

Discovery and recognition of motion primitives in human activities

Marta Sanzari¹, Valsamis Ntouskos¹, Fiora Pirri¹,

¹ Dipartimento di Ingegneria Informatica Automatica e Gestionale, University of Rome 'Sapienza', Alcor LAB

Abstract

We present a novel framework for the automatic discovery and recognition of motion primitives in videos of human activities. Given the 3D pose of a human in a video, human motion primitives are discovered by optimizing the ‘motion flux’, a quantity which captures the motion variation of a group of skeletal joints. A normalization of the primitives is proposed in order to make them invariant with respect to a subject anatomical variations and data sampling rate. The discovered primitives are unknown and unlabeled and are unsupervisedly collected into classes via a hierarchical non-parametric Bayes mixture model. Once classes are determined and labeled they are further analyzed for establishing models for recognizing discovered primitives. Each primitive model is defined by a set of learned parameters. Given new video data and given the estimated pose of the subject appearing on the video, the motion is segmented into primitives, which are recognized with a probability given according to the parameters of the learned models. Using our framework we build a publicly available dataset of human motion primitives, using sequences taken from well-known motion capture datasets. We expect that our framework, by providing an objective way for discovering and categorizing human motion, will be a useful tool in numerous research fields including video analysis, human inspired motion generation, learning by demonstration, intuitive human-robot interaction, and human behavior analysis.

1 Introduction

Activity recognition is widely acknowledged as a core topic in computer vision, witness the huge amount of research done in recent years spanning a wide number of applications from sport to cinema, from human robot interaction to security and rehabilitation.

Activity recognition has evolved from earlier focus on action recognition and gesture recognition. The main difference being that activity recognition is completely general as it concerns any kind of human activity, which can last few seconds or minutes or hours, from daily activities such as cooking, self-care, talking at the phone, cleaning a room, up to sports or recreation such as playing basketball or fishing. Nowadays there are a number of publicly available datasets dedicated to the collection of any kind of human activity, likewise a number of challenges (see for example the ActivityNet challenge [1]).

On the other hand, the interest in motion primitives is due to the fact that they are essential for deploying an activity. Think about sport activities, or cooking, or performing arts, which require to purposefully select a specific sequences of movements. Likewise daily activities such as cleaning, or cooking, or washing the dishes or preparing the table require precise motion sequences to accomplish the task. Indeed, the compositional nature of human

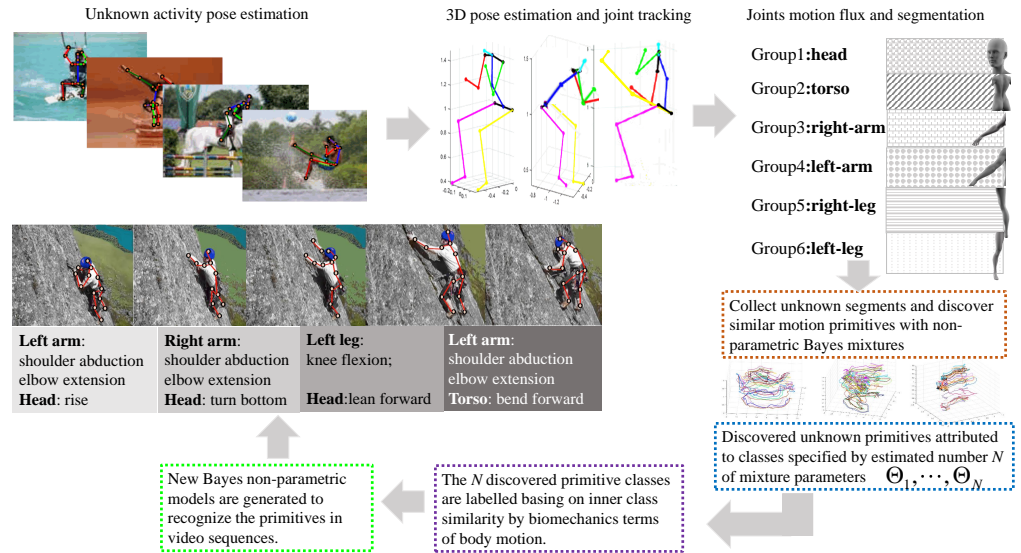


Fig 1. The above schema presents the proposed framework and the process to obtain from video sequences the discovered motion primitives.

activities, under body and kinematics constraints, has attracted the interest of many research areas such as in computer vision [2, 3], in neurophysiology [4, 5], in sports and rehabilitation [6], and in biomechanics [7] and in robotics [8, 9, 10].

The goal of this work is to automatically discover the start and end points where primitives of 6 identified body parts occur throughout the course of an activity, and recognize each of the occurred primitives. The idea is that these primitives sort out a non-complete set of human movements, which combined together can form a wide range of human activities, in so providing a compositional approach to the analysis of human activities.

The steps of the proposed method are as follows. Given a video of a human activity both the 2D pose and 3D pose of the human are estimated (see [11], and also [12]). Once the 3D poses of the joints of interest are determined, we compute the *motion flux*. The motion flux method provides a model from first principles for human motion primitives, and it effectively discovers where primitives begin and end on human activity motion trajectories.

Motion primitives discovered by the motion flux are unknown: they are segments of motion about which only the involved specific body part is known. These primitives are collected into classes by a non-parametric Bayes model, namely the Dirichlet process mixture model (DPM), which gives the freedom to not choose the number of mixture components. By suitably eliminating very small clusters it turns out that discovered primitives can be collected into 69 classes (see Fig. 12). For each of them the mixture model returns a parameter set identifying the precise primitive class. We label the computed parameters with terms taken from the biomechanics of human motion, by inspecting only a representative primitive for each discovered class. Out of these generated classes we form a new layer of the hierarchical model, to generate the parameters for each class, further used for primitives recognition. Under this last models each primitive category is approximated by a DPM with a number of components mirroring the inner idiosyncratic behavior of each primitive class.

Motion primitives classification is finalized by providing a label for each primitive. Namely, given an activity (possibly unknown) and an unknown primitive discovered by motion flux, we find the model the primitive belongs to, hence the primitive is labeled by that model.

Experiments show that the motion flux is a good model for segmenting the motion of body parts. Likewise, the unsupervised non-parametric model provides both a good classification of similar motion primitives and a good estimation of primitive labels, as shown in the results

(see Section 6). The approach therefore is quite general and it turns out to be very useful to any researcher who would like to explore the compositional nature of any activity, using both the proposed method and the motion primitives dataset provided.

To the best of our knowledge just few works, among which we recall [2, 3], have faced the problem of discovering motion primitives in video activities or motion capture (MoCap) sequences, quantitatively evaluating the ability to recognize them.

Despite the lack of works on motion primitives we show that they are quite an expressive *language* for ascertaining specific human behaviors. To prove that, in a final application for video surveillance, described in Section 7, we show that motion primitives can play a compelling role in detecting distinct classes of dangerous activities. In particular, we show that dangerous activities can be detected with off-the-shelf classifiers, once motion primitives have been extracted in the videos. Comparisons with state of the art results prove the relevance of motion primitives in discovering specific behaviors, since motion primitives embed significant time-space features easily usable for classification.

The contributions of the work, schematically shown in Fig. 1 are the followings:

1. We introduce the motion flux method to discover motion primitives, relying on the variation of the velocity of a group of joints.
2. We introduce a hierarchical model for the classification and recognition of the unlabeled primitives, discovered by the motion flux.
3. We show a relevant application of human motion primitives for video surveillance.
4. We created a new dataset of human motion primitives from three public MoCap datasets ([13], [14], [15]).

2 Related work

Human motion primitives are investigated in several research areas, from neurophysiology to vision to robotics and biomechanics. Clearly, any methodology has to deal with the vision process, and many of the earliest more relevant approaches to human motion highlighted that understanding human motion requires view independent representations [16, 17] and that a fine grained analysis of the motion field is paramount to identify primitives of motion. In early days this required a massive effort in visual analysis [18] to obtain the poses, the low level features, and segmentation. Nowadays, scientific and technological advances have made it possible to exploit several methods to measure human motion, such as the availability of a number of MoCap databases [13, 15, 19], see for a review [20]. Furthermore recent findings result in methods that can deliver 3D human poses from videos if not even from single frames [21, 11, 22, 12]. Since then 3D MoCap data have been widely used to study and understand human motion, see for example [23, 24, 25] in which Gaussian Process Latent Variable Models or Dirichlet processes are used to classify actions, or [26] in which a non-parametric Bayesian approach is used to generate behaviors for body parts and classify actions based on these behaviors. In [27] temporal segmentation of collaborative activities is examined, or in [28] different descriptors are exploited to achieve arm-hand action recognition.

Neurophysiology Neurophysiology studies on motion primitives [29, 4, 30, 31, 32, 33] are based on the idea that kinetic energy and muscular activity are optimized in order to conserve energy. In these works it has been observed that curvature and velocity of joint motion are related. Earliest works such as Lacquaniti et al. [34] proposed a relation between curvature and angular velocity. In particular, using their notation, letting C be the curvature and A the angular velocity, they called the equation $A = KC^{\frac{2}{3}}$ the Two-Thirds Power law, valid for certain class of two-dimensional movements. Viviani and Schneider [35] formulated an extension of this law, relating the radius of curvature R at any point s along the trajectory with

the corresponding tangential velocity V , in their notation:

$$V(s) = K(s) \left(\frac{R(s)}{1 + \alpha R(s)} \right)^\beta \quad (1)$$

where the constants $\alpha \geq 0$, $K(s) \geq 0$ and β has a value close to $= \frac{1}{3}$. An equivalent Power law for trajectories in 3D space is introduced by [36] and it is called the curvature-torsion power law and is defined as $\nu = \alpha \kappa^\beta |\tau|^\gamma$, where κ is the curvature of the trajectory, τ the torsion, ν the spatial movement speed, β and γ are constants.

Computer Vision The interpretation of motion primitives as simple individual actions or gestures is often purported, in any case they are related to segmentation of videos and 3D motion capture data. Many approaches explore video sequences segmentation to align similar action behaviors [37] or for spatio-temporal annotation as in [38]. Lu et al. [39] propose to use a hierarchical Markov Random Field model to automatically segment human action boundaries in videos. Similarly, [40] develop a motion capture segmentation method. Besides these works, only [41, 2, 3, 42] have targeted motion primitives, to the best of our knowledge. [41] focuses on 2D primitives for drawing, on the other hand [2] does not consider 3D data and generate the motion field considering Lukas-Kanade optical flow for which Gaussian mixture models are learned. None of these approaches provide quantitative results for motion primitives, but only for action primitives, which makes their method not directly comparable with ours. [3, 42] use 3D data and explicitly mention motion primitives, providing quantitative results. The authors account for the velocity field via optical flow basing the recognition of motion primitives on harmonic motion context descriptors. Since [3] deal only with upper torso gestures we compare with them only the primitives they mention. In [42] the authors achieve motion primitives segmentation from wrist trajectories of sign language gestures, obtaining unsupervised segmentation with Bayesian Binning. Again here no comparison for motion primitives discovery or recognition is possible as original data are not available.

Robotics In robotics the paradigm of transferring human motion primitives to robot movements is paramount for imitation learning and, more recently to implement human-robot collaboration [43]. A good amount of research in robotics has approached primitives in terms of Dynamic Movement Primitives (DMP) [43] to model elementary motor behaviors as attractor systems, representing them with differential equations. Typical applications are learning by imitation or learning from demonstration [44, 45, 46, 47], learning task specifications [48], modeling interaction primitives [8]. Motion primitives are represented either via Hidden Markov models or Gaussian Mixture Models (GMM). [49] present an approach based on HMM for imitation learning of arm movements, and [50] model arm motion primitives via GMM.

It is apparent that in most of the approaches motion primitives are only observed and modeled, instead we are able to learn and model them using respectively the *motion flux* quantity and a hierarchical model. The main contribution of our work is indeed the introduction of a new ability for a robot to automatically discover motion primitives observing 3D joints raw pose data. The outcome of our approach is also a motion primitives dataset not requiring human manual operation.

Our view of motion primitive shares the hypothesis of energy minimality during motion, fostered by neurophysiology, likewise the idea to characterize movements using the proper geometric properties of the skeleton joints space motion. However, for primitive discovery, we go beyond these approaches capturing the variation of the velocity of a group of joints using this as the baseline for computing the change in motion by maximizing the motion flux.

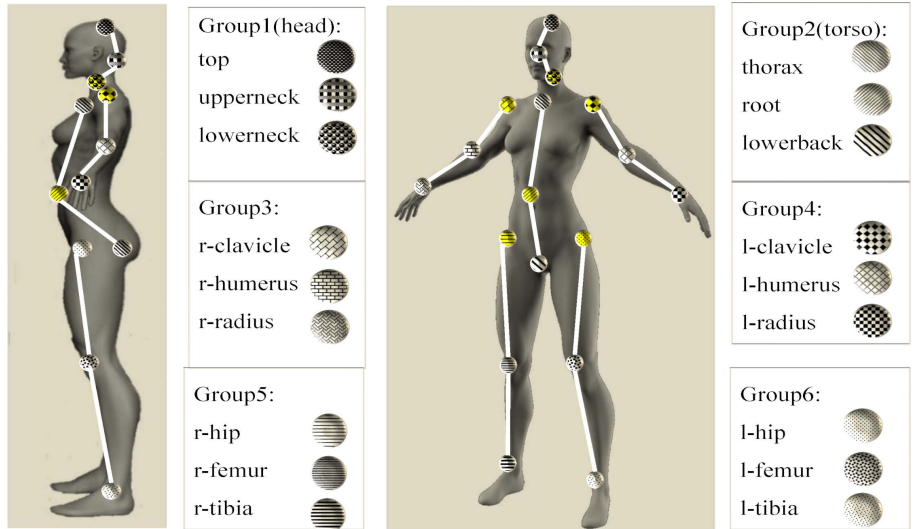


Fig 2. The six groups partitioning the human body with respect to motion primitives are shown, together with the joints specifying each group and the skeleton hierarchy inside each group: joints in yellow are the *parent joints* in the skeleton hierarchy.

3 Preliminaries

The 3D pose of a subject, as she appears in each frame of a video presenting a human activity, is inferred according to the method introduced in [11]. Other methods for inferring the 3D pose of a subject are available, we refer in particular, to the method introduced by [12], which improves [11] in accuracy.

3D pose data for a single subject are given by the joints configuration. Joints are associated with the subject skeleton as shown in Fig. 2 and are expressed via transformation matrices \mathcal{T} in $SE(3)$:

$$\mathcal{T} = \begin{bmatrix} R & \mathbf{d} \\ \mathbf{0}^{1 \times 3} & 1 \end{bmatrix} \quad (2)$$

Here $R \in SO(3)$ is the rotation matrix, and $\mathbf{d} \in \mathbb{R}^3$ is the translation vector. $\mathcal{T} \in SE(3)$ has 6 DOF and it is used to describe the pose of the moving body with respect to the world inertial frame. $SO(3)$ and $SE(3)$ are Lie groups and their identity elements are the 3×3 and 4×4 identity matrices, respectively. We consider an ordered list $\mathcal{J} = \{j_1^1, j_2^1, \dots, j_{K-1}^m, j_K^m\}$ of $K = 18$ joints forming the skeleton hierarchy, as shown in Fig. 2, with $m = 1, \dots, 6$ being the groups each joint belongs to. The 6 groups G_1, \dots, G_6 we consider here correspond to head, torso, right and left arm, right and left leg.

Each joint j_i^m , $i = 1, \dots, 18$, belonging to a group G_m , $m = 1, \dots, 6$, has one parent joint $j_i^{m,*}$, which is the joint of the group closest to the root joint $root = j_4^2 \in \mathcal{J}$ and it belongs to the group G_2 , the torso. Parent joints for each group are illustrated in yellow on the woman body in the left of Fig. 2, they are in the order $(j_3^1, j_4^2, j_7^3, j_{10}^4, j_{13}^5, j_{16}^6)$.

A MoCap sequence of length N is formed by a sequence of frames of poses. Each frame of poses is defined by a set of transformations $\{\mathcal{T}_{i,m}^k \in SE(3) : k = 1, \dots, N, m = 1, \dots, 6\}$ involving all joints $j_i^m \in \mathcal{J}$, $i = 1, \dots, 18$, according to the skeleton hierarchy. Given a MoCap sequence of length N , for each frame k the pose of each joint is *root-sequence* normalized, to ensure pose invariance with respect to a common reference system of the whole skeleton. Let $\mathcal{T}_{i,m}^k$ be the pose of the joint j_i^m , according to the skeleton hierarchy, at frame k

in the sequence, and let $j_i^{m,*}$ be the parent node of j_i^m , then the *root-sequence* normalization is defined as follows:

$$\hat{\mathcal{T}}_{i,m}^k = \left((\mathcal{T}_{root,2}^1)^{-1} \mathcal{T}_{j_i^{m,*},m}^1 \right) \left((\mathcal{T}_{j_i^{m,*},m}^k)^{-1} \mathcal{T}_{i,m}^k \right). \quad (3)$$

Here $(\mathcal{T}_{root,2})$ is the transformation of the root node, which is the joint j_4^2 belonging to the group G_2 , the torso. Equation (3) says that the pose $\mathcal{T}_{i,m}^k$ of joint $j_i^m \in G_m$ at frame k is *root-sequence* normalized if obtained by a sequence of transformations seeing first a transformation with respect to its parent node $(\mathcal{T}_{j_i^{m,*},m}^k)^{-1}$, at frame k , and then with respect to the transformation of the parent node with respect to the root node, taken at the initial frame of the sequence. In Fig. 3 are shown joints position data for each skeleton group after *sequence-root* normalization for all sequences in the dataset. More details on the skeleton structure and its transformations can be found in [26, 11].

4 Motion Primitive Discovery

We are considering now the problem of discovering motion primitives within a motion sequence displaying an activity in a video. We begin by providing the definition of a joint trajectory on which the temporal analysis is performed.

Definition 4.1 (Joint Trajectory). The trajectory of a joint j is given by the path followed by the skeletal joint j in a given interval of time $I = [t_1, t_2]$. Formally:

$$\xi_j : I \subset \mathbb{R} \mapsto \mathbb{R}^3, \quad (4)$$

Based on the definition above, motion primitives correspond to segments of the joint trajectories of a group G . We identify motion primitives as trajectory segments where the variation of the velocity of the joints is maximal and where the endpoints of the segment correspond to stationary poses of the subject [51].

Preprocessing To overcome problems related to the finite sampling frequency of the poses in the data, we compute smooth versions of the joint trajectories by cubic spline interpolation. This interpolation provides a continuous-time trajectory for all the joints of the group with smooth velocity and continuous acceleration, satisfying natural constraints of human motion.

Motion Flux The motion flux captures the variation of the velocity of a group with respect to its rest pose. The total variation of the joint group velocity is evaluated along a direction \mathbf{g} that corresponds to stationary poses of the group. For groups 1, 3 and 4 this direction is defined by the segment connecting the ‘lowerneck’ and ‘upperneck’ joints while for groups 2, 5 and 6 by the segment connecting the ‘root’ with the ‘lowerback’ joints.

Definition 4.2 (Motion Flux). Let $G = \{j_1, \dots, j_K\}$ be a group consisting of K joints and \mathbf{v}_j the velocity of joint $j \in G$. The *motion flux* with respect to the time interval $I = [t_1, t_2]$ is defined as

$$\Phi(t_2, t_1) \doteq \sum_{j \in G} \int_{t_1}^{t_2} |\dot{\mathbf{v}}_j(t) \cdot \mathbf{g}| dt. \quad (5)$$

Discovery In order to discover a motion primitive, we identify a time interval between two time instances (endpoints) where the group velocity is minimal while the motion flux within the interval is maximal. This is done by performing an optimization based on the motion flux of a group G , as defined in eq. (5). More specifically, the time interval of a motion primitive is identified by maximizing the following energy-like function:

$$P(\rho; t_0) \doteq \Phi(\rho, t_0) - \frac{\beta_v}{2} \sum_{j \in G} (\|\mathbf{v}_j(\rho)\|^2 + \|\mathbf{v}_j(t_0)\|^2) + \beta_s \sum_{j \in G} (s_j(\rho) - s_j(t_0)), \quad (6)$$

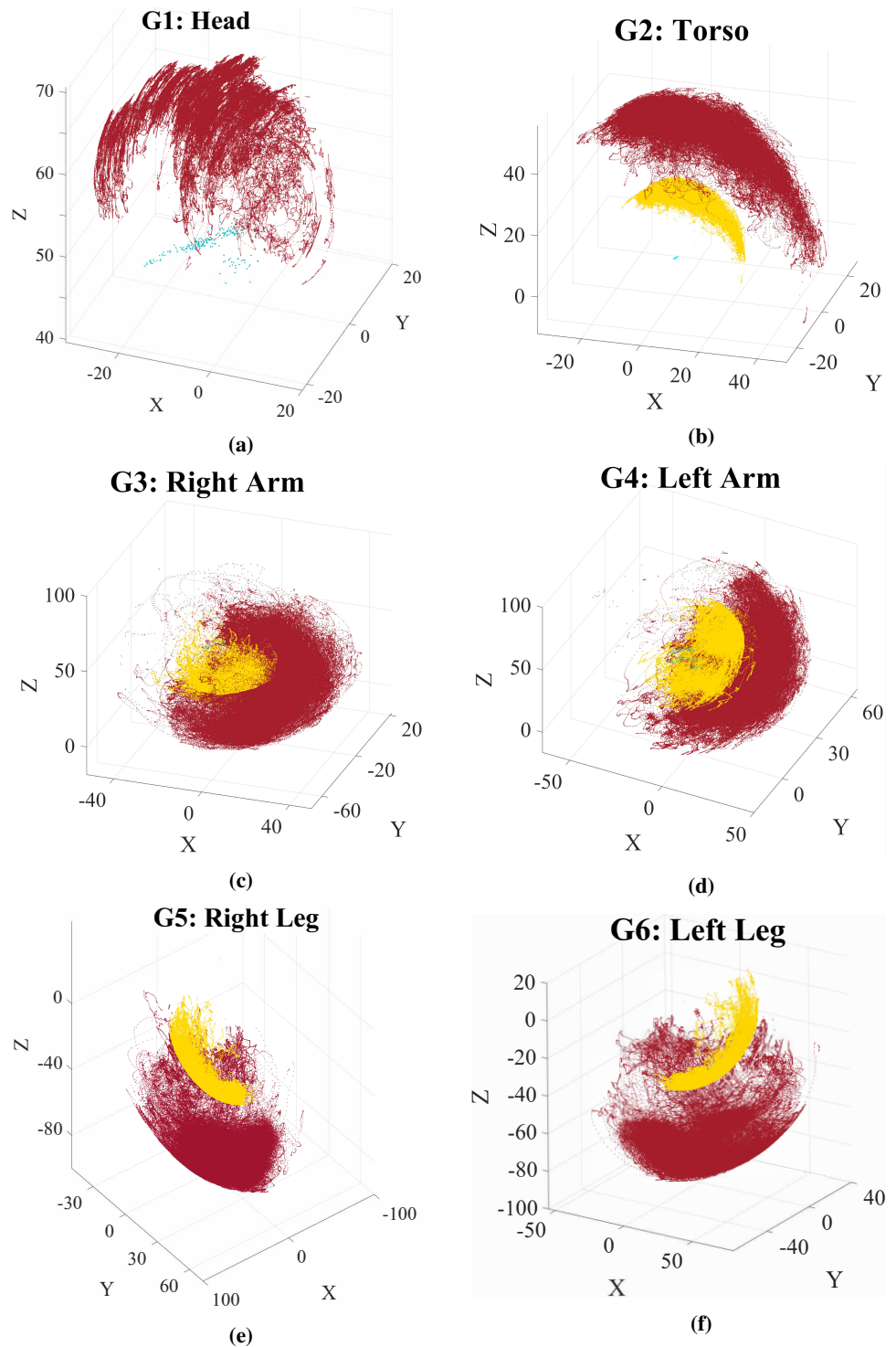
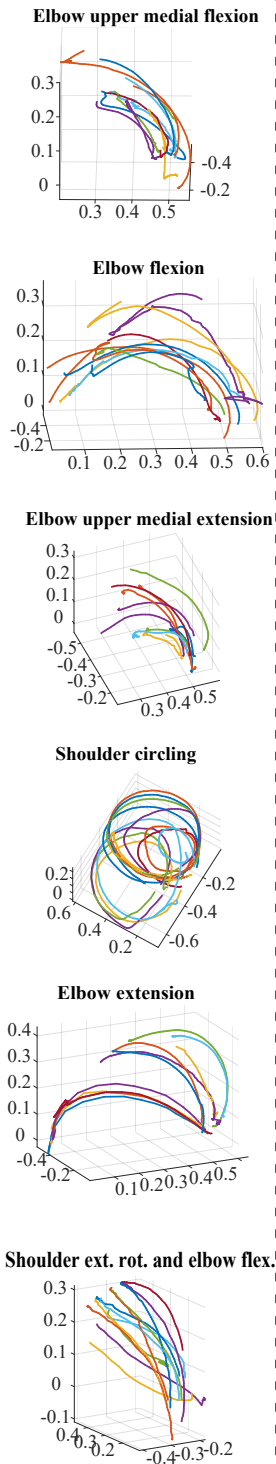


Fig 3. Sequences of joint positions, for each skeleton group, after the *root-sequence* normalization described in Section 3. Position data are in cm. The green points show the most internal group joint data (e.g. the hip for the leg); the yellow points show the intermediate group joint data (e.g. the knee for the leg); the red points show the most external group joint data (e.g. the ankle for the leg). The joints data are collected from the datasets described in Section 6.

Primitive Model Training



Primitive Discovery and Recognition

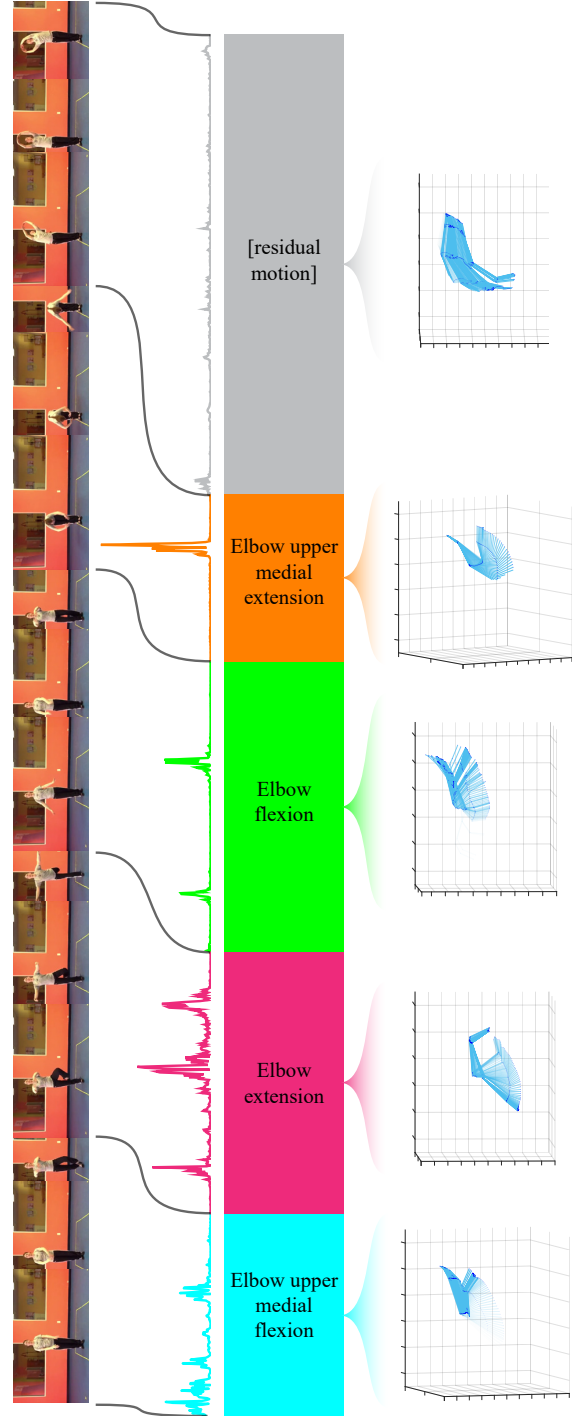


Fig 4. Overview of motion primitive discovery and recognition framework. The top section shows primitives of the group ‘Arm’ from six different categories. Primitives are discovered by maximizing the *motion flux* energy function, presented here above the colored bar, though deprived of velocity and length components. These sets of primitives are used to train the hierarchical models for each category. Primitives are then recognized according to the learned models. The recognized motion primitive categories are depicted with different colors. At the bottom, the group motion in the corresponding interval is shown.

where $s_j(t) = \int_0^t \|\dot{\xi}_j(\tau)\| d\tau$ is the arc length function of ξ_j . The last term of eq. (6) is a regularizer based on the length of the trajectory segment, introduced in order to avoid excessively long primitives. The hyper-parameter β_v acts as penalizer associated to the soft-constraint on the stationarity of the poses at the start and end of the primitive, while β_s controls the strength of the regularization on the primitive length. Both β_v and β_s depend on the scaling of the data and the sampling rate of the joint trajectories.

Given a starting time instant t_0 , a motion primitive is extracted by identifying the time instant ρ , which corresponds to a local maximum of (6). The optimality condition of (6) gives:

$$\sum_{j \in G} \left(|\dot{\mathbf{v}}_j(\rho) \cdot \mathbf{g}| - \beta_v \frac{\dot{\mathbf{v}}_j(\rho)}{\|\mathbf{v}_j(\rho)\|} - \beta_s \|\dot{\xi}_j(\rho)\| \right) = 0. \quad (7)$$

Given the one-dimensional nature of the problem, finding the zeros of (7) and verifying whether they correspond to local maxima of (6) is trivial.

Based on the previous we provide a formal definition of a motion primitive.

Definition 4.3 (Motion Primitive). A motion primitive of a group of joints G is defined by the trajectory segments of all joints $j \in G$ corresponding to a common temporal interval $I = [t_{start}, t_{end}] \subset \mathbb{R}$ such that $P(t_{start}; t_{end}) > P(\rho; t_{start}) \forall \rho \in (t_{start}, t_{end})$. Namely

$$\gamma_G^I = \{\xi_{j_1}(t), \dots, \xi_{j_K}(t)\} \text{ for } t \in [t_{start}, t_{end}]. \quad (8)$$

Primitive discovery in an activity A set of primitives is extracted from an entire sequence of an activity ς by sequentially finding the time instances which maximize (6).

Let t_0 and t_{seq} denote the starting and ending instances of the sequence, respectively. Let also

$$t_n = \arg \max_{\rho \in [t_{n-1}, t_{seq}]} P(\rho; t_{n-1}), \quad (9)$$

and $\mathcal{I}_\varsigma = \{[t_{n-1}, t_n] \mid n \in \mathbb{N} \text{ and } t_n \leq t_{seq}\}$ the set of time intervals defining successive motion primitives in the sequence. The set of motion primitives discovered in the entire sequence ς is given by

$$\Gamma_G^\varsigma = \{\gamma_G^I \mid I \in \mathcal{I}_\varsigma\}. \quad (10)$$

As noted in the introduction, and also shown in Figure 5, there is a significant motion variation across subjects, activities and sampling rates. For example, for the upper limbs it is known that the range of motion varies from person to person and is influenced by gait speed [52]. This is in turn influenced by the specific task, and determining ranges of motion is still a research topic [53] (for a review on range of motions for upper limbs, see [52]). This makes analysis and recognition of motion primitives taken from different datasets, activities and subjects problematic. To induce invariance with respect to these factors we apply anatomical normalization.

More specifically, the main source of variation of the primitives is due to the anatomical differences among the subjects. To remove the influence of these differences on the primitives we consider a scaling factor k_G based on the length ℓ_G of the limb defined by group G , namely $k_G = 1/\ell_G$. Hence, given a primitive γ_G^I we scale the trajectory of each joint by the constant k_G . By applying the anatomical normalization to the entire collection of motion primitives for group G discovered across all sequences of a dataset \mathcal{D} we obtain the set of motion primitive of the group, namely

$$\Gamma_G = \{\Gamma_G^\varsigma \mid \varsigma \in \mathcal{D}\}. \quad (11)$$

In Section 6 we provide a quantitative evaluation of the normalization effectiveness, together with a comparison with additional normalization candidates.

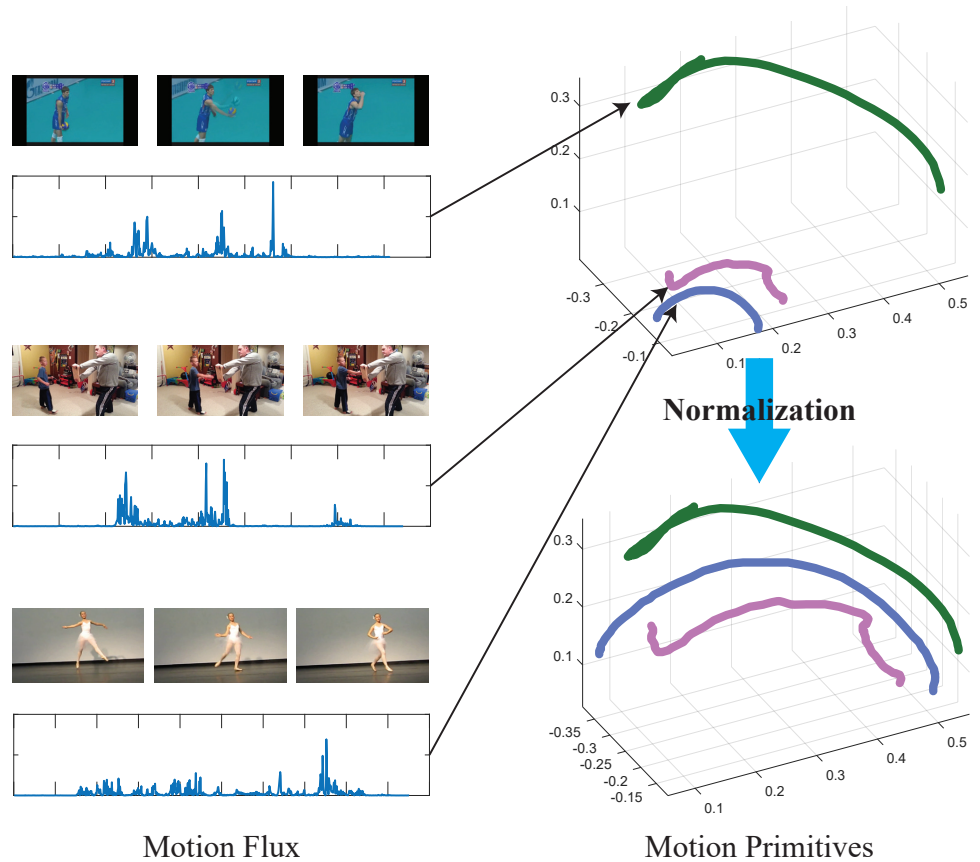


Fig 5. Left: Motion flux of three motion primitives of group G_3 labeled as ‘Elbow Flexion’, discovered from video sequences taken from the ActivityNet dataset. **Right:** Motion primitives before and after the normalization, for clarity only the curve of the out most joint is shown.(Best seen on screen, zoomed-in)

5 Motion Primitive Recognition

In the previous section we have shown that for each group of joints $G_m, m = 1, \dots, 6$, the motion flux obtains the interval $I = [t_{start}, t_{end}]$ matching the joint trajectory of a sequence in so determining a primitive as a path $\gamma_{G_m}^I : I \subset \mathbb{R} \mapsto \mathbb{R}^9$, given a video sequence of a human activity. Here \mathbb{R}^9 is due to the path being related to the 3 joints of each group G_m , as indicated in Fig. 2. We have also seen that the path is normalized by the link length of a limb, to limit variations due to bodies dissimilarities. For clarity from now on we shall denote each primitive with γ unless the context requires to add superscripts and subscripts, and in general subscripts and superscripts are local to this section, also we shall refer to the group a primitive or trajectory belongs to both with G_m and more in general with G .

We expect that the following facts will be true of the discovered motion primitives:

1. Each primitive of motion is independent of the gender, (adult) age, and body structure, under normalization.
2. Each primitive of motion can be characterized independently of the specific activity, hence the same primitive can occur in several activities (see Section 6 for a distribution of discovered primitives in a set of activities).

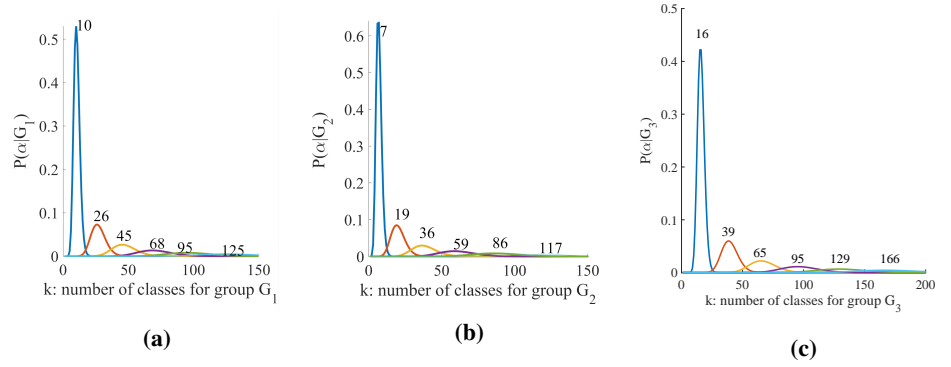


Fig 6. Number k of components for groups G_1 , G_2 and G_3 . Values of k are computed adjusting α so as to maximize the posterior $p(\alpha, G_m)$, given the data, namely the sampled primitives in the groups.

3. The motion flux ensures that each unknown segmented primitive belongs to a class such that: the number of classes is finite and the set of classes can be mapped onto a subset of motion primitives defined in biomechanics (see e.g table 1.1 of [54]).

To show experimentally the above results we shall introduce a hierarchical classification. The hierarchical classification first partitions the primitives of each group into classes. Once the classes are generated a class representative is chosen and inspected to assign a label to the class. We show that the classes correspond to a significant subset of the motion primitives defined in biomechanics, thus ensuring a proper partition. Each class is then further partitioned into subclasses to comply with the inner diversification of each class of primitives. This last classification is further used for recognition of unknown discovered primitives.

Primitive recognition is used to both test experimentally the three above results of the introduced motion flux method and for applications where discovering and recognition of primitives of human motion is relevant (see for example [55]).

5.1 Solving primitive classes

We describe in the following the method leading to the generation of all the primitive classes illustrated in Fig 12.

We consider three MoCap datasets [15, 13, 14] guaranteeing the ground truth for the human pose and segment the activities according to the motion flux method, described in the previous section. Let Γ_G be the set of primitives collected for group G according to equation (11). Let $\gamma_\nu \in \Gamma_G$, $\nu = 1, \dots, S$, with S the number of primitives in Γ_G , $\gamma_\nu = (\xi_{j_1}^\nu, \xi_{j_2}^\nu, \xi_{j_3}^\nu)$ is formed by the trajectories of the joints in G . Out of these trajectories we choose the one of the most external joint (see Figure 2) that we indicate with ξ_E^ν . We order these trajectories, each designating a primitive in group G , with an enumeration $\langle \Gamma_G \setminus \xi_E \rangle_{\nu=1}^S$, S the number of discovered primitives for group G . Note that we can arbitrarily enumerate the primitives of a group, restricted to a single joint, though they are unlabeled and unknown, and this is what the first model should solve.

At this step, model generation amounts to find the classes of primitives for each group G , taking the trajectories ξ_E^ν in the enumeration $\langle \Gamma_G \setminus \xi_E \rangle_{\nu=1}^S$ as observations.

Feature Vectors Given a trajectory ξ_E^ν , with ν the index in the enumeration $\langle \Gamma_G \setminus \xi_E \rangle_{\nu=1}^S$, a feature vector is obtained by first computing curvature $\kappa(s(t))$ and torsion $\tau(s(t))$ on the trajectory ξ_E^ν , where $s(t)$ indicates the arc length as already defined in Section 4 for trajectories. Then we take three contiguous points $(x_{i-1}, y_{i-1}, z_{i-1}), \dots, (x_{i+1}, y_{i+1}, z_{i+1})$ on the trajectory $\hat{\xi}_E^\nu$ decimated by a factor of 5 [56], keeping the curvature and torsion of the

sampled points, after decimation. We choose curvature and torsion as they suffice to specify a 3D curve up to a rigid transformation. The formed feature vector is indicated by \mathcal{F}_i , where the index i is the index of the middle point (x_i, y_i, z_i) , it is of size 17×1 and it is defined as follows:

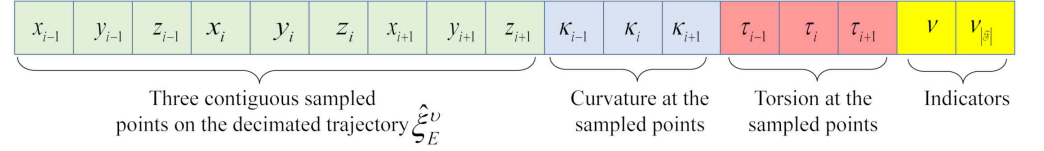


Fig 7. Transposed feature vector of 3 contiguous sampled points on the decimated trajectory.

The last two elements $\nu, \nu_{|\mathcal{F}_i|} \in \mathbb{R}$ of \mathcal{F}_i are indicators. Namely, the indicator ν is the index, in the enumeration $\langle \Gamma_G \setminus \xi_E \rangle_{\nu=1}^S$, identifying the trajectory the 3 points belong to, the three points are the first 6 element of the feature vector. On the other hand, the indicator $\nu_{|\mathcal{F}_i|}$ specifies the number of features vectors the decimated trajectory $\hat{\xi}_E^\nu$ is decomposed into, here $|\cdot|$ indicates the cardinality; These two indicators, allow to recover the path a feature vector belongs to, and are normalized and denormalized as follows. Let \mathbb{F}_G^ξ be the set of all feature vectors for the trajectories in $\langle \Gamma_G \setminus \xi_E \rangle_{\nu=1}^S$, and let their number be W . Accordingly, let $\nu_{|\mathcal{F}|} = (\nu_{|\mathcal{F}_1|}, \dots, \nu_{|\mathcal{F}_W|})$, then the normalization and denormalization for the element $\nu_{|\mathcal{F}_i|}$ (and similarly for ν) is defined as follows, with g indicating the denormalization:

$$\hat{\nu}_{|\mathcal{F}_i|} = \frac{\nu_{|\mathcal{F}_i|} - \min(\nu_{|\mathcal{F}|})}{\max(\nu_{|\mathcal{F}|}) - \min(\nu_{|\mathcal{F}|})} \quad (12)$$

$$g(\hat{\nu}_{|\mathcal{F}_i|}) = \hat{\nu}_{|\mathcal{F}_i|}(\max(\nu_{|\mathcal{F}|}) - \min(\nu_{|\mathcal{F}|})) + \min(\nu_{|\mathcal{F}|})$$

Generation of the primitives classes Given the feature vectors for each trajectory in the enumeration $\langle \Gamma_G \setminus \xi_E \rangle_{\nu=1}^S$, the goal is to cluster them and return a cluster for each class of primitives. Since we do not even know the number of classes the primitives should be partitioned into, a good generative model to approximate the distribution of the observations is the Dirichlet process mixture (DPM) [57, 58]. The Dirichlet process assigns probability measures to the set of measurable partitions of the data space. This induces in the limit a finite mixture since, by the discreteness of the distributions sampled from the process, parameters have positive probability to take the same value, in so realizing components of the mixture. Here we assume that feature vectors in the data space are realizations of normal distributions with a conjugate prior. Namely the variables have precision priors following the Wishart distribution and location parameters prior following the normal distribution. The Dirichlet mixture model is based on the definition of a Dirichlet process $\Pi(\cdot, \cdot)$ with $\Pi \sim DP(H, \alpha)$ (D being the Dirichlet distribution), where H is the base distribution and α the precision parameter of the process (see [59]). In the Dirichlet process mixture the value of the precision α of the underlying Dirichlet process influences the number of classes generated by the model.

For determining the number of classes for each group G we estimate the posterior $P(\alpha|G)$, of the precision parameter α according to a mixture of two gamma distributions, as described in [60], choosing the best value. This is a rather complex simulation process since it requires different initializations of the parameters of the gamma distribution for α within the estimation of the parameters of the DPM, for each group G . Here the parameters of the DPM are estimated according to [61]. Distributions of α for the groups G_1, G_2 and G_3 , according to different simulation processes, are given in Fig. 6 where the number of components k for the maximum values of each distribution, are indicated. Finally the DPM returns the parameters of the components (for each group G) given the feature vector \mathcal{F}_i , as:

$$\Theta_G = \langle k, \{\Theta_w \mid \Theta_w = (\pi_w, \mu_w, \Sigma_w), w = 1, \dots, k\} \rangle, k \geq 1. \quad (13)$$

$$p(\mathcal{F}_i | \Theta_G) = \sum_{w=1}^k \pi_w \mathcal{N}(\mathcal{F}_i | \mu_w, \Sigma_w).$$

Note that the number of components k is unknown and estimated by the DPM, hence it is one of the parameters for each group. The parameters μ_w and Σ_w are the mean vector and covariance matrix of the w -th Gaussian component of the mixture, indicated by \mathcal{N} , and π_w is the w -th weight of the mixture, with $\sum_w \pi_w = 1$. Hence, $p(\mathcal{F}_i|\Theta_G)$ is the probability of the feature vector \mathcal{F}_i , given the parameters Θ_G .

We expect that each $\Theta_w \in \Theta_G$ indicates the parameters of a component C_w^G , collecting primitives of the same type, in group G . In other words, we expect that two feature vectors, say $\mathcal{F}_p, \mathcal{F}_q$, of group G , belong to the same component if their likelihood are both maximized by the same parameters $\Theta_w \in \Theta_G$.

Assigning primitives to classes The classification returns, for each group G_m , the number k of components indicated in Fig. 12, say $k = 10$ for G_1, G_5, G_6 , $k = 7$ for G_2 and $k = 16$ for G_3, G_4 , also thanks to the specification of the α parameter, as highlighted above (see Fig. 6). Components are formed by features vectors. To retrieve the trajectories and generate a corresponding class of primitives, ready to be labeled, we use the normalized indicators placed in position 16th and 17th of the feature vector (Fig 7) and the denormalization function g . Let $C_w^{G_m}$ be a component of the mixture of the group G_m , identified by parameters $\Theta_w \in \Theta_{G_m}$. Algorithm 1 shows how to compute the class of primitives:

```

Input: Component  $C_w^{G_m}$  of DPM
Output: Class  $\mathcal{L}_w^{G_m}$  of primitives
Initialize  $U_{\xi_E}^\nu = \emptyset, \nu = 1, \dots, S$ ,  $S$  number of primitives in  $\Gamma_{G_m}$ 
foreach Feature vector  $\mathcal{F}_i$  in  $C_w^{G_m}$  do
    compute  $g(\nu)$  and associate it with the trajectory  $\xi_E^\nu$ ;
     $U_{\xi_E}^\nu = \{\mathcal{F}_i\} \cup U_{\xi_E}^\nu$ ;
    compute  $g(\nu_{|\mathcal{F}|})$ , number of feature vectors the trajectory  $\xi_E^\nu$  is decomposed into;
end
if  $|U_{\xi_E}^\nu| \geq 0.8g(\nu_{|\mathcal{F}|})$  then
    find the primitive  $\gamma_\nu \in \Gamma_{G_m}$  designated by  $\xi_E^\nu$ 
    assign the pair  $(\gamma_\nu, \Theta_w)$  to  $\mathcal{L}_w^{G_m}$ 
end
return Class  $\mathcal{L}_w^{G_m}$ .

```

Algorithm 1: Obtaining classes of primitives from DPM components. Here $|\cdot|$ indicates cardinality.

At this point we have generated the classes $\mathcal{L}_w^{G_m}, w = 1, \dots, k, k \in \{7, 10, 16\}$ of primitive for each group G_m . To label the classes we proceed as follows. Let $p(\gamma_\nu|\Theta_w) = 1/g(\nu_{|\mathcal{F}|}) \sum_i p(\mathcal{F}_i|\Theta_w)\delta(\mathcal{F}_i)$, where $\delta(\mathcal{F}_i) = 1$ if $\mathcal{F}_i \in U_{\xi_E}^\nu$ and 0 otherwise. For each class $\mathcal{L}_w^{G_m}$ the class representative is the primitive maximizing $p(\gamma_\nu|\Theta_w)$. The representative primitive is observed and labeled by inspection, according to the nomenclature given in biomechanics, see [54]. The same label is assigned to the class $\mathcal{L}_w^{G_m}$, without need to inspect all other primitives assigned to the class.

Average Hausdorff distances between each primitive in a class and its class representative, for each class in group G_2 , are given in Table 1, where classes for G_2 are enumerated according to the labels illustrated in Fig. 12. Note that in Table 1 R_w is the class representative, so $R_w \in \mathcal{L}_w^{G_m}, w = 1, \dots, 7; \forall \xi_E \in R_w$ abbreviates $\forall \xi_E \in \mathcal{L}_w^{G_2}, \xi_E \neq R_w$. Finally, \mathcal{L}_w abbreviates $\mathcal{L}_w^{G_2}$. Note that distances with elements of other classes are obviously not considered, hence the dashes in other classes columns.

Table 1. Average Hausdorff distance to each class representative in G_2

	R_1	R_2	R_3	R_4	R_5	R_6	R_7
$\forall \xi_{E \setminus R_1} \in \mathcal{L}_1$	0.121	-	-	-	-	-	-
$\forall \xi_{E \setminus R_2} \in \mathcal{L}_2$	-	0.173	-	-	-	-	-
$\forall \xi_{E \setminus R_3} \in \mathcal{L}_3$	-	-	0.144	-	-	-	-
$\forall \xi_{E \setminus R_4} \in \mathcal{L}_4$	-	-	-	0.112	-	-	-
$\forall \xi_{E \setminus R_5} \in \mathcal{L}_5$	-	-	-	-	0.081	-	-
$\forall \xi_{E \setminus R_6} \in \mathcal{L}_6$	-	-	-	-	-	0.142	-
$\forall \xi_{E \setminus R_7} \in \mathcal{L}_7$	-	-	-	-	-	-	0.114

5.2 Models for recognition

The recognition problem is stated as follows. Given an unlabeled primitive γ_u , for group G_m obtained by segmenting an activity (from any dataset) with the motion flux method, γ_u is labeled by the label of class $\mathcal{L}_w^{G_m}$, if:

$$p(\gamma_u | \Theta_w) > p(\gamma_u | \Theta_i), \quad \forall i, i \neq w \quad (14)$$

We found experimentally that relying on the same parameters used for finding the classes of primitives, described in the previous sub-section, does not lead to optimal results. In fact, recomputing a DPM model for each class and introducing a loss function on the set of hypotheses, computed by thresholding the best classes, leads to an improvement up to the 20% in the recognition of an unknown primitive.

To this end we compute a DPM for each class $\mathcal{L}_w^{G_m}$ using as observations the primitives collected in the class, by Algorithm 1. Therefore the generated DPM model \mathcal{M}_w for each class $\mathcal{L}_w^{G_m}$ is made by a number of components with parameters $\Theta_w = \{\Theta_{w_1}, \dots, \Theta_{w_\rho}\}$, with ρ varying according to the components generated for class $\mathcal{L}_w^{G_m}$. The number of components mirrors the idiosyncratic behavior of each class of primitives, therefore ρ varies for each class $\mathcal{L}_w^{G_m}$. To generate these DPM models we use all the three trajectories of the primitives $\gamma \in \mathcal{L}_w^{G_m}$, and for each of them we use the same decimation and feature vector as shown in Fig. 7.

Given the refined classification, the recognition problem, at this point, is stated as follows. Let $\gamma_u = (\xi_{u_1}, \xi_{u_2}, \xi_{u_3})$ be an unknown primitive, of a specific group G , and let $\{\mathcal{F}_{u_1}, \dots, \mathcal{F}_{u_q}\}$ be the set of features the three trajectories are decomposed into. Then $\gamma_u \in \mathcal{L}_w^{G_m}$, hence is labeled by the label of this class, if:

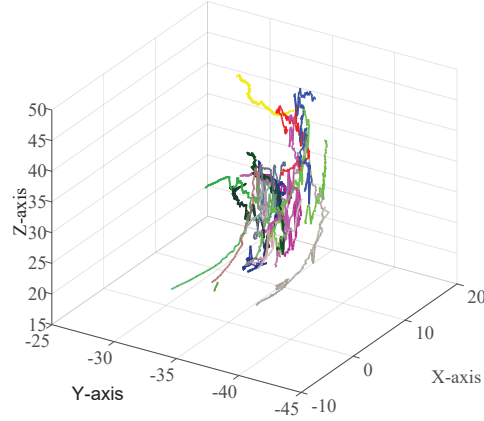
$$p(\mathcal{F}_{u_1}, \dots, \mathcal{F}_{u_q} | \Theta_w) = \sum_{j=1}^{\rho} \pi_j \prod_{n=1}^q p(\mathcal{F}_{u_n} | \Theta_{w_j}) > p(\mathcal{F}_{u_1}, \dots, \mathcal{F}_{u_q} | \Theta_h) = \sum_{j=1}^{\rho'} \pi'_j \prod_{n=1}^q p(\mathcal{F}_{u_n} | \Theta_{h_j}) \quad (15)$$

for any parameter set Θ_h associated with a class $\mathcal{L}_h^{G_m}$ of the group G_m . Here π_j and π'_j are the mixture weights, with $\sum_j \pi_j = 1$ and ρ, ρ' indicate the number of components of the chosen models. For example, the model of class $\mathcal{L}_w^{G_2}$, with $w = 1$, will have a set of parameters $\Theta_w = \{\Theta_{w_1}, \dots, \Theta_{w_\rho}\}$, while the model of class $\mathcal{L}_{w'}^{G_2}$, with $w' = 3$, will have a set of parameters $\Theta_{w'} = \{\Theta_{w'_1}, \dots, \Theta_{w'_{\rho'}}\}$, with $w_\rho \neq w'_{\rho'}$.

This formulation is much more flexible than (14), also because it computes the class label by considering all the components and therefore it does not care whether the features are scattered amid components, and does not need to reconstruct the whole trajectories as was done for generating the classes of primitives. Furthermore, under this refined classification we can improve (15) considering a geometric measure to reinforce the statistics measure in the choice of the class label for γ_u .

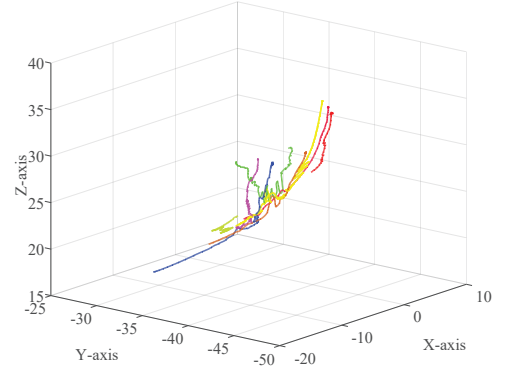
More precisely, let us form a set of hypotheses for an unknown primitive with feature set

Component 1 of DPM model for Elbow Flexion



(a)

Component 4 of DPM model for Shoulder Abduction



(b)

Fig 8. Manifold generated by a component of the DPM model for Elbow Flexion on the left and from a component of Shoulder Abduction on the right.

$\{\mathcal{F}_{u_1}, \dots, \mathcal{F}_{u_q}\}$ as follows (we are still assuming a specific group G_m):

$$\mathbb{H} = \{ \langle C_{w_j}, \Theta_{w_j} \rangle \mid \prod_{n=1}^q p(\mathcal{F}_{u_n} | \Theta_{w_j}) > \eta, \langle C_{w_j}, \Theta_{w_j} \rangle \in \mathcal{M}_w, w = 1, \dots, k \} \quad (16)$$

Namely C_{w_j} is a component of the DPM \mathcal{M}_w , with $w = 1, \dots, k$, k the number of classes in group G_m , and $j = 1, \dots, \rho$, such that the associated parameter Θ_{w_j} makes the joint probability of the features, the primitive is decomposed into, greater than a threshold η . This means that we are collecting in \mathbb{H} those components coming from all the models of group G_m , whose joint probability of the feature set of the unknown primitives γ_u forms an hypotheses set, or a set from which we can select the correct label to assign to γ_u .

The advantage of the hypotheses set is that we delay the decision of choosing the labeled class for the unknown primitive to further evidence, which we collect by using geometric measures. The role of these geometric measures is essentially to evaluate the similarity between the curve segments coming out from the features of γ_u and those coming from the observations which are indexed in the components in \mathbb{H} . In the following we succinctly describe the new geometric features, which are computed as follows, both for the features of the unknown primitive γ_u and for the features coming from the observations indexed in C_{w_j} . Let us consider any pair $\langle C_{w_j}, \Theta_{w_j} \rangle \in \mathbb{H}$, by definition (16), C_{w_j} indexes features $\{\mathcal{F}_{\nu_1}, \dots, \mathcal{F}_{\nu_s}\}$, s varying according to the specific component C_{w_j} . For each of these features we consider the points of the trajectory ξ^ν , recovered from the decimated trajectory $\hat{\xi}^\nu$, between $(x_{i-1}, y_{i-1}, z_{i-1})$ and $(x_{i+1}, y_{i+1}, z_{i+1})$. Let us consider these curve segments, which we combine whenever they occur in sequence in C_{w_j} and call any of these curve segments \mathbf{y} . In particular, the collection of these segments in C_{w_j} is called the manifold of C_{w_j} , denoted $man(C_{w_j})$, and the collection of segments generated from the features of γ_u is denoted $man(\gamma_u)$, examples are given in Fig. 8.

We compute for each \mathbf{y} both in $man(C_{w_j})$ and in $man(\gamma_u)$ the tangent \mathbf{t} , normal \mathbf{n} and binormal \mathbf{b} vectors. Based on these vectors, we compute the ruled surface $\mathcal{R} = \frac{\mathbf{n} \times \mathbf{n}'}{\|\mathbf{n} \times \mathbf{n}'\|}$, where \mathbf{n}' is the derivative of \mathbf{n} . The ruled surface forms a ribbon of tangent planes to the curve segment \mathbf{y} . In particular, let us distinguish the curve segments in $man(\gamma_u)$ denoting them \mathbf{y}_u . We compute the distances between any curve segment $\mathbf{y} \in man(C_{w_j})$ and $\mathbf{y}_u \in man(\gamma_u)$ as the distance between the projection \mathbf{y}_π of \mathbf{y} on the ruled surface tangent to \mathbf{y} , and the *closest point* \mathbf{q} of \mathbf{y}_u to \mathbf{y}_π . We denote this distance $\delta(\mathbf{y}_u, \mathbf{y})$. We consider also the distance between

the Frenet frames at closest points \mathbf{q} of \mathbf{y}_u and point \mathbf{q}' of \mathbf{y}_π denoted F_R and computed as follows: $F_R(\mathbf{q}, \mathbf{q}') = \text{trace}((\mathcal{I} - R_{\mathbf{q}, \mathbf{q}'})(\mathcal{I} - R_{\mathbf{q}, \mathbf{q}'})^\top)$, with \mathcal{I} the identity matrix and $R_{\mathbf{q}, \mathbf{q}'}$ the rotation, in the direction from \mathbf{q} to \mathbf{q}' . Then the cost of a component C_{w_j} in \mathbb{H} , given an unknown primitive γ_u , with feature set $\{\mathcal{F}_{u_1}, \dots, \mathcal{F}_{u_q}\}$, is defined as:

$$\text{Cost}(C_{w_j} \in \mathbb{H} | \gamma_u) = \max\{\delta(\mathbf{y}_u, \mathbf{y}) + F_R(\mathbf{q}, \mathbf{q}') | \mathbf{y}_u \in \text{man}(\gamma_u) \text{ and } \mathbf{y} \in \text{man}(C_{w_j})\} \quad (17)$$

Note that both $\delta(\mathbf{y}_u, \mathbf{y})$ and $F_R(\mathbf{q}, \mathbf{q}')$ were both computed looking at the minimum distance between a considered curve segment and the projection on the ruled surface of the other curve segment. Hence the component minimizing the above cost and maximizing the probability in (15) will indicate the class label, since its related parameter indicates exactly a component of one of the classes $\mathcal{L}_w^{G_m}$. Note that if in (15) η is taken to be equal to $\max(\prod_{n=1}^q p(\mathcal{F}_{u_n} | \Theta_{w_j}))$ then \mathbb{H} would have only a single element $\langle C_{w_j}, \Theta_{w_j} \rangle$. Hence to find the correct label for γ_u we push η as high as possible using the above cost. More precisely, the component of the class $\mathcal{L}_w^{G_m}$ which should label the unknown primitive γ_u is computed as follows:

$$C^* = \arg \min_{C_{w_j}} \sup_{\eta} \{\text{Cost}(C_{w_j} \in \mathbb{H} | \gamma_u) | \prod_{n=1}^q p(\mathcal{F}_{u_n} | \Theta_{w_j}) > \eta\} \quad (18)$$

To conclude this section we can note that the computation of the hierarchical model that first generates the primitive classes and then uses these generated sets to estimate model parameters to be used in the recognition of an unknown primitive, has an exponential cost, in the dimension of the features and in the size of the observations. However using the computed models to recognize an unknown primitive is $\mathcal{O}(n^2 \log n)$ where n is the size of γ_u , since all the curve segments in the models can be precomputed together with the models. Results on both the primitive generation and on recognition are given in the next section.

6 Experiments

In this section we evaluate the proposed framework for discover and classification of human motion primitives. For all the evaluations we consider three reference MoCap public datasets [15, 13, 14].

First we evaluate the accuracy of the motion primitives discovered using the motion flux, further we evaluate the accuracy of the classification and recognition. Additionally, we examine the distribution of recognized primitives with respect to the type of performed activity on the ActivityNet dataset [1]. Finally, we address the dataset of human motion primitives we have created, which consists of the primitives discovered on the three reference MoCap datasets using the motion flux, and the DPM models established for each primitive category.

6.1 Reference Datasets

The datasets we consider for the evaluation of the motion flux are the Human3.6M dataset (H3.6M) [13], the CMU Graphics Lab MoCap database (CMU) [14] and the KIT Whole-Body Human Motion Database (KIT-WB) [15]. The sampling rates used in these datasets are 50Hz for H3.6M, 60/120Hz for CMU and 100Hz for KIT-WB. In order to have the same sampling rate for all sequences we have transformed all of them to 50Hz. The pose of the joints specified in Fig. 2 are extracted for each frame of the sequences as described in the preliminaries, considering the ground-truth 3D poses. For KIT-WB the trajectories of the joints are computed from the marker positions taken from the C3D files. We considered 40 activities from the three reference datasets. Fig. 9 shows the total number of motion primitives discovered for the five most general activities according to the ActivityNet taxonomy based on the motion flux for each group G_m . Table 2 shows the total number of motion primitives discovered from the three datasets.

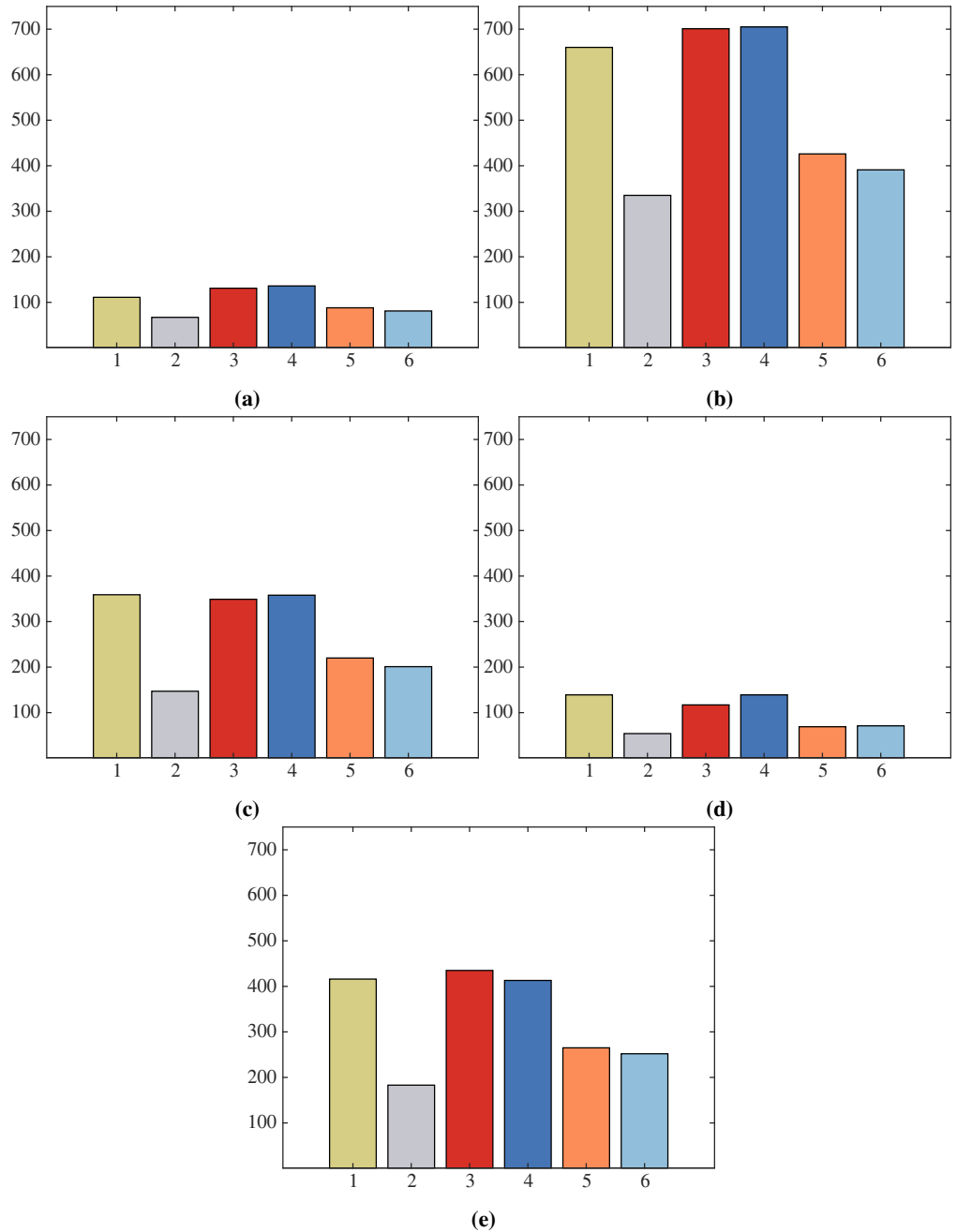


Fig 9. Total number of discovered primitives for each group for the five most general categories of the ActivityNet dataset. Clock-wise from top-left: *Eating and drinking Activities*; *Sports, Exercise, and Recreation*; *Socializing, Relaxing, and Leisure*; *Personal Care*; *Household Activities*. Each color corresponds to a different group following the convention of Fig. 12. Note: Axes scale is shared among the plots.

6.2 Motion Primitive Discovery

To evaluate the accuracy of primitive discovery based on the motion flux, we created a baseline relying on a synthetic dataset of motion primitives. This was necessary to mitigate the difficulty in measuring accuracy, due to the lack of a ground truth.

Table 2. Total number of unlabeled primitives discovered for each group using the motion flux on the reference datasets

	G1	G2	G3	G4	G5	G6
Total	1665	759	1773	1703	1152	1015

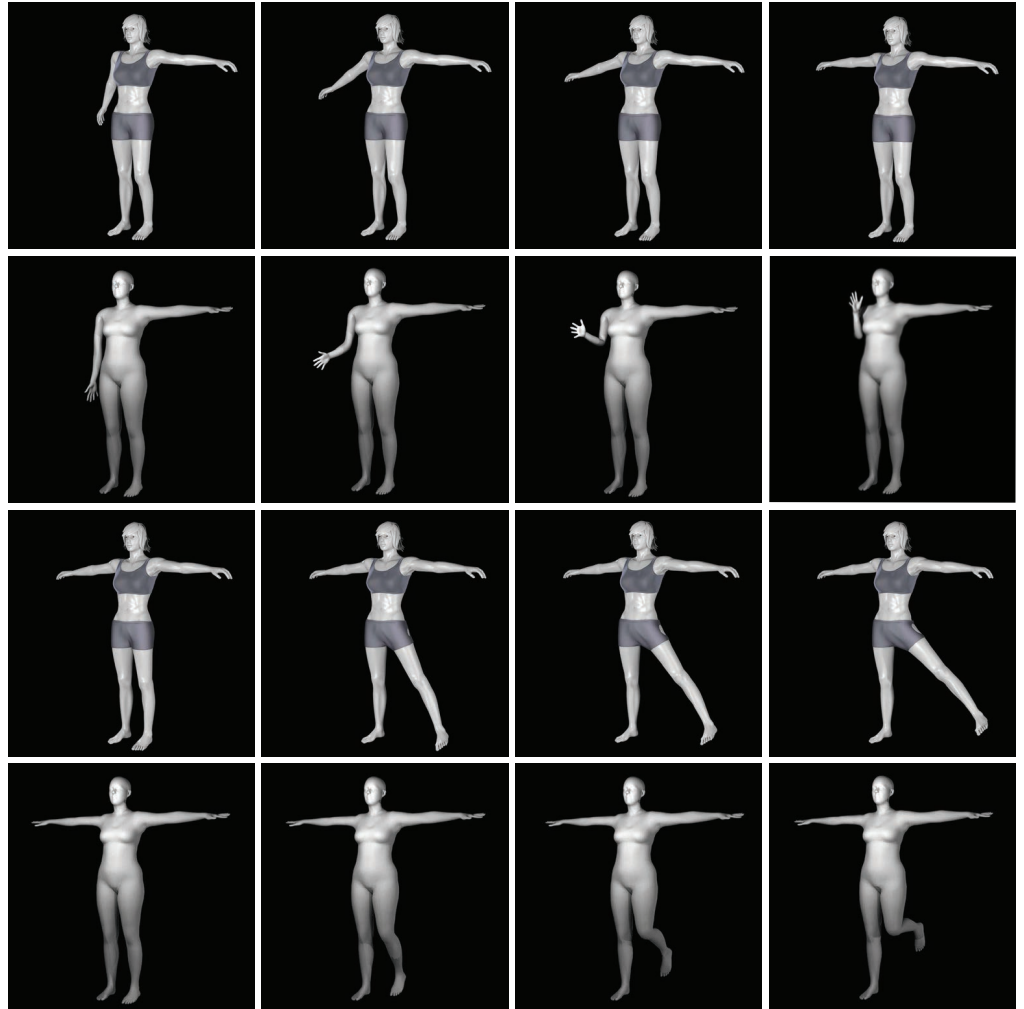


Fig 10. Example of synthetic motion primitive, specifically right arm Shoulder Abduction (first row) and Elbow Flexion (second row), left leg Hip abduction (third row) and Knee Flexion (fourth row). For each synthetic motion primitive the four imaged poses match four representative poses extracted from the animation of the aforementioned primitive.

The synthetic dataset of motion primitives we created is formed by animations of 3D human models for each of the 69 primitive classes discovered in Sec. 5. The human models were downloaded from the dataset provided by [62] or acquired from [63, 64]. To obtain further characters the shapes of the human models were randomly modified taking care of human height and limb length limits.

Animations of the characters were produced moving the skeleton joints belonging to the 3D human models from a start pose to an end pose representing the primitives. Specifically, for each primitive of each skeleton group the animation was generated in Maya or Blender (depending on the 3D human model format) moving the group joints according to angles, gait speed and limbs proportions as described in [52, 53, 54, 55].

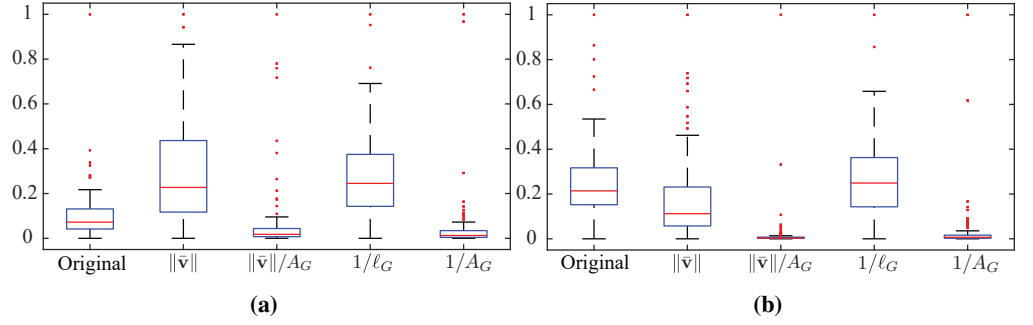


Fig 11. Arc length distribution of original and scaled primitives of a specific category for group G_1 (left) and G_4 (right). The first box in each box plot, corresponds to the original arc length distribution, the next four are the arc length distributions obtained scaling the primitives original data using the detailed scaling factors. Each box indicates the inner 50th percentile of the trajectory data, top and bottom of the box are the 25th and 75th percentiles, the whiskers extend to the most extreme data points not considered outliers, crosses are the outliers.

The dataset reference skeleton, see Fig. 2 is matched with the 3D human mesh models by fitting the joint poses of the synthetic data to the reference skeleton, basing on MoSh [65, 66]. Examples of synthetic motion primitives, namely the primitives Shoulder abduction and Elbow flexion for the right arm, and Hip abduction and Knee flexion for the left leg, are illustrated in Fig. 10, where for each primitive four representative poses extracted from the animations are shown.

The baseline for evaluating accuracy was created generating 4500 random length sequences of synthetic motion primitives placing them one after another in a random order. Between two consecutive primitives a transition phase from the end pose of the preceding one to the beginning pose of the subsequent one was added.

With this procedure we know precisely the endpoints of each primitive.

Then we applied the ‘motion flux’ method described in Sec. 4 to the 3D joints trajectories extracted from the automatically generated sequences and collected the end points of the discovered primitives.

We use the Mean Absolute Error (MAE) and Root Mean Square Error (RMSE) metrics to assess the accuracy of the collected endpoints with respect to the known end points in the generated sequences. Let S be the total number of generated sequences. Let $\{\hat{e}_{i,s}\}_{i=1}^{N_G^{(s)}}$ be the i -th automatically discovered endpoint based on the motion flux for the generated sequence $s = \{1, \dots, S\}$, with $N_G^{(s)}$ the number of primitives for Group G and sequence s . Denoting $\{\bar{e}_{i,s}\}_{i=1}^{N_G^{(s)}}$ the i -th endpoint in the generated sequence s , the MAE and RMSE metrics are defined as follows:

$$MAE = \frac{1}{S} \sum_{s=1}^S \frac{\sum_{i=1}^{N_G^{(s)}} |\bar{e}_{i,s} - \hat{e}_{i,s}|}{N_G^{(s)}}, \quad RMSE = \sqrt{\frac{1}{S} \sum_{s=1}^S \frac{\sum_{i=1}^{N_G^{(s)}} (\bar{e}_{i,s} - \hat{e}_{i,s})^2}{N_G^{(s)}}}.$$

Results shown in Table 3 prove that the proposed method discovers motion primitives quite accurately, since the endpoints are close to those of the automatically generated sequences.

Table 3. Accuracy of discovered primitive endpoints (in number of frames)

	G1	G2	G3	G4	G5	G6	Overall
MAE	2.8	3.2	2.9	3.4	3.6	4.1	3.3
RMSE	3.7	4.2	4.1	4.6	4.8	5.2	4.4

Furthermore, to evaluate the effects of the normalization in Fig. 11 we show the arc length

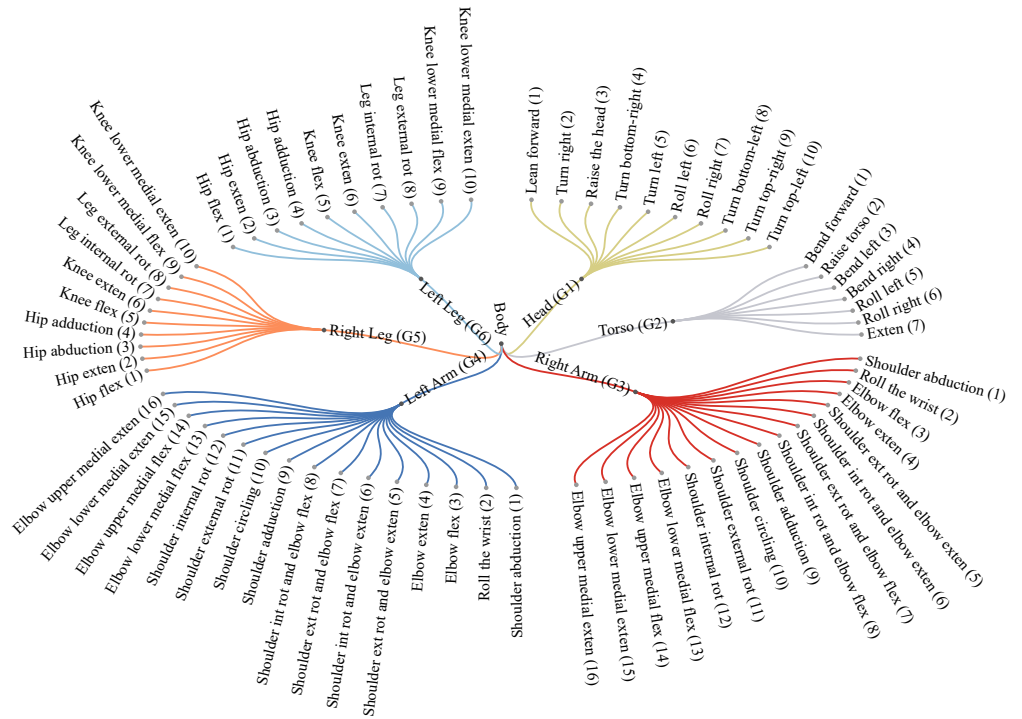


Fig 12. Diagram showing the motion primitives of each group. Abbreviation *ext* stands for external, *int* for internal, *rot* for rotation, *exten* for extension, and *flex* for flexion.

distribution of motion primitives with and without normalization, as well as considering different normalization constants.

For comparison we consider alternative normalization constants based on anatomical properties and execution style. Specifically, we consider normalization based on the average velocity along $\gamma \in \Gamma_G$, denoted as $\|\bar{v}\|$, and based on the area A_G covered by group G during its motion. The first is related to the execution speed of the motion and the sampling rate of the data, while the latter is considering anatomical differences among the subjects.

In Fig. 11 the first box in each plot corresponds to the original distribution and the following boxes correspond to the distributions resulting by scaling the original one with $\|\bar{v}\|$, $\|\bar{v}\|/A_G$, $1/\ell_G$, and $1/A_G$, respectively. We note that normalizing the primitives based on the inverse of the limb length, i.e. ℓ_G , consistently results to an arc length distribution closer to the normal, minimizing the number of outliers indicated by red crosses in the figure. This result is consistent across different activities and groups justifying the choice of $k_G = 1/\ell_G$ for anatomical normalization.

6.3 Motion Primitive Classification and Recognition

As discussed in Section 5, the set of primitive categories for each group is generated by a DPM model given the collection of discovered primitives as observations. In this way a total of 69 types of primitives were identified, each described by the distribution parameters. By inspecting a representative primitives for each category, we observed that they correspond to a subset of motion primitives defined in biomechanics. Therefore we generated new DPM models to obtain parameters and corresponding labels for each category. The labeled collection of motion primitives is depicted in Fig. 12.

To evaluate the coherence of the generated classes we performed 10 cycles of random sampling, with a rate of 10% at each cycle, of the primitives in each class and verified the class

consistency. Only $\sim 2\%$ of the primitives were not correctly classified, according to the label assigned to the class.

For the recognition we adopted the protocol P2 used for pose estimation (see [11, 67]) using one specific subject for testing. Table 4 presents the average accuracy of the recognition for each group, as well as an ablation study with respect to the components of the cost function used in eq. (18). Fig. 13 shows the corresponding confusion matrices. The results suggest that the DPM classification together with the proposed recognition method capture the main characteristics of each motion primitive category.

Finally, we evaluate the recognition accuracy by considering the same sequences though computing the subject’s pose directly from the video frames using [11]. The corresponding results are shown in parentheses in the last column of Table 4. We note that the recognition accuracy decreases in average just by 4% by using the estimated pose.

Table 4. Primitive recognition accuracy and ablation study

Group	Projection on tangent plane	Frenet frame rotation	Torsion	Curvature	All
G1	0.82	0.80	0.70	0.72	0.84 (0.82)
G2	0.85	0.82	0.75	0.75	0.86 (0.84)
G3	0.80	0.80	0.73	0.74	0.82 (0.78)
G4	0.80	0.79	0.75	0.77	0.83 (0.76)
G5	0.87	0.86	0.72	0.72	0.88 (0.81)
G6	0.86	0.86	0.71	0.73	0.88 (0.82)
Average	0.83	0.82	0.73	0.76	0.85 (0.81)

6.4 Primitives in Activities

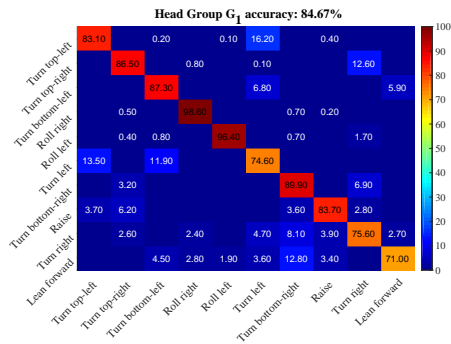
We examine the distribution of discovered motion primitives with respect to the activities been performed by the subjects. We perform our analysis on the sequences of the ActivityNet dataset. More specifically we use the 3D pose estimation algorithm of [11] on the video sequences of ActivityNet. We then extract motion primitives using the motion flux and perform recognition based on the extracted poses. We consider only the segments of the videos labeled with a corresponding activity. Additionally, we use only the segments were a single subject is detected and at least the upper body is visible. Fig. 14 display the distribution of the motion primitives for the five most general activities according to the ActivityNet taxonomy.

6.5 Motion Primitives Dataset

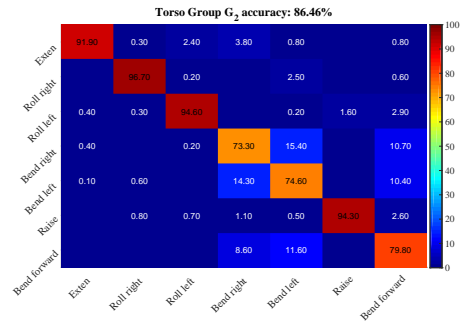
The dataset of annotated motion primitives extracted from the MoCap sequences of H3.6M [13], CMU [14] and KIT-WB [15] has been made publicly available at <https://github.com/MotionPrimitives/MotionPrimitives>. The dataset provides the start and end frames of each motion primitive together with the corresponding label as well as a reference to the MoCap sequence from which the motion primitive has been extracted.

6.6 Comparisons with state of the art on motion primitive recognition

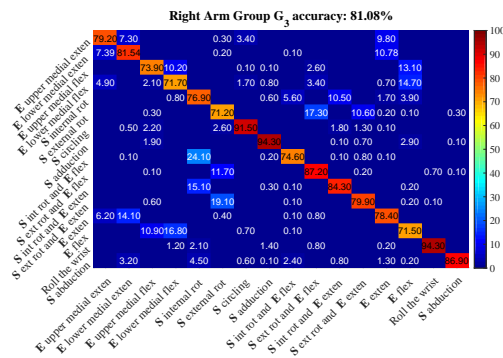
We consider here the results of [3], so far the only work providing quantitative results on human motion primitives, as far as we know. Here performance is evaluated for 4 actions of the arms (gestures), namely *Point right*, *Raise arm*, *Clap* and *Wave*. The authors perform two tests, one without noise in the start and end frames of the primitives and one where the primitives are affected by noise. In the noise-free case their overall accuracy is 94.4% while in



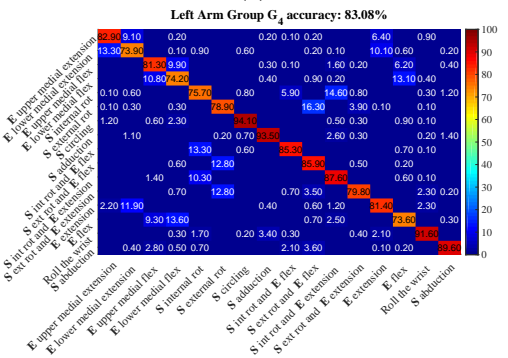
(a)



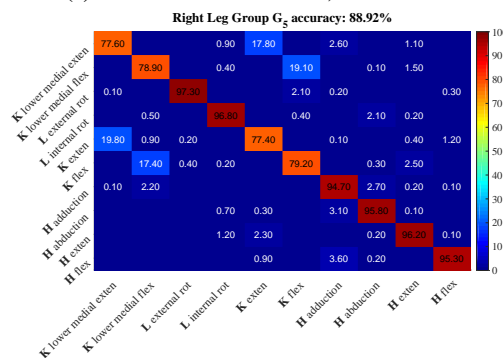
(b)



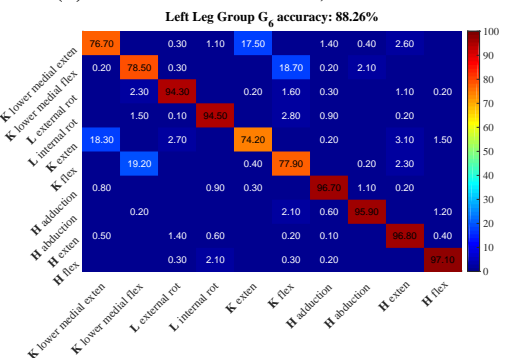
(c) S indicates the shoulder, E the elbow.



(d) S indicates the shoulder, E the elbow.



(e) K indicates the knee, L the leg and H the hip



(f) K indicates the knee, L the leg and H the hip

Fig 13. Confusion matrices for motion primitive recognition. The matrices for G1 and G2 are shown at the top, G3 and G4 at the middle, while G5 and G6 are shown at the bottom.

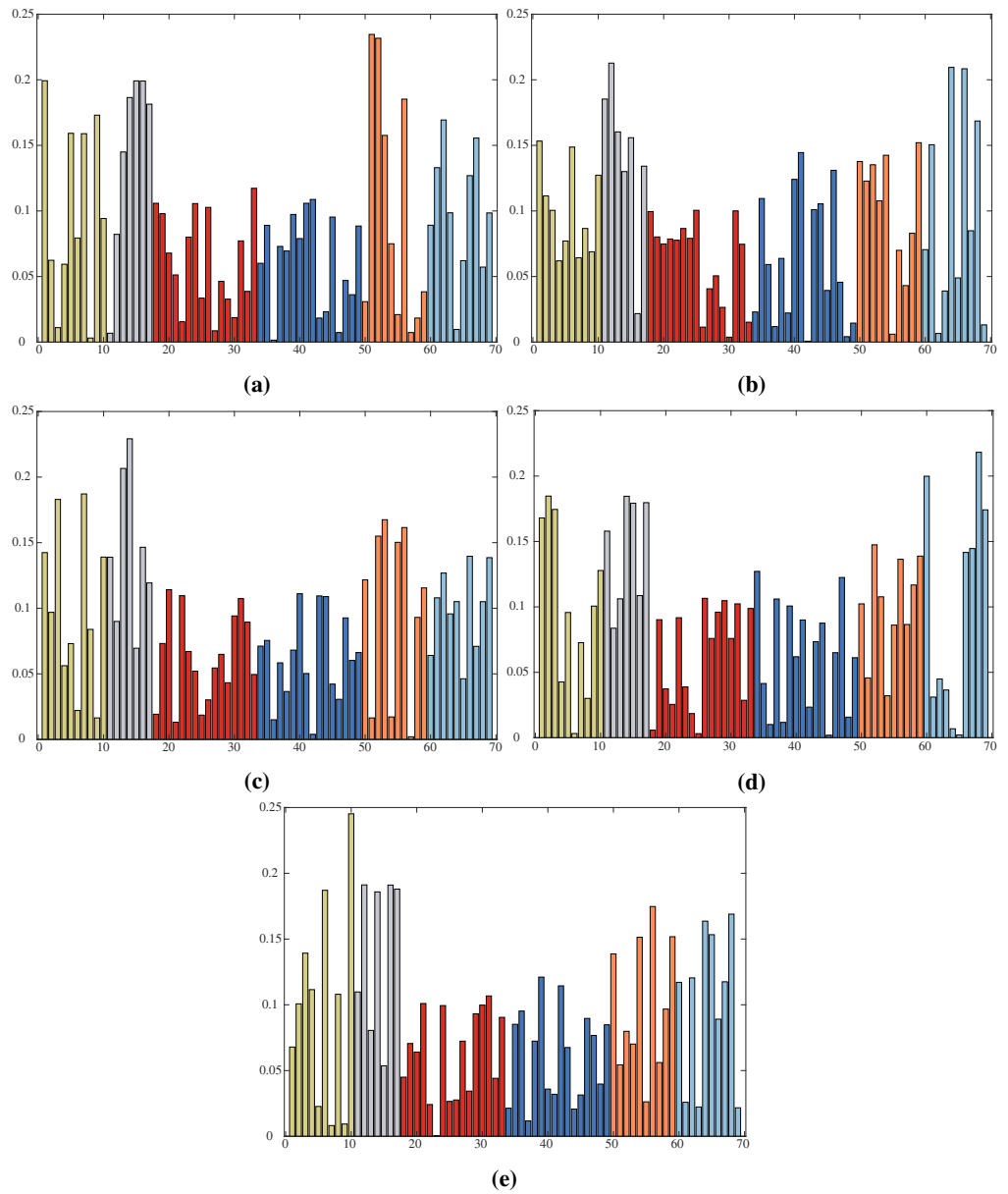


Fig 14. Distribution of the 69 primitives for the five most general categories of the ActivityNet dataset. Clock-wise from top-left: *Eating and drinking Activities*; *Sports, Exercise, and Recreation*; *Socializing, Relaxing, and Leisure*; *Personal Care*; *Household Activities*. Each color corresponds to a different group following the convention of Fig. 12.

the presence of noise the accuracy is 86.9%. Our results are not immediately comparable with the ones of [3] since we use public datasets (see above §6.1, while they have built their own dataset, which is not publicly available. Furthermore, we have obtained by our classification process 16 primitives for each arm which are in accordance with biomechanics primitives. This notwithstanding, we mapped their 22 primitives, denoted by the letters A, \dots, V to our defined primitives of the groups of *Left arm* and *Right arm* (see Table 5). To maintain the use of public datasets we have extracted videos from our reference datasets (see above §6.1) to obtain the 4 above mentioned gestures from 10 different subjects. Hence, we have computed the motion primitives recognition accuracy on these video sets, to compare with [3]. The results are shown in Table 5.

Table 5. Comparison with the 22 motion primitives of [3]

		Shoulder abd.	Shoulder add.	Elbow ext.	Elbow flex.	Shoulder Int. Rot. and elbow flex.	Shoulder Ext. Rot. and elbow ext.	Elbow Upper med. flex.	Elbow Up- per med ext.
A,B,C	Point right	92.3	96.8						
D,E,F		(89.6)	(93.5)						
		82.5							
G,H,I	Raise arm			84.5	77.5				
J,K,L				(81.4)	(73.6)				
				87.5					
M,N,O	Clap					91.7	89.2		
P,Q,R						(87.6)	(85.9)		
						90.0			
S,T	Wave							85.4	87.7
U,V								(81.3)	(82.9)
								87.5	

In Table 5 the capital letters in the first column indicate the primitives in the language of [3]. In the second column are listed the actions formed by the primitives indicated in the first column. In the first row are indicated the primitive taken from our biomechanics language, which we mapped on the [3] primitives. Results are on the diagonal, in gray the results of [3]. We have indicated in parentheses the values illustrated in the confusion matrices. While the values in the confusion matrices were mean precision averages over all experiments for all actions in all the considered datasets, here the results are with respect to an amount of videos comparable to the experiments of [3], hence they are significantly better for the indicated primitives. Despite the results are not quite comparable since we have measured our results on public databases, and in 3D, we can observe that our approach outperforms in all but one case the results in [3].

6.7 Discussion

The results show that our framework discovers and recognizes motion primitives with high accuracy with respect to the manually defined baseline while providing competitive results with respect to [3], the only work, to the best of our knowledge, providing quantitative results on similarly defined motion primitives.

Additionally, given the importance of studying human motion in a wide spectrum of research fields, ranging from robotics to bioscience, we believe that the human motion primitives dataset will be particularly useful in exploring new ideas and for enriching knowledge in these areas.

7 An application of the motion primitives model to surveillance videos

In this section we show how to set up an experiment by using motion primitives. In particular, the application we have chosen is the detection in surveillance videos of dangerous human behaviors. To set up the experiment we consider videos of anomalous and dangerous behaviors, and prove that idiosyncratic primitives, among those identified in Figure 12, appear to characterize these behaviors. The application is quite interesting because it highlights how the combination of primitives allows to detect specific human behaviors. On the one side the motion primitives are used for detection and on the other side they can be used also for characterizing classes of actions or classes of activities.

7.1 Related works and datasets on abnormal behaviors

There is a significant amount of literature on *abnormality* detection in surveillance videos. Only few of them, though, are concerned with dangerous behaviors. These methods can be further divided into those detecting dangerous crowd behaviors, in which the individual motion is superseded by large flows as in [68, 69, 70, 71], and those detecting closer dangerous human behaviors.

Among the latter there are methods focusing on fights [72], methods specialized on violence [73, 74, 75, 76], on aggressive behaviors [77], and on crime [78]. A review on methods for detecting abnormal behaviors, taking into account some of the above mentioned ones, and also discussing available datasets, is provided in [79].

In the last years, also due to the above studies, a number of datasets have been created from real surveillance videos, or from movies repositories. The most used ones are *UCSD Anomaly* [80], *Avenue Dataset* [81], the *Behave* [82] dataset, the *Violent Flows* dataset [71], the *Hockey Fight Dataset* [83], the *Movies Fight Dataset* from [83] too and, finally, the recent *UCF-crime* introduced by [78]. To these datasets some authors, studying abnormal behaviors in surveillance videos, have added specific activities from *UCF101* [84].

To detect dangerous behaviors we considered four of the above datasets most suitable for the task of analyzing human behaviors with small groups of subjects. The first dataset is the *Hockey Fight Dataset* provided by [83], which is formed by 1000 clips of actions from hockey games of the National Hockey League (NHL). A second dataset, also introduced by [83] is the *Movies Fight dataset*, which is composed of 200 video clips obtained from action movies, 100 of which show a fight. Videos in both these datasets are untrimmed but divided in those where there are fights and those where there are no fights. The third dataset is the *UCF-Crime dataset* introduced by [78]. This dataset is formed by 1900 untrimmed surveillance videos of 13 realworld anomalies, including *abuse, arrest, arson, assault, road accident, burglary, explosion, fighting, robbery, shooting, stealing, shoplifting, and vandalism*, and normal videos. These videos have varying length from 30 sec. up to several minutes. In a number of these videos, like explosion and road accident, no human behavior is observable. Among the others there are a number of videos not including human behaviors. Therefore we have chosen a subset of all the UCF-crime dataset for both training and testing. In particular, we have chosen *abuse, arrest, assault, burglary, fighting, robbery, shooting, stealing, and vandalism*. Finally we have taken videos from *UCF101* dataset, which includes 101 human activities.

Given the above selected datasets we aim at showing that once the primitives are computed an off-the-shelf classifier can be used to detect specific behaviors, in this case the dangerous ones.

The method we propose requires to compute the primitives on a selected training set, separating the untrimmed videos with dangerous behaviors from the normal ones, as described below, and then training a non-linear kernel SVM on the two datasets, as illustrated in §7.3. The trained classifier is then tested on the test sets and results are reported in §7.4, comparing

with state of the art approaches.

The main idea we want to convey here is that once primitives are computed all the relevant features for distinguishing a behavior are embedded in the primitive category of the specific group (see §7.4) and therefore the classifier has to deal just with them and not with other features such as poses, images, time and tracking, in so alleviating the classifier burden and allowing to deal with state of the art classifiers. Furthermore, the primitive parameters, used to estimate the primitive classes, are no more needed for the further classification of behaviors. This is the main advantage of human motion primitives modeling, namely their effectiveness in characterizing specific behaviors.

7.2 Primitives computation

For primitives computation we collected all the videos from hokey and fight-movie datasets, we collected from the UCF-crime dataset the videos from *abuse*, *arrest*, *assault*, *burglary*, *fighting*, *robbery*, *shooting*, *stealing*, and *vandalism*. Finally, from UCF101 we collected 276 videos from the datasets *Punch* and *SumoWrestling* and further 276 videos from other sports, randomly chosen as in [72]. The total number of videos collected is 3050 for primitive computation, as illustrated in the following table:

Table 6. Datasets for primitive computation in dangerous behaviors detection

	Hockey		Fight-Movies		UCF-crime		UCF101	
	Danger.	Normal	Danger.	Normal	Danger.	Normal	Danger.	Normal
Video sets	500	500	100	100	650	650	276	276
Training	70%	70%	70%	70%	70%	70%	100%	70%
Test	30%	30%	30%	30%	30%	30%	0%	30%

To compute the primitives for each subject from a small group of people appearing in a frame of a video, we have fitted 3D poses basing on the SMPL model [62] of *human mesh recovery* (HMR) [85]. HMR recovers together with joints and pose also a full 3D mesh from a single image (see Figures 15 and 16), and it is accurate enough to estimate multiple subject poses in a single frame.

Having more than a subject requires to track each subject pose across frames, in order to compute the motion primitives for each of them. To this end we used the joints given by SMPL model in world frame, for the following body joints (see the preliminary Section 3): left and right *hip*, left and right *clavicle* (called shoulder in HMR), and the *head*. These joints are well suited for tracking since they have slower motion with respect to other body parts. Tracking amounts to find the rotations and translations amid all the bodies appearing in two consecutive frames, and identifying the rotation and translation pertaining to each subject across the two frames. Consider two consecutive frames indexed by t and $t+1$, and let $\mathcal{J}^{(t)} = \{j_1^{(t)}, \dots, j_5^{(t)}\}$ and $\mathcal{J}^{(t+1)} = \{j_1^{(t+1)}, \dots, j_5^{(t+1)}\}$ be the joints in world frame of the above mentioned body components, where joint subscripts indicate in the order left and right hip, left and right clavicle and head. We first find the translation \mathbf{d} and rotation R between any two set of joints appearing in the frames t and $t+1$ (see also Section 3):

$$(R, \mathbf{d}) = \arg \min_{R \in SO(3), \mathbf{d} \in \mathbb{R}^3} \sum_{i=1}^5 w_i \|(R j_i^{(t)} + \mathbf{d}) - j_i^{(t+1)}\|^2 \quad (19)$$

With $w_i > 0$ weights for each pair of joints in (t) and $(t+1)$. Let $\hat{\mathcal{J}} = (\sum_{i=1}^5 w_i j_i) / \sum_{i=1}^5 w_i$ be the weighted centroids of the set of joints \mathcal{J} . The minimization in (19) is solved by computing the singular value decomposition $U \Sigma V^\top$ of the covariance matrix $\bar{\mathcal{J}}^{(t)} W (\bar{\mathcal{J}}^{(t+1)})^\top$ of the normalized joints $\bar{\mathcal{J}}^{(t)}, \bar{\mathcal{J}}^{(t+1)}$, obtained by subtracting the weighted centroid to each joints set. Here W is the diagonal matrix of the

weights w_i . Let H be the diagonal matrix $\text{diag}(\mathbf{1}, \det(VU^\top))$, then the rotations and translations between sets of joints are found as:

$$R = VHU^\top \quad \text{and} \quad \mathbf{d} = \hat{\mathcal{J}}^{(t+1)} - R\hat{\mathcal{J}}^{(t)} \quad (20)$$

Finally, once we have obtained the rotation matrices and the translation vectors between the sets of considered joints of all the fitted skeletons, from frame t to frame $t + 1$, we can track each individual skeleton S_k . A skeleton $S_k^{(t+1)}$ belongs to the same subject fitted by skeleton $S_k^{(t)}$, at frame t , if the rotation R_k and translation \mathbf{d}_k , obtained according to eq. (20) between the chosen joints $\mathcal{J}^{(t)}$ of $S_k^{(t)}$ and $\mathcal{J}^{(t+1)}$ of $S_k^{(t+1)}$, satisfy

$$(R_k, \mathbf{d}_k) = \arg \min_{R_k \in SO(3), \mathbf{d}_k \in \mathbb{R}^3, k=1:s} \|\mathcal{J}^{(t+1)} - ((\mathcal{J}^{(t)} R_k)^\top + \mathbf{d}_k)^\top\|_F \quad (21)$$

With $\|\cdot\|_F$ the Frobenious norm and $s = N_S! / ((N_S - 2)!2!)$, with N_S the common number of fitted skeletons S in both frame t and $t + 1$.

Once the skeletons are tracked we can compute the unknown primitives from the flux (see Section 4) as paths $\gamma_{G_m}^T : I \subset \mathbb{R} \mapsto \mathbb{R}^9$, for each group G_m , with I the time interval, specified by the frame sequence, and scale it as described in Section 4. We can then use the parameters Θ learned with the recognition model, detailed in §5.2, to assign a label $\mathcal{L}_w^{G_m}$ to each primitive segmented by the motion flux as precised in eq. (18). Namely, we find the model identified by the parameter Θ_w , which maximizes the probability of the primitive under consideration. We recall that for each group G_m , $m = 1, \dots, 6$ there are q models with $q \in \{7, 10, 16\}$ (see the primitives representation in Figure 12).

Our model of motion primitives relies significantly on the accuracy of the 3D pose estimation. We have chosen the model HMR [85] based on SMPL [62], in place of [26, 12], since it is most recent and highly accurate. Still not all the videos chosen obtain a reasonable fitting, therefore after skeleton fitting and tracking a number of videos from UCF-crime have been removed from the considered set.

7.3 Training a non-linear binary classifier

All the computed primitives are labeled by their name (e.g. *Elbow flex*), according to the recognition model, as specified above. A set of primitives for a given video is formed as follows. Primitive names are embedded into real numbers $r \sim \text{Unif}(0, 1)$, such that for each primitive name there is a precise real number. Given frame t for each skeleton appearing in the frame we form a vector of dimension 6×1 , where the 6 elements are the corresponding embedded primitive names occurring at frame t . Let $\gamma_{G_m}^{(t)}$ denote the primitive of the body group G_m , and u the mapping of the primitive name to the real number:

$$\mathbf{x}_j^{(t)} = (u(\gamma_{G_1}^{(t)}), u(\gamma_{G_2}^{(t)}), \dots, u(\gamma_{G_6}^{(t)}))^\top \quad (22)$$

Where j indicates the j -th skeleton appearing in frame t . Note that t and j are actually indicated just for forming the training set, to select from all the gathered vectors \mathbf{x} those that have changing primitives. Namely, for training, from the set of all vectors in each frame, we have retained only those vectors in which at least one primitive changes, for each recorded skeleton.

For training we have selected videos for both dangerous behaviors and normal behaviors, thus labeling them with 1 for dangerous and -1 for normal behaviors, as follows. We selected 70% of fighting and 70% of not fighting from both hockey and fight movies; from UCF101 we have selected all videos in *Punch* and *SumoWrestling*, getting 276 videos and further 276 videos randomly from sport activities. For UCF-crime we proceeded as follows. We have selected the videos from all the crime activities specified above with time length less than

60sec. and cropped the first and last 10sec., in order to do a weak supervised training, namely, as in [78] we have not trimmed the video. Thus we obtained 173 videos for abnormal activities and we selected 173 videos from the normal activities. The total number of videos for training is 1634 videos. All the remaining video with computed primitives have been used for testing.

The resulting data structure is:

$$\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\} \text{ with } \mathbf{x} \in \mathbb{R}^6, y \in \{-1, 1 \mid -1 \text{ if normal, } 1 \text{ if dangerous}\} \quad (23)$$

The SVM [86] is a popular classification method computing, for two non-separable classes, the classifier:

$$\begin{aligned} f(\mathbf{x}) &= (\sum_{i=1}^n y_i \alpha_i K(\mathbf{x}_i, \mathbf{x}) + b) \\ \hat{y} &= \text{sgn}(f(\mathbf{x})) \end{aligned} \quad (24)$$

where K is the kernel function $\varphi(\mathbf{x}_i)^\top \varphi(\mathbf{x}_j)$ with φ the feature map, here we considered the RBF kernel $\exp(-\eta \|\mathbf{x}_i - \mathbf{x}_j\|_{\ell_2}^2)$, with η a tunable parameter. Classification is obtained by solving the constrained optimization problem:

$$\max_{\alpha} \frac{1}{2} \alpha^\top \Omega \alpha - \mathbf{e}^\top \alpha \quad \text{subject to } \mathbf{y}^\top \alpha = 0, 0 \leq \alpha_i \leq \lambda \quad (25)$$

Here Ω is a square $n \times n$ positive semidefinite matrix, with $\omega_{i,j} = y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$, \mathbf{e} is a vector of ones, the non zero α_i define the support vectors, and λ is the regularization parameter of the primal optimization problem $\min_{\mathbf{w}, b, \xi} \frac{1}{2} \mathbf{w} \mathbf{w}^\top + \lambda \sum_{i=1}^n \xi_i$ [87]. To obtain posterior probabilities we applied the Platt scaling [88], proposing a sigmoid model to fit a posterior on the SVM output:

$$P(y = 1 | f(\mathbf{x})) = \frac{1}{1 + \exp(Af(\mathbf{x}) + B)} \quad (26)$$

Here the parameters A and B are fitted by solving the maximum likelihood problem:

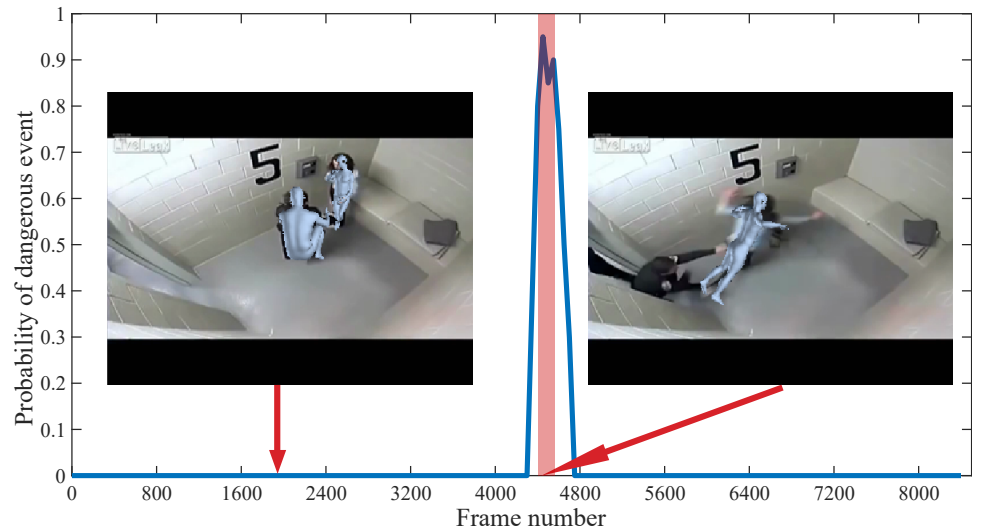
$$\min_{z=(A,B)} F(z) = - \sum_{i=1}^n (t_i \log(p_i) + (1 - t_i) \log(1 - p_i)) \quad (27)$$

Using as prior the number of positive N_+ and negative N_- examples in the training data, with $p_i = P(y = 1 | f(\mathbf{x}_i))$, $t_i = (N_+ + 1)/(N_+ + 2)$ if $y_i = 1$ and $1/(N_- + 2)$ if $y_i = -1$. See also [89] for an improved algorithm with respect to [88].

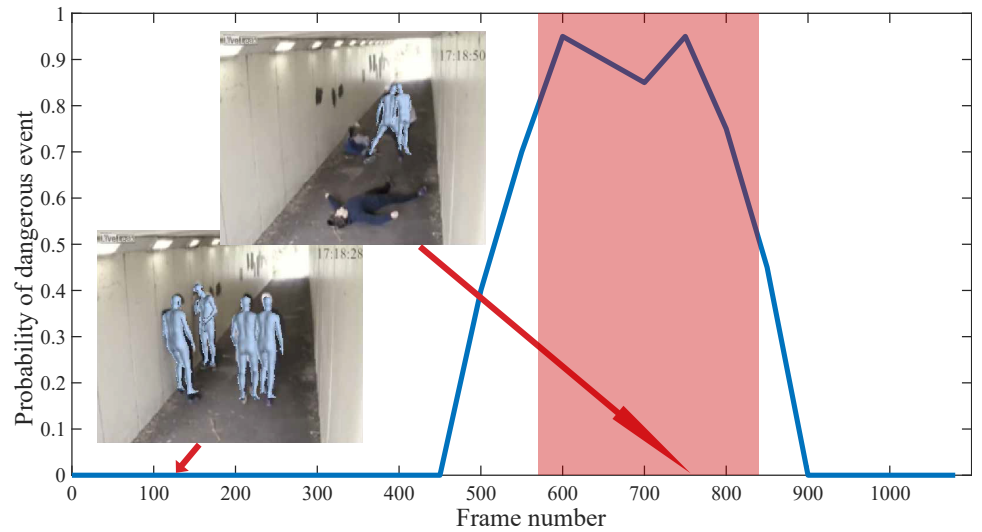
To obtain the probability that at a given frame t a dangerous event occurs we compute the average response to the primitives of each subject which has been detected. More precisely, let s be the number of subjects in frame t for which the primitives are computed, then the observation $\mathbf{x}^{(t)} = (\mathbf{x}_1^{(t)}, \dots, \mathbf{x}_s^{(t)})$. Given $\mathbf{x}^{(t)}$, and assuming that the SVM scores for each $\mathbf{x}_i^{(t)}$ are independent, we can define the probability that a dangerous event Y is occurring at t , in a surveillance video, as the expectation:

$$P(Y | \mathbf{x}^{(t)}) = \sum_{i=1}^s p(\hat{y}_i^{(t)} | \mathbf{x}_i^{(t)}) P(y_i = 1 | f(\mathbf{x}_i^{(t)})) \quad (28)$$

Here $p(\hat{y}_i^{(t)} | \mathbf{x}_i^{(t)})$ is computed by remapping the scores to $[0, 1]$ such that $\sum_{i=1}^s p(\hat{y}_i^{(t)} | \mathbf{x}_i^{(t)}) = 1$. Testing has been done on the videos on which the primitives have been precomputed, and the results are shown together with comparisons with the state of the art in §7.4. Note that the method is not yet suitable for online detection of dangerous behaviors, still it can be advanced to online detection, by lifting the computation of the flux with motion anticipation.



(a)



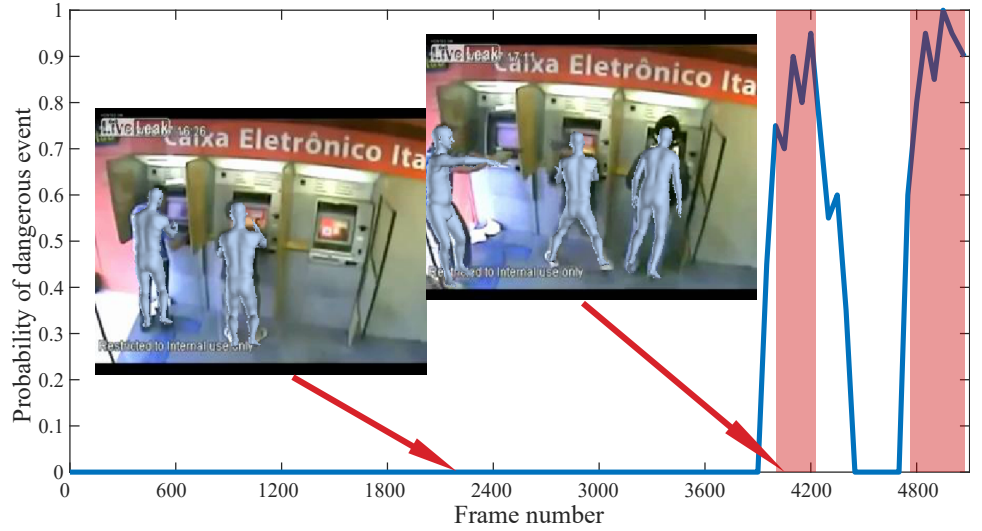
(b)

Fig 15. Results of the proposed method on videos from UCF-Crime dataset. From top: *Abuse*, *Fighting*. Colored window shows ground truth anomalous region.

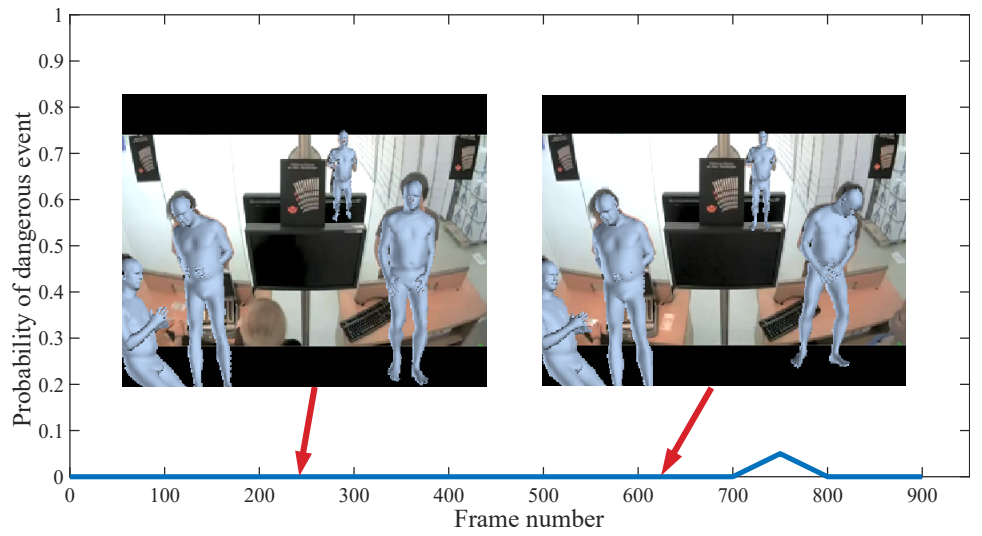
7.4 Results and comparisons with the state of the art

We discuss now the results achieved by our method for abnormal behavior detection based on human motion primitives. Figure 15 shows some qualitative results of dangerous behaviors detection in four videos. Three videos correspond to crime activities, namely *Abuse*, *Fighting* and *Shooting*, while the last displays a normal activity. The curve plotted in the graphs provides for each frame the probability that a dangerous event is occurring, according to eq. (28). The highlighted region corresponds to the interval where a crime activity occurs. From this graphs it is evident that the crime activity detection follows closely the ground truth. For each example we also show two representative frames overlaid with the human meshes identified by HMR. Similarly, Figure 17 shows some representative examples of fitted human meshes for videos taken from Hockey and Movie Fights datasets.

Fig. 19 presents the ROC curves of the proposed method for the four datasets considered,



(a)



(b)

Fig 16. Results of the proposed method on videos from UCF-Crime dataset. From top: *Shooting, Normal*. Colored window shows ground truth anomalous region.

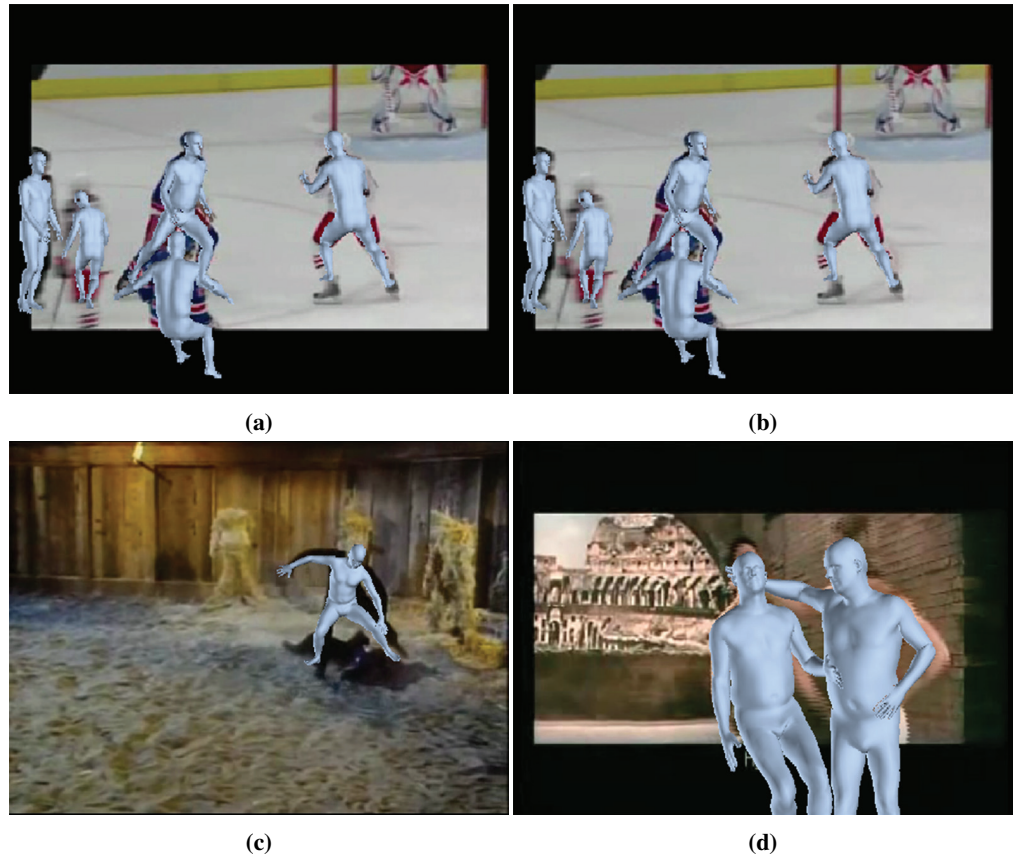


Fig 17. Instances of videos with human meshes fitted using HMR from Hockey and Movies datasets [83].

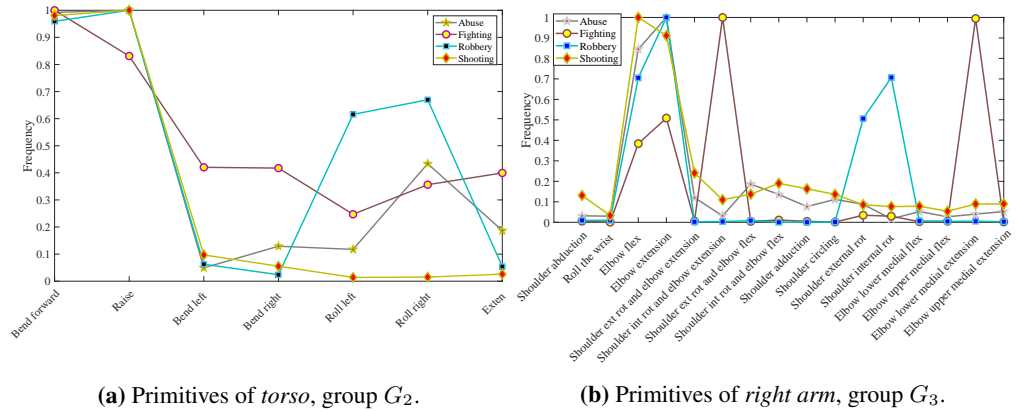


Fig 18. Frequency graphs of the occurrences of primitives for groups G_2 (torso) and G_3 (right arm) in the videos of *Abuse*, *Fighting*, *Robbery*, and *Shooting* of the dataset UCF-crime.

namely UCF-Crime, UCF101, Hockey Fights and Movie Fights. The corresponding values of the area under curve (AUC) are 76.15%, 91.92%, 98.44% and 98.77%, respectively. Table 8 presents the mean accuracy, its standard deviation and the area under the receiver-operating-characteristic (ROC) curve of our method in comparison with other state-of-the-art methods. The results of the other methods are taken from [72]. We observe that our method achieves better performance on the Hockey Fights and Movies Fights datasets

Table 7. AUC comparison with state-of-the-art methods on the UCF-Crime dataset.

Method	Binary classifier	Hasan et al. [90]	Lu et al. [91]	[78]	[78] w. constraints	Ours
AUC	50.0	50.6	65.51	74.44	75.41	76.15

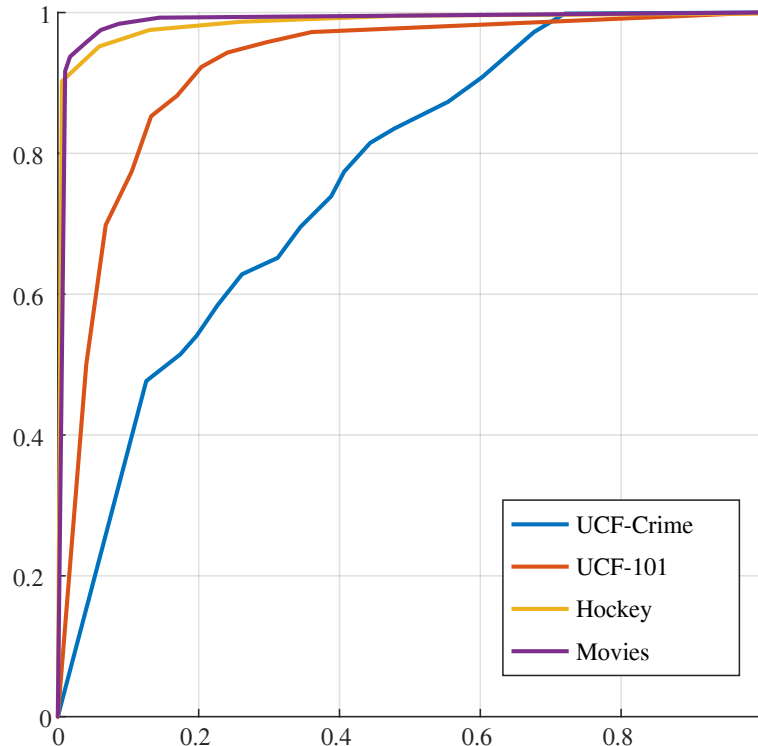


Fig 19. ROC curves of the proposed method for UFC-Crime, UFC101, Hockey and Movies datasets.

while it has very similar performance with the best performing method on the UCF101 dataset.

Additionally, in Figure 18 we present the frequency graphs of primitive occurrences for groups G2 and G3, for the crime activities *Abuse*, *Fighting*, *Robbery*, and *Shooting*. The graphs show that each type of activity manifests itself by a different combination of idiosyncratic motions of the limbs. This fact can be used to achieve finer grained categorization of the crime activities, however, we do not examine further this possibility in this work.

Figure 19 presents the ROC curves of the proposed method for the four datasets considered, namely UCF-Crime, UCF101, Hockey Fights and Movie Fights. The corresponding values of the area under curve (AUC) are 76.15%, 91.92%, 98.44% and 98.77%, respectively. Table 8 presents the mean accuracy, its standard deviation and the area under the receiver-operating-characteristic (ROC) curve of our method in comparison with other state-of-the-art methods. The results of the other methods are taken from [72]. We observe that our method achieves better performance on the Hockey Fights and Movies Fights datasets while it has very similar performance with the best performing method on the UCF101 dataset.

Finally, Table 7 gives a comparison of the results achieved by our method on the UCF-Crime dataset in comparison with results from other state-of-the-art methods as reported in [78]. In this case we have to highlight that our results are not directly comparable with the ones reported in [78] as we restrict our analysis on videos where human subjects are visible. Nevertheless, the results indicate that also on this database the proposed method is able to achieve state-of-the-art performance on crime activity detection.

Table 8. Comparison with state-of-the-art methods on the datasets Movies, UCF101 and Hockey.

Method	Classifier	Datasets		
		Movies	Hockey	UCF101
BoW (STIP)	SVM	82.3±0.9/0.88	88.50.2/0.95	72.51.5/0.74
	AdaBoost	75.30.83/0.83	87.1±0.2/0.93	63.1±1.9/0.68
	RF	97.7±0.5/0.99	96.5±0.2/0.99	87.3±0.8/0.94
BoW (MoSIFT)	SVM	63.4±1.6/0.72	83.9±0.6/0.93	81.3± 1/0.86
	AdaBoost	65.3±2.1/0.72	86.9±1.6/0.96	52.8±3.6/0.62
	RF	75.1±1.6/0.81	96.7±0.7/0.99	86.3±0.8/0.93
ViF	SVM	96.7±0.3/0.98	82.3±0.2/0.91	77.7±2.16/0.87
	AdaBoost	92.8±0.4/0.97	82.2±0.4/0.91	78.4±1.7/0.86
	RF	88.9±1.2/0.97	82.4±0.6/0.9	77±1.2/0.85
LMP	SVM	84.4±0.8/0.92	75.9±0.3/0.84	65.9±1.5/0.74
	AdaBoost	81.5±2.1/0.86	76.5±0.9/0.82	67.1±1/0.71
	RF	92±1/0.96	77.7±0.6/0.85	71.4±1.6/0.78
[75]	SVM	85.4±9.3/0.74	90.1±0/0.95	93.4±6.1/0.94
	AdaBoost	98.9±0.22/0.99	90.1±0/0.90	92.8±6.2/0.94
	RF	90.4±3.1/0.99	61.5±6.8/0.96	64.8±15.9/0.93
[72] v1	SVM	87.9±1/0.97	70.8±0.4/0.75	72.1±0.9/0.78
	AdaBoost	81.8±0.5/0.82	70.7±0.2/0.7	71.7±0.9/0.72
	RF	97.7±0.4/0.98	79.3±0.5/0.88	74.8±1.5/0.83
[72] v2	SVM	87.2±0.7/0.97	72.5±0.5/0.76	71.2±0.7/0.78
	AdaBoost	81.7±0.2/0.82	71.7±0.3/0.72	71±0.8/0.72
	RF	97.8±0.4/0.97	82.4±0.6/0.9	79.5±0.9/0.85
Ours	SVM	99.1±0.3/0.99	97.2±0.8/0.98	93.3±2.1/0.92

8 Conclusions

We presented a framework for automatically discovering and recognizing human motion primitives from video sequences based on the motion of groups of joints of a subject. To this end the motion flux is introduced which captures the variation of the velocity of the joints within a specific interval. Motion primitives are discovered by identifying intervals between rest instances that maximize the motion flux. The unlabeled discovered primitives have been separated into different categories using a non-parametric Bayesian mixture model.

We experimentally show that each primitive category naturally corresponds to movements described using biomechanical terms. Models of each primitive category are built which are then used for primitive recognition in new sequences. The results show that the proposed method is able to robustly discover and recognize motion primitives from videos, by using state-of-the-art methods for estimating the 3D pose of the subject of interest. Additionally, the results suggest that the motion primitives categories are highly discriminative for characterizing the activity been performed by the subject.

Finally, a dataset of motion primitives is made publicly available to further encourage result reproducibility and benchmarking of methods dealing with the discovery and recognition of human motion primitives.

9 Acknowledgments

This research is supported by European Union’s Horizon 2020 Research and Innovation programme under grant agreement No 643950, project SecondHand

References

1. Ghanem B, Niebles JC, Snoek C, Heilbron FC, Alwassel H, Khrisna R, et al. ActivityNet Challenge 2017 Summary. arXiv:171008011. 2017.
2. Yang Y, Saleemi I, Shah M. Discovering Motion Primitives for Unsupervised Grouping and One-Shot Learning of Human Actions, Gestures, and Expressions. TPAMI. 2013;35(7).
3. Holte MB, Moeslund TB, Fihl P. View-invariant gesture recognition using 3D optical flow and harmonic motion context. *Comp Vis and Im Underst*. 2010;114(12):1353–1361.
4. Flash T, Hochner B. Motor primitives in vertebrates and invertebrates. *Curr Op in Neurob*. 2005;15(6):660–666.
5. Polyakov F. Affine differential geometry and smoothness maximization as tools for identifying geometric movement primitives. *Biological cybernetics*. 2017;111(1):5–24.
6. Ting LH, Chiel HJ, Trumbower RD, Allen JL, McKay JL, Hackney ME, et al. Neuromechanical principles underlying movement modularity and their implications for rehabilitation. *Neuron*. 2015;86(1):38–54.
7. Hogan N, Sternad D. Dynamic primitives of motor behavior. *Biological cybernetics*. 2012; p. 1–13.
8. Amor HB, Neumann G, Kamthe S, Kroemer O, Peters J. Interaction primitives for human-robot cooperation tasks. In: ICRA; 2014. p. 2831–2837.
9. Moro FL, Tsagarakis NG, Caldwell DG. On the kinematic Motion Primitives (kMPs)—theory and application. *Frontiers in neurorobotics*. 2012;6.
10. Azad P, Asfour T, Dillmann R. Toward an unified representation for imitation of human motion on humanoids. In: *Robotics and Automation*; 2007. p. 2558–2563.
11. Sanzari M, Ntouskos V, Pirri F. Bayesian Image Based 3D Pose Estimation. In: ECCV. vol. 8; 2016. p. 566–582.
12. Tome D, Russell C, Agapito L. Lifting from the deep: Convolutional 3d pose estimation from a single image. *CVPR 2017 Proceedings*. 2017; p. 2500–2509.
13. Ionescu C, Papava D, Olaru V, Sminchisescu C. Human3.6M: Large Scale Datasets and Predictive Methods for 3D Human Sensing in Natural Environments. TPAMI. 2014;36(7):1325–1339.
14. CMU Mocap Database. <http://mocap.cs.cmu.edu/>.
15. Mandery C, Terlemez O, Do M, Vahrenkamp N, Asfour T. The KIT whole-body human motion database. In: ICAR; 2015. p. 329–336.
16. Weinland D, Ronfard R, Boyer E. Automatic Discovery of Action Taxonomies from Multiple Views. In: CVPR. vol. 2; 2006. p. 1639–1645.
17. Li Y, Fermuller C, Aloimonos Y, Ji H. Learning shift-invariant sparse representation of actions. In: CVPR; 2010. p. 2630–2637.

18. Turaga P, Chellappa R, Subrahmanian VS, Udrea O. Machine Recognition of Human Activities: A Survey. *Trans on Circuits and Systems for Video Technology*. 2008;18(11):1473–1488.
19. Sigal L, Balan AO, Black MJ. HumanEva: Synchronized Video and Motion Capture Dataset and Baseline Algorithm for Evaluation of Articulated Human Motion. *IJCV*. 2009;87(1):4.
20. Moeslund TB, Hilton A, Krüger V. A survey of advances in vision-based human motion capture and analysis. *Computer vision and image understanding*. 2006;104(2):90–126.
21. Akhter I, Black MJ. Pose-conditioned joint angle limits for 3D human pose reconstruction. In: *CVPR*; 2015. p. 1446–1455.
22. Zhou X, Zhu M, Leonardos S, Derpanis K, Daniilidis K. Sparseness Meets Deepness: 3D Human Pose Estimation from Monocular Video. In: *CVPR*; 2016.
23. Ntouskos V, Papadakis P, Pirri F. Discriminative Sequence Back-Constrained GP-LVM for MOCAP Based Action Recognition. In: *Proceedings of the 2nd International Conference on Pattern Recognition Applications and Methods*; 2013. p. 87–96.
24. Ntouskos V, Papadakis P, Pirri F. Probabilistic Discriminative Dimensionality Reduction for Pose-Based Action Recognition. In: *Pattern Recognition Applications and Methods*. vol. 318 of *Advances in Intelligent Systems and Computing*; 2015. p. 137–152.
25. Pirri F, Pizzoli M. Inference about Actions: Levesques view on action ability and Dirichlet processes. In: Lakemeyer G, McIlraith SA, editors. *Knowing, Reasoning, and Acting Essays in Honour of Hector J. Levesque*; 2011.
26. Natola F, Ntouskos V, Sanzari M, Pirri F. Bayesian non-parametric inference for manifold based MoCap representation. In: *Proceedings of the IEEE International Conference on Computer Vision*; 2015. p. 4606–4614.
27. Natola F, Ntouskos V, Pirri F. Collaborative Activities Understanding from 3D Data. *Doctoral Consortium on Pattern Recognition Applications and Methods (DCPRAM)*. 2015.
28. Fanello S, Gori I, Pirri F. Arm-Hand Behaviours Modelling: From Attention to Imitation. In: *Advances in Visual Computing*; 2010. p. 616–627.
29. Bizzi E, Mussa-Ivaldi FA. Toward a Neurobiology of Coordinate Transformations. In: *The Cognitive Neurosciences*; 1995. p. 495–506.
30. Flash T, Meirovitch Y, Barliya A. Models of human movement: Trajectory planning and inverse kinematics studies. *RAS*. 2013;61(4):330–339.
31. Viviani P, Flash T. Minimum-jerk, two-thirds power law, and isochrony: converging approaches to movement planning. *J of Exp Psy: Human Perception and Performance*. 1995;21(1):32.
32. Flash T, Handzel AA. Affine differential geometry analysis of human arm movements. *Bio Cyb*. 2007;96(6):577–601.
33. Biess A, Liebermann DG, Flash T. A Computational Model for Redundant Human Three-Dimensional Pointing Movements: Integration of Independent Spatial and Temporal Motor Plans Simplifies Movement Dynamics. *J Neuroscience*. 2007;27(48):13045–13064.

34. Lacquaniti F, Terzuolo C, Viviani P. The law relating the kinematic and figural aspects of drawing movements. *Acta Psychologica*. 1983;54(13).
35. Viviani P, Schneider R. A developmental study of the relationship between geometry and kinematics in drawing movements. *J of Experimental Psychology: Human Perception and Performance*. 1991;17(1).
36. Maoz U, Flash T. Spatial constant equi-affine speed and motion perception. *J of Neurophysiology*. 2014;111(2):336–349.
37. Gong D, Medioni G, Zhao X. Structured Time Series Analysis for Human Action Segmentation and Recognition. *TPAMI*. 2014;36(7):1414–1427.
38. Lillo I, Niebles JC, Soto A. A Hierarchical Pose-Based Approach to Complex Action Understanding Using Dictionaries of Actionlets and Motion Poselets. In: *CVPR*; 2016.
39. Lu J, Xu R, Corso JJ. Human action segmentation with hierarchical supervoxel consistency. In: *CVPR*; 2015. p. 3762–3771.
40. Bouchard D, Badler N. In: *Semantic Segmentation of Motion Capture Using Laban Movement Analysis*. Springer; 2007. p. 37–44.
41. Vecchio DD, Murray RM, Perona P. Decomposition of human motion into dynamics-based primitives with application to drawing tasks. *Automatica*. 2003;39(12):2085–2098.
42. Endres D, Meirovitch Y, Flash T, Giese MA. Segmenting sign language into motor primitives with Bayesian binning. *Frontiers in computational neuroscience*. 2013;7.
43. Ijspeert AJ, Nakanishi J, Hoffmann H, Pastor P, Schaal S. Dynamical movement primitives: learning attractor models for motor behaviors. *Neural computation*. 2013;25(2):328–373.
44. Gams A, Petri T, Do M, Nemeč B, Morimoto J, Asfour T, et al. Adaptation and coaching of periodic motion primitives through physical and visual interaction. *RAS*. 2016;75:340–351.
45. Pastor P, Hoffmann H, Asfour T, Schaal S. Learning and generalization of motor skills by learning from demonstration. In: *ICRA*; 2009. p. 763–768.
46. Kober J, Peters JR. Policy search for motor primitives in robotics. In: *Adv. in neural inf. proc. systems*; 2009. p. 849–856.
47. Park DH, Hoffmann H, Pastor P, Schaal S. Movement reproduction and obstacle avoidance with dynamic movement primitives and potential fields. In: *ICHR*; 2008. p. 91–98.
48. Ureche ALP, Umezawa K, Nakamura Y, Billard A. Task Parameterization Using Continuous Constraints Extracted From Human Demonstrations. *IEEE Trans Robot*. 2015.
49. Asfour T, Gyarfas F, Azad P, Dillmann R. Imitation Learning of Dual-Arm Manipulation Tasks in Humanoid Robots. In: *International Conference on Humanoid Robots*; 2006. p. 40–47.
50. Luo R, Berenson D. A framework for unsupervised online human reaching motion recognition and early prediction. In: *IROS*; 2015. p. 2426–2433.

51. Marr D, Vaina L. Representation and Recognition of the Movements of Shapes. *Proceedings of the Royal Society of London B: Biological Sciences*. 1982;214(1197):501–524.
52. de los Reyes-Guzmán A, Dimbwadyo-Terrer I, Trincado-Alonso F, Monasterio-Huelin F, Torricelli D, Gil-Agudo A. Quantitative assessment based on kinematic measures of functional impairments during upper extremity movements: A review. *Clinical Biomechanics*. 2014;29(7):719–727.
53. Gates DH, Walters LS, Cowley J, Wilken JM, Resnik L. Range of motion requirements for upper-limb activities of daily living. *American J of Occupational Therapy*. 2016;70(1).
54. Hamill J, Knutzen KM. *Biomechanical basis of human movement*. Lippincott Williams & Wilkins; 2006.
55. Abernethy B. *Biophysical foundations of human movement*. Human Kinetics; 2013.
56. Alt H, Guibas LJ. Discrete geometric shapes: Matching, interpolation, and approximation. In: *Handbook of computational geometry*. Elsevier; 2000. p. 121–153.
57. Ferguson TS. A Bayesian analysis of some nonparametric problems. *Ann Stat*. 1973; p. 209–230.
58. Antoniak CE. Mixtures of Dirichlet processes with applications to Bayesian nonparametric problems. *Ann Stat*. 1974; p. 1152–1174.
59. Teh YW. Dirichlet process. In: *Encyclopedia of machine learning*. Springer; 2011. p. 280–287.
60. West M. Hyperparameter estimation in Dirichlet process mixture models. *Duke University ISDS Discussion Paper# 92-A03*; 1992.
61. Jain S, Neal RM. A split-merge Markov chain Monte Carlo procedure for the Dirichlet process mixture model. *J of Comp and Graph Statistics*. 2004.
62. Loper M, Mahmood N, Romero J, Pons-Moll G, Black MJ. SMPL: A skinned multi-person linear model. *ACM Transactions on Graphics (TOG)*. 2015;34(6):248.
63. <https://www.turbosquid.com/>.
64. <https://renderpeople.com/3d-people/>.
65. Loper M, Mahmood N, Black MJ. MoSh: Motion and shape capture from sparse markers. *ACM Transactions on Graphics (TOG)*. 2014;33(6):220.
66. Varol G, Romero J, Martin X, Mahmood N, Black MJ, Laptev I, et al. Learning from synthetic humans. In: *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2017)*; 2017.
67. Tekin B, Rozantsev A, Lepetit V, Fua P. Direct prediction of 3d body poses from motion compensated sequences. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*; 2016. p. 991–1000.
68. Mohammadi S, Perina A, Kiani H, Murino V. Angry crowds: Detecting violent events in videos. In: *European Conference on Computer Vision*. Springer; 2016. p. 3–18.
69. Mohammadi S, Kiani H, Perina A, Murino V. Violence detection in crowded scenes using substantial derivative. In: *2015 12th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)*. IEEE; 2015. p. 1–6.

70. Mousavi H, Mohammadi S, Perina A, Chellali R, Mur V. Analyzing tracklets for the detection of abnormal crowd behavior. In: 2015 IEEE Winter Conference on Applications of Computer Vision (WACV). IEEE; 2015. p. 148–155.
71. Hassner T, Itcher Y, Kliper-Gross O. Violent flows: Real-time detection of violent crowd behavior. In: Computer Vision and Pattern Recognition Workshops (CVPRW), 2012 IEEE Computer Society Conference on. IEEE; 2012. p. 1–6.
72. Gracia IS, Suarez OD, Garcia GB, Kim TK. Fast fight detection. PLoS one. 2015;10(4):e0120448.
73. Zhou P, Ding Q, Luo H, Hou X. Violence detection in surveillance video using low-level features. PLoS one. 2018;13(10):e0203668.
74. Gao Y, Liu H, Sun X, Wang C, Liu Y. Violence detection using oriented violent flows. Image and vision computing. 2016;48:37–41.
75. Deniz O, Serrano I, Bueno G, Kim TK. Fast violence detection in video. In: Computer Vision Theory and Applications (VISAPP), 2014 International Conference on. vol. 2. IEEE; 2014. p. 478–485.
76. Xu L, Gong C, Yang J, Wu Q, Yao L. Violent video detection based on MoSIFT feature and sparse coding. In: ICASSP; 2014. p. 3538–3542.
77. Kooij JF, Liem M, Krijnders JD, Andringa TC, Gavrilu DM. Multi-modal human aggression detection. Computer Vision and Image Understanding. 2016;144:106–120.
78. Sultani W, Chen C, Shah M. Real-world Anomaly Detection in Surveillance Videos. Center for Research in Computer Vision (CRCV), University of Central Florida (UCF). 2018.
79. Mabrouk AB, Zagrouba E. Abnormal behavior recognition for intelligent video surveillance systems: A review. Expert Systems with Applications. 2018;91:480–491.
80. Mahadevan V, Li W, Bhalodia V, Vasconcelos N. Anomaly detection in crowded scenes. In: Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on. IEEE; 2010. p. 1975–1981.
81. Lu C, Shi J, Jia J. Abnormal event detection at 150 fps in matlab. In: Proceedings of the IEEE international conference on computer vision; 2013. p. 2720–2727.
82. Blunsden S, Fisher R. The BEHAVE video dataset: ground truthed video for multi-person behavior classification. Annals of the BMVA. 2010;4(1-12):4.
83. Nievas EB, Suarez OD, Garca GB, Sukthankar R. Violence detection in video using computer vision techniques. In: International conference on Computer analysis of images and patterns. Springer; 2011. p. 332–339.
84. Soomro K, Zamir A, Shah M. UCF101-Action Recognition Data Set; 2012.
85. Kanazawa A, Black MJ, Jacobs DW, Malik J. End-to-end recovery of human shape and pose. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR); 2018.
86. Vapnik V. The nature of statistical learning theory. Springer science & business media; 2013.
87. Scholkopf B, Smola AJ. Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond. Cambridge, MA, USA: MIT Press; 2001.

88. Platt J, et al. Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. *Advances in large margin classifiers*. 1999;10(3):61–74.
89. Lin HT, Lin CJ, Weng RC. A note on Platts probabilistic outputs for support vector machines. *Machine learning*. 2007;68(3):267–276.
90. Hasan M, Choi J, Neumann J, Roy-Chowdhury AK, Davis LS. Learning Temporal Regularity in Video Sequences. In: *CVPR*; 2016.
91. Lu C, Shi J, Jia J. Abnormal Event Detection at 150 FPS in MATLAB. In: *ICCV*; 2013.