

Real-time marker-less multi-person 3D pose estimation in RGB-Depth camera networks

Marco Carraro¹, Matteo Munaro¹, Jeff Burke² and Emanuele Menegatti¹

Abstract—This paper proposes a novel system to estimate and track the 3D poses of multiple persons in calibrated RGB-Depth camera networks. The multi-view 3D pose of each person is computed by a central node which receives the single-view outcomes from each camera of the network. Each single-view outcome is computed by using a CNN for 2D pose estimation and extending the resulting skeletons to 3D by means of the sensor depth. The proposed system is marker-less, multi-person, independent of background and does not make any assumption on people appearance and initial pose. The system provides real-time outcomes, thus being perfectly suited for applications requiring user interaction. Experimental results show the effectiveness of this work with respect to a baseline multi-view approach in different scenarios. To foster research and applications based on this work, we released the source code in OpenPTrack, an open source project for RGB-D people tracking.

I. INTRODUCTION

The human body pose is rich of information. Many algorithms and applications, such as Action Recognition [1], [2], [3], People Re-identification [4], Human-Computer-Interaction (HCI) [5] and Industrial Robotics [6], [7], [8] rely on this type of data. The recent availability of smart cameras [9], [10], [11] and affordable RGB-Depth sensors as the first and second generation Microsoft Kinect, allow to estimate and track body poses in a cost-efficient way. However, using a single sensor is often not reliable enough because of occlusions and Field-of-View (FOV) limitations. For this reason, a common solution is to take advantage of camera networks. Nowadays, the most reliable way to perform human Body Pose Estimation (BPE) is to use marker-based motion capture systems. These systems show great results in terms of accuracy (less than 1mm), but they are very expensive and require the users to wear many markers, thus imposing heavy limitations to their diffusion. Moreover, these systems usually require offline computations in complicated scenarios with many markers and people, while the system we propose provides immediate results. A real-time response is usually needed in security applications, where person actions should be detected in time, or in industrial applications, where human motion is predicted to prevent collisions with robots in shared workspaces. Aimed by those reasons, the research on marker-less motion capture systems has been particularly active in recent years.

¹Marco Carraro, Matteo Munaro and Emanuele Menegatti are with the Intelligent Autonomous Systems Laboratory (IAS-Lab), Department of Information Engineering (DEI), University of Padova, Via Ognissanti 72, 35129, Padova, Italy. {marco.carraro, emg}@dei.unipd.it

²Jeff Burke is with REMAP in the School of Theater, Film and Television at UCLA, Los Angeles, California, USA 90095

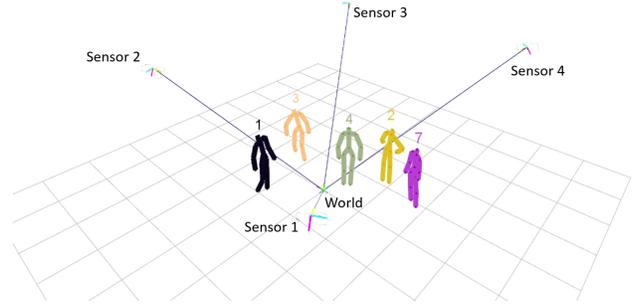


Fig. 1: The output provided by the system we are proposing. In this example, five persons are seen from a network composed of four Microsoft Kinect v2.

In this work, we propose a novel system to estimate the 3D human body pose in real-time. To the best of our knowledge, this is the first open-source and real-time solution to the multi-view, multi-person 3D body pose estimation problem. Figure 1 depicts our system output. The system relies on the feed of multiple RGB-D sensors (from 1 to N) placed in the scene and on an extrinsic calibration of the network. In this work, this calibration is performed with the `calibration_toolkit` [12]¹. The multi-view poses are obtained by fusing the single view outcomes of each detector, that runs a state-of-the-art 2D body pose estimator [13], [14] and extend it to 3D by means of the sensor depth. The contribution of the paper is two-fold: i) we propose a novel system to fuse and update 3D body poses of multiple persons in the scene and ii) we enriched a state-of-the-art single-view 2D pose estimation algorithm to provide 3D poses. As a further contribution, the code of the project has been released as open-source as part of the OpenPTrack [15], [16] repository. The proposed system is:

- *multi-view*: The fused poses are computed taking into account the different poses of the single-view detectors;
- *asynchronous*: The fusion algorithm does not require the different sensors to be synchronous or have the same frame rate. This allows the user to choose the detector computing node accordingly to his needs and possibilities;
- *multi-person*: The system does not make any assumption on the number of persons in the scene. The overhead due to the different number of persons is negligible;

¹https://github.com/iaslab-unipd/calibration_toolkit

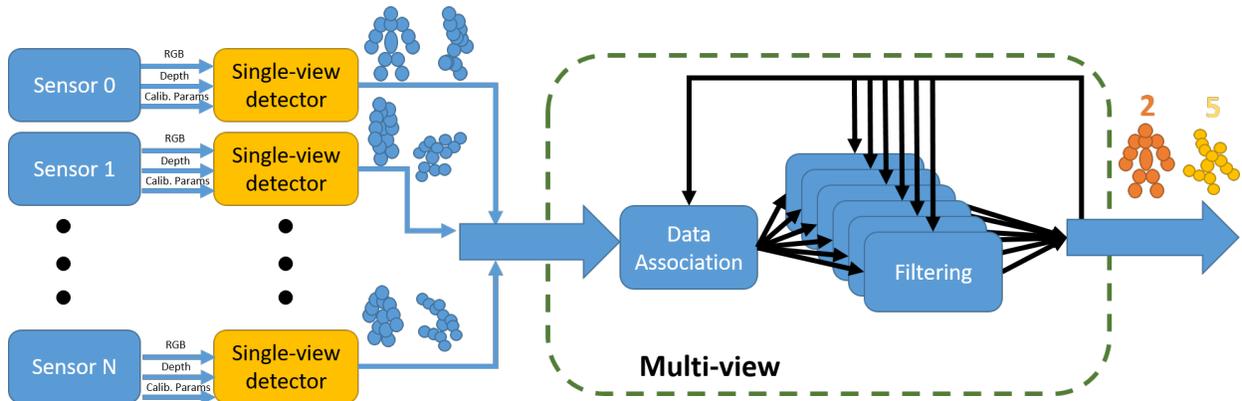


Fig. 2: The system overview. The camera network is composed of several RGB-D sensors (from 1 to N). Each single-view detector takes the RGB and Depth images as input and computes the 3D skeletons of the people in the scene as the output using the calibration parameters K . The information is then sent to the multi-view central node which is in charge of computing the final pose estimation for each person in the scene. First, a data association is performed to determine which pose detection is belonging to which pose track, then a filtering step is performed to update the pose track given the detection.

- *scalable*: No assumptions are made on the number or positions of the cameras. The only request is an offline one-time extrinsic calibration of the network;
- *real-time*: The final pose framerate is linear to the number of cameras in the network. In our experiments, a single-camera network can provide from 5 fps to 15 fps depending on the Graphical Processing Unit (GPU) exploited by the detector. The final framerate of a camera network composed of k nodes is the sum of their single-view framerate;
- *low-cost*: The system relies on affordable low-cost RGB-D sensors controlled by consumer GPU-enabled computers. No specific hardware is required.

The remainder of the paper is organized as follows: in Section II we review the literature regarding human BPE from single and multiple views, while Section III describes our system and the approach used to solve the problem. In Section IV experimental results are presented, and, finally in Section V we present our final conclusions.

II. RELATED WORK

A. Single-view body pose estimation

Since a long time, there have been a great interest about single-view human BPE, in particular for gaming purposes or avatar animation. Recently, the advent of affordable RGB-D sensors boosted the research in this and other Computer Vision fields. Shotton et al. [17] proposed the skeletal tracking system licensed by Microsoft used by the XBOX console with the first-generation Kinect. This approach used a random forest classifier to classify the different pixels as belonging to the different body parts. This work inspired an open-source approach that was released by Buys et al. [18]. This same work was then improved by adding the OpenPTrack people detector module as a preprocessing step [19]. Still, the performance of the detector remained

very poor for non frontal persons. In these last years, many challenging Computer Vision problems have been finally resolved by using *Convolutional Neural Networks* (CNNs) solutions. Also single-view BPE has seen a great benefit from these techniques [20], [21], [22], [14]. The impressive pose estimation quality provided by those solution is usually paid in terms of computational time. Nevertheless, this limitation is going to be leveraged with newer network layouts and Graphical Processing Units (GPU) architectures, as proved by some recent works [22], [14]. In particular, the work of Cao et. al [14] was one of the first to implement a CNN solution to solve people BPE in real-time using a bottom-up approach. The authors were able to compute 2D poses for all the people in the scene with a single forward pass of their CNN. This work has been adopted here as part of our single-view detectors.

B. Multi-view body pose estimation

Multiple views can be exploited to be more robust against occlusions, self-occlusions and FOV limitations. In [23] a Convolutional Neural Network (CNN) approach is proposed to estimate the body poses of people by using a low number of cameras also in outdoor scenarios. The solution combines a generative and discriminative approach, since they use a CNN to compute the poses which are driven by an underlying model. For this reason, the collaboration of the users is required for the initialization phase. In our previous work [19], we solved the single-person human BPE by fusing the data of the different sensors and by applying an improved version of [18] to a virtual depth image of the frontalized person. In this way, the skeletonization is only performed once, on the virtual depth map of the person in frontal pose. In [24], a 3D model is registered to the point clouds of two Kinects. The work provides very accurate results, but it is computationally expensive and not scalable to multiple persons. The authors of [25] proposed a pure geometric approach to infer the

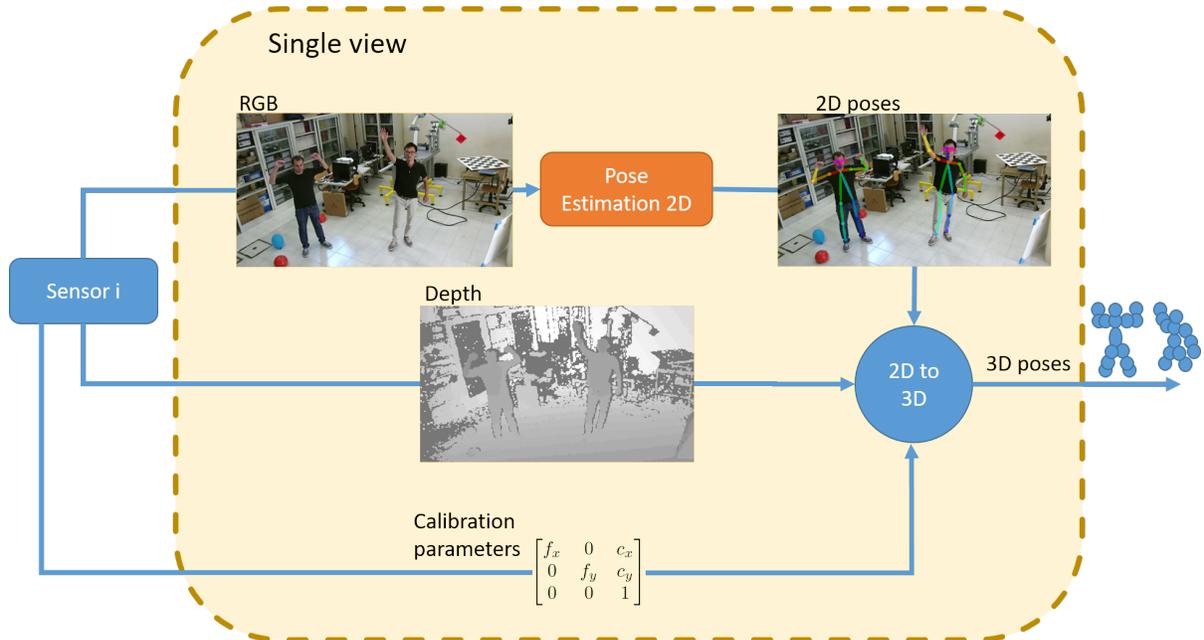


Fig. 3: The single-view pipeline followed for each sensor. At each new frame composed of a color image (RGB), a depth image and the calibration parameters, the 3D pose of each person in the scene is computed from the 2D one. Then, the results are sent to the central computer which will compute the multi-view result.

multi-view pose from a synchronous set of 2D single-view skeletons obtained using [26]. The third dimension is computed by imposing a set of algebraic constraints from the triangulation of the multiple views. The final skeleton is then computed by solving a least square error method. While the method is computationally promising (skeleton computed in 1s per set of synchronized images with an unoptimized version of the code), it does not scale with the number of persons in the scene. In [27] a system composed of common RGB cameras and RGB-D sensors are used together to record a dance motion performed by a user. The fusion method is obtained by selecting the best skeleton match between the different ones obtained by using a probabilistic approach with a particle filter. The system performs well enough for its goal, but it does not scale to multiple people and requires an expensive setup. In [28] the skeletons obtained from the single images are enriched with a 3D model computed with the visual hull technique. In [29] two orthogonal Kinects are used to improve the single-view outcome of both sensors. They used a constrained optimization framework with the bone lengths as hard constraints. While the work provides a real-time solution and there are no hard assumption on the Kinect positions, it was tested just with one person and two orthogonal Kinect sensors. Similarly to many recent works [25], [28], [27], we use a single-view state-of-the-art body pose estimator, but we augment this result with 3D data and we then combine the multiple views to improve the overall quality.

III. SYSTEM DESIGN

Figure 2 shows an overview of the proposed system. It can be split into two parts: i) the single view, which

is the same for each sensor and it is executed locally and ii) the multi-view part which is executed just by the master computer. In the single-view part (see Figure 3), each detector estimates the 2D body pose of each person in the scene using an open-source state-of-the-art single-view body pose estimator. In this work, we use the OpenPose²[13], [14] library, but the overall system is totally independent of the single-view algorithm used. The last operation made by the detector is to compute the 3D positions of each joint returned by OpenPose. This fusion is done by exploiting the depth information coming from the RGB-D sensor used. The 3D skeleton is then sent to the master computer for the fusion phase. This is done by means of multiple Unscented Kalman Filters used on the detection feeds, as explained in Section III-C.

A. Camera Network setup

The camera network can be composed of several RGB-D sensors. In order to know the relative position of each camera, we calibrate the system using a solution similar to our previous works [16], [15]. From this passage we fix a common *world* reference frame \mathcal{W} and we obtain a transformation $\mathcal{T}_C^{\mathcal{W}}$, for each camera C in the network, which transforms points in the camera coordinate system to the *world* reference system.

B. Single-view Estimation of 3D Poses

Each node in the network is composed of an RGB-D sensor and a computer to elaborate the images. Let ${}^{\text{ri}}\mathfrak{F} =$

²<https://github.com/CMU-Perceptual-Computing-Lab/openpose>

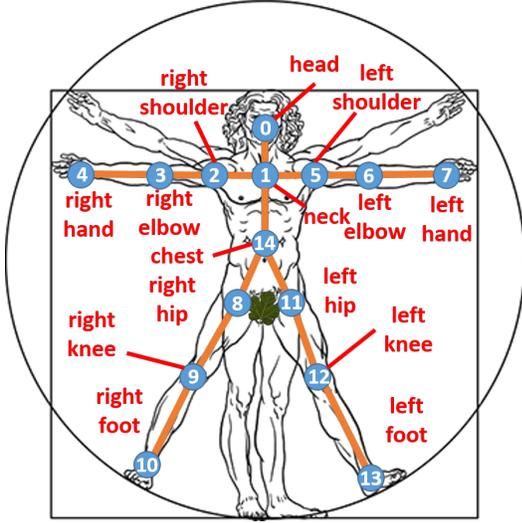


Fig. 4: The human model used in this work.

$\{^R C, ^R D\}$ be a frame captured by the detector R and composed of the color image C and the depth image D all in the R reference frame. The color and depth images in \mathfrak{F} are considered as synchronized. We then apply *OpenPose* to $^R C$ obtaining the raw two dimensional skeletons $\overline{\mathfrak{S}} = \{\overline{S}_0, \overline{S}_1, \dots, \overline{S}_k\}$. Each $S = \{j_i | 0 \leq i \leq m\} \in \overline{\mathfrak{S}}$ is a set of 2D joints which follows the human model depicted in Figure 4. The goal of the single-view detector is to transform $\overline{\mathfrak{S}}$ in the set of skeletons $\widehat{\mathfrak{S}} = \{\widehat{S}_0, \widehat{S}_1, \dots, \widehat{S}_k\}$ where each $\widehat{S} \in \widehat{\mathfrak{S}}$ is a three dimensional skeleton. Given the RGB image I , let's consider a point $p = (x_p, y_p) \in I$ and its corresponding depth $d = proj(x_p, y_p)$. Considering (f_x, f_y) and (c_x, c_y) respectively the focal point and the optical center of the sensor, the relationship to compute the 3D point $P_R = (X_R, Y_R, Z_R)$ in the camera reference system R is explained in Equation 1.

$$p = \begin{bmatrix} x_p \\ y_p \\ d \end{bmatrix} = \begin{bmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} X_R \\ Y_R \\ Z_R \end{bmatrix} = K P_R \quad (1)$$

Since the depth data is potentially noisy or missing, we compute the depth d associated to the point $p = (x_p, y_p)$ by applying a median to the set $\mathcal{D}(p)$, as shown in Equations 2, 3.

$$\mathcal{D}(p = (x_p, y_p)) = \{(x, y) | \|(x, y) - (x_p, y_p)\| < \epsilon\} \quad (2)$$

$$d = \phi(p) = \text{median}\{proj(x, y) | (x, y) \in \mathcal{D}(p)\} \quad (3)$$

Given $\overline{\mathfrak{S}}$, we then proceed to the calculation of $\widehat{\mathfrak{S}}$ as shown in Equation 4.

$$\forall 0 \leq j < k, \quad \overline{S}_j = \{\overline{j}_i = (x_i, y_i) | 0 \leq i < m\} \in \overline{\mathfrak{S}},$$

$$\widehat{S}_j = \left\{ \widehat{j}_i = \begin{bmatrix} |K^{-1}(\overline{j}_i)|_x \\ |K^{-1}(\overline{j}_i)|_y \\ \phi(\overline{j}_i) \end{bmatrix}, 0 \leq i < m \right\} \in \widehat{\mathfrak{S}} \quad (4)$$

Algorithm 1 The algorithm performed by the master computer to decide the association between the different skeletons in a detection and the current tracks.

INPUT:

- ${}^W \widehat{\mathfrak{S}}_i = \{S_0, S_1, \dots, S_{k-1}\}$ - a new detection set from sensor i in the world reference frame
- $\mathfrak{T} = \{T_0, T_1, \dots, T_{l-1}\}$ - the current set of tracked persons pose.
- ϵ - maximum distance for a detection to be considered for the association

OUTPUT:

- $\mathcal{M} = \{(S_i, T_j) \in {}^W \widehat{\mathfrak{S}}_i \times \mathfrak{T}\}$ - the association between the pose tracked and the new observations
- $\mathcal{N} \subseteq {}^W \widehat{\mathfrak{S}}_i$ - the detections without an association. They will initialize a new track.
- $\mathfrak{T}_o \subseteq \mathfrak{T}$ - the tracks without an associated observations. They will be considered for removal

```

1: procedure DATA_ASSOCIATION( ${}^W \widehat{\mathfrak{S}}_i, \mathfrak{T}, \epsilon$ )
2:    $\mathfrak{T}_o \leftarrow \emptyset$ 
3:    $C \leftarrow \mathbf{0}_{k \times l}$ 
4:   for each  $T_i \in \mathfrak{T}$  do
5:     for each  $S_j \in {}^W \widehat{\mathfrak{S}}_i$  do
6:        $x_t(j) \leftarrow \text{centroid}(S_j)$ 
7:        $z_t(i, j) \leftarrow *v$  that  $T_i$  would have if  $S_j$  were
         associated to it*
8:        $\widehat{z}_{t|t-1}(i) \leftarrow *$ prediction step of  $\mathcal{K}_{im}*$ 
9:        $\Sigma_t(i) \leftarrow \Sigma_t(\mathcal{K}_{im})$ 
10:       $\tilde{z}_t(i, j) = z_k(i, j) - \widehat{z}_{t|t-1}(i)$ 
11:       $C_{ij} \leftarrow \tilde{z}_t^T(i, j) \cdot \Sigma_t(i)^{-1} \cdot \tilde{z}_t(i, j)$ 
12:       $X \leftarrow \text{solve\_Munkres}(C)$ 
13:      for  $i \in [0, l - 1]$  do
14:        for  $j \in [i + 1, k - 1]$  do
15:          if  $X_{ij} == 1$  and  $C_{ij} < \epsilon$  then
16:             $\mathcal{M} \leftarrow \mathcal{M} \cup \{(S_j, T_i)\}$ 
17:            * update  $\mathcal{K}_{im}$  with  $S_j$  *
18:       $\mathcal{N} \leftarrow \{S_i | \nexists T_j, (S_i, T_j) \in \mathcal{M}\}$ 
19:       $\mathfrak{T}_o \leftarrow \{T_i | \nexists S_j, (S_j, T_i) \in \mathcal{M}\}$ 
20:      return  $\mathcal{M}, \mathcal{N}, \mathfrak{T}_o$ 

```

C. Multi-view fusion of 3D poses

The master computer is in charge of fusing the different information it is receiving from the single-view detectors in the network. One of the common limitations in motion capture systems is the necessity to have synchronized cameras. Moreover, off-the-shelves RGB-D sensors, such as the Microsoft Kinect v2, do not have the possibility to trigger the image acquisition. In order to overcome this limitation, our solution merges the different data streams asynchronously. This allows the system to work also with other RGB-D sensors or other low-cost embedded machine. At time t , the master computer maintains a set of tracks $\mathfrak{T} = \{T_0, T_1, \dots, T_l\}$ where each pose tracked T_i is composed of the set of states of m different Kalman Filters, one per

		r-shoulder	r-elbow	r-wrist	l-shoulder	l-elbow	l-wrist	r-hip	r-knee	r-ankle	l-hip	l-knee	l-ankle
single-camera network	MAF ₃₀	>100	>100	>100	>100	>100	>100	>100	>100	>100	>100	>100	>100
	MAF ₄₀	>100	>100	>100	>100	>100	>100	>100	>100	>100	>100	>100	>100
	Ours	54.9 ± 58.6	42.4 ± 47.4	42.4 ± 40.0	77.7 ± 74.4	79.1 ± 82.7	70.0 ± 61.8	51.7 ± 43.7	54.5 ± 31.0	63.3 ± 34.2	97.8 ± 30.3	57.5 ± 38.9	69.2 ± 37.6
2-camera network	MAF ₃₀	62.0 ± 33.0	62.9 ± 32.0	63.1 ± 34.5	83.3 ± 33.4	85.8 ± 37.8	94.8 ± 45.4	76.4 ± 30.6	75.9 ± 27.4	88.3 ± 35.6	>100	85.4 ± 35.5	93.3 ± 37.0
	MAF ₄₀	83.7 ± 41.8	84.0 ± 40.9	83.1 ± 43.7	>100	>100	>100	99.2 ± 40.4	96.3 ± 38.0	>100	>100	>100	>100
	Ours	20.7 ± 17.2	21.0 ± 17.5	24.3 ± 17.5	32.1 ± 23.0	33.4 ± 26.3	39.8 ± 35.1	22.4 ± 16.7	42.8 ± 17.2	59.7 ± 28.6	98.3 ± 21.2	39.9 ± 18.3	58.6 ± 27.1
4-camera network	MAF ₃₀	28.7 ± 16.4	31.0 ± 16.9	32.2 ± 22.5	41.5 ± 17.9	39.9 ± 19.6	44.7 ± 29.5	40.2 ± 15.0	48.7 ± 12.8	58.6 ± 21.2	94.1 ± 26.1	52.1 ± 17.8	57.8 ± 27.9
	MAF ₄₀	38.4 ± 21.2	40.8 ± 21.7	41.6 ± 26.3	53.0 ± 23.2	52.7 ± 24.6	57.6 ± 33.1	50.7 ± 19.4	56.2 ± 16.7	66.0 ± 24.5	96.6 ± 30.8	61.2 ± 23.1	67.5 ± 31.6
	Ours	22.7 ± 18.9	21.3 ± 18.5	26.3 ± 19.9	22.5 ± 22.1	26.7 ± 25.9	31.8 ± 29.7	23.9 ± 18.0	46.5 ± 19.7	55.9 ± 25.1	95.4 ± 22.0	45.1 ± 20.5	49.1 ± 25.2

TABLE I: The results of the experiments. Each number represents the mean and the standard deviation of the reprojection error on a reference camera (see Equation 5)

each joint, i.e: $T_i = \{\mathcal{S}(\mathcal{K}_{i0}), \mathcal{S}(\mathcal{K}_{i1}), \dots, \mathcal{S}(\mathcal{K}_{im})\}$. The additional Kalman Filter \mathcal{K}_{im} is maintained for the data association algorithm. At time $t+1$, it may arrive a detection $\widehat{\mathcal{S}}_i = \{\widehat{S}_0, \widehat{S}_1, \dots, \widehat{S}_k\}$ from the sensor i of the network. The master computer first refers the detection to the common *world* coordinate system \mathcal{W} (see Section III-A):

$$\mathcal{W}\widehat{\mathcal{S}}_i = \mathcal{T}_i^{\mathcal{W}} \cdot \widehat{\mathcal{S}}_i = \{\mathcal{T}_i^{\mathcal{W}} \cdot S_j \mid \forall S_j \in \widehat{\mathcal{S}}_i\}$$

Then, it associates the different skeletons in $\mathcal{W}\widehat{\mathcal{S}}_i$ as new observations for the different tracks in \mathcal{T} if they belong to them or initializes new tracks if some of the skeletons do not belong to any $T_i \in \mathcal{T}$. At this stage, the system also decides if a track is old and has to be removed from \mathcal{T} . This step is important to prevent \mathcal{T} to grow big causing time computing problems with systems which are running for hours. We refer to this phase as *data association*. Algorithm 1 shows how it is performed. The data association is done by considering the centroid of each skeleton S contained in the detection $\mathcal{W}\widehat{\mathcal{S}}_i$. The centroid is calculated as the chest joint $j_{14} \in S$, if this is valid, otherwise it is replaced with a weighted mean of the neighbor joints. Lines [6-9] of Algorithm 1 refers to the calculation of a cost associated to the case if the detection pose S_j would be associated to the track T_i . To calculate this, we consider the Mahalanobis distance between the likelihood vector at time t $\tilde{z}_t(i, j)$ and $\Sigma_t(\mathcal{K}_{i, x_t})$: the covariance matrix of the Kalman filter associated to the centroid of T_i . At this point, computing the optimal association between tracks and detections is the same as solving the Hungarian algorithm associated to the cost matrix C ; Line 11 refers to the use of the Munkres algorithm which efficiently computes the optimal matrix X with a 1 on the associated couples. Nevertheless, this algorithm does not consider a maximum distance between tracks and detections. Thus, it may happen that a couple is wrongly associated in the optimal assignment. For this reason, when inserting the couples in \mathcal{M} , we check also if the cost of the couple in the initial cost matrix C is below a threshold.

Once solved the data association problem, we can assign the tracks ID to the different skeletons. Indeed, we know which are the detection at the current time t belonging to the tracks in the system and, additionally, we know also which tracks need to be created (i.e. new detections with no associated track) and the tracks to consider for the removal. Let n be the number of people in the scene, we used a set of Unscented Kalman Filters $\mathcal{R} = \{\mathcal{K}_{ij}, 0 \leq i < n, 0 \leq j \leq m\}$ where the generic $\mathcal{K}_{ij} \in \mathcal{R}$ is in charge of computing the new position of the joint j of the person i at time t ,

given the new detection received from one of the detectors at time t and the prediction of the filter \mathcal{K}_{ij} computed from the previous position at time $t-1$ of the same joint j .

The state of each Kalman Filter \mathcal{K}_{ij} is dimensioned with the three dimensional position of the joint j . We used as motion model a constant velocity model, since it is good to predict joint movements in the small temporal space between two good detections of that joint.

IV. EXPERIMENTS

The algorithm described in this paper does not require any synchronization between the cameras in the networks. This fact makes particularly difficult to find a fair comparison between our proposed system and other state-of-the-art works. Thus, in order to provide useful indication on how our system performs, we recorded and manually annotated a set of RGB-D frames while a person was freely moving in the field-of-view of a 4-sensors camera network. We compare our algorithm with a baseline method called MAF (Moving Average Filter), in which the outcome of the generic joint i at time t is computed as an average of the last k frames. In order to be as fair as possible, we fixed $k \geq 30$ to provide comparable results in terms of smoothness. We also demonstrated the effectiveness of the multi-view fusion by comparing our results with the poses obtained by considering just one and two cameras of the same network. In this comparison, we report the average reprojection error with respect to one of the cameras, C_0 . Equation 5 shows how this error is calculated with ${}^{\mathcal{W}}P$ as the generic joint expressed in the world reference system and p^* as the corresponding ground truth :

$$e_{\text{repr}} = |p^* - K \cdot \mathcal{T}_{\mathcal{W}}^{C_0} \cdot {}^{\mathcal{W}}P| \quad (5)$$

Table I shows the results we achieved. As depicted, the proposed method outperforms the baseline in all the cases: single-view, 2-camera network and 4-camera network. In the first two cases (single and 2-camera network) the improvement is from 50% to 60%, while, when multiple views are available, it is from 18% to 32%. It is also interesting to note that the most noisy joints are the ones relative to the legs as confirmed by other state-of-the-art works [14], [20], [21].

A. Implementation Details

The system has been implemented and tested with Ubuntu 14.04 and Ubuntu 16.04 operating system using the Robot Operating System (ROS) [30] middleware. The code is entirely written in C++ using the Eigen, OpenCV and PCL libraries.

V. CONCLUSIONS AND FUTURE WORKS

In this paper we presented a framework to compute the 3D body pose of each person in a RGB-D camera network using only its extrinsic calibration as a prior. The system does not make any assumption on the number of cameras, on the number of persons in the scene, on their initial poses or clothes and does not require the cameras to be synchronous. In our experimental setup we demonstrated the validity of our system over both single-view and multi-view approaches. In order to provide the best service to the Computer Vision community and to provide also a future baseline method to other researchers, we released the source code under the BSD license as part of the OpenPTrack library³. As future works, we plan to add a human dynamic model to guide the prediction of the Kalman Filters to further improve the performance achievable by our system (in particular for the lower joints) and to further validate the proposed system on a new RGB-Depth dataset annotated with the ground truth of the single links of the persons' body pose. The ground truth will be provided by a marker based commercial motion capture system.

ACKNOWLEDGEMENT

This work was partially supported by U.S. National Science Foundation award IIS-1629302

REFERENCES

- [1] F. Han, X. Yang, C. Reardon, Y. Zhang, and H. Zhang, "Simultaneous feature and body-part learning for real-time robot awareness of human behaviors," 2017.
- [2] M. Zanfir, M. Leordeanu, and C. Sminchisescu, "The moving pose: An efficient 3d kinematics descriptor for low-latency action recognition and detection," in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2752–2759, 2013.
- [3] C. Wang, Y. Wang, and A. L. Yuille, "An approach to pose-based action recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 915–922, 2013.
- [4] S. Ghidoni and M. Munaro, "A multi-viewpoint feature-based re-identification system driven by skeleton keypoints," *Robotics and Autonomous Systems*, vol. 90, pp. 45–54, 2017.
- [5] A. Jaimes and N. Sebe, "Multimodal human-computer interaction: A survey," *Computer vision and image understanding*, vol. 108, no. 1, pp. 116–134, 2007.
- [6] C. Morato, K. N. Kaipa, B. Zhao, and S. K. Gupta, "Toward safe human robot collaboration by using multiple kinects based real-time human tracking," *Journal of Computing and Information Science in Engineering*, vol. 14, no. 1, p. 011006, 2014.
- [7] S. Michieletto, F. Stival, F. Castelli, M. Khosravi, A. Landini, S. Ellero, R. LandÄs, N. Boscolo, S. Tonello, B. Varaticeanu, C. Nicolescu, and E. Pagello, "Flexicoil: Flexible robotized coils winding for electric machines manufacturing industry," in *ICRA workshop on Industry of the future: Collaborative, Connected, Cognitive*, 2017.
- [8] F. Stival, S. Michieletto, and E. Pagello, "How to deploy a wire with a robotic platform: Learning from human visual demonstrations," in *FAIM 2017*, 2017.
- [9] Z. Zivkovic, "Wireless smart camera network for real-time human 3d pose reconstruction," *Computer Vision and Image Understanding*, vol. 114, no. 11, pp. 1215–1222, 2010.
- [10] M. Carraro, M. Munaro, and E. Menegatti, "A powerful and cost-efficient human perception system for camera networks and mobile robotics," in *International Conference on Intelligent Autonomous Systems*, pp. 485–497, Springer, Cham, 2016.
- [11] M. Carraro, M. Munaro, and E. Menegatti, "Cost-efficient rgb-d smart camera for people detection and tracking," *Journal of Electronic Imaging*, vol. 25, no. 4, pp. 041007–041007, 2016.
- [12] F. Basso, R. Levorato, and E. Menegatti, "Online calibration for networks of cameras and depth sensors," in *OMNIVIS: The 12th Workshop on Non-classical Cameras, Camera Networks and Omnidirectional Vision-2014 IEEE International Conference on Robotics and Automation (ICRA 2014)*, 2014.
- [13] S.-E. Wei, V. Ramakrishna, T. Kanade, and Y. Sheikh, "Convolutional pose machines," in *CVPR*, 2016.
- [14] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh, "Realtime multi-person 2d pose estimation using part affinity fields," in *CVPR*, 2017.
- [15] M. Munaro, A. Horn, R. Illum, J. Burke, and R. B. Rusu, "Opentrack: People tracking for heterogeneous networks of color-depth cameras," in *IAS-13 Workshop Proceedings: 1st Intl. Workshop on 3D Robot Perception with Point Cloud Library*, pp. 235–247, 2014.
- [16] M. Munaro, F. Basso, and E. Menegatti, "Opentrack: Open source multi-camera calibration and people tracking for rgb-d camera networks," *Robotics and Autonomous Systems*, vol. 75, pp. 525–538, 2016.
- [17] J. Shotton, T. Sharp, A. Kipman, A. Fitzgibbon, M. Finocchio, A. Blake, M. Cook, and R. Moore, "Real-time human pose recognition in parts from single depth images," *Communications of the ACM*, vol. 56, no. 1, pp. 116–124, 2013.
- [18] K. Buys, C. Cagniard, A. Baksheev, T. De Laet, J. De Schutter, and C. Pantofaru, "An adaptable system for rgb-d based human body detection and pose estimation," *Journal of visual communication and image representation*, vol. 25, no. 1, pp. 39–52, 2014.
- [19] M. Carraro, M. Munaro, A. Roitberg, and E. Menegatti, "Improved skeleton estimation by means of depth data fusion from multiple depth cameras," in *International Conference on Intelligent Autonomous Systems*, pp. 1155–1167, Springer, Cham, 2016.
- [20] E. Insafutdinov, L. Pishchulin, B. Andres, M. Andriluka, and B. Schiele, "Deepcut: A deeper, stronger, and faster multi-person pose estimation model," in *European Conference on Computer Vision*, pp. 34–50, Springer, 2016.
- [21] L. Pishchulin, E. Insafutdinov, S. Tang, B. Andres, M. Andriluka, P. V. Gehler, and B. Schiele, "Deepcut: Joint subset partition and labeling for multi person pose estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4929–4937, 2016.
- [22] J. Carreira, P. Agrawal, K. Fragkiadaki, and J. Malik, "Human pose estimation with iterative error feedback," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [23] A. Elhayek, E. de Aguiar, A. Jain, J. Thompson, L. Pishchulin, M. Andriluka, C. Bregler, B. Schiele, and C. Theobalt, "Marconi: Convnet-based marker-less motion capture in outdoor and indoor scenes," *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 3, pp. 501–514, 2017.
- [24] Z. Gao, Y. Yu, Y. Zhou, and S. Du, "Leveraging two kinect sensors for accurate full-body motion capture," *Sensors*, vol. 15, no. 9, pp. 24297–24317, 2015.
- [25] M. Lora, S. Ghidoni, M. Munaro, and E. Menegatti, "A geometric approach to multiple viewpoint human body pose estimation," in *Mobile Robots (ECMR), 2015 European Conference on*, pp. 1–6, IEEE, 2015.
- [26] Y. Yang and D. Ramanan, "Articulated human detection with flexible mixtures of parts," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 12, pp. 2878–2890, 2013.
- [27] Y. Kim, "Dance motion capture and composition using multiple rgb and depth sensors," *International Journal of Distributed Sensor Networks*, vol. 13, no. 2, p. 1550147717696083, 2017.
- [28] A. Kanaujia, N. Haering, G. Taylor, and C. Bregler, "3d human pose and shape estimation from multi-view imagery," in *Computer Vision and Pattern Recognition Workshops (CVPRW), 2011 IEEE Computer Society Conference on*, pp. 49–56, IEEE, 2011.
- [29] K.-Y. Yeung, T.-H. Kwok, and C. C. Wang, "Improved skeleton tracking by duplex kinects: a practical approach for real-time applications," *Journal of Computing and Information Science in Engineering*, vol. 13, no. 4, p. 041007, 2013.
- [30] M. Quigley, K. Conley, B. Gerkey, J. Faust, T. Foote, J. Leibs, R. Wheeler, and A. Y. Ng, "Ros: an open-source robot operating system," in *ICRA workshop on open source software*, vol. 3, p. 5, Kobe, 2009.

³https://github.com/marketto89/open_pttrack