

Tracking the articulated motion of the human body with two RGBD cameras

Damien Michel^a, Costas Panagiotakis^b, Antonis A. Argyros^{a,c}

^a*Institute of Computer Science, FORTH*

^b*Department of Business Administration (Agios Nikolaos), TEI of Crete*

^c*Computer Science Department, University of Crete*

Abstract

We present a model-based, top-down solution to the problem of tracking the 3D position, orientation and full articulation of the human body from markerless visual observations obtained by two synchronized RGBD cameras. Inspired by recent advances to the problem of model-based hand tracking [15], we treat human body tracking as an optimization problem that is solved using stochastic optimization techniques. We show that the proposed approach outperforms in accuracy state of the art methods that rely on a single RGBD camera. Thus, for applications that require increased accuracy and can afford the extra complexity introduced by the second sensor, the proposed approach constitutes a viable solution to the problem of markerless human motion tracking. Our findings are supported by an extensive quantitative evaluation of the method that has been performed on a publicly available data set that is annotated with ground truth.

Keywords: markerless human motion capture, 3D human tracking, 3D pose estimation, articulated object tracking, 3D reconstruction

Email address: argyros@ics.forth.gr (Antonis A. Argyros)

URL: <http://www.ics.forth.gr/~argyros> (Antonis A. Argyros)

1. Introduction

The estimation of the articulated motion of the human body is very important to a number of real world applications, ranging from surveillance to game design and human computer interaction. It is considered to be a difficult problem because of its high dimensionality and the variability of the tracked person regarding appearance, body dimensions, etc. A number of practical approaches simplify or even avoid these problems by using special hardware or by interfering with the subject and/or the environment by means of visual markers or full body suits [25]. However, unobtrusive, markerless tracking is definitely preferable since it does not interfere with the environment, the subject and its actions. The methods that use markerless visual data as their only input fall into two basic categories, the top-down and the bottom-up ones. Top-down approaches can provide accurate, physically plausible solutions at the cost of a high computational complexity. Bottom-up methods are typically faster, but rely on a discrete set of training poses whose selection determines the accuracy of the obtained results.

In this paper, we propose a model-based, top-down solution to the problem of tracking the 3D pose and articulation of a human body. This is formulated as a optimization problem that minimizes the discrepancy between the 3D occupancy of hypothesized instances of a human body model and the volume reconstructed from the observations. The input to the method comes from two wide baseline, extrinsically calibrated, off-the-shelf RGBD sensors [27] whose depth maps are fused to give rise to the required volumetric representation of the human body. The required volumetric representation can also be obtained by computing the visual hull of a human body figure through standard techniques [24] employing a network of multiple, conventional cameras. Nevertheless, the setup of two RGBD cameras is preferable due to its lower cost and complexity.

Optimization is performed based on an a stochastic method (Particle

Swarm Optimization - PSO) [9]. We demonstrate experimentally the accuracy achieved by the baseline PSO (*bPSO*) optimization method that borrows directly from recent advances on the problem of tracking the articulated motion of the human hand [15]. We also propose a new variant called perturbed PSO (*pPSO*) which systematically perturbs the solutions provided by *bPSO*. We demonstrate that *pPSO* outperforms *bPSO*. We also compare both *bPSO* and *pPSO* with a widely employed method [17] which we will refer to as *OpenNI* for estimating the human skeleton based on a single RGBD camera. Experimental results show that compared to *OpenNI*, the proposed *pPSO* method provides more accurate results. Thus, in applications where increased accuracy is worth the extra complexity introduced by the second sensor, *pPSO* is the preferred choice. On the other hand, both *bPSO* and *pPSO*, being model-based tracking approaches, require knowledge of the parameters of the human body and its 3D pose in the first frame of a sequence. To address these practical problems, we also propose and evaluate another variant, called *HYBRID*, which combines *pPSO* and *OpenNI*, aiming at combining the merits of both in a single method.

The rest of the paper is organized as follows. Section 2 reviews existing approaches to the problem of markerless human motion capture and tracking. Section 3 describes the proposed approaches, by detailing the human body model employed, the observation model, the objective function used to compare hypotheses and observations as well as the optimization methods used to minimize it. It also presents how *pPSO* and *OpenNI* are combined into the *HYBRID* method. Section 4 presents the experimental evaluation of the proposed method in a standard dataset that is annotated with ground truth. Finally, Section 5 summarizes the paper by drawing the main conclusions from this research.

2. Related work

Because of its high theoretical and practical interest, human motion capture based on vision has been the theme of numerous research efforts. The complete review of these works is beyond the scope of this paper. The interested reader is referred to [11, 19] where extended surveys are provided. More recently, Chen et al. [2] surveyed methods for human motion estimation based on depth cameras.

Most commercial solutions to the problem of human motion capture make use of special markers that are placed on carefully selected (e.g., joints) points of the subject’s body. In this paper we are interested in markerless motion capture techniques because, being unobtrusive, present obvious practical advantages over the marker-based solutions.

Markerless human motion capture techniques may be classified into two broad classes, the *bottom-up* and the *top-down* ones. Bottom up methods [22, 1, 20, 18, 21] extract a set of features from the input images, and try to map them to the human pose space. This is achieved with a learning process that involves a typically large database of known poses that cover as much as possible the whole human poses search space. The type of descriptors employed, the mapping method and the actual poses database are the factors determining the accuracy and efficiency of these methods. Due to their nature, most of their computing time is spend on the offline processes of database creation and mapping, while the online performance is rather good.

Top-down approaches [4, 6, 5, 26, 3, 28] use a fully articulated model of the human body and try to estimate the joints angles that would make the appearance of this model fit best the visual input. The model is usually made of a base skeleton and an attached surface. In some methods, complex surface deformations are allowed [6]. Having defined a model of the human body, different pose hypotheses can be formed. A typical top-down method consists

86 of generating hypotheses and comparing them to the input visual data. The
 87 comparison is performed based on an objective function that measures the
 88 discrepancy between a pose hypothesis and the actual observations. The
 89 minimization of this objective function determines the pose that best explains
 90 the available observations. Typically, this is formulated as an optimization
 91 problem that amounts to the exploration of a very high dimensional search
 92 space. Kinematic constraints based on physiological data are often applied
 93 to the model, excluding non realistic poses and reducing significantly that
 94 search space. Constraining not only the pose but also the motion itself can
 95 further help reducing the complexity, for example with Kalman filters [10].
 96 However, this means a reduced generality and the necessity to build and learn
 97 human motion models.

98 The main advantage of top-down methods is their flexibility. The em-
 99 ployed model can be changed easily, and the whole search space can be
 100 explored without any form of training. The price to pay for this flexibility
 101 is the computational cost of the online process. Due to their generative na-
 102 ture, most of the computational work needs to be performed online. Two
 103 more shortcomings is the requirement for knowing the body model param-
 104 eters of each individual and the requirement of providing an initial pose to be
 105 tracked.

106 Instead of trying to estimate the full body model in a single step, a
 107 variety of methods first identify body parts. Then, they either report them
 108 as the solution or they further combine them into a full model [20, 21].
 109 As in the case of hand tracking and according to the related categorization
 110 of Oikonomidis et al. [16], we can identify *disjoint evidence methods* and
 111 *joint evidence methods* [4, 6, 5, 26, 3, 28]. Joint evidence methods handle
 112 effortlessly collisions, self occlusions and all part interactions while disjoint
 113 evidence methods have to handle them explicitly.

114 This paper presents a model-based, top-down pose estimation method

115 that employs a single hypothesis. Furthermore, it is a joint-evidence method.
 116 The 3D body pose recovery is treated as a minimization problem whose ob-
 117 jective function quantifies the discrepancy between the 3D structure and
 118 appearance of hypothesized 3D body model instances, and visual observa-
 119 tions of a human body. Observations come from two off-the-shelf Kinect
 120 sensors. Optimization is performed through a variant of PSO tailored to
 121 the needs of the specific problem. Zhang et al. [28] proposed a solution to
 122 articulated human motion tracking that is also based on two RGBD sensors.
 123 Their approach differs in the observation model that is being used and in
 124 the employed optimization technique. Other versions of PSO have been em-
 125 ployed in the past for human body pose tracking [26, 3, 12], as well as for
 126 multicamera-based and RGBD camera-based hand pose estimation [14, 15].
 127 For example in [26] body-pose hypotheses are used to render silhouettes that
 128 are compared with respective observations. They adopt a hierarchical ap-
 129 proach to the problem and employ a PSO variant to solve it. Their approach
 130 differs from our methodology both in the observation model and in the opti-
 131 mization strategy. In particular, we propose and present a novel PSO variant
 132 (*pPSO*). We also present another variant, called *HYBRID*, which combines
 133 the *pPSO* with the *OpenNI* method and lifts the requirements for knowing
 134 the body model parameters and the initial body pose.

135 An extensive evaluation on a standard, publicly available data set anno-
 136 tated with ground truth shows that *pPSO* outperforms the baseline (*bPSO*)
 137 optimization method and is more accurate than an extensively used, bottom
 138 up *OpenNI* approach [17] to the same problem. Finally, the *HYBRID* ap-
 139 proach performs slightly worse compared to *pPSO* with respect to accuracy.
 140 This is because *pPSO* operates with more accurate, manually derived human
 141 body models, while *HYBRID* estimated them automatically, but with some
 142 error. Nevertheless, *HYBRID* is far more practical as it avoids cumbersome
 143 initialization processes.

144 3. Tracking human body articulations

145 The input to the proposed method is a volumetric representation of the
146 human body (Figure 1(e)). This can be obtained by two RGBD sensors that
147 are configured in a wide baseline setup, or by computing the visual hull of the
148 human body based on a number of conventional RGB cameras. The first op-
149 tion is preferable because it involves fewer cameras and because the resulting
150 volumetric representation describes more accurately a 3D shape compared
151 to its visual hull. The depth information also facilitates the segmentation of
152 the human figure from the rest of the environment.

153 The adopted 3D human model comprises of a set of appropriately as-
154 sembled geometric primitives. Each body pose is represented as a vector
155 of 35 parameters. Body articulation tracking is formulated as the problem
156 of estimating the 35 body model parameters that minimize the discrepancy
157 between the body hypotheses and the actual observations. To quantify this
158 discrepancy, a representation of the volume occupied by a given body model
159 is produced and compared to the volumetric representation generated by the
160 two RGBD cameras. An appropriate objective function is thus formulated
161 and a variant of PSO is employed to search for the optimal body configura-
162 tion. The result of this optimization process is the output of the method for
163 the given frame. Temporal continuity is exploited to track the body articu-
164 lation in a sequence of frames. The remainder of this section describes these
165 algorithmic steps in more detail.

166 3.1. Observing a human

167 At a certain moment in time, the input to the method is a set of two
168 640×480 depth images of a human, as provided by two intrinsically and ex-
169 trinsically calibrated RGBD sensors [27]. Figure 1(a), (b) and (c), (d), show
170 the RGB and depth information acquired by two such sensors. Foreground
171 is segmented through change detection that is performed on the depth in-

172 formation. More specifically, depth views of the environment without and
 173 with the human are available. Image points that exhibit pixelwise depth
 174 differences that exceed a certain threshold are detected and attributed to the
 175 scene foreground. The threshold used in this process is determined based
 176 on a study of the depth error estimation of the Kinect [23]. The resulting
 177 largest foreground blob in each depth image is kept for further consideration.
 178 A conservative estimation of the human spatial extend is performed by ap-
 179 plying a closing morphological operator to these blobs with a circular mask
 180 of radius $r = 1$. Due to sensor limitations, the depth of some points that
 181 lie within the detected foreground is unknown. However, it is necessary to
 182 give at least an approximate depth value to these points in order to produce
 183 a correct 3D reconstruction that is needed for further processing. Thus, the
 184 depth at such points is set equal to the median of the non-null depths of
 185 points within a radius of 2 pixels. Averaging instead of median filtering was
 186 also tested, giving rise to negligible differences in accuracy. The depth values
 187 of the background pixels is set to infinity.

188 A 3D space of $150 \times 150 \times 150$ voxels is then considered. Each voxel is a
 189 cube with side equal to $s_v = 15mm$ resulting in a volumetric representation
 190 of a 3D space of $2.25 \times 2.25 \times 2.25$ meters. The center of this space is
 191 set equal to the mean position of the 3D points located onto the largest
 192 foreground blob of one of the two RGBD cameras. Each voxel of this space
 193 is set to 1 representing the initial assumption that the whole voxel space
 194 is fully occupied by the human figure. Then, the depth values of the two
 195 extrinsically calibrated RGBD cameras are used to carve out voxels that are
 196 not occupied. More specifically, for each 3D voxel v we compute its Euclidean
 197 distance d from an RGBD camera and compare this to the distance \hat{d} that is
 198 estimated from the depth values provided by that camera. If $d < \hat{d}$ then this
 199 voxel should be carved out and takes the value of 0. This test is performed
 200 for both RGBD views. At the end of this process, voxels with a value of 1

201 provide the volumetric representation o_v of the human. An example of such
202 a representation is shown in Figure 1(e).

203 We also compute the outer surface o_s of o_v . On o_s , we apply a 3D distance
204 transform using a spherical kernel of a radius equal to 7 voxels. In the
205 resulting map o_{sd} , voxels that also belong to o_s have a value of 1, voxels that
206 are more than 7 voxels apart from any surface voxel have a value of 0 and
207 the rest of the voxels have a value from 1 to 0 that is inversely proportional
208 to their distance (0 to 7 voxels) to the closest voxel of o_s .

209 The observation model $o = \{o_v, o_{sd}\}$ that feeds the rest of the process
210 consists of o_v and o_{sd} .

211 3.2. Modeling a human

212 The employed human model consists of a main body, two legs, two arms
213 and the head (Figure 1(f)). The main body is modeled with two articulated
214 elliptic cylinders and three ellipsoids for the caps and the junction. The head
215 is made of one cylinder and a sphere. Each arm consists of three spheres and
216 two truncated cones, while a leg has two such cones, two spheres for the knee
217 and the ankle, respectively, and one ellipsoid for the foot. In Figure 1(f), the
218 human model is depicted with color-coded geometric primitives (yellow for
219 elliptic cylinders, red for ellipsoids, green for spheres and blue for truncated
220 cones). For the *bPSO* and the *pPSO* variants, the parameters of these primi-
221 tives (lengths, radii, etc) are manually set. For the *HYBRID* approach, these
222 are automatically estimated based on the output of the *OpenNI* method.

223 The kinematics of each arm is modelled using six parameters encoding
224 angles. Two parameters determine the shoulder position with respect to the
225 torso, three parameters the upper arm with respect to the shoulder and one
226 parameter the elbow with respect to the upper arm. Six parameters are also
227 used for a leg, three for the root, one for the knee and two for the ankle. Two
228 parameters are used for the head, and three parameters for the articulation

229 between the torso and the hip. The global position of the body is represented
 230 using a fixed point on the hip. The global orientation is parametrized using
 231 Euler angles. The above parametrization encodes a 35 degrees of freedom
 232 (DOFs) human model with each DOF represented by a single parameter.

233 3.3. Evaluating a human hypothesis

234 Having defined the parametric 3D model of a human, the goal is to esti-
 235 mate the model parameters that are most compatible to the visual observa-
 236 tions (Section 3.1). To do so, given a human pose hypothesis h , a volumetric
 237 representation h_v of the human model at pose h is generated through graph-
 238 ics rendering. The volume h_v is rendered in a voxel space with identical
 239 characteristics to those of o_v . A distance function $D_v(h_v, o_v)$ is defined as
 240 follows:

$$D_v(o_v, h_v) = 1 - \frac{2 \sum(o_v \wedge h_v)}{\sum(o_v \wedge h_v) + \sum(o_v \vee h_v)}. \quad (1)$$

241 Intuitively, D_v quantifies the volumetric discrepancy between the observation
 242 volume o_v and the hypothesis volume h_v . In Eq.(1), symbols \wedge and \vee denote
 243 logical operations between the binary values of corresponding voxels and
 244 summations are over the set of all voxels. When the volumes h_v and o_v are
 245 disjoint, the quotient in Eq.(1) is equal to 0. If these volumes are identical
 246 and coincide, the quotient is equal to 1. Thus, D_v is equal to 0 if volumes
 247 coincide and 1 if they are totally disjoint.

248 Besides volumetric discrepancy, we also compute a surface alignment dis-
 249 crepancy. To define this, we first compute the outer surfaces h_s of the vol-
 250 umetric representation h_v of the hypothesis h . Then, the surface alignment
 251 discrepancy $D_s(o_{sd}, h_s)$ is defined as:

$$D_s(o_{sd}, h_s) = 1 - \frac{1}{n_p} \sum(o_{sd} \cdot h_s). \quad (2)$$

252 In Eq.(2), the sum is over all voxels and \cdot denotes standard multiplication of
 253 the values of the corresponding maps. o_{sd} is defined as in Section 3.1. Thus,

254 D_s takes a value of 0 if the surface of the hypothesis coincides perfectly with
 255 that of the observation.

256 Given the distance functions D_v and D_s , it is now possible to define the
 257 function $E(o, h)$ that measures the discrepancy between the observation o
 258 and a given body pose hypothesis h :

$$E(o, h) = D_v(o_v, h_v) + D_s(o_{sd}, h_s). \quad (3)$$

259 The minimization of $E(o, h)$ with respect to h yields the body pose that best
 260 (as quantified by the objective function) explains the observations. The next
 261 section details how this minimization is actually achieved.

262 It should be noted that the reconstructed volume of the subject is a
 263 superset of his/her actual volume. This is a direct effect of the fact that 3D
 264 reconstruction is performed with a space-carving-like method, which cannot
 265 handle occlusions. As an example, assume that a human is observed from
 266 sideways and the volume between his arm and torso cannot be reconstructed,
 267 as the arm occludes this space. Having two RGBD cameras in a wide baseline
 268 configuration minimizes this type of effects but does not eliminate them. This
 269 is exactly why the objective function of the optimization process employs
 270 the term $D_s(o_{sd}, h_s)$ that is related to the coverage of the surface of the
 271 reconstruction. As the 3D reconstruction is a superset of the actual volume,
 272 solutions that “float” in the 3D reconstructed space are equally good in terms
 273 of volume coverage. For the case of the previous example, these are all arm
 274 configurations that occupy some of the reconstructed volume. From those,
 275 the additional surface coverage term selects the one that best matches also
 276 the visible surface of the 3D reconstruction. If the 3D reconstruction did not
 277 suffer from these extra, non-carved voxels, then the volume term would be
 278 enough to guide optimization and the surface term would not introduce any
 279 further useful constraints.

280 3.4. Particle Swarm Optimization

281 Particle Swarm Optimization (PSO) is a popular optimization algorithm
282 that was introduced by Kennedy and Eberhart in [8, 9]. PSO looks for the
283 optimum of an objective function employing a population of entities that
284 evolve according to rules that emulate social interaction.

285 Central to PSO are the notions of *particles* and *generations*. A particle
286 holds a position/candidate solution in the parametric space where the search
287 is performed. Each particle can estimate the fitness of its position by eval-
288 uating the objective function at that point. Each particle is aware of the
289 position at which it has achieved its own best objective function value. It
290 also knows the global best position that has ever been achieved by any of the
291 rest of the particles. Two forces are defined that attract a particle to these
292 two positions. The particles evolve themselves by moving in the search space
293 under the previously described forces in iterations called generations. The
294 details of this process are provided in [15].

295 It has been observed that given enough particles and generations, the
296 swarm reaches the global minimum of the objective function. The required
297 number of particles and generations is problem-dependent and, thus, experi-
298 mentally identified. A number of studies have shown that PSO is very compe-
299 tent in optimizing complex, multidimensional, multimodal, non-differentiable
300 objective functions. The product of the number of particles to that of gener-
301 ations determines the computational requirements of the optimization pro-
302 cess. This is because this product represents the number of objective function
303 evaluations that constitutes the most computationally demanding part of the
304 algorithm.

305 Typically, the particles are initialized at random positions and zero veloc-
306 ities. Each dimension of the multidimensional parameter space is bounded in
307 some range. As in [15], if during the position update a particle has a velocity
308 that forces it to move to a point p_o outside the bounds of the parameter

space, that particle effectively moves to the point p_b inside the bounds that minimizes the distance $|p_o - p_b|$.

3.5. Baseline PSO (bPSO)

In this work, PSO operates on the 35-dimensional 3D body pose parameter space. This also implies that the intrinsic human model parameters (lengths and radii of the primitives of the human model) need to be known in advance. The objective function to be optimized (i.e., minimized) is $E(O, h)$ (Eq. 3) and the population is a set of candidate 3D body poses hypothesized for a single frame. Thus, the process of tracking a human requires the solution of a sequence of optimization problems, one for each acquired frame. By exploiting temporal continuity, the solution over frame F_t is used to generate the initial population for the optimization problem for frame F_{t+1} . More specifically, the first member of the population h_{ref} for frame F_{t+1} is the solution for frame F_t . This implies that for the first frame F_0 , a human body configuration close to the actual one needs to be provided. The rest of the population consists of perturbations of h_{ref} . The variance of these perturbations is experimentally determined and depends on the characteristics of the observed motion and the image acquisition frame rate. The optimization for frame F_{t+1} is executed for a fixed amount of generations. After all generations have evolved, the best hypothesis h_{best} constitutes the solution for time step $t + 1$.

3.6. Perturbed PSO (pPSO)

It has been verified experimentally that *bPSO* is competent in estimating the 6D global pose of the human body. However, the estimation of the 29 remaining parameters that are related to limb angles is not equally satisfactory. The swarm often gets stuck to local minima. To overcome this problem and to increase accuracy, we propose a PSO variant which we call

336 *pPSO* that performs systematic perturbations/randomization on the articu-
 337 lation parameters. More specifically, the human body model is decomposed
 338 into seven branches, as shown in Figure 2. Each branch consists of a set
 339 b_p of primitives and has a set b_d of internal articulation parameters. *pPSO*
 340 operates exactly as *bPSO* for a percentage of its generations. This percent-
 341 age has been identified experimentally to be 40%. After those generations,
 342 each particle is perturbed in a very specific way. First, one branch is ran-
 343 domly selected. Then, only the parameters of this branch are perturbed
 344 and replicated in the global particle representation. Additionally, the local
 345 (particle-dependent) best position for this particle is reset to the new particle
 346 position. After each and every particle is perturbed in this way, all particles
 347 are left to interact as in the *bPSO* scheme for g_p generations. In all reported
 348 experiments, the value of g_p was set to 6 generations. The process is repeated
 349 until the rest 60% of the PSO generations are lapsed.

350 Two different perturbation strategies have been identified and tested. In
 351 the first one, samples are drawn from a uniform distribution in the range
 352 of minimum/maximum allowed values of the corresponding parameter. In
 353 the second case, samples are drawn from a Gaussian distribution centered
 354 at the particle’s previous position and with a standard deviation equal to
 355 one sixth of the range of the corresponding uniform distribution. It has
 356 been verified experimentally that both result in the same tracking error but
 357 the error variance is slightly higher for the case of Gaussian perturbation.
 358 In fact, Gaussian perturbation performs better on slow actions, but worse
 359 on fast actions. This can be explained by the fact that Gaussian sampling
 360 performs a more local search in the parameter space compared to the uniform
 361 sampling, making it more difficult to recover from track losses.

362 Special care should be taken when a particle is perturbed with respect to
 363 its torso or hip branches. As shown in Figure 2, these two branches are not
 364 leafs in the kinematics hierarchy. Thus, the perturbation of these branches

365 affects the parameters of the rest of the branches, too. For this reason, as
 366 soon as these branches are perturbed, the global human kinematics model
 367 consistency needs to be enforced. This is achieved by employing inverse
 368 kinematics. Consequently, a perturbation on the torso or the hip will in fact
 369 influence most, if not all the 35 parameters.

370 The particular scheme for perturbing particles/candidate solutions is jus-
 371 tified by the study of the morphology of the human body and the objective
 372 function of the optimization problem. The human torso accounts for most
 373 of the body’s volume and, therefore, for the largest part of the objective
 374 function. Fine tuning a solution requires checking alternative configurations
 375 of the human limbs that are much smaller in size and less influential to the
 376 objective function. Thus, a targeted particle perturbation that affects only
 377 a branch at a time gives more chances to the algorithm to explore the true
 378 minimum of the objective function.

379 Another reason why perturbation proves valuable stems from the coarse-
 380 ness of the employed model. Consider, for example, the case in which the
 381 arms of a subject are stretched straight and turned around at the shoulder
 382 joint for 90 degrees. In this particular case, the method will probably loose
 383 tracking of the roll around the shoulder joint because this motion does not
 384 produce some significant, observable difference in the volume occupied and
 385 the surface covered by the subject. Due to the perturbation step of *pPSO*,
 386 several solutions relatively far from the computed one will be tested. This
 387 prevents *pPSO* from getting trapped in local minima and enables the effective
 388 tracking of subsequent, unambiguous motions.

389 3.7. *Hybrid human body pose tracking (HYBRID)*

390 As stated in sections 3.5 and 3.6, both *bPSO* and *pPSO* require

- 391 • Knowledge of the human body shape parameters (i.e., lengths, radii of
 392 the geometric primitives comprising the human body model).

- 393 • A coarse estimation of the human body pose for the first frame of a
394 sequence.

395 These requirements hinder the practical exploitation of these algorithms be-
396 cause their fulfilment is associated with considerable effort. In order to alle-
397 viate this problem, we capitalize on the *OpenNI* appearance-based method
398 for human skeleton estimation [17] to come up with a new variant, which we
399 call *HYBRID* and which operates as follows.

400 At a first stage, given an input RGBD sequence of an articulating human,
401 the *OpenNI* method is employed to estimate the articulation. The result of
402 this process is twofold:

- 403 1. The lengths of the parts of a skeletal model of the human body for each
404 frame of the sequence.
- 405 2. An estimation of the human body pose for each frame of the sequence.

406 Based on (1), we compute a human body model of constant shape pa-
407 rameters that consists of primitives that are compatible to those estimated
408 by the *OpenNI* method. More specifically, on top of the 35 mobilities of
409 our model, 9 parameters (constant for each sequence) are added to be able
410 to fit the human body model to a particular subject. These are the upper
411 body length (UBL), the lower body length (LBL), the shoulders-neck dis-
412 tance (SND), the head neck distance (HND), the legs-hip distance (LHD),
413 the back-arm length (BAL), the forearm length (FAL), the back-leg length
414 (BLL) and the front-leg length (FLL). Table 2 presents ground truth values
415 as well as estimations of these parameters for a number of subjects. The
416 parenthesis next to the name of each parameter refers to the corresponding
417 body segment(s) in Figure 2.

418 For a certain sequence, the human skeletons estimated by *OpenNI* provide
419 the 3D positions of human body joints and extremities. For each valid frame,
420 the aforementioned distances are calculated (9 parameters, 15 distances, since

some appear twice for the two arms and legs). The median value for each parameter across all the sequence is selected as the representative one. Other dimensions (i.e, radii of primitives) are set accordingly, based on anatomical studies.

These parameters are then used to define the human body model and inverse kinematics fits the model to the *OpenNI* solution for the first frame. Then, *pPSO* is employed to track the derived human model. Moreover, the solution suggested for each frame by *OpenNI* identifies a particle for *pPSO*. This has the additional advantage that in case of a tracking loss from *pPSO*, tracking can be recovered by considering the fairly accurate *OpenNI* recommendation as a candidate solution.

4. Experimental evaluation

The experimental evaluation of the proposed method was based on the Berkeley Multimodal Human Action Database (MHAD) [13]. This dataset features 12 human subjects.

Figure 3 shows one frame of each subject. From this figure it can be verified that the employed data set includes subjects of considerable variability with respect to age, size and body types. This is also shown quantitatively in Table 2.

The subjects of the dataset perform 11 different activities (01-jumping, 02-jumping jacks, 03-bending, 04-punching, 05-waving two hands, 06-waving one hand, 07-clapping, 08-throwing, 09-sit down/stand up, 10-sit down and 11-stand up). In each sequence, each activity is repeated several times. A motion capture system has been used to provide ground truth information regarding the position of all joints in all sequences. Additionally, the activities are recorded with a multicamera setup consisting of several conventional cameras as well as by two extrinsically calibrated Kinect sensors. In all experiments reported in this paper, the RGBD data provided by the two

449 Kinect sensors feed the proposed methods. The resulting tracking results are
 450 compared against the ground truth resulting from the motion capture data.

451 To quantify the accuracy in body pose estimation, we adopt the metric
 452 used in [7]. More specifically, the distance between a set of corresponding
 453 3D points in the ground truth and in the estimated body model is measured.
 454 Each such point (four per leg, three per arm and one for the head) is marked
 455 in Figure 2 with a red “x”. The average of all these distances over all the
 456 frames of the sequence constitutes the resulting error estimate Δ . Another
 457 metric reports the percentage of these distances that are within some pre-
 458 defined threshold A_t . We will refer to this metric as the accuracy in human
 459 body pose estimation. A_t was set to $10cm$ for all experiments. For example,
 460 an accuracy of $A_t = 70\%$ for a sequence means that 70% of the joints were
 461 estimated at positions that are within less than $10cm$ from the ground truth,
 462 in all frames.

463 Several experiments were carried out to assess quantitatively and qualita-
 464 tively the accuracy and the performance of the proposed human articulation
 465 tracking method. The goal of the first experiment was to assess the error in
 466 joints position estimation as a function of the computational budget devoted
 467 to PSO. To do so, we choose one of the sequences of the MHAD dataset that
 468 consists of 80 consecutive human poses showing a human performing activity
 469 02 (jumping jacks). The rationale for selecting this particular activity and
 470 sequence is that (a) it is executed in high speed and (b) it involves the whole
 471 body, so all body model parameters change values as a function of time.
 472 Thus, it is expected that this sequence constitutes a worst case scenario, at
 473 least among activities represented in the specific dataset.

474 Figure 4 illustrates the error Δ in joints position estimation as a function
 475 of the *pPSO* parameters (number of generations and particles per genera-
 476 tion). As explained in Section 3, the product of these parameters determines
 477 the computational budget of the proposed methodology, as it accounts for

478 the number of objective function evaluations. The horizontal axis of the plot
 479 denotes the number of PSO generations. Each plot of the graph corresponds
 480 to a different number of particles per generation. Each point in each plot
 481 is the average of the error Δ for 5 runs of an experiment with the specific
 482 parameters. A first observation is that Δ decreases monotonically as the
 483 number of generations increase. Additionally, as the particles per generation
 484 increase, the resulting error decreases. Nevertheless, employing more than 65
 485 generations and more than 200 particles results in a reduction of the error Δ
 486 that is disproportionally low compared to the increase in the required com-
 487 putational budget. For this reason, 200 particles evolving in 65 generations
 488 was retained in all further experiments.

489 The second experiment aimed at evaluating the performance of the meth-
 490 ods across different human subjects. All twelve sequences showing the twelve
 491 different subjects performing the same activity (activity 04, boxing) were
 492 considered. *bPSO* and *pPSO* require knowledge of the parameters of the
 493 human body models as well as body configuration parameters for the first
 494 frame of a sequence. In our experiments, these subject-specific model pa-
 495 rameters and initial model configurations were estimated manually for *bPSO*
 496 and *pPSO*, and automatically for *HYBRID*, based on the results obtained by
 497 the *OpenNI* method (see Section 3.7). Additionally, the *pPSO*, *bPSO* and
 498 *HYBRID* methods were assigned exactly the same computational budget.

499 Figure 5 illustrates the error Δ and the accuracy of the *pPSO* and the
 500 *HYBRID* methods. For the purposes of comparative evaluation, errors and
 501 accuracies are also provided for *bPSO* and for the *OpenNI* skeleton estima-
 502 tion method [17]. Table 1 summarizes the individual errors and accuracies
 503 shown in Figure 5. It can be verified that the *pPSO* method outperforms
 504 all other methods in all aspects (average error, standard deviation of er-
 505 ror and accuracy). *HYBRID* outperforms the *bPSO* and *OpenNI*, showing
 506 that model-based optimization improves the guess made by *OpenNI*. Still,

Method	Mean Δ	Std. Δ	Accuracy (%)
<i>OpenNI</i>	52.9	49.5	87.3
<i>bPSO</i>	62.2	69.5	82.3
<i>pPSO</i>	41.8	33.1	94.4
<i>HYBRID</i>	45.5	44.1	92.5

Table 1: Comparison of *pPSO* and *HYBRID* with the baseline PSO method (*bPSO*) and the *OpenNI* method for the case of different humans performing the same action (boxing).

507 *HYBRID* does not outperform *pPSO*. The reason for this is that *HYBRID*
508 operates on automatically estimated human body models that are less ac-
509 curate compared to the ones on which *pPSO* operates (identified manually).
510 Table 2 shows characteristic distances and metrics regarding the 12 human
511 subjects as measured manually (columns G) and as estimated in the *HY-*
512 *BRID* algorithm (columns H). It can be verified that the body dimensions
513 estimated by *HYBRID* deviate considerably from the actual, ground truth
514 measurements.

515 One interesting question that arises is how *HYBRID* would perform if
516 it was provided with the manually estimated human body models on which
517 *pPSO* operates. It turns out that in this case, mean Δ is 40.9, the standard
518 deviation of Δ is 32.6 and the accuracy is 94.7%. Thus, *HYBRID* outperforms
519 marginally *pPSO* in that case. The slight difference in performance between
520 *HYBRID* and *pPSO* is explained by the fact that no track loss occurred
521 during *pPSO* tracking from which *HYBRID* could recover.

522 In a third experiment, the goal was to assess the proposed method with
523 respect to different activities. For that purpose, the evaluation was performed
524 on image sequences showing a single subject performing the eleven different

Subject	S01		S02		S03		S04		S05		S06	
Metric	G	H	G	H	G	H	G	H	G	H	G	H
UBL (HI)	26	19	30	21	33	21	29	20	32	22	28	20
LBL (IJ)	15	19	17	21	18	21	17	20	19	22	17	20
SND (CH, C'H)	19	15	19	15	17	14	15	15	17	16	19	14
HND (GH)	20	25	20	25	20	25	20	21	20	26	20	20
LHD (FJ, F'J)	10	9	11	9	10	8	9	9	9	9	9	8
BAL (BC, B'C')	24	25	28	27	31	28	24	23	26	28	26	26
FAL (AB, A'B')	23	26	25	31	26	32	24	25	25	31	24	27
BLL (EF, E'F')	36	41	43	47	44	47	37	39	42	45	42	44
FLL (DE, D'E')	42	37	48	42	47	43	41	35	45	42	45	41

(a)

Subject	S07		S08		S09		S10		S11		S12	
Metric	G	H	G	H	G	H	G	H	G	H	G	H
UBL (HI)	25	20	30	20	27	20	27	21	28	20	24	20
LBL (IJ)	15	20	18	20	15	20	16	21	16	20	21	20
SND (CH, C'H)	17	17	18	15	15	13	17	14	17	15	18	14
HND (GH)	20	20	20	21	20	17	20	24	20	25	20	24
LHD (FJ, F'J)	8	7	9	9	8	7	8	8	8	9	9	8
BAL (BC, B'C')	22	25	24	26	23	22	27	25	26	27	25	25
FAL (AB, A'B')	22	24	24	26	23	27	24	29	24	28	22	27
BLL (EF, E'F')	35	39	39	40	35	42	41	43	41	44	38	43
FLL (DE, D'E')	41	35	43	39	41	37	43	40	44	41	41	39

(b)

Table 2: Characteristic measures of the shape of the human subjects of the MHAD dataset, (a) subjects 01-06, (b) subjects 07-12. Columns (G) are the manually measured, ground truth values and columns (H) the one estimated by the *HYBRID* method. The parenthesis next to the name of each measure refers to the corresponding body segment(s) in Figure 2.

Method	Mean Δ	Std. Δ	Accuracy (%)
<i>OpenNI</i>	54.5	46.2	86.3
<i>bPSO</i>	50.6	48.8	89.5
<i>pPSO</i>	39.3	27.3	96.3
<i>HYBRID</i>	42.8	25.2	96.3

Table 3: Comparison of *pPSO* and *HYBRID* with the baseline PSO method (*bPSO*) and the *OpenNI* method for the case of all actions performed by the same subject (subject 09).

activities. Figure 6 illustrates the obtained results in a way analogous to that of Figure 5. Again, *pPSO* outperforms the rest of the methods with respect to the mean error Δ . It should also be noted that for actions like bending (action 03) and sit-down/stand-up (action 09) that exhibit considerable self- and body-object occlusions, the proposed method performs considerably better. In this experiment, the *HYBRID* method has the smallest error variance and equal accuracy to that of *pPSO*.

We again tested the performance of *HYBRID* in the case that it is fed with the manually estimated human body model for that subject. It turns out that mean Δ is 37.9, the standard deviation of Δ is 22.9 and the accuracy is 98.2%. Thus, *HYBRID* outperforms *pPSO* in that case, showing that, given accurate models, the combination of the bottom-up *OpenNI* method with the top-down *pPSO* method improves the tracking performance.

Finally, Table 4 summarizes all performed experiments. It can be verified that *pPSO* achieves a significant reduction in mean error and error variance compared to the rest of the methods as well as a significant increase in accuracy. If *HYBRID* is employed on the manually estimated human body models, then mean Δ becomes 39.7, the standard deviation of Δ becomes 28.4

Method	Mean Δ	Std. Δ	Accuracy (%)
<i>OpenNI</i>	54.5	46.2	86.3
<i>bPSO</i>	50.6	48.8	89.5
<i>pPSO</i>	39.3	27.3	96.3
<i>HYBRID</i>	44.5	35.4	94.2

Table 4: Aggregate comparison of *pPSO* and *HYBRID* with the baseline PSO method (*bPSO*) and the *OpenNI* method for all tested sequences.

and the accuracy becomes 96.3%. So, its performance is considerably improved, but it does not outperform that of *pPSO*.

Figure 7 shows characteristics snapshots of the MHAD dataset and the skeletons that have been extracted by the *pPSO*, *OpenNI* and *bPSO* methods superimposed on the RGB frame of one of the two employed RGBD sensors. Finally, Figure 8 provides additional characteristic examples of the solutions provided by the *pPSO* method. A much more complete qualitative assessment of the performance of the proposed method can be performed based on the supplementary material accompanying this paper which is available at <http://youtu.be/n5irgHVuFwc>. It should be noted that no temporal smoothing has been performed between successive frames.

The proposed method runs on a computer equipped with a 8-core Intel i7 950 CPU, 4 GBs RAM. On this system, the average computing time for our non-optimized CPU-only implementation is 20 sec/frame. However, all involved computations are inherently data parallel and tailored for a GPU implementation. This is also evidenced by the real-time performance (20 fps) that is achievable by GPU implementations of similar approaches for the case of 3D hand tracking [15].

561 5. Discussion

562 We proposed a model-based method for tracking the articulated motion
563 of the human body using a volumetric 3D representation that is built by
564 fusing the depth measurements provided by two calibrated RGBD sensors.
565 The proposed method follows a hypothesize-and-test approach that casts
566 the articulated motion tracking problem into a search problem in a high-
567 dimensional space. Searching is performed with a stochastic optimization
568 technique, called PSO, resulting in a baseline implementation called *bPSO*.
569 We also proposed a perturbation scheme that is applied on top of the *bPSO*
570 solutions that results in the *pPSO* method. Finally, in order to raise the
571 practical difficulties the limitations of *pPSO* with respect to its need for
572 tailored human models and initialization in the first frame, we proposed the
573 *HYBRID* method that combines *pPSO* with *OpenNI*. A series of experiments
574 performed on a ground-truth-annotated data set demonstrated quantitatively
575 and qualitatively that *pPSO* outperforms in error and accuracy the rest of
576 the methods. This is even more striking in the challenging cases where the
577 body configuration exhibits significant self occlusions. Thus, in situations
578 where small error and high accuracy is more important than the burden and
579 the overhead of using a second RGBD sensor, the proposed *pPSO* marker-
580 less human articulations tracking method constitutes an attractive approach.
581 *HYBRID* performs worse than *pPSO* because of the less accurate (yet auto-
582 matic) estimation of the human body models. Still, the fact that *HYBRID*
583 is fully automatic, is a significant advantage that, depending on application,
584 might be more important than its lacking accuracy. In fact, as demonstrated
585 experimentally, if *HYBRID* is given the chance to operate on accurate (non-
586 automatically extracted) human body models, it performs comparably and,
587 in some cases, better compared to *pPSO*.

588 Acknowledgements

589 This work was partially funded by the European Commission under con-
590 tract FP7-IST-288146 HOBbit and by the European Union (European So-
591 cial Fund - ESF) and Greek national funds through the Operational Pro-
592 gram “Education and Lifelong Learning” of the National Strategic Refer-
593 ence Framework (NSRF) - Research Funding Project: THALIS-UOA- ERA-
594 SITECHNIS MIS 375435.

- 595 [1] Bisacco, A., Ming-Hsuan, Y., Soatto, S., 2007. Fast human pose es-
596 timation using appearance and motion via multi-dimensional boosting
597 regression. In: IEEE Computer Vision and Pattern Recognition.
- 598 [2] Chen, L., Wei, H., Ferryman, J., 2013. A survey of human motion anal-
599 ysis using depth imagery. Pattern Recognition Letters 34 (15), 1995 –
600 2006.
- 601 [3] Corazza, S., Mundermann, L., Gambaretto, E., Ferrigno, G., Andriac-
602 chi, T., 2010. Markerless motion capture through visual hull, articulated
603 icp and subject specific model generation. International Journal of Com-
604 puter Vision 87 (1-2), 156–169.
- 605 [4] Deutscher, J., Reid, I., 2005. Articulated body motion capture by
606 stochastic search. International Journal of Computer Vision 61 (2), 185–
607 205.
- 608 [5] Gall, J., Rosenhahn, B., Brox, T., Seidel, H.-P., 2010. Optimization and
609 filtering for human motion capture. International Journal of Computer
610 Vision 87 (1-2), 75–92.
- 611 [6] Gall, J., Stoll, C., de Aguiar, E., Theobalt, C., Rosenhahn, B., Seidel,
612 H. P., 2009. Motion capture using joint skeleton tracking and surface

estimation. In: IEEE Computer Vision and Pattern Recognition. pp. 1746–1753.

[7] Hamer, H., Schindler, K., Koller-Meier, E., Van Gool, L., 2009. Tracking a Hand Manipulating an Object. In: IEEE International Conference on Computer Vision.

[8] Kennedy, J., Eberhart, R., Jan. 1995. Particle Swarm Optimization. In: International Conference on Neural Networks. Vol. 4. IEEE, pp. 1942–1948.

[9] Kennedy, J., Eberhart, R., Yuhui, S., 2001. Swarm intelligence. Morgan Kaufmann.

[10] Mikic, I., Trivedi, M., Hunter, E., Cosman, P., 2003. Human body model acquisition and tracking using voxel data. International Journal of Computer Vision 53 (3), 199–223.

[11] Moeslund, T. B., Hilton, A., Kru, V., 2006. A Survey of Advances in Vision-based Human Motion Capture and Analysis. Computer Vision and Image Understanding 104, 90–126.

[12] Mussi, L., Ivekovic, S., Cagnoni, S., 2010. Markerless articulated human body tracking from multi-view video with gpu-pso. In: Tempesti, G., Tyrrell, A., Miller, J. (Eds.), Evolvable Systems: From Biology to Hardware. Vol. 6274 of Lecture Notes in Computer Science. pp. 97–108.

[13] Ofli, F., Chaudhry, R., Kurillo, G., Vidal, R., Bajcsy, R., 2013. Berkeley MHAD: A Comprehensive Multimodal Human Action Database. In: IEEE Workshop on Applications on Computer Vision (WACV).

[14] Oikonomidis, I., Kyriazis, N., Argyros, A. A., 2010. Markerless and Efficient 26-DOF Hand Pose Recovery. Asian Conference on Computer Vision, 744–757.

- 639 [15] Oikonomidis, I., Kyriazis, N., Argyros, A. A., 2011. Efficient Model-
640 based 3D Tracking of Hand Articulations using Kinect. In: British Ma-
641 chine Vision Conference. Dundee, UK.
- 642 [16] Oikonomidis, I., Kyriazis, N., Argyros, A. A., 2011. Full DOF Track-
643 ing of a Hand Interacting with an Object by Modeling Occlusions and
644 Physical Constraints. In: IEEE International Conference on Computer
645 Vision.
- 646 [17] OpenNI, November 2010. OpenNI User Guide. OpenNI organization,
647 last viewed 19-01-2011 11:32.
648 URL <http://www.openni.org/documentation>
- 649 [18] Pons-Moll, G., Leal-Taixe, L., Truong, T., Rosenhahn, B., 2011. Effi-
650 cient and robust shape matching for model based human motion cap-
651 ture. In: Mester, R., Felsberg, M. (Eds.), Pattern Recognition. Vol. 6835
652 of Lecture Notes in Computer Science. Springer Berlin Heidelberg, pp.
653 416–425.
- 654 [19] Poppe, R., 2007. Vision-based human motion analysis: An overview.
655 Computer Vision and Image Understanding 108 (1-2), 4 – 18, special
656 Issue on Vision for Human-Computer Interaction.
- 657 [20] Shotton, J., Fitzgibbon, A., Cook, M., Sharp, T., Finocchio, M., Moore,
658 R., Kipman, A., Blake, A., 2011. Real-Time Human Pose Recognition
659 in Parts from Single Depth Images.
- 660 [21] Sigal, L., Isard, M., Haussecker, H., Black, M., 2012. Loose-limbed peo-
661 ple: Estimating 3d human pose and motion using non-parametric belief
662 propagation. International Journal of Computer Vision 98 (1), 15–48.

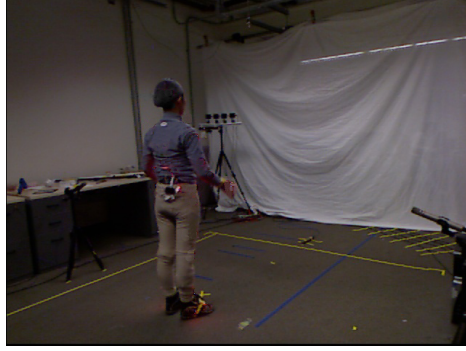
- 663 [22] Sminchisescu, C., Kanaujia, A., Li, Z., Metaxas, D., 2005. Discrimi-
 664 native density propagation for 3d human motion estimation. In: IEEE
 665 Computer Vision and Pattern Recognition. Vol. 1. pp. 390–397 vol. 1.
- 666 [23] Smisek, J., Jancosek, M., Pajdla, T., 2011. 3d with kinect. In: IEEE
 667 ICCV Workshops. pp. 1154–1160.
- 668 [24] Tzevanidis, K., Zabulis, X., Sarmis, T., Koutlemanis, P., Kyriazis, N.,
 669 Argyros, A., 2010. From multiple views to textured 3d meshes: a gpu-
 670 powered approach. In: ECCV Workshops 2010. pp. 5–11.
- 671 [25] Vicon, 2013. Vicon: Motion capture systems.
 672 URL <http://www.vicon.com>
- 673 [26] Vijay, J., Trucco, E., Ivekovic, S., 2010. Markerless human articulated
 674 tracking using hierarchical particle swarm optimisation. Image and Vi-
 675 sion Computing 28 (11), 1530–1547.
- 676 [27] Wilson, J. L., 2010. Microsoft kinect for xbox 360. PC Magazine Com-
 677 munications.
- 678 [28] Zhang, L., Sturm, J., Cremers, D., Lee, D., Oct. 2012. Real-time human
 679 motion tracking using multiple depth cameras. In: Proceedings of the
 680 International Conference on Intelligent Robot Systems (IROS).



(a)



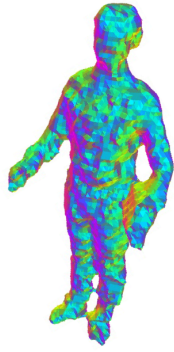
(b)



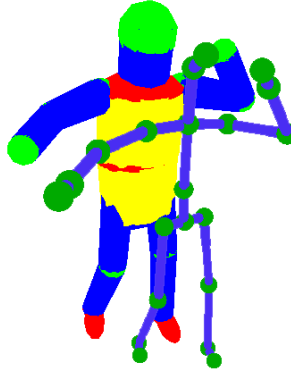
(c)



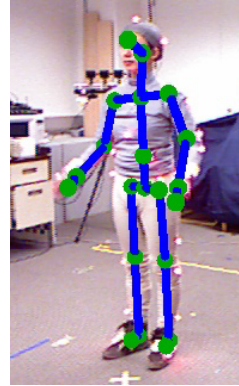
(d)



(e)



(f)



(g)

Figure 1: Graphical illustration of the proposed method. Two RGB frames ((a), (c)) and the corresponding depth maps ((b), (d)). The volume (e) occupied by the person is reconstructed using the depth maps. The proposed method fits the employed human body model (f) to this volume, recovering the body articulation (g).

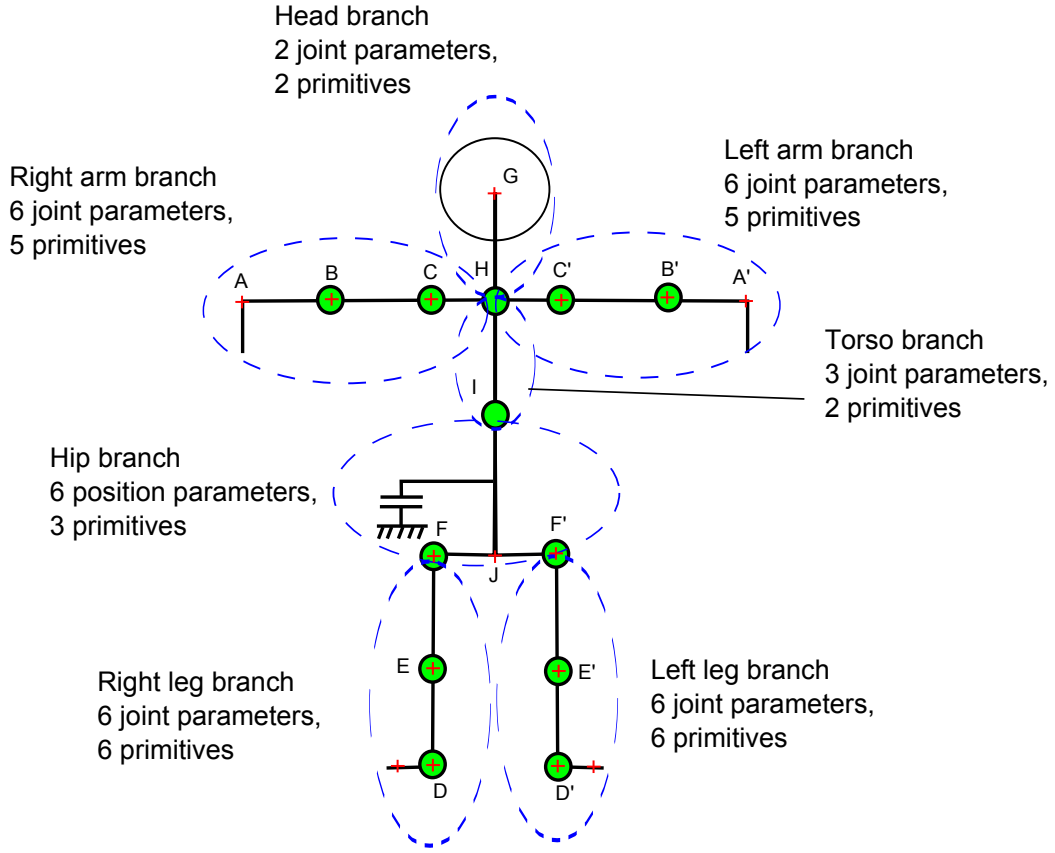


Figure 2: Definition of human body branches. The perturbation of the torso and hip branches results in changes in the parameters of their child branches. Model points with a red "x" denote joints whose 3D position is taken into account in defining the tracking error in the quantitative experimental evaluation of the method.



Figure 3: The twelve subjects of the MHAD dataset.

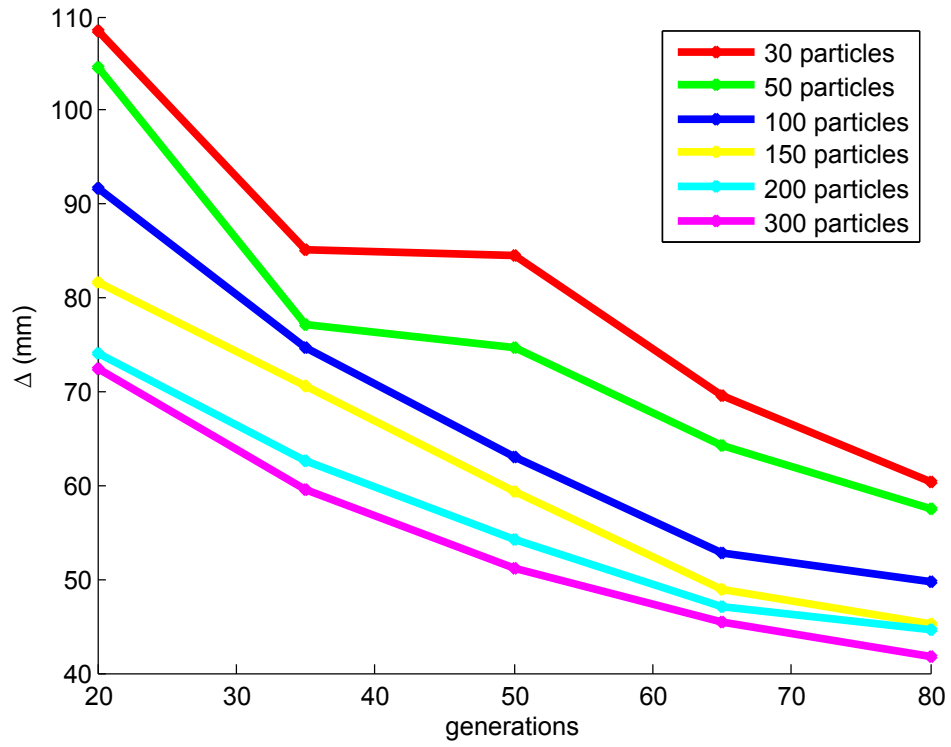
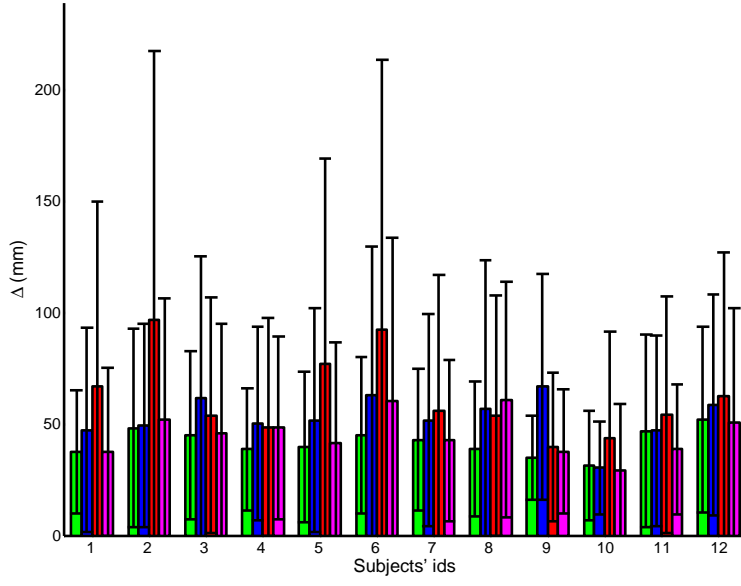
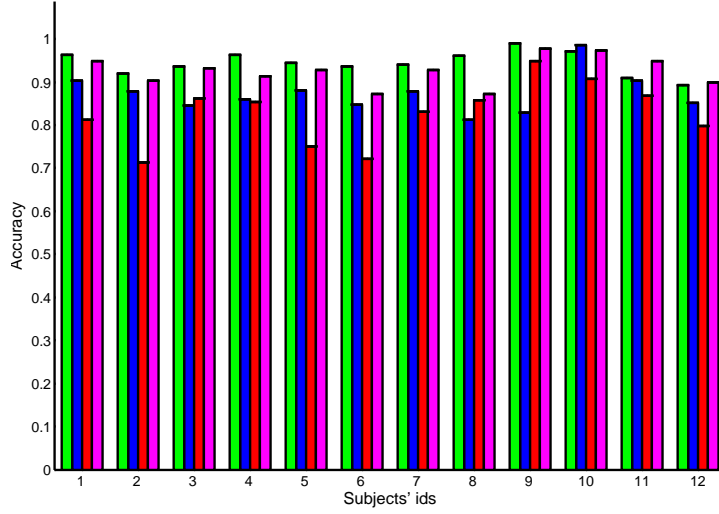


Figure 4: Quantitative evaluation of the performance of the method with respect to the *pPSO* parameters.

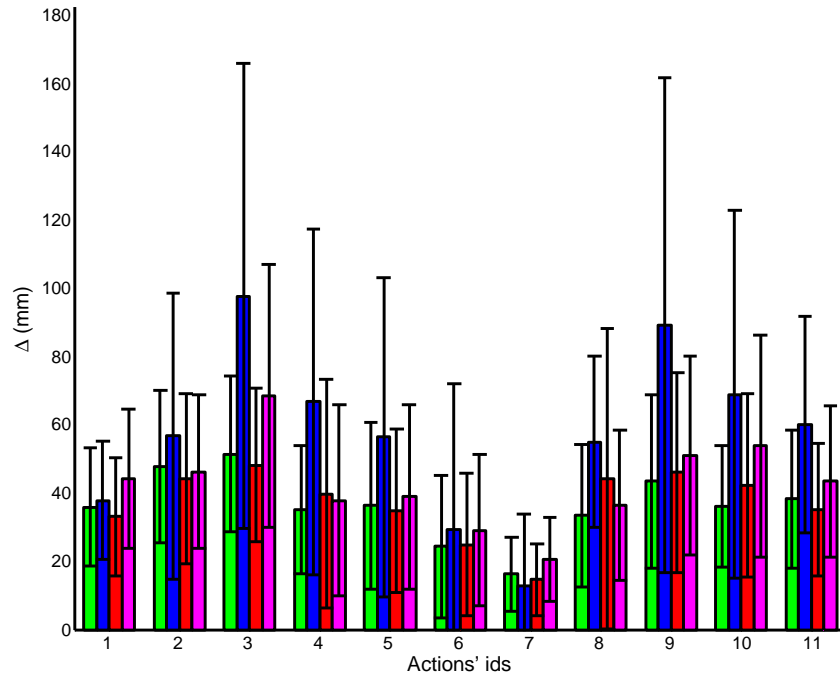


(a)

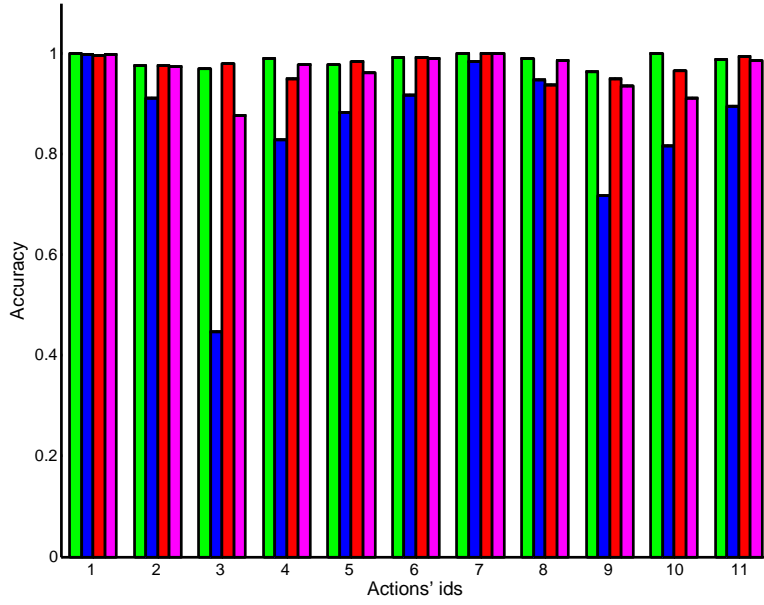


(b)

Figure 5: Quantitative evaluation of the method applied to 12 subjects performing the same action (boxing). (a) Error Δ and variances, (b) accuracy for the proposed method (*pPSO*, green bars), baseline method (*bPSO*, red bars), *OpenNI* human skeleton estimation (blue bars) and *HYBRID* (purple bars).



(a)



(b)

Figure 6: Quantitative evaluation of the method applied to 11 actions performed by the same person. (a) Error Δ and variances, (b) accuracy for the proposed method (green bars) *OpenNI* human skeleton estimation (blue bars) and *HYBRID* method (purple bars).

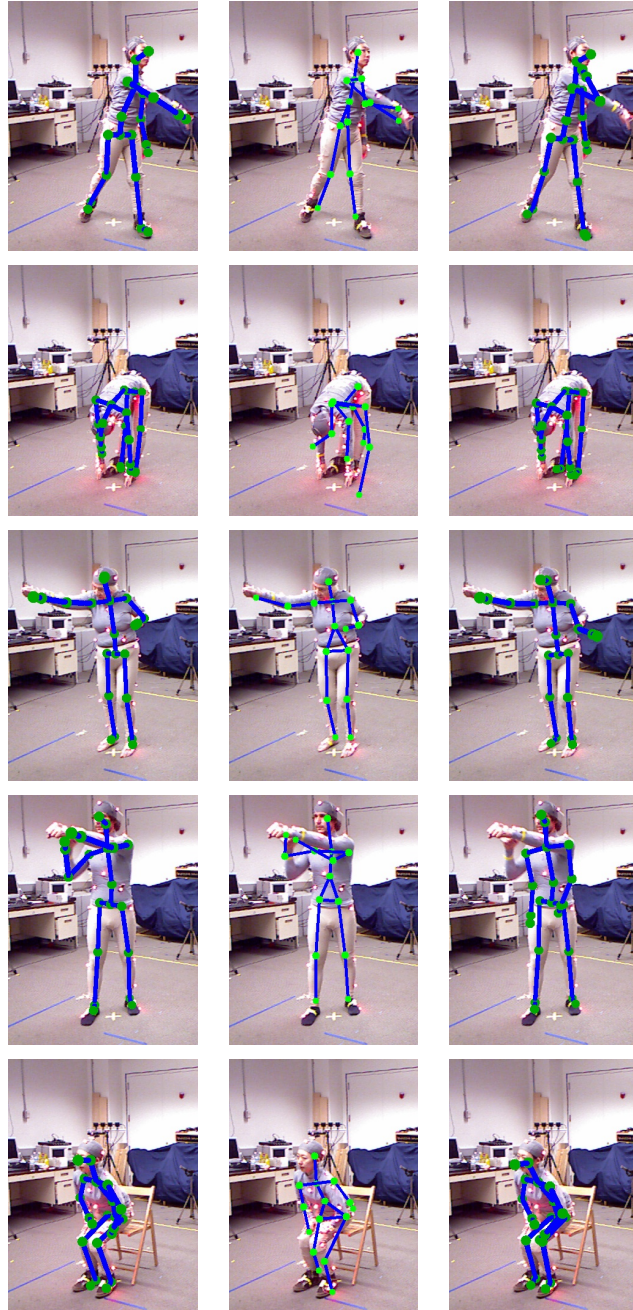


Figure 7: Qualitative comparison of the *pPSO* (left), OpenNI (middle) and *bPSO* (right) methods based on characteristic frames of the MHAD dataset.



Figure 8: Various configurations on different subjects evaluated by the method.