

A Multi-view RGB-D Approach for Human Pose Estimation in Operating Rooms

Abdolrahim Kadkhodamohammadi¹, Afshin Gangi^{1,2}, Michel de Mathelin¹, Nicolas Padoy¹

¹ICube, University of Strasbourg, CNRS, IHU Strasbourg, France

²Radiology Department, University Hospital of Strasbourg, France

{kadkhodamohammad, gangi, demathelin, npadoy}@unistra.fr

Abstract

Many approaches have been proposed for human pose estimation in single and multi-view RGB images. However, some environments, such as the operating room, are still very challenging for state-of-the-art RGB methods. In this paper, we propose an approach for multi-view 3D human pose estimation from RGB-D images and demonstrate the benefits of using the additional depth channel for pose refinement beyond its use for the generation of improved features. The proposed method permits the joint detection and estimation of the poses without knowing a priori the number of persons present in the scene. We evaluate this approach on a novel multi-view RGB-D dataset acquired during live surgeries and annotated with ground truth 3D poses.

1. Introduction

Recovering the configuration of human body parts, which is also referred to as Human Pose Estimation (HPE), can benefit a wide variety of applications such as video surveillance and behavior analysis [16] and human computer interactions [20]. Vision-based HPE has been actively researched over the years and promising results have been reported on various challenging datasets recorded in common indoor and outdoor scenes [28, 41, 33, 3, 40, 19].

Even though state-of-the-art models such as [19, 41] achieve impressive results on standard computer vision datasets, our experiments show that they do not necessarily generalize well to special environments like operating rooms (ORs). The quantitative results presented in Section 3 on data recorded during real surgeries show that there is still a large margin for improvement. These results are also in agreement with [23], who reported that the Kinect skeleton tracker [30], which has been successfully used in the game industry, does not generalize well to the OR environment. The main issues are incorrect background subtraction and the mixing of the body parts belonging to different persons. Our supplementary video presents qualitatively how the aforementioned approaches perform on our OR data.

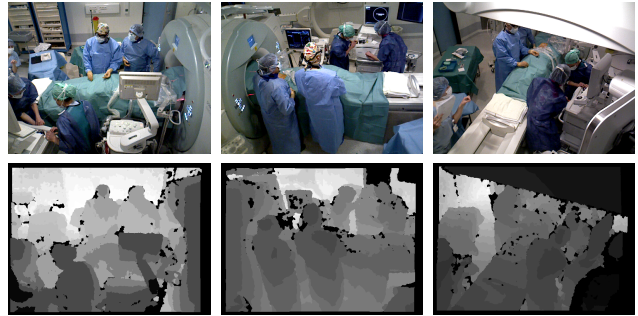


Figure 1. Synchronized pairs of color and depth images from a novel multi-view operating room dataset. The images are recorded during live surgeries using a three-view RGB-D camera system.

We believe that the drop in performance is caused by the inherent visual challenges present in such an environment. Figure 1 shows an operating room during real surgeries: one can notice that the visual appearances of many different surfaces are very similar, which makes it difficult to distinguish persons from the background. People are also wearing loose and textureless gowns that make it hard to discriminate body parts. In addition, the camera positioning possibilities are very limited due to the objects and equipment that need to be frequently displaced in the room, such as ceiling mounted articulated arms, screens and the respiratory tower. Furthermore, performing a surgery requires the collaboration of multiple people, which increases the risk of occlusions.

In [23], we have proposed an approach based on pictorial structures for clinician detection and pose estimation. The approach relies on a pair of color and depth images captured from a single viewpoint using a low-cost RGB-D sensor similar to the Microsoft Kinect camera [30]. That work has demonstrated that the combination of color and depth information along with the use of depth information to model 3D constraints between neighboring body parts greatly improves the pose estimation results. Following the findings of [23], we propose in this paper to use a multi-view system based on RGB-D cameras. The system uses three cameras in order to capture the environment from complementary

views and to reduce the risk of occlusions in such a cluttered environment. We show that the advantages of using depth maps in such a multi-view approach go beyond the mere generation of improved appearance features.

In a multi-view RGB system, correspondences across views are traditionally established by relying on appearance similarity and triangulation [6, 14, 2, 7], which is unreliable in OR environments containing many surfaces that are visually similar. Instead, the depth data enables us to back-project points to 3D and is not affected by the visual appearance of the surfaces in the scene [23]. It also enables us to back-project points that are only visible in one view, while in multi-view RGB systems, points should be visible in at least two views.

Current multi-view human pose estimation approaches have been proposed either for single-person scenarios [14, 8, 18, 2] or for multi-person scenarios in which the number of persons is known in advance [25, 6, 7]. The approach proposed in this work makes no assumption about the number of persons in the scene. To this end, our approach first processes each view separately to detect putative skeletons. Next, *a priori* information about the environment, modeled using random forests, is applied to filter spurious skeletons. The resulting skeletons are then merged across views¹. Finally, a novel energy function is optimized to incorporate evidence across views and update initial part positions directly in 3D.

Our single-view RGB-D pose estimation approach extends 3D Pictorial Structures (PS) [23] by incorporating Convolutional Neural Networks (ConvNets) for the part detection [19]. ConvNets have recently enjoyed a great success in solving many vision-based tasks including human pose estimation [36, 34, 29, 19, 27]. They are capable of learning strong detectors that can incorporate a wide image context through deep network architectures with large receptive fields [38]. Mutual spatial constraints among body parts are however not explicitly modeled, even though they are essential to guarantee joint consistency in the predicted body configuration, especially in multi-person and cluttered environments such as ORs. Therefore, we use a deep ConvNet-based part detector constructed for RGB-D data in conjunction with a 3D pairwise dependency model to enforce body kinematic constraints directly in 3D. This is in contrast to current methods that rely on 2D displacement [40, 21, 35] or visual similarities [19] among body joints. Enforcing body kinematic constraints in 3D is crucial to reliably estimate body part configurations of different individuals who are close to each other in the projected 2D image and are visually similar.

Incorrect detections and occlusions can however result in spurious skeleton candidates in each view that can mislead the multi-view merging algorithm. We argue that in a

specific environment like the operating room, *a priori* information about the room should be leveraged to identify spurious candidates. Therefore, we also propose a method to learn a prior on the 3D body kinematic and room layout constraints. This prior, based on random forest, is used to recognize and remove skeletons with unlikely 3D shapes or positions. Relying directly on high level 3D skeleton information enables the model to better explore the *a priori* information of the OR and to build a stronger prior compared to traditional pose priors that are based on the displacement or visual similarity among parts [40, 19, 23].

This paper investigates multi-person multi-view pose estimation using RGB-D data and makes the following contributions. First, we extend the 3D pictorial structures of [23] to use a ConvNet body part detector on RGB-D images. Second, we propose a random forest based method to automatically learn a prior to incorporate *a priori* information about the environment. Third, we introduce a novel multi-view energy formulation to estimate 3D body configurations by leveraging depth data and reasoning across all views. Finally, we have evaluated the approach on a novel multi-view OR dataset generated from several days of recordings during *live surgeries*.

Related work. Approaches using RGB-D data for human pose estimation are mostly based on background subtraction [4, 30], which is a very challenging task in the OR. An exception is [9], which relies on random forests and color to cluster image pixels into body parts and then estimate body poses. However, the performance of the system degrades dramatically due to occlusions and clutter, which mislead the pixel classifier and the color-based clustering algorithm. The multi-view RGB-D systems that have been proposed also rely on background subtraction and usually address simple scenarios in laboratory setups [39]. For this reason, we focus in the following description on part-based approaches (using RGB images) that do not rely on background subtraction, since they are at the core of our method. Such approaches represent the human body as a set of body parts and model body kinematic constraints using a deformation model. The pictorial structures framework is the main part-based approach and has driven much of the progress in the field since its introduction in [13].

In single-view pictorial structures, exact inference was made tractable by the seminal work of [12], which has been extended in different ways by either constructing a stronger body part detector [10, 41] or by improving the deformation model [28, 23]. Most recently, with the availability of large training sets and high computational power, state-of-the-art results are obtained by using deep convolutional neural networks as body part detectors [34, 19, 26].

Single-view multiple human pose estimation is often performed in two steps: person detection followed by pose estimation [15, 32]. But, when people are in close proximity

¹We assume that the extrinsic parameters of the cameras are known.

to each other, *e.g.* in the operating room, a body part can be assigned to more than one person due to weak body part detections or occlusions. [19], which is built on [27], has proposed a multi-person pose estimation approach based on integer linear programming to jointly detect people and estimate their body part configurations. In this approach, part detection is performed using deep ConvNets and interpart pairwise constraints are enforced based on 2D displacement and appearance similarity between body parts. But, image-based pairwise constraints are not very discriminative in operating rooms since people are wearing textureless gowns with similar colors (see Figure 1). Instead, [23] uses more discriminative 3D pairwise constraints. We follow a similar formulation to [23] using 3D pairwise constraints, but include a deep ConvNet-based detector in our model instead of a support vector machine with engineered visual features.

Most work on multi-view human pose estimation focuses on single-person scenarios in controlled environments to reduce the ambiguity of data association by relying on background subtraction or exemplar-based approaches [14, 18, 31]. In order to be robust to cluttered background, part-based approaches are used to generate part hypotheses per view and then estimate the body configuration in 2D [1, 2] or in 3D [8].

Recently, a PS approach has been proposed to estimate the poses of multiple persons in a multi-view setup [5]. The body pose estimation is performed in 3D by relying on 2D view appearance cues and on multi-view cues computed through triangulation. This approach has been extended in [6] to also include temporal cues. However, all these methods require prior knowledge of the number of persons present and have been evaluated on scenarios recorded in controlled laboratory environments that include people in upright poses only. In contrast, our method detects the number of persons and does not rely on triangulation that might not always be reliable in complex and cluttered environments. Moreover, we have evaluated our approach on a dataset recorded during real surgeries, which contains many visual challenges and where people exhibit a much wider range of articulations compared to those used in the aforementioned works.

A multi-view clinician pose estimation approach has been proposed in [7]. Background subtraction and tracking over the entire sequence are used to find the trajectories of the persons and to localize them using bounding boxes. Then, the approach presented in [5] and a ConvNet-based RGB part detector are used to estimate the pose of each individual given the bounding boxes. To evaluate this approach, two sequences containing a constant number of individuals have been recorded during two medical procedures simulated by actors. In contrast, our approach relies on a multi-view set of images from a single time-step. Moreover, our dataset has been generated from four days of

recordings during live surgeries, which is more challenging due to: (1) a high variation in number of persons per frame and (2) the presence of many movable objects in the scene, which makes it difficult to compute the foreground.

2. Method

We start this section by recapitulating the clinician pose estimator of [23] and then present the different components that lead to our multi-view RGB-D approach.

2.1. Single-view body pose estimator

In [23], we have presented a pictorial structures model to estimate body configurations on RGB-D images. This model represents the body as a set of n joints and learns multiple mixtures of parts to be robust to appearance changes. The model uses ten body joints to indicate upper-body poses, since lower body parts are often occluded in operating rooms. A body configuration is specified by a pair (l, t) , where $l = \{l_1 \dots l_n\}$ indicates the 2D positions of the body joints and t_i belongs to a set of m possible mixture types $t = \{t_1 \dots t_n\}$ for each body joint. The pose estimation is defined as an energy minimization over a tree-structured graph $G = (V, E)$, whose nodes are the body joints and whose edges indicate dependencies between joints. The body joint dependencies are defined following the human body skeleton. Given a pair of aligned color and depth images denoted by I and D , respectively, the score associated with a body configuration (l, t) is defined as:

$$S(I, D, l, t) = \sum_{i \in V} \phi(I, D, l_i) + \sum_{ij \in E} w_{ij}^{t_i, t_j} \cdot \psi(D, l_i, l_j) + \sum_{i \in V} b_i^{t_i} + \sum_{ij \in E} b_{ij}^{t_i, t_j}, \quad (1)$$

where the first term is the appearance model, which is also referred to as the part detector, and the second term is the deformation model that enforces pairwise dependencies between body joints. The last two terms are part type compatibility score functions, where $b_i^{t_i}$ captures the score of assigning a particular mixture type to part i and $b_{ij}^{t_i, t_j}$ is the score associated with the co-occurrence of a particular pair of part types. The compatibility score functions and $w_{ij}^{t_i, t_j}$ are the model parameters. These parameters are learned using a structured support vector machine formulation.

The part detector assigns a confidence score for placing the body joint i at image location l_i . [23] relies on hand-crafted features, namely Histogram of Oriented Gradients (HOG) and Histogram of Depth Differences (HDD). In this work, we compute the part detection scores using the Deep ConvNet model presented in Section 2.1.1.

The deformation model is parametrized by w_{ij} and $\psi(D, l_i, l_j)$. The weights w_{ij} encode the deformations

between pairs of joints. The relative displacement of joint i w.r.t. joint j is captured by $\psi(D, l_i, l_j) = [|d_{3D}|, dc, dc^2, dr, dr^2]^T$, where $|d_{3D}|$ is the absolute 3D Euclidean distance between the joints and (dc, dr) are the relative displacements along columns and rows of the image. The 3D joint positions are computed by back-projecting 2D joints into 3D using the depth image. This term enables the model to incorporate more reliable 3D body part lengths for assembling body joints, which is important in order to discriminate between detections on the surface of a person and detections on the background.

Estimating the body configuration in this model corresponds to finding the optimal body joint positions and mixture types given color and depth images. This process, called inference, computes $(l^*, t^*) = \operatorname{argmax}_{l,t} S(I, D, l, t)$. It is performed in 3D using an efficient algorithm that makes exact inference tractable. For more information, we refer the reader to [23].

2.1.1 ConvNet-based RGB-D body part detector

Motivated by the great success of deep convolutional neural networks in recent years [29, 17, 34, 19, 26], we propose to use RGB-D body part detectors based on deep ConvNets in order to automatically learn feature representations instead of relying on engineered feature representations such as HOG or HDD. To this end, we build on the very deep residual network [17], which has recently been used for part detection and shown promising results [19]. The body part detection is formulated as a multi-label classification problem, where a set of n scores is generated at each image location to denote the probability of part presence. The scores are obtained by using sigmoid activation functions on the output neurons. We adapt the network to learn body part detectors for pairs of color and depth images. We change the input layer to accept four dimensional data (*i.e.* three color channels and depth channel). We also change the `res3d_pose` layer to generate part score maps for ten upper-body parts instead of the fourteen full body parts. During pose estimation, we use the ConvNet-based body part detector to predict confidence scores for all parts at every image locations. Hereafter, we refer to this HPE model as *Deep3DPS*.

Fine-tuning. We initialize the network from the pre-trained model of [19], which is trained on the *MPII Human Pose* dataset. We fine-tune the network on the single-view clinician pose dataset from [23] that consists of 1451 RGB-D frames including 1991 persons using the Caffe framework [22]. We scale the images down to 85% and use a batch size of two. Similarly to [19], we generate target training score maps for all body joints by assigning the positive label 1 for all image locations within 15 pixels to the ground truth location and negative label 0 otherwise. During training, we

use all positive samples and keep at most three times more negative samples. The network is trained with cross entropy loss and stochastic gradient descent for 50k iterations. The initial learning rate is set to 5×10^{-5} for the adapted layers and 5×10^{-6} for the rest. This yields the best results in our experiments. In [19], the network is trained for three tasks: body part detection, location refinement, which is the relative row and column displacement from a scoremap location to the ground truth, and regression to other parts. However, training for the last two tasks did not yield any performance improvement during our experiments. We therefore only train for the body part detection task.

2.2. Random forests based prior

To design a robust method, we believe that it is essential to include priors specific to the environment. Even though a general body kinematic prior is included in the pose estimation model through pairwise constraints, it cannot be guaranteed that these constraints are always properly enforced due to the high complexity of the pose estimation model that predicts human poses directly from image pixel values. In addition, this prior only captures body kinematic constraints and does not incorporate *a priori* information about the environment. In an environment like the OR, constraints such as possible human poses and possible locations can also be used to improve the reliability of the method. Such constraints cannot be easily handcrafted. Furthermore, including them in the pose estimation model would need higher-order dependency terms. Adding such terms would increase the number of model parameters and, more importantly, dramatically increase the complexity of the inference algorithm. We therefore propose to automatically learn the prior, which we formulate as a binary classification problem that takes a skeleton estimated by the single-view detector as input and outputs whether this skeleton corresponds to a spurious detection or not.

We base our approach on Random Forests (RF), which are an ensemble of decision trees consisting of two types of nodes: split and leaf nodes. In each split node, a decision function forwards samples to one of the branches until they finally reach a leaf node containing a prediction function. In our case, we use RF with binary trees and the mean over all predictions to aggregate the votes across all trees. The trees are learned automatically given a labeled training set, which we construct using the skeletons estimated by our single-view pose estimator on a set of images for which ground-truth is available. The detected skeletons are compared to the ground truth using the probability of correct keypoints (PCK) metric, which is commonly used for evaluation in multiple-person pose estimation [41, 23, 27]. We label a detected skeleton as positive if the head, neck, and left and right shoulders are correctly localized according to PCK.

For RF training, we propose to combine various features

computed from the 3D skeletons, which are all expressed in the common room reference frame. The reference coordinate system is chosen w.r.t. the operating table in default position, which makes the prior generalizable to other ORs. This enables our prior to encode two types of information: room layout and possible clinician poses. Certain parts of the room, such as the floor or the ceiling, are for instance not expected to have clinicians or certain body parts. Thus, as first set of features, we use the positions of the 3D body parts to enable the RF to build an internal representation of their spatial occupancy probability. To capture the set of possible human poses in the OR, we include a second set of features, namely the relative 3D displacements between all pairs of body joints. The prior also serves to verify 3D part lengths and exclude incorrect skeletons that may occur due to weak detections and foreground/background confusions. As third feature, we include the detection score of the individual skeleton to incorporate detection confidence. To enable our prior to better encode high-level information, we use the RF method in a multi-layer scheme, referred to as auto-context in the machine learning literature [37]. A multi-layer model is learned, where the first RF layer is constructed using only the three aforementioned types of features, while the other layers use another extra feature that is the classification confidence generated by the previous RF layer.

2.3. Multi-view human pose estimation

2.3.1 Multi-view fusion

The objective of the multi-view fusion is to combine the 3D skeletons across all views. For a given *frame*, defined as a set of RGB-D images recorded from all cameras at the same time step, detections from all views are first put in a set. The two closest skeletons that do not originate from the same view are then merged. This procedure is iterated until no pair of merging candidates is left in the set, where the condition for merging two skeletons is that the distance between their heads and the distance between their necks are both smaller than a constant T_s . Since the left/right side labels of the individual detections are not always reliable, to ensure a consistent merging of the 3D joints we use the 3D positions of the shoulders to find the correct association between the two skeletons. Finally, for all skeletons resulting from a merging step, the left and right side labels are set based on a majority vote among the supporting skeletons. If a merged skeleton originates from only two supporting skeletons, which do not agree on the side label, we set the side according to the skeleton with highest confidence.

As a result, we obtain a set of initial 3D skeletons generated from skeletons coming from one or more views. Then, a new multi-view energy function, presented next, is used to drive the body parts towards their optimal 3D locations by jointly optimizing over all views.

2.3.2 Multi-view RGB-D Optimization

We formulate our multi-view RGB-D approach as an energy minimization over the same graph G as in Section 2.1 and define the energy function $E(\Delta)$ over the graph as:

$$E(\Delta) = \sum_{i \in V} \left(\lambda_1 \cdot \Phi^{conf}(\delta_i) + \lambda_2 \cdot \Phi^{depth}(\delta_i) \right) + \sum_{(i,j) \in E} \Psi_{i,j}(\delta_i, \delta_j), \quad (2)$$

where λ_1 and λ_2 are weighting coefficients, $\Delta = \{\delta_1 \dots \delta_n\}$ is a set of displacement labels for all body parts, $\delta_i \in \mathbb{R}^3$ is a 3D displacement offset for part i , $\Phi(\cdot)$ are the unary potentials and $\Psi_{i,j}(\delta_i, \delta_j)$ is a pairwise dependency term enforcing body physical constraints.

The first term in (2) incorporates part detection confidence scores computed by the ConvNet part detector. Given the list of all views *views*, we define:

$$\Phi^{conf}(\delta_i) = \sum_{v \in views} conf\left(proj(P(\delta_i), v)\right), \quad (3)$$

where $P(\delta_i)$ is the 3D position of part i displaced by an offset δ_i and $proj(p_{3D}, v)$ projects the 3D point p_{3D} . In order to provide a smooth cost function, we compute the distance transforms of the deep ConvNet score maps using the generalized distance transform algorithm [12]. We find that this transformation is necessary to avoid local minima. $conf(p_{2D}) \in [0..1]$ is the value of the distance transform of the score map of part i at location p_{2D} . The second term is defined as:

$$\Phi^{depth}(\delta_i) = \sum_{v \in views} \left| D\left(proj(P(\delta_i), v)\right) - Z(P(\delta_i), v) \right|, \quad (4)$$

where $D(p_{2D})$ is the depth value at image location p_{2D} , $Z(p_{3D}, v)$ is the z value of the 3D point p_{3D} in the coordinate system of the view and $|\cdot|$ is the absolute value operator. To reduce the effect of the noise present in the depth image, we smooth the depth image with a median filter of size 7×7 px. This term quantifies the distance between the displaced 3D joint and the surfaces captured by the depth cameras. Therefore, it can help to avoid placing parts in ghost 3D locations that do not correspond to any surface in the scene. These two unary terms incorporate multi-view cues, where the RGB-D ConvNet is used to include image evidence and depth is used to integrate a reprojection cost across all views.

The pairwise term is used to enforce kinematic constraints, namely body part lengths between pairs of joints. Let $\Psi_{i,j}$ be defined as:

$$\Psi_{i,j}(\delta_i, \delta_j) = \left| \|P(\delta_i) - P(\delta_j)\| - \mu_{i,j} \right|, \quad (5)$$



Figure 2. Multi-view examples illustrating the results of the RF-based prior. Accepted skeletons are shown in orange and rejected skeletons in purple.

where $\|\cdot\|$ is the \mathcal{L}_2 -norm and $\mu_{i,j}$ is the average distance between joints i and j , *i.e.* average part length. The average part lengths are computed over the entire training dataset. Note that since the body part lengths are relatively constant in 3D, it is here not needed to learn person-specific average part lengths.

Inference. In order to recover 3D body part configurations, we need to perform inference in 3D. This problem corresponds to optimizing the energy function in Eq. (2). Note that using the optimization algorithm of [23] would require to construct a 3D state space that includes all 2D positions back-projected to 3D (amounting to the number of views multiplied by the size of the images) augmented with extra nodes for each back-projected node to account for occlusions. Such a large state space would degenerate the performance and slow down the inference. Similarly, the inference approach from [8] would limit us to use simple binary pairwise terms. Instead, we perform discrete optimization using the *fast-PD* algorithm [24], which casts the optimization problem in an integer programming framework and exploits solutions from both primal and dual problems for efficiency. To perform the optimization, we define a set of discrete displacement labels \mathcal{L} for each body joint by sampling densely from a cube centered at the initial joint position. The sampling function is parametrized by (k, s) , where k is the number of samples along each 3D direction and s is the step size between the samples. We perform the optimization iteratively by starting with a coarse label set that covers a large 3D space. At the end of each iteration, we update the part positions based on the displacement labels and then generate a finer label set for the next iteration.

3. Experimental results

Datasets. We have generated a novel multi-view RGB-D dataset, illustrated in Figure 1, by recording all activities in an operating room for four days. For quantitative analysis, the 3D upper body poses of 1378 clinicians have been manually annotated in 741 multi-view frames that are evenly

Setting	Head	Shld	Elbow	Wrist	Hip	Avg
Deep3DPS (DeeperNet)	89.6	56.5	50.6	54.3	42.9	58.8
Deep3DPS (RGB)	93.7	74.9	69.6	71.8	66.6	75.3
Deep3DPS (Depth)	91.0	75.0	69.1	68.0	63.2	73.2
Deep3DPS (RGBD)	93.4	77.0	71.5	73.7	69.1	76.9
+Auxiliary tasks	91.4	72.1	64.9	68.4	63.5	72.1
3DPS (IHOG+HDD) [23]	90.8	74.2	62.2	63.4	57.5	69.6
Insafutdinov et al. [19] ²	91.1	53.7	47.5	50.1	38.4	56.2
Yang and Ramanan [41] ³	30.4	35.2	19.6	24.3	16.7	25.2

Table 1. Pose estimation results of several single-view approaches using PCK metric.

distributed across the dataset. All clinicians who have more than 50% of their upper-body parts visible in at least one view have been annotated in these frames. The annotations are performed using a tool that displays a 3D point cloud reconstructed from all three views as well as the corresponding individual 2D images. This tool allows the user to move the body joints either in the 2D views or in the 3D point cloud. Whenever a joint is moved in 2D, the 3D position of the corresponding joint in the average 3D skeleton is updated using the depth map and then reprojected back to all views. Thus, the annotator can verify the correctness of the annotated skeletons using both 3D visualization and 2D re-projections across the views.

In order to have a fair comparison with our method in [23], the single-view dataset of [23] is used for training all 3DPS single-view pose estimation models and for fine-tuning the network. For all models, evaluation is performed on the new multi-view dataset. As the multi-view dataset is used for random forest training, to evaluate the model, a 4-fold leave-one-out cross-validation is performed, where one folds is used for testing and the rest for training. The evaluation reports the average results of the cross-validation.

3.1. Single-view pose estimation

Table 1 reports the performance of different models on the multi-view dataset using the *PCK* metric [41, 23, 19]. All 3DPS models have been trained on the single-view dataset used in [23]. The 3DPS method using the pre-trained network of [19] as body part detector, referred to as *DeeperNet*, achieves a better performance compared to the full *DeeperCut* approach from [19] that estimates the body poses via a joint optimization across all people. These results indicate that in an environment with many visually similar surfaces, a 3D deformation model, even with tree-structured graph, is more reliable than a fully connected deformation model which relies on appearance and 2D displacement constraints. Fine-tuning the network on

²We use the model that was made publicly available by the authors.

³The model is trained on the Buffy dataset [11] using the public implementation of the approach.



Figure 3. Examples of multi-view pose estimation results. Each row shows a multi-view frame. The 3D skeletons obtained after multi-view energy optimization are projected to the views.

the single-view dataset significantly improves the results (*Deep3DPS (RGB)*: 75.3% vs. 58.8% PCK), as it allows the network to adapt its representation for learning a better encoder for such an environment. We have also trained the network to detect body parts using only depth data, *Deep3DPS (Depth)*, which achieves competitive results. The best performance is obtained when the network relies on both color and depth images: the resulting model, called *Deep3DPS (RGB-D)*, is therefore used as single-view pose estimator during the rest of the experiments. But, we observe that on this data, training the network for the auxiliary tasks suggested in [19], namely location refinement and regression to other parts, degrades the performance. We believe that this is due to both a much smaller training set and to the strong foreshortening of the body parts because of the top views of the cameras.

As baseline, we evaluate the performance of the best 3DPS model from [23], which relies on a 3D deformation model similar to our approach, but with handcrafted color and depth features. Our best model improves the performance over this baseline by $\sim 7\%$ on the same experimental setup. This highlights the benefits of deep ConvNets in constructing more discriminative body part detectors by automatically learning feature representations and also incorporating a wider context. Evaluation of state-of-the-art RGB models [19, 41] trained on common computer vision datasets shows that they do not generalize to the OR environment due to both loose clinical clothes and the presence of many visually similar surfaces.

3.2. Random forest based prior

The *Deep3DPS (RGB-D)* model is applied to detect skeletons in each view of the multi-view dataset separately. The skeletons are back-projected into 3D and transformed into a common reference frame. We use these 3D skeletons to train our prior, as explained in section 2.2. We also augment the data by flipping the skeletons to exchange the left-right body parts. Due to the small size of the training set,

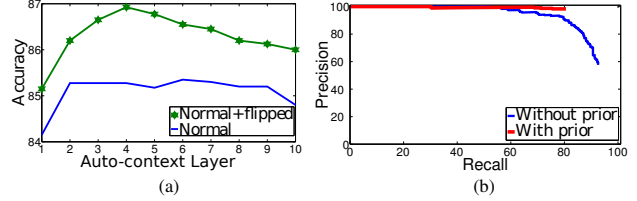


Figure 4. (a) Accuracy of the RF-based prior in detecting spurious skeletons. (b) Precision-recall curves for 3D clinician detection.

we learn 100 shallow trees with a maximum depth of 10. Figure 4(a) shows the detection accuracy of the RF-based method in distinguishing spurious detections. The results show that the method always detects valid skeletons with an accuracy superior to 84%. One observes that augmenting the training set by flipping the skeletons consistently improves the results by enabling the forest to learn a richer prior model that is not confused by the noisy side detections. The results also indicate that auto-context enhances the performance up to the fourth iteration and then tends to overfit. We therefore use the output of the RF trained on the augmented training set at the fourth layer to identify spurious skeletons during the remaining evaluations. Figure 2 illustrates the results of the approach on sample frames from the multi-view dataset and also on a few frames recorded in a different room from totally different viewpoints. It can be seen that the method has correctly identified spurious skeletons in both datasets and generalizes well. The proper generalization is due to the fact that the reference frame is defined on the floor at the center of the operating table, the main element in any operating room.

3.3. Multi-view 3D detection and pose estimation

We set T_s to 30 cm to avoid merging skeletons across persons who are close to each other. We evaluate 3D clinician detection using the precision-recall curves. A detection is accepted as a true positive if the distance between the ground-truth and the detection is below 30 cm for both the head and neck. We use the fusion algorithm described in Section 2.3.1 to generate a set of 3D skeleton candidate per frame. Figure 4(b) shows the 3D clinician detection results after multi-view fusion with and without the RF-based prior. The high precision obtained by our method when the OR prior is used indicates the high quality of the generated skeletons.

To optimize part positions based on multi-view cues, we generate four label sets $\{(k, s) : (3, 50), (5, 10), (7, 2), (7, 1)\}$, where the step sizes are in centimeter. We solve the optimization in four iterations by going from a large and coarse search space towards a small and fine search space, which allows us to more efficiently explore the 3D space. The parameters used in all experiments are $\lambda_1 = 2$ and $\lambda_2 = 0.5$, that are

Part name	One view			Two views			Three views		
	initial	after opt.	opt.- Φ^{depth}	initial	after opt.	opt.- Φ^{depth}	initial	after opt.	opt.- Φ^{depth}
Head	7 ± 4	7 ± 4	7 ± 4	6 ± 3	6 ± 3	6 ± 3	5 ± 2	5 ± 2	5 ± 2
Neck	7 ± 4	7 ± 4	7 ± 4	5 ± 3	5 ± 3	5 ± 3	4 ± 2	4 ± 2	4 ± 2
Shld	25 ± 25	19 ± 16	21 ± 18	22 ± 16	15 ± 10	19 ± 13	14 ± 14	10 ± 7	12 ± 9
Hip	28 ± 22	27 ± 19	28 ± 20	24 ± 13	23 ± 13	24 ± 14	18 ± 10	17 ± 9	18 ± 10
Elbow	31 ± 22	27 ± 19	30 ± 21	30 ± 18	23 ± 15	27 ± 18	19 ± 14	16 ± 11	18 ± 14
Wrist	42 ± 34	32 ± 21	35 ± 24	34 ± 22	25 ± 16	28 ± 18	24 ± 18	18 ± 13	20 ± 15
avg†	32 ± 26	26 ± 19	29 ± 21	28 ± 17	22 ± 14	25 ± 16	19 ± 14	15 ± 10	17 ± 12

Table 2. Mean and standard deviation of 3D part localization error in centimeter. The results are presented as a function of the number of supporting views used to generate the initial 3D skeletons (distribution: 1 view: 30%; 2 views: 43%; 3 views: 27%). † The average is computed for all parts except the head and neck since they are not included in the optimization. See Section 3.3 for details.

selected using grid search over a set of 50 frames from the multi-view dataset. The mean and standard deviation (STD) of the 3D Euclidean distances between the predicted body part positions and the ground-truth positions are used to evaluate 3D body part localizations.

In Table 2, we present the evaluation results for multi-view body part localization as a function of the number of supporting views. Please note that since the head and neck localization errors are close to the expected error in low-cost RGB-D cameras, we do not update these two joints during our optimization. This table presents localization errors for the initial 3D skeletons obtained by the fusion algorithm and the error after performing the multi-view optimization. One can notice that the proposed multi-view fusion method correctly associates skeletons across views by consistently reducing the localization errors as the number of supporting views increases. However, we observe that if we ignore the left and right labels of the detections and assign the label based on shoulder distances with ground truth, the localization errors decrease by ~ 10 cm for skeletons with one or two supporting views and ~ 3 cm for skeletons with three supporting views. These results indicate that the side detection in individual views is not very reliable. But, if a person is detected in all views, the proposed voting algorithm can make a more reliable prediction. The multi-view optimization significantly reduces the localization error for skeletons with any number of supporting views. Interestingly, the optimization improves the results even for skeletons with one supporting view by properly incorporating the depth-based reprojection costs and detection confidences. The deep RGB-D part detector is the main driver of the optimization. To evaluate the effect of the depth-based reprojection cost, we also report the results without this term in column ‘opt.- Φ^{depth} ’. The drop in performance highlights its importance. The 2D projections of 3D poses obtained using the proposed multi-view optimization are shown for a few frames in Figure 3.

4. Conclusions

In this paper, we propose a multi-view RGB-D approach for detecting and estimating the body part positions of medical staff in 3D. A ConvNet-based body part detector combined with a 3D pairwise deformation model is used to recover body poses in each view. A method based on multi-layer random forests is then proposed to automatically learn *a priori* information about the OR and remove spurious detections per view, which allows us to reliably detect the body poses of persons in the scene. Then, these detections are back-projected to 3D and merged across views. Finally, a novel optimization function is introduced to update the part positions by relying jointly on the body part confidence maps, depth data and multi-view cues. The method has been quantitatively evaluated on a new multi-view dataset acquired during live surgeries. Experimental results show significant improvements over state-of-the-art methods for the task of single-view pose estimation in multi-person scenarios, indicating the benefit of combining deep part detectors and 3D pairwise constraints in building robust models. The multi-view formulation also achieves very promising results showing the benefits of the deep ConvNet detector and of depth data for correctly driving parts towards their optimal locations. To the best of our knowledge, this is the first multi-view approach that performs both human detection and pose estimation in a real scenario without any prior knowledge on the number of persons present, as well as the first multi-view RGB-D approach presented for pose estimation.

Acknowledgements

This work was supported by French state funds managed by the ANR within the Investissements d’Avenir program under references ANR-11-LABX-0004 (Labex CAMI), ANR-10-IDEX-0002-02 (IdEx Unistra), ANR-10-IAHU-02 (IHU Strasbourg) and ANR-11-INBS-0006 (FLI).

References

- [1] S. Amin, M. Andriluka, M. Rohrbach, and B. Schiele. Multi-view pictorial structures for 3d human pose estimation. In *British Machine Vision Conference (BMVC)*, September 2013.
- [2] S. Amin, P. Müller, A. Bulling, and M. Andriluka. Test-time adaptation for 3d human pose estimation. In X. Jiang, J. Hornegger, and R. Koch, editors, *Pattern Recognition*, volume 8753 of *Lecture Notes in Computer Science*, pages 253–264. Springer International Publishing, 2014.
- [3] M. Andriluka, L. Pishchulin, P. Gehler, and B. Schiele. 2D human pose estimation: New benchmark and state of the art analysis. In *2014 IEEE Conference on Computer Vision and Pattern Recognition*, pages 3686–3693, June 2014.
- [4] A. Baak, M. Müller, G. Bharaj, H. Seidel, and C. Theobalt. A data-driven approach for real-time full body pose reconstruction from a depth camera. In *Proceedings of the International Conference on Computer Vision (ICCV)*, pages 1092–1099. IEEE, 2011.
- [5] V. Belagiannis, S. Amin, M. Andriluka, B. Schiele, N. Navab, and S. Ilic. 3d pictorial structures for multiple human pose estimation. In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*, pages 1669–1676. IEEE, 2014.
- [6] V. Belagiannis, X. Wang, B. Schiele, P. Fua, S. Ilic, and N. Navab. Multiple human pose estimation with temporally consistent 3D pictorial structures. In *ChaLearn Looking at People Workshop, European Conference on Computer Vision (ECCV2014)*. IEEE, September 2014.
- [7] V. Belagiannis, X. Wang, H. B. B. Shitrit, K. Hashimoto, R. Stauder, Y. Aoki, M. Krantzfelder, A. Schneider, P. Fua, S. Ilic, H. Feussner, and N. Navab. Parsing human skeletons in an operating room. *Machine Vision and Applications*, pages 1–12, 2016.
- [8] M. Burenius, J. Sullivan, and S. Carlsson. 3d pictorial structures for multiple view articulated pose estimation. In *CVPR*, pages 3618–3625. IEEE, 2013.
- [9] K. Buys, C. Cagniard, A. Baksheev, T. D. Laet, J. D. Schutter, and C. Pantofaru. An adaptable system for rgb-d based human body detection and pose estimation. *Journal of Visual Communication and Image Representation*, 2013.
- [10] M. Eichner and V. Ferrari. Better appearance models for pictorial structures. In *British Machine Vision Conference, BMVC 2009, London, UK, September 7-10, 2009. Proceedings*, pages 1–11, 2009.
- [11] M. Eichner, M. Marin-Jimenez, A. Zisserman, and V. Ferrari. 2d articulated human pose estimation and retrieval in (almost) unconstrained still images. *International Journal of Computer Vision*, 99(2):190–214, 2012.
- [12] P. F. Felzenszwalb and D. P. Huttenlocher. Pictorial structures for object recognition. *International Journal of Computer Vision*, 61(1):55–79, 2005.
- [13] M. A. Fischler and R. A. Elschlager. The representation and matching of pictorial structures. *IEEE Transactions on Computers*, C-22(1):67–92, Jan 1973.
- [14] J. Gall, B. Rosenhahn, T. Brox, and H.-P. Seidel. Optimization and filtering for human motion capture. *International Journal of Computer Vision*, 87(1):75–92, 2010.
- [15] G. Gkioxari, B. Hariharan, R. Girshick, and J. Malik. Using k-poselets for detecting people and localizing their keypoints. In *Proceedings of the 2014 IEEE Conference on Computer Vision and Pattern Recognition, CVPR '14*, pages 3582–3589. IEEE, 2014.
- [16] D. Gowsikhaa, S. Abirami, and R. Baskaran. Automated human behavior analysis from surveillance videos: a survey. *Artificial Intelligence Review*, 42(4):747–765, 2014.
- [17] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *CoRR*, abs/1512.03385, 2015.
- [18] M. Hofmann and D. M. Gavrilu. Multi-view 3D human pose estimation in complex environment. *International Journal of Computer Vision*, 96(1):103–124, 2011.
- [19] E. Insafutdinov, L. Pishchulin, B. Andres, M. Andriluka, and B. Schiele. *DeeperCut: A Deeper, Stronger, and Faster Multi-person Pose Estimation Model*, pages 34–50. Springer, Cham, 2016.
- [20] A. Jaimes and N. Sebe. Multimodal human-computer interaction: A survey. *Computer Vision and Image Understanding*, 108(12):116 – 134, 2007. Special Issue on Vision for Human-Computer Interaction.
- [21] A. Jain, J. Tompson, M. Andriluka, G. W. Taylor, and C. Bregler. Learning human pose estimation features with convolutional networks. In *International Conference on Learning Representations (ICLR)*, April 2014.
- [22] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. *arXiv preprint arXiv:1408.5093*, 2014.
- [23] A. Kadkhodamohammadi, A. Gangi, M. de Mathelin, and N. Padoy. Articulated clinician detection using 3D pictorial structures on RGB-D data. *Medical Image Analysis*, 35:215 – 224, 2017.
- [24] N. Komodakis, G. Tziritas, and N. Paragios. Performance vs computational efficiency for optimizing single and dynamic mrfs: Setting the state of the art with primal-dual strategies. *Computer Vision and Image Understanding*, 112(1):14 – 29, 2008. Special Issue on Discrete Optimization in Computer Vision.
- [25] X. Luo, B. Berendsen, R. T. Tan, and R. C. Veltkamp. Human pose estimation for multiple persons based on volume reconstruction. In *Pattern Recognition (ICPR), 2010 20th International Conference on*, pages 3591–3594, Aug 2010.
- [26] A. Newell, K. Yang, and J. Deng. *Stacked Hourglass Networks for Human Pose Estimation*, pages 483–499. Springer, 2016.
- [27] L. Pishchulin, E. Insafutdinov, S. Tang, B. Andres, M. Andriluka, P. Gehler, and B. Schiele. DeepCut: Joint subset partition and labeling for multi person pose estimation. In *CVPR*, 2016.
- [28] B. Sapp, D. Weiss, and B. Taskar. Parsing human motion with stretchable models. In *CVPR*, pages 1281–1288. IEEE, 2011.
- [29] J. Schmidhuber. Deep learning in neural networks: An overview. *Neural Networks*, 61:85 – 117, 2015.

- [30] J. Shotton, R. Girshick, A. Fitzgibbon, T. Sharp, M. Cook, M. Finocchio, R. Moore, P. Kohli, A. Criminisi, A. Kipman, and A. Blake. Efficient human pose estimation from single depth images. *PAMI*, 35(12):2821 – 2840, 2012.
- [31] L. Sigal, A. O. Balan, and M. J. Black. Humaneva: Synchronized video and motion capture dataset and baseline algorithm for evaluation of articulated human motion. *International Journal of Computer Vision*, 87(1):4–27, 2009.
- [32] M. Sun and S. Savarese. Articulated part-based model for joint object detection and pose estimation. In *Proceedings of the 2011 International Conference on Computer Vision, ICCV '11*, pages 723–730. IEEE, 2011.
- [33] R. Tokola, W. Choi, and S. Savarese. Breaking the chain: liberation from the temporal markov assumption for tracking human poses. In *Proceedings of the International Conference on Computer Vision*, 2013.
- [34] J. Tompson, R. Goroshin, A. Jain, Y. LeCun, and C. Bregler. Efficient object localization using convolutional networks. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 648–656, June 2015.
- [35] J. Tompson, A. Jain, Y. LeCun, and C. Bregler. Joint training of a convolutional network and a graphical model for human pose estimation. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 1799–1807. Curran Associates, Inc., 2014.
- [36] A. Toshev and C. Szegedy. DeepPose: Human pose estimation via deep neural networks. In *CVPR*, 2014.
- [37] Z. Tu and X. Bai. Auto-context and its application to high-level vision tasks and 3D brain image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(10):1744–1757, 2010.
- [38] S.-E. Wei, V. Ramakrishna, T. Kanade, and Y. Sheikh. Convolutional pose machines. In *CVPR*, 2016.
- [39] W. Xu, P. c. Su, and S. c. S. Cheung. Human pose estimation using two rgb-d sensors. In *2016 IEEE International Conference on Image Processing (ICIP)*, pages 1279–1283, Sept 2016.
- [40] W. Yang, W. Ouyang, H. Li, and X. Wang. End-to-end learning of deformable mixture of parts and deep convolutional neural networks for human pose estimation. In *CVPR*. IEEE, 2016.
- [41] Y. Yang and D. Ramanan. Articulated human detection with flexible mixtures of parts. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(12):2878–2890, 2013.