

# Real-time 3D human pose and motion reconstruction from monocular RGB videos

Anastasios Yiannakides<sup>1,2</sup> | Andreas Aristidou<sup>1,2</sup>  | Yiorgos Chrysanthou<sup>1,2</sup>

<sup>1</sup>Department of Computer Science,  
University of Cyprus AND RISE Research  
Center, Nicosia, Cyprus

<sup>2</sup>Department of Computer Science,  
University of Cyprus

## Correspondence

Anastasios Yiannakides, Department of  
Computer Science, University of Cyprus,  
Nicosia 1678, Cyprus.  
Email: tasyiann@gmail.com

## Funding information

RESTART 2016-2020 Programmes for  
Technological Development and  
Innovation, Grant/Award Number:  
P2P/JPICH\_DH/0417/0052; European  
Union's Horizon 2020 Research and  
Innovation Programme, Grant/Award  
Number: 739578, RISE-Call:  
H2020-WIDESPREAD-01-2016-2017-T

## Abstract

Real-time three-dimensional (3D) pose estimation is of high interest in interactive applications, virtual reality, activity recognition, and most importantly, in the growing gaming industry. In this work, we present a method that captures and reconstructs the 3D skeletal pose and motion articulation of multiple characters using a monocular RGB camera. Our method deals with this challenging, but useful, task by taking advantage of the recent development in deep learning that allows two-dimensional (2D) pose estimation of multiple characters and the increasing availability of motion capture data. We fit 2D estimated poses, extracted from a single camera via OpenPose, with a 2D multiview joint projections database that is associated with their 3D motion representations. We then retrieve the 3D body pose of the tracked character, ensuring throughout that the reconstructed movements are natural, satisfy the model constraints, are within a feasible set, and are temporally smooth without jitters. We demonstrate the performance of our method in several examples, including human locomotion, simultaneously capturing of multiple characters, and motion reconstruction from different camera views.

## KEYWORDS

monocular video, motion reconstruction, estimation openPose, 3D pose

## 1 | INTRODUCTION

Motion capture (mocap) is the technological process used for acquiring three-dimensional (3D) position and orientation information of a moving object. Despite recent advances in mocap technology, which allows high-quality acquisition and portrayal of movements, common capture systems are cost demanding and require a setup of capturing devices that is not accessible to the general public, especially for home use. Recently, scholars have turned their attention in using more affordable technologies such as RGB or depth cameras. Capturing and reconstructing the motion of multiple characters using single, monocular cameras have a wide spectrum of applications in surveillance, human-computer interaction, activity recognition, behavioral analysis, and virtual reality. It is also of high interest in the growing gaming industry, where users require cost-effective and easy configurable systems for real-time human-machine interaction. Real-time pose estimation and 3D motion reconstruction using data only from a single RGB camera is, however, a challenging task. This is because skeletal pose estimation, using such sparse data, is a highly underconstrained problem. A human pose is usually represented and parameterized by a set of joint positions and rotations, where its heterogeneous, dynamic, and highly versatile nature, as well as the changes in camera viewpoint, makes motion reconstruction a difficult job.

Most papers in the literature estimate the human pose in two-dimensional (2D) for one or multiple characters by localizing joint keypoints in pixel space<sup>1,2</sup> or by extracting the shape silhouette and then retrieving the closest neighbor from a database.<sup>3–5</sup> More recently, and with the advent of deep learning (DL), the community has been moving to learning-based discriminative methods, where the effectiveness of 2D human pose estimation has greatly improved.<sup>6–8</sup> The 3D skeletal reconstruction, though, is a much harder problem.<sup>9,10</sup> Even though there are methods that are effective at 3D pose reconstruction, they are usually not real-time implementable and suffer from depth and scale ambiguities. In addition, the reconstructed motion is temporally inconsistent and unstable because they treat each frame independently and do not employ bone length constraints. The big challenge in this context is to learn rich features to encode depth, spatial and temporal relation of the body parts so as to ensure smooth motion reconstruction.<sup>11,12</sup> While some recent methods run at high frame (e.g., other works<sup>11,13</sup>), they are still unsuitable for use in closely interactive characters or crowds. This is because they require a tracked bounding box for each person, and thus, they can only reconstruct the motion of a single person at a time.

The novel contribution of this paper, in principle, is that it fits 2D deep estimated poses, taken from a single, monocular camera, with the 2D multiview joint projections of 3D motion data, to retrieve the 3D body pose of the tracked character. Our method takes advantage of the recent advances in deep and convolutional networks that allow 2D pose estimation of multiple characters and the large (and increasing) availability of mocap data. More particularly, our method infers the 3D human poses in real time using only data from a single video stream. To deal with the limitations of the prior work such as the bone length constraints violations, the simultaneously capturing of multiple characters, and the temporal consistency of the reconstructed skeletons, we generate a database with numerous 2D projections by rotating a small angle at a time, the yaw axis of 3D skeletons. Then, we match the input 2D poses (which are extracted from a single video stream using the OpenPose network<sup>8</sup>) with the projections on the database and retrieve the best 3D skeleton pose that is temporally consistent to the skeleton of the previous frames, producing natural and smooth motion.

Our approach is capable of estimating the posture of multiple characters in real time, whereas the reconstructed pose is always within a natural and feasible set. The performance of our method in reconstructing 3D articulated motion from 2D poses is demonstrated in several examples, including video streams taken from different points of view and using multiple characters in the scene.

## 2 | RELATED WORK

There has been a growing demand in recent years for realistic 3D animation in media and entertainment, as well as for research and training purposes. 3D motion acquisition can be roughly categorized as being achieved either using *marker-based* or *marker-less* mocap systems.

**Marker-based** systems are generally producing high-quality, realistic animations and are mainly used by the movies and entertainment industries. They use fiducial markers, which are attached near each joint to identify motion, and they provide real-time acquisition of labeled or algorithmically tracked data to reconstruct the subjects' articulated motion. More specifically, they utilize data captured from special markers (passive and active) to triangulate the 3D position of a subject inferred from a number of high-speed cameras.<sup>14,15</sup> The main strengths of optical systems are their high acquisition accuracy, speed, and high sample rate capture. However, optical hardware is expensive, intrusive by nature, and lacks portability; calibration is needed for precise application of the markers, whereas extensive and tedious manual effort is required to recover the misapplied or (self-)occluded markers.<sup>16,17</sup>

More recently, a number of autonomous systems have been developed that use a variety of sensors such as inertial measurement units (IMUs).<sup>18–20</sup> Inertial mocap technology uses a number of gyroscopes and accelerometers to measure rotational rates, while these rotations are translated to a skeleton model. Inertial systems do not require external cameras to capture the sensors, and thus, they are portable and functional in outdoor environments. Nevertheless, marker-based systems are costly and, thus, not suitable for home use and, more particularly, for interactive applications or gaming.

In general, there is a tendency to reduce the required equipment and the objects attached to the body for motion tracking and reconstruction. One way to deal with sparse data is to model the trajectories of the end effectors and then apply kinematics models to fill in the gaps.<sup>21,22</sup> The most popular way, though, is the use of **marker-less** systems (or vision systems) that are becoming more and more popular because they are easy to set up, have low cost, and are less intrusive, meaning that subjects are not required to wear special equipment for tracking. Usually, the subject's silhouette is captured from a single or multiple angles using a number of vision or RGB-depth cameras.<sup>23,24</sup> A voxel representation of the body is extracted over time, while animation is achieved by fitting a skeleton into the 3D model; see other works for example.<sup>25–29</sup>

These approaches can be broadly classified into discriminative, generative, and hybrid approaches. Generative methods reconstruct human pose by fitting a template model to the observed data.<sup>30,31</sup> Discriminative methods infer mode-to-depth correspondences and cluster pixels to hypothesize body joint positions, fit a model, and then track the skeleton.<sup>32,33</sup> Hybrid methods use a combination of the aforementioned techniques to achieve higher accuracy.<sup>34</sup> Other works also include methods for motion capturing using an architecture of multiple color-depth sensors to better deal with occlusions or unobserved view angles.<sup>35,36</sup>

Skeletal pose estimation from a single camera is a much harder problem; the problem has been tackled by localizing the joint keypoint positions in pixel space<sup>1,2</sup> or by extracting the shape silhouette and then retrieving the closest neighbor from a database.<sup>3–5</sup> For a recent overview of 3D pose estimations, refer to the work of Sarafianos et al.<sup>37</sup>

More recently, DL has brought revolutionary advances in computer vision and graphics, including efficient methods for pose estimation. Learning-based discriminative methods are quite popular for the real-time 2D pose estimation of single or multiple characters.<sup>6–8,38</sup> More recently, a number of methods have tackled a much harder problem, that of 3D skeletal estimation from single color images or videos.<sup>9,10,13,39–41</sup> A common limitation of the discriminative methods in 3D pose reconstruction, though, is that they typically run off-line, and because they do not take into consideration the temporal consistency of motion (they reconstruct the 3D joint positions on per image), the generated motion is oscillating. Moreover, they are restricted to the viewpoints learned from the training data,<sup>42</sup> while they suffer from depth and scale ambiguities. The problem of different camera views can be dealt by training the network to predict the same pose using images from multiple views.<sup>43,44</sup> Another major limitation is that the reconstructed 3D pose is not assigned on a character model (3D joint positions are estimated independently) to enforce kinematic constraints, resulting in temporal bone length violations. The methods in the works of Mehta et al.<sup>11</sup> and Martinez et al. (VNect)<sup>12</sup> animate a modeled character, work in real time, and produce smooth animations in terms of their temporal coherence. However, because they require each character to be tracked by a bounding box, they only reconstruct single-person skeletons at a time, making them unsuitable for closely interacting characters. More recently, an enormous effort has been devoted to deep and convolutional methods that map all human pixels of an RGB image to 3D surface of the human body.<sup>45–48</sup>

Our method overcomes most of the prior work limitations, including the 3D capturing of multiple characters, the skeletal model constraints, and the production of smooth animation for the articulated character.

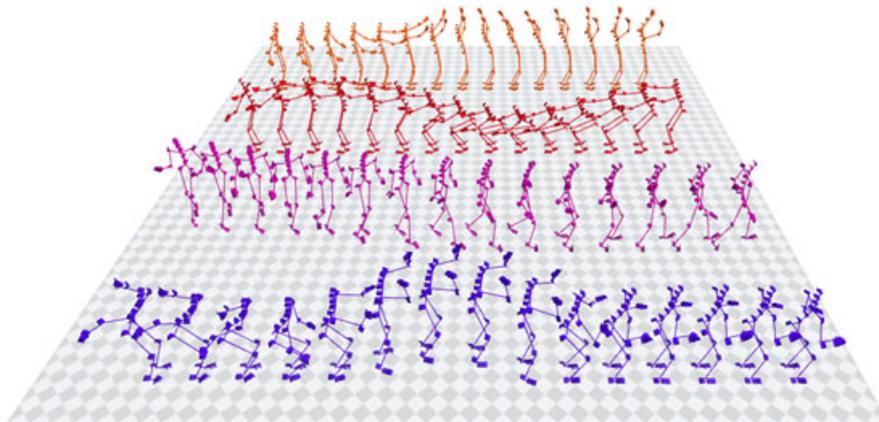
### 3 | METHOD OVERVIEW

Our approach for motion reconstruction can be decomposed into two main parts: (a) the preprocessing step, where a 2D pose database is defined by mocap data projections, and its entries are associated with 3D skeletons, and (b) the run-time step, where 2D joint positions extracted from a video are matched to the 2D pose projections database in order to recover the 3D skeleton and reconstruct the motion.

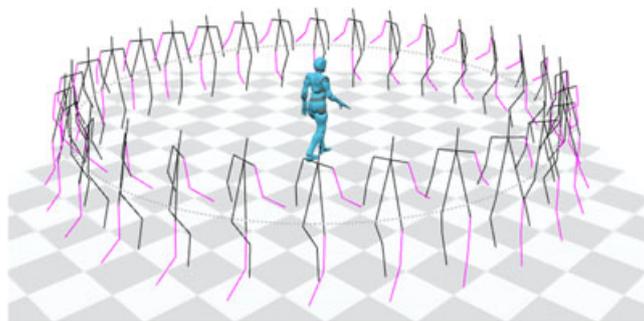
More specifically, as a first step, we retarget motion data taken from an online mocap database<sup>49</sup> into a universal skeleton format, and then compute the 2D projections of the 3D skeleton from many different views. The 2D projections are associated with their 3D skeletons and their viewing information, and stored in a database. The 2D projections are thereafter grouped into mutually exclusive clusters, and for each cluster, we allocate a representative pose. Then, in the second stage, for an input frame taken from a video stream, we extract in real time its 2D pose using the OpenPose network.<sup>8</sup> The 2D pose is then rescaled so as to be consistent with the 2D projections stored in the database. We select the  $k$ -best matches from the database to the input 2D pose and retrieve their 3D skeletons. To achieve smooth reconstruction of the articulated motion, we select the 3D skeleton from the  $k$ -matched projections that is temporally consistent to the previous frame.

### 4 | MOTION DATABASE

The first step of our method, in a preprocessing time, is to define a pose database that will be used to retrieve the 3D pose estimates. We used a small data set  $\mathfrak{D}$  of mocap data in the biovision hierarchy format (.bvh), taken from the Carnegie Mellon University (CMU) mocap library<sup>49</sup> (see Figure 1). The human poses in these CMU 3D data are represented by  $n = 30$  joint positions and rotations. In our work, we use 2D poses that are estimated from an input monocular video using OpenPose<sup>8</sup> (see Section 5 for more details) that are represented by  $m = 14$  joint locations (see Figure 1). Thus, in order to have a uniform and comparable skeleton, we retarget the CMU .bvh data to a 3D skeleton that its projection in T-pose



**FIGURE 1** A number of skeleton sequences used in our data set  $\mathfrak{D}$ , taken from the Carnegie Mellon University (CMU) motion capture library



**FIGURE 2** Multiview skeleton projection. For each frame, we rotate the three-dimensional (3D) skeleton by an angle  $\theta$  at a time on the yaw axis and create two-dimensional (2D) projections

matches the 2D skeleton (in T-pose) returned by OpenPose (see Figure 2 for the new skeleton); the 2D pose projections are then scaled so as their bounding box remains constant over time. This step is crucial because we will compare pixelwise the joint locations of the 2D skeletons joint by joint.

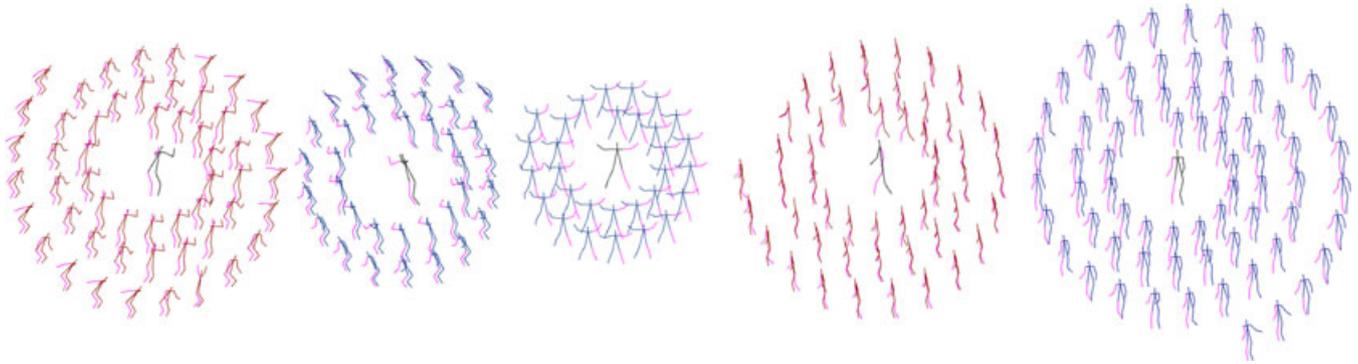
In order to make our method invariant to the camera view and robust against fine-grained pose variations, for each frame, we compute the 2D projections of the 3D skeleton rotated by an angle  $\theta$  at a time ( $\theta = 12^\circ$ ) on the yaw axis, resulting in total to 30 projections per frame. The main idea is to assign for each 3D pose multiple 2D projection representation from many different camera viewpoints and is illustrated in Figure 2. Each 2D projection is then associated with its 3D skeleton (pointed to its frame on the .bvh animation) by indexing the corresponding frame on the animation and its rotation angle. Note that, in this work, we place the camera at 1.5 meters above earth plane (common height for video recording), and similarly, we project the 3D skeletons assuming that the camera is at the same height.

To allow fast skeleton indexing and retrieval, we position the 2D pose projections into  $d$ -dimensional space using multidimensional scaling (MDS)<sup>50</sup>; we tried different dimensionalities,  $\mathbb{R}^d$ , and the quality of embedding seems to be equally well for  $2 \leq d \leq 5$ . The distance metric used to compute the distance between the 2D skeletal projections is defined as

$$\text{dist}_{ij} = \sum_{k=1}^m d(\mathbf{p}_k, \mathbf{q}_k), \quad (1)$$

where  $m$  is the number of joints and  $\mathbf{p}_k, \mathbf{q}_k \in \mathbb{R}^2$  are the joint positions of the two skeletons  $i$  and  $j$ . The term  $d(\mathbf{p}_k, \mathbf{q}_k)$  represents the Euclidean norm between the two joint positions.

2D projection representations in  $\mathbb{R}^d$  are then grouped into mutually exclusive clusters. We use the  $k$ -means clustering algorithm to create clusters that are separated by similar characteristics. For each cluster, we also define a representative pose that is the one closest to the centroid of the cluster. Figure 3 portrays the 2D poses of five selected clusters and their representative skeletons. As can be seen, our database is rich enough and each pose has a lot of repetitions with a big variation of different poses.



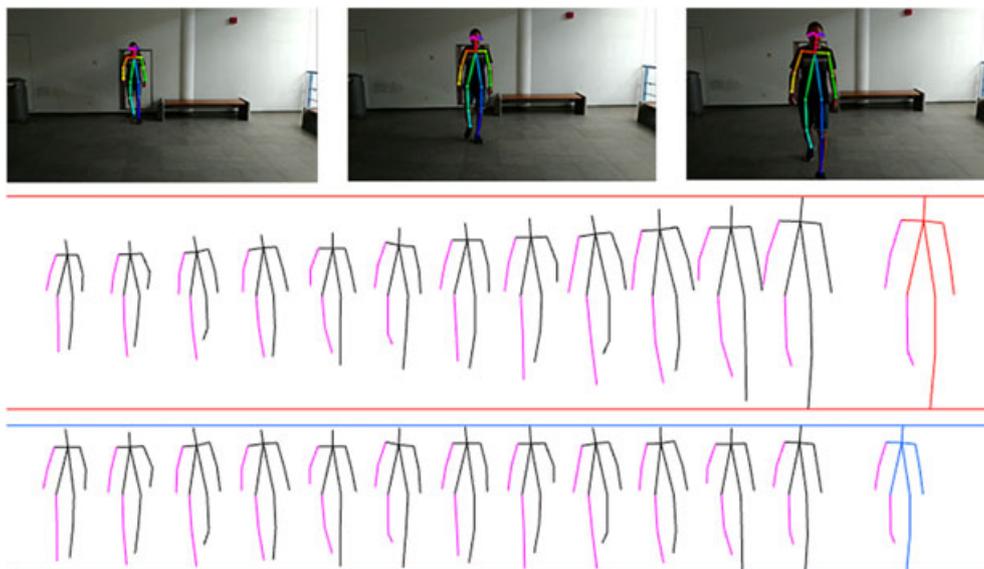
**FIGURE 3** Five selected clusters with their representatives. The projected poses with the smallest distance to the centroid are placed nearest to the center, and as we move further from the center, the poses are less correlated to the centroid

## 5 | MOTION RECONSTRUCTION

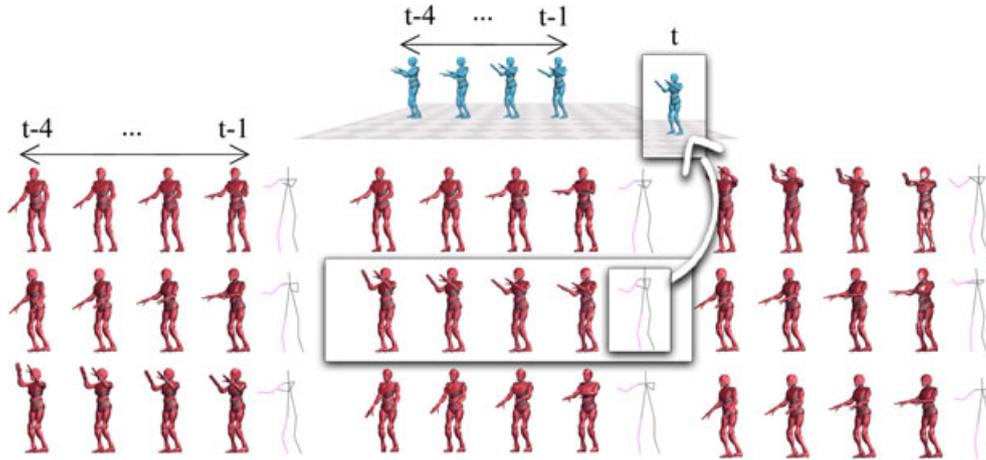
For an input video stream, we use OpenPose to estimate, in real time, the 2D pose of the characters. OpenPose is a bottom-up approach that uses Part Affinity Fields, a set of 2D vector fields that encodes the location and orientation of limbs over the image domain, offering robustness to early commitment. It offers state-of-the-art accuracy; it returns the 2D key points of the joints for multiple articulated characters at the same time, but most importantly, it decouples runtime complexity from the number of people in the image, achieving low computational cost.

In this project, we assume that the character's skeleton is fully visible at the camera. Because the 2D joint locations are inferred from a video, the size of the estimated pose is depended on the size of the tracked character (e.g., adult or child), and its depth (e.g., how far or near is to the camera). To deal with this size ambiguity and be consistent with the 2D projections stored in the database, we introduce a scaling step. For each OpenPose 2D pose, we export its bounding box and rescale it so as the height of the bounding box remains constant over time. Figure 4 shows a sequence of 2D poses of a walking character moving closer to the camera (the size of its pose increases over time), as inferred by OpenPose and its rescaled pose that its size remains constant over time.

Having the rescaled 2D pose from OpenPose, we then need to search and retrieve the nearest pose to the entry from the projections database, and associate it with its corresponding 3D skeleton. However, one of the main challenges to deal with is that multiple 3D poses may correspond to the same 2D pose after projection. This introduces severe ambiguities in 3D pose estimation. Moreover, applying per-frame pose estimation on a motion sequence does not ensure temporal consistency of motion, and small pose inaccuracies lead to temporal jitter and motion oscillations. We dealt with these



**FIGURE 4** Scaling the OpenPose two-dimensional (2D) pose estimation to ensure comparison consistency



**FIGURE 5** To ensure temporal consistency in motion, for each candidate skeleton, we compare its  $w$  previous frames in the original animation (shown in red) with the  $w$  previously reconstructed skeletons (shown in blue) and then select the candidate skeleton that its previous frames return the smallest average distance (highlighted box)

challenges in two steps. First, we compare the entry pose with the cluster representations and then select the  $i = 5$  closest clusters with the smallest Euclidean distance (see Equation 1) between the input pose and their representative. Note that searching on a single cluster is not a good idea because there is no guarantee that the closest representative contains the closest projection for our entry. For each of those  $i$  clusters, we retrieve the  $j = 10$  nearest poses projections. Then, by taking advantage of the property that motion is locally linear, we retrieve the 3D skeleton that is temporally more consistent to the reconstructed poses of the previous frames. More particularly, for each of the  $l = i \times j$  selected nearest pose projections, we extract their associated 3D skeleton (their temporal location is pointed to the .bvh animation). Thereafter, we compare the 3D skeletons of their  $w$  previous frames (we empirically conclude that  $w = 4$  is enough for retrieving the most temporally consistent skeleton), in the .bvh animation, with the  $w$  previously reconstructed skeletons (each of the  $w$  skeletons is weighted differently), and select the one that its  $w$  previous frames have the smallest average distance (see Figure 5 for a visual explanation). To compute the distances between the 3D skeletons, we use the Lee et al.<sup>51</sup> distance metric, which is the sum of the difference in rotation between joints. The distance between two skeletons is defined as

$$\text{dist}_{ij}^2 = \sum_{k=1}^m \left\| \log \left( q_{j,k}^{-1} q_{i,k} \right) \right\|^2, \quad (2)$$

where  $m$  is the number of joints and  $q_{i,k}, q_{j,k} \in \mathbb{S}^3$  are the complex forms of the quaternion for the  $k$ th joint of the two skeletons  $i$  and  $j$ , respectively. The log-norm term  $\| \log(q_{j,k}^{-1} q_{i,k}) \|^2$  represents the geodesic norm in quaternion space, which yields the distance from  $q_{i,k}$  to  $q_{j,k}$  on  $\mathbb{S}^3$ . The retrieved 3D skeleton is placed at the current frame, in the animation, after it is rotated by an angle  $\theta$  in the yaw axis.

Finally, apart from retrieving the pose that is more related to the previous frames, in order to remove unwanted oscillation, we additionally smooth the reconstructed 3D motion using real-time filtering, for example, the Savitzky–Golay<sup>52</sup> or the 1 € filter.<sup>53</sup> The filtering is applied on the joint rotations (represented in Euler angles) of the 3D skeleton. Motion data are further edited, again in real time, using FBBIK<sup>54</sup> to avoid common synthesis artifacts such as foot sliding and floor penetration. It is important to note here that, by retrieving existing 3D skeletons, we ensure that bone length are constant over time, a common limitation of prior work in 2D and 3D pose estimation. The data were transformed into .bvh format that includes absolute root position and orientation of the relative joint angles.

## 6 | RESULTS AND DISCUSSION

In this section, we demonstrate the performance of our method in reconstructing 3D skeletons from RGB video sequences. We first provide the implementation details and its computational performance, and then present several experiments that illustrate the effectiveness of our work.

## 6.1 | Implementation details

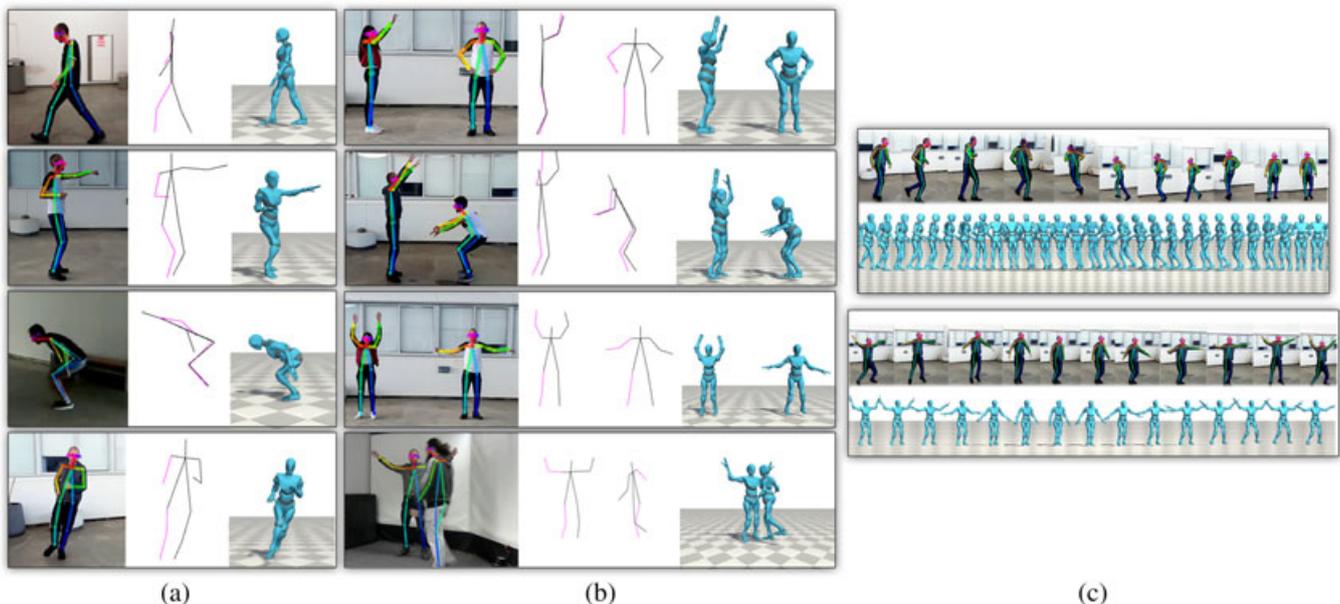
We have implemented our system in the Unity game engine using C#. All experiments were run on a six-core PC with Intel i7-6850 K at 3.6GHz, 32 GB of RAM, and nVIDIA Titan XP GPU. We created a 2D projection database using the 3D mocap data taken from the CMU mocap library. These data were originally sampled at 120 frames per second (fps), but because human motion is locally linear and in order to reduce the computational cost, we resample it to 24 fps without much loss of the temporal information (see the work of Forbes et al.<sup>55</sup>). In our experiments, we only used a small data set of different actions, in total 5 min of motion (72,000 frames). Because for each frame we extract 30 projections (rotated by an angle  $\theta = 12^\circ$  on the yaw axis), our 2D projection database  $\mathfrak{D}$  consists in total with 216,000 pose projections. Selecting the right number of clusters for organizing  $\mathfrak{D}$  is very important because it will allow fast indexing and searching of poses, but it will also guarantee the quality of motion retrieval. We empirically concluded that our clusters are compact for roughly  $K = 500$ ; we ended up at this number of clusters by increasing their number until the mean distance between the skeletons of each cluster and its centroid do not change more than 1%.

Our method requires approximately 12 hr to compute the distance matrix between the poses and to create the pose embedding, 14 hr for dimensionality reduction and 15 minutes for the pose clustering (in MATLAB R2018b). Although this process is time consuming, it is only required once at a preprocessing time. Once it is done, 3D motion retrieval is fast and performed in real time (see the supplementary video for a live demonstration). There are several factors that affect the runtime, but also the quality of the reconstructed animation, for example, the size of the database (bulkier  $\mathfrak{D}$  means larger variety of movements, thus better reconstruction quality, but it requires higher computational cost; this is mainly because of the complexity in matrix computation and MDS), the rotation angle for the projections (again, more frequent projections will increase the quality at the cost of higher computational cost).

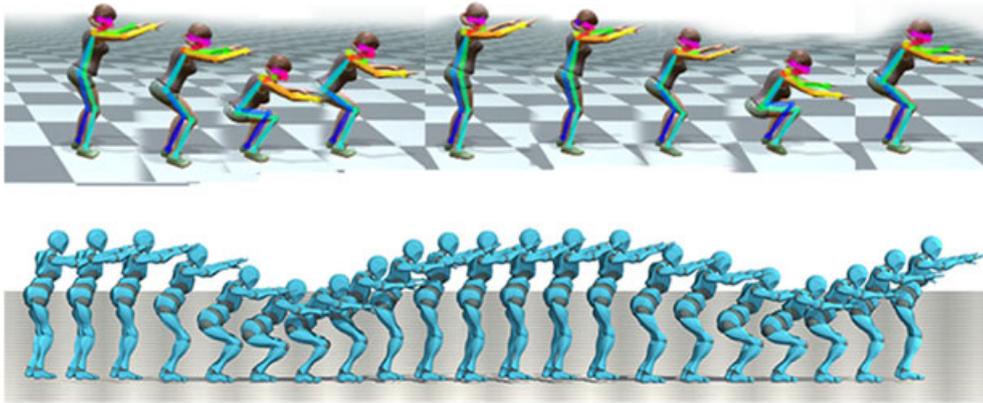
## 6.2 | Evaluation

We conducted several experiments to evaluate the performance of our method. Figure 6a visualizes the 3D pose reconstruction on several motion examples, including walking, boxing, jumping, and running. Figure 6b illustrates the corresponding results for multiple, closely interactive characters. Similarly, Figure 6c shows the effectiveness of our method in tracking and capturing motion that is temporally coherent. It can be observed that the 3D skeletons of the articulated motion are smoothly reconstructed over time.

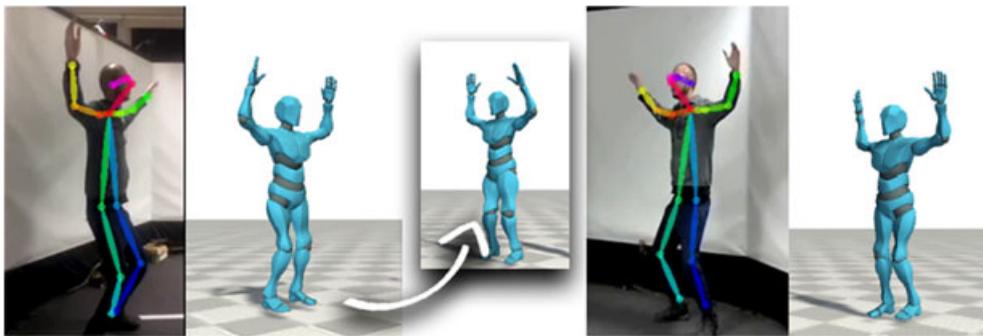
To quantitatively evaluate the performance of our method, we animated a .bvh file (that is not part of our database) using a virtual character and then reconstruct its articulated motion (see Figure 7). We then computed the difference



**FIGURE 6** Three-dimensional (3D) skeleton reconstruction of various different frames for (a) single character and (b) multiple interactive characters in the scene. The RGB video frames are shown on the left side of the figure with the 2D pose estimations overlaid, as returned by OpenPose; in the middle is the best match from our 2D pose projections database to the input pose; and on the right is our 3D skeleton reconstructions. (c) Smooth motion reconstruction in a sequence of skeletons



**FIGURE 7** Three-dimensional (3D) motion reconstruction on virtually animated motion capture (mocap) data



**FIGURE 8** Three-dimensional (3D) pose reconstruction of the same action recorded from different angles. The arrow shows the same pose but rotated to be aligned and visually comparable with the pose from the other view

between the poses of the original and reconstructed characters; the average Euclidean distance between their joints is only 9.4 cm (for a character with height 175 cm). Please refer to our supplementary video for an overlay comparison of the two skeletons. We further evaluated the effectiveness of our method in pose reconstruction by recording the same action from two different viewing angles and then reconstructing their 3D motion. Again, we computed the average Euclidean distance between the two 3D skeletons per joint, that is, 12.6 cm at recordings with different angle views approximately  $135^\circ$  and 6.8 cm for  $75^\circ$ . Similarly, the corresponding distance for the filtered VNect reconstruction is 16.4 cm and 15.2 cm, respectively. Note that, for the comparison between the two poses, we discard the translation and rotation of the root joint so as to fairly compare the two 3D poses. Figure 8 illustrates the two 3D poses from different angles. For a side-by-side animated comparison of the two 3D skeletons, please refer to our accompanied video.

Results demonstrate the effectiveness of our method in reconstructing, at real time, the articulated motion in various different actions, in video streams taken from different points of view, virtually generated videos, and using single or multiple characters in the scene. At the same time, our method ensures the smoothness and temporal consistency of motion, and that the character's bone length constraints are not violated.

## 7 | CONCLUSIONS, LIMITATIONS, AND FUTURE WORK

We have presented a method that estimates, in real time, the 3D pose of a human character using a single, monocular camera and reconstructs its articulated motion. Our method is capable of tracking and reconstructing multiple 3D human postures at the same time, ensuring that the estimated 3D poses satisfy the character's bone constraints and are always within a natural and feasible set. Common prior work limitations were efficiently dealt such as the temporal consistency of poses and the production of smooth and linear motion. We have evaluated the performance of our method in several examples, including a large variety of locomotion with single and multiple characters in the scene, on data taken from

different points of view, artificially generated data, and so forth; our results demonstrate the efficiency of the proposed approach.

Our method has some limitations. First, the accuracy of the 2D joint estimation can greatly affect our method's 3D estimation performance. This is more obvious at the hand and feet joints, which are more vulnerable to noise. There are two ways to overcome this limitation: (a) one way is to use a weighted distance metric and enforce less influence (see Equation 2) on these joints, or (b) to use an architecture of two or more cameras for a better 2D pose estimation. Moreover, in our work, we assume that the character's skeleton is fully visible at the camera; however, this is not always possible due to occlusions from other characters or objects in the scene. Future work will see the introduction of a pose recovery method that selects the most appropriate skeleton, from the database, and matches with the sparse data.

A second limitation lies on the camera's projection angle, which must match the view angle of the camera. A possible way to make our system invariant to different camera viewpoints is to further extend the database with additional projections on different axis, different height (view), or camera parameters. 3D pose interpolations will be employed to match the 3D keypoints so as to enable 3D pose detection with different camera configurations. Nevertheless, learning the camera viewpoint to improve 3D pose estimation is currently an active and challenging topic in computer vision.<sup>56</sup>

Another limitation of our method is that the results are highly related to the training data. A larger variety of movements will increase the variety of reconstructed actions, but at the same time, it will increase the computational cost. Finally, because there is no calibration of the scene, our method is not trained to account for global translation and rotation. In future work, scene calibration and root translations will be added.

## ACKNOWLEDGEMENTS

This work has been supported through the RESTART 2016-2020 Programmes for Technological Development and Innovation by the Cyprus Research Promotion Foundation, with protocol number P2P/JPICH\_DH/0417/0052. It has also been partly supported by the project that has received funding from the European Union's Horizon 2020 Research and Innovation Programme under Grant Agreement 739578 (RISE-Call: H2020-WIDESPREAD-01-2016-2017-TeamingPhase2) and the Government of the Republic of Cyprus through the Directorate General for European Programmes, Coordination and Development.

## ORCID

Andreas Aristidou  <https://orcid.org/0000-0001-7754-0791>

## REFERENCES

1. Bourdev L, Malik J. Poselets: body part detectors trained using 3D human pose annotations. Proceedings of the IEEE 12th International Conference on Computer Vision (ICCV '09); Sep 29–Oct 2; Kyoto, Japan. Washington, DC: IEEE Computer Society; 2009. p. 1365–1372.
2. Felzenszwalb PF, Girshick B, McAllester D, Ramanan D. Object detection with discriminatively trained part-based models. *IEEE Trans Pattern Anal Mach Intell.* 2010;32(9):1627–1645.
3. Wang C, Wang Y, Lin Z, Yuille AL, Gao W. Robust estimation of 3D human poses from a single image. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR '14); 2014 Jun 23–28; Columbus, OH. Washington, DC: IEEE Computer Society; 2014. p. 2369–2376.
4. Akhter I, Black MJ. Pose-conditioned joint angle limits for 3D human pose reconstruction. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR'15); 2015 Jun 7–12; Boston, MA. Washington, DC: IEEE Computer Society; 2015.
5. Atrevi DF, Vivet D, Emile B, Duculty F. 3D human poses estimation from a single 2D silhouette. Proceedings of the International Conference on Computer Vision Theory and Applications (VISSAP '16); 2016 Feb 27–29; Rome, Italy. p. 361–369.
6. Newell A, Yang K, Deng J. Stacked hourglass networks for human pose estimation. Proceedings of the European Conference on Computer Vision (ECCV '16); 2016 Oct 11–14; Amsterdam, The Netherlands. Berlin, Germany: Springer; 2016.
7. Wei SE, Ramakrishna V, Kanade T, Sheikh Y. Convolutional pose machines. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR '16); 2016 27–30; Las Vegas, NV. Washington, DC: IEEE Computer Society; 2016. p. 4724–4732.
8. Cao Z, Hidalgo G, Simon T, Wei S-E, Sheikh Y. OpenPose: realtime multi-person 2D pose estimation using part affinity fields. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR '18); 2018 Jun 18–22; Salt Lake City, UT. Washington, DC: IEEE Computer Society; 2018.
9. Bogo F, Kanazawa A, Lassner C, Gehler P, Romero J, Black MJ. Keep it SMPL: automatic estimation of 3D human pose and shape from a single image. Computer vision: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part V. Berlin, Germany: Springer; 2016. p. 561–578.

10. Pavlakos G, Zhu L, Zhou X, Daniilidis K. Learning to estimate 3D human pose and shape from a single color image. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR '18)*; 2018 Jun 18–22; Salt Lake City, UT. Washington, DC: IEEE Computer Society; 2018.
11. Mehta D, Sridhar S, Sotnychenko O, et al. Vnect: real-time 3D human pose estimation with a single RGB camera. *ACM Trans Graph*. 2017;36(4). Article No. 4.
12. Martinez J, Hossain R, Romero J, Little JJ. A simple yet effective baseline for 3D human pose estimation. *Proceedings of the IEEE International Conference on Computer Vision (ICCV '17)*; 2017 Oct 22–29; Venice, Italy. Piscataway, NJ: IEEE; 2017. p. 2659–2668.
13. Wang K, Lin L, Jiang C, Qian C, Wei P. 3D human pose machines with self-supervised learning. *IEEE Trans Pattern Anal Mach Intell*. 2019.
14. PhaseSpace Inc. Optical MoCap systems. 2019. <http://www.phasespace.com>
15. Vicon. Vicon motion capture systems. 2019. <http://www.vicon.com>
16. Aristidou A, Lasenby J. Real-time marker prediction and CoR estimation in optical motion capture. *Vis Comput*. 2013;29(1):7–26.
17. Aristidou A, Cohen-Or D, Hodgins JK, Shamir A. Self-similarity analysis for motion capture cleaning. *Comput Graph Forum*. 2018;37(2):297–309.
18. Xsens Technologies BV. Motion capture systems. 2019. <http://www.xsens.com>
19. Tautges J, Zinke A, Krüger B, et al. Motion reconstruction using sparse accelerometer data. *ACM Trans Graph*. 2011;30(3). Article No. 18.
20. Slyper R, Hodgins JK. Action capture with accelerometers. *Proceedings of the 2008 ACM SIGGRAPH/Eurographics Symposium on Computer Animation (SCA '08)*; 2008 Jul 7–9; Dublin, Ireland. Aire-la-Ville, Switzerland: Eurographics Association; 2008. p. 193–199.
21. Aristidou A, Chrysanthou Y, Lasenby J. Extending FABRIK with model constraints. *Comput Animat Virtual Worlds*. 2016;27(1):35–57.
22. Carreno-Medrano P, Gibet S, Marteau P-F. From expressive end-effector trajectories to expressive bodily motions. *Proceedings of the 29th International Conference on Computer Animation and Social Agents (CASA '16)*; 2016 May 23–25; Geneva, Switzerland. New York, NY: ACM; 2016. p. 157–163.
23. Zimmermann C, Welschehold T, Dornhege C, Burgard W, Brox T. 3D human pose estimation in RGBD images for robotic task learning. *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA '18)*; 2018 May 21–25; Brisbane, Australia. Piscataway, NJ: IEEE; 2018.
24. Biswas A, Admoni H, Steinfeld A. Fast on-board 3D torso pose recovery and forecasting. *Proceedings of the International Conference on Robotics and Automation (ICRA '19)*; 2019 May 20–24; Montreal, Canada. Piscataway, NJ: IEEE; 2019.
25. Cheung GKM, Baker S, Kanade T. Shape-from-silhouette of articulated objects and its use for human body kinematics estimation and motion capture. *Proceedings of the 2003 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR '03)*; 2003 Jun 18–20; Madison, WI. Washington, DC: IEEE Computer Society; 2003. p. 77–84.
26. de Aguiar E, Stoll C, Theobalt C, Ahmed N, Seidel H-P, Thrun S. Performance capture from sparse multi-view video. *ACM Trans Graph*. 2008;27(3). Article No. 98.
27. Vlastic D, Baran I, Matusik W, Popović J. Articulated mesh animation from multi-view silhouettes. *ACM Trans Graph*. 2008;27(3). Article No. 97.
28. Gall J, Yao A, Van Gool L. 2D action recognition serves 3D human pose estimation. In: *Computer vision: 11th European Conference on Computer Vision, Heraklion, Crete, Greece, September 5–11, 2010, Proceedings, Part III*. Berlin, Germany: Springer-Verlag; 2010. p. 425–438.
29. Liu Y, Gall J, Stoll C, Dai Q, Seidel H-P, Theobalt C. Markerless motion capture of multiple characters using multiview image segmentation. *IEEE Trans Pattern Anal Mach Intell*. 2013;35(11):2720–2735.
30. Wei X, Zhang P, Chai J. Accurate realtime full-body motion capture using a single depth camera. *ACM Trans Graph*. 2012;31(6). Article No. 188.
31. Ye M, Yang R. Real-time simultaneous pose and shape estimation for articulated objects using a single depth camera. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR '14)*; 2014 Jun 23–28; Columbus, OH. Washington, DC: IEEE Computer Society; 2014. p. 2353–2360.
32. Shotton J, Sharp T, Kipman A, et al. Real-time human pose recognition in parts from single depth images. *Commun ACM*. 2013;56(1):116–124.
33. Sharp T. The vitruvian manifold: inferring dense correspondences for one-shot human pose estimation. *Proceedings of the 2012 IEEE Conference on Computer Vision and Pattern Recognition (CVPR '12)*; 2012 Jun 16–21; Providence, RI. Washington, DC: IEEE Computer Society; 2012. p. 103–110.
34. Baak A, Müller M, Bharaj G, Seidel HP, Theobalt C. A data-driven approach for real-time full body pose reconstruction from a depth camera. *Proceedings of the International Conference on Computer Vision (ICCV'11)*; 6–13; Barcelona, Spain. Washington, DC: IEEE Computer Society; 2011. p. 1092–1099.
35. Vlastic D, Peers P, Baran I, et al. Dynamic shape capture using multi-view photometric stereo. *ACM Trans Graph* 2009;28(5). Article No. 174.
36. Berger K, Ruhl K, Schroeder Y, Bruemmer C, Scholz A, Magnor M. Markerless motion capture using multiple color-depth sensors. In: *Eisert P, Hornegger J, Polthier K, editors. Vision, modeling, and visualization*. Aire-la-Ville, Switzerland: Eurographics Association; 2011. p. 317–324.
37. Sarafianos N, Boteanu B, Ionescu B, Kakadiaris IA. 3D human pose estimation: a review of the literature and analysis of covariates. *Comput Vis Image Underst*. 2016;152:1–20.
38. Papandreou G, Zhu T, Kanazawa N, et al. Towards accurate multi-person pose estimation in the wild. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR '17)*; 2017 Jul 21–26; Honolulu, HI. Washington, DC: IEEE Computer Society; 2017.

39. Zhou X, Zhu M, Pavlakos G, Leonardos S, Derpanis KG, Daniilidis K. MonoCap: monocular human motion capture using a CNN coupled with a geometric prior. *IEEE Trans Pattern Anal Mach Intell.* 2018;41:901–914.
40. Yang W, Ouyang W, Wang X, Ren JSJ, Li H, Wang X. 3D human pose estimation in the wild by adversarial learning. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR '18)*; 2018 Jun 18–23; Salt Lake City, UT. Washington, DC: IEEE Computer Society; 2018. p. 5255–5264.
41. Pavlo D, Feichtenhofer C, Grangier D, Auli M. 3D human pose estimation in video with temporal convolutions and semi-supervised training. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR '19)*; 2019 Jun 16–20; Long Beach, CA. Washington, DC: IEEE Computer Society; 2019.
42. Simo-Serra E, Quattoni A, Torras C, Moreno-Noguer F. A joint model for 2D and 3D pose estimation from a single image. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR '13)*; 2013 Jun 23–28; Portland, OR. Washington, DC: IEEE Computer Society; 2013. p. 3634–3641.
43. Rhodin H, Meyer F, Spörri J, et al. Learning monocular 3D human pose estimation from multi-view images. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR '18)*; 2018 Jun 23–28; Salt Lake City, UT. Washington, DC: IEEE Computer Society; 2018. p. 8437–8446.
44. Rhodin H, Salzmann M, Fua P. Unsupervised geometry-aware representation for 3D human pose estimation. In: *Computer vision: 15th European Conference, Munich, Germany, September 8–14, 2018, Proceedings, Part X.* Berlin, Germany: Springer-Verlag; 2018. 765–782.
45. Lassner C, Romero J, Kiefel M, Bogo F, Black MJ, Gehler PV. Unite the people: closing the loop between 3D and 2D human representations. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR '17)*; 2017 Jul 21–26; Honolulu, HI. Washington, DC: IEEE Computer Society; 2017. p. 4704–4713.
46. Kanazawa A, Black J. M, Jacobs DW, Malik J. End-to-end recovery of human shape and pose. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR'18)*; Jun 18–23; Salt Lake City, UT. Washington, DC: IEEE Computer Society; 2018. p. 7122–7131.
47. Xu W, Chatterjee A, Zollhöfer M., et al. Monoperfcap: human performance capture from monocular video. *ACM Trans Graph.* 2018;37(2). Article No. 27
48. Güler RA, Neverova N, Kokkinos I. Densepose: dense human pose estimation in the wild. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR'18)*; Jun 18–23; Salt Lake City, UT. Washington, DC: IEEE Computer Society; 2018. p. 7297–7306.
49. Carnegie Mellon University. MoCap Library. 2019. <http://mocap.cs.cmu.edu/>
50. Seber GAF. *Multivariate observations.* Hoboken, NJ: John Wiley & Sons; 1984.
51. Lee J, Chai J, Reitsma PSA, Hodgins JK, Pollard NS. Interactive control of avatars animated with human motion data. *ACM Trans Graph.* 2002;21(3):491–500.
52. Orfanidis SJ. *Introduction to signal processing.* Upper Saddle River, NJ: Prentice-Hall; 1995.
53. Casiez G, Roussel N, Vogel D. 1 € filter: a simple speed-based low-pass filter for noisy input in interactive systems. *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (CHI '12)*; 2012 May 5–10; Austin, TX. New York, NY: ACM; 2012. p. 2527–2530.
54. Root-Motion. FINAL-IK. 2019. <http://root-motion.com/>
55. Forbes K, Fiume E. An efficient search algorithm for motion data using weighted PCA. *Proceedings of the 2005 ACM SIGGRAPH/Eurographics Symposium on Computer Animation*; 2015 Jul 29–31; New York, NY: ACM; 2005. p. 67–76.
56. Ghezghieh MF, Kasturi R, Sarkar S. Learning camera viewpoint using CNN to improve 3D body pose estimation. *Proceedings of the Fourth International Conference on 3D Vision (3DV'16)*; 2016 Oct 25–28; Stanford, CA. Piscataway, NJ: IEEE; 2016.

## AUTHOR BIOGRAPHIES



**Anastasios Yiannakides** is a bachelor degree student at the Department of Computer Science, University of Cyprus. Since 2018, he joined the Graphics and Hypermedia lab, and works on projects related to 3D character animation. His interests are focus on motion synthesis and skeletal reconstruction, and involves deep and convolutional learning.



**Andreas Aristidou** is a senior post-doc researcher at the Graphics and Hypermedia lab, University of Cyprus. He had been a Cambridge European Trust fellow at the University of Cambridge, where he obtained his PhD (2011). Andreas has a BSc in Informatics and Telecommunications from the National and Kapodistrian University of Athens (2005) and he is an honor graduate of Kings College London (2006). He worked as a research fellow at the Shandong University (China), and IDC Herzliya (Israel), and participated in a number of EU funded projects. His main interests are focused on character animation, motion analysis,

synthesis, and classification, and involve motion capture, inverse kinematics, deep and reinforcement learning, and applications of Conformal Geometric Algebra in graphics.



**Yiorgos Chrysanthou** is a Professor at the Computer Science Department of the University of Cyprus where he is heading the Graphics and Hypermedia lab. He is also the Research Director of the newly established Centre of Excellence on Interactive Media, Smart Systems and emerging Technologies (RISE). Yiorgos was educated in the UK (Queen Mary College, University of London) and worked for several years as a research fellow and a lecturer at University College London. He has published over 80 papers in journals and international conferences and served as the local or overall coordinator of 27 research projects, related to 3D graphics, virtual reality and applications. His research interests lie in the general area of 3D Computer Graphics, recently focusing more on computer animation, algorithms for real-time AR and VR rendering and reconstruction of urban environments.

### SUPPORTING INFORMATION

Additional supporting information may be found online in the Supporting Information section at the end of the article.

**How to cite this article:** Yiannakides A, Aristidou A, Chrysanthou Y. Real-time 3D human pose and motion reconstruction from monocular RGB videos. *Comput Anim Virtual Worlds*. 2019;e1887. <https://doi.org/10.1002/cav.1887>