. In: CVPR. (2020) 6040–6049
28. Pavlakos, G., Zhou, X., Derpanis, K.G., Daniilidis, K.: Coarse-to-fine volumetric prediction for single-image 3d human pose. In: CVPR, IEEE (2017) 1263–1272
29. Zhou, X., Huang, Q., Sun, X., Xue, X., Wei, Y.: Towards 3d human pose estimation in the wild: a weakly-supervised approach. In: ICCV. (2017)
30. Rogez, G., Weinzaepfel, P., Schmid, C.: Lcr-net++: Multi-person 2d and 3d pose detection in natural images. TPAMI (2019)
31. Kreiss, S., Bertoni, L., Alahi, A.: Pifpaf: Composite fields for human pose estimation. In: CVPR. (2019) 11977–11986
32. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. In: NIPS. (2015) 91–99
33. Redmon, J., Farhadi, A.: Yolov3: An incremental improvement. arXiv preprint arXiv:1804.02767 (2018)
34. Zhou, X., Wang, D., Krähenbühl, P.: Objects as points. arXiv preprint arXiv:1904.07850 (2019)
35. Moon, G., Yong Chang, J., Mu Lee, K.: V2v-posenet: Voxel-to-voxel prediction network for accurate 3d hand and human pose estimation from a single depth map. In: CVPR. (2018) 5079–5088
36. Yan, Y., Mao, Y., Li, B.: Second: Sparsely embedded convolutional detection. Sensors **18**(10) (2018) 3337
37. Xiang, D., Joo, H., Sheikh, Y.: Monocular total capture: Posing face, body, and hands in the wild. In: CVPR. (2019) 10965–10974
38. Pishchulin, L., Insafutdinov, E., Tang, S., Andres, B., Andriluka, M., Gehler, P.V., Schiele, B.: Deepcut: Joint subset partition and labeling for multi person pose estimation. In: CVPR. (2016) 4929–4937
39. Ci, H., Wang, C., Ma, X., Wang, Y.: Optimizing network structure for 3d human pose estimation. In: ICCV. (2019)
40. Ershadi-Nasab, S., Noury, E., Kasaei, S., Sanaei, E.: Multiple human 3d pose estimation from multiview images. Multimedia Tools and Applications **77**(12) (2018) 15573–15601
41. Ionescu, C., Papava, D., Olaru, V., Sminchisescu, C.: Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. T-PAMI **36**(7) (2014) 1325–1339