# 3D real-time human action recognition using a spline interpolation approach

Enjie Ghorbel, Rémi Boutteau, Jacques Boonaert, Xavier Savatier, Stéphane Lecoeuche

# 3D real-time human action recognition using a spline interpolation approach

Enjie Ghorbel[1][2], Rémi Boutteau[1], Jacques Boonaert[2], Xavier Savatier[1] and Stéphane Lecoeuche[2]

[1] Institut de recherche en systémes électroniques embarqués (IRSEEM-ESIGELEC)
e-mail: enjie.ghorbel@esigelec.fr, remi.boutteau@esigelec.fr, xavier.savatier@esigelec.fr
[2] Unité de recherche en informatique et automatique (URIA-Ecole des Mines de Douai)
e-mail: jacques.boonaert@mines-douai.fr, stephane.lecoeuche@mines-douai.fr

*Abstract*—**This paper presents a novel descriptor based on skeleton information provided by RGB-D videos for human action recognition. These features are obtained, considering the motion as continuous trajectories of skeleton joints. With the discrete information of skeleton joints position, a cubic-spline interpolation is applied to joints position, velocity and acceleration components. The training and classification steps are done using a linear SVM. In the literature, many human motion descriptors based on RGB-D cameras had already been proposed with good accuracy, but with a high computational time. The main interest of this proposed approach is its ability to calculate human motion descriptors with a low computation cost while such a descriptor leads to an acceptable accuracy of recognition. Thus, this approach can be adapted to human computer interaction applications. For the purpose of validation, we apply our method to the challenging benchmark MSR-Action3D and introduce a new indicator which is the ratio between accuracy and execution time per descriptor. Using this criterion, we show that our algorithm outperforms the state-of-art methods in terms of the combined information of rapidity and accuracy.**

*Keywords*—**action recognition, depth cameras, skeleton representation, cubic-spline interpolation.**

## I. INTRODUCTION

Nowadays, action recognition becomes an important field of research, due to its various applications, e.g. Video surveillance, Human Computer Interaction, etc. Generally, the action recognition pipeline can be divided into two main steps: the motion descriptors extraction and the actions classification. The accuracy of classification heavily depends on the chosen features. Indeed, the more discriminative the descriptor, the higher the accuracy. Based on classical videos (RGB videos), a wide range of human motion descriptors were proposed [1]. These descriptors are calculated using color properties.

With the recent advances on imaging, a new type of acquisition devices was introduced: Red Green Blue-Depth cameras (RGB-D cameras), e.g. Kinect. These cameras provide two types of videos: RGB videos and depth videos. A third modality can be extracted with libraries such as Kinect SDK [2] and OpenNI[3]: the human skeleton sequence. Even if the depth-based descriptors are generally more accurate than skeleton-based ones, this last representation remains a good alternative, especially when a real-time behavior is required such as Human Computer Interaction applications (HCI). The emergence of these new modalities, depth videos and skeleton

sequences, has allowed researchers to propose new more discriminative descriptors: every year, various skeleton and depth descriptors leading to a higher accuracy are introduced, but their computational cost do not seem to be a priority, so that they can not be used if real time performances are required. As stated before, in the area of Human Computer Interaction, a low-cost descriptor with an acceptable accuracy is needed.

In this paper, we propose a real-time descriptor based on skeleton joints. We calculate our descriptor using a spline interpolation of joints position, joints velocity and joints acceleration. To reduce the anthropometric variability, we perform a skeleton normalization as a pre-processing step. Subsequently, a sampling proportional to the videos length is applied in order to overcome their temporal variability. The classification step is realized using a linear SVM classifier [4]. Our principal contributions include the following three parts : first, the introduction of a new low-cost human motion descriptor and second, the benchmarking of the state-of-the-art available descriptors in terms of execution time and accuracy. Third, an adaptive sampling is proposed according to the videos length. This paper is organized in 5 sections. Section 2 represents the state-of-art related to action recognition using RGB-D cameras, followed by section 3 which describes the proposed approach. Then, Section 4 details the experimental settings and the obtained results. Finally, section 5 resumes this paper and exposes our future work.

## II. RELATED WORK

Recently, various descriptors based on RGB-D cameras have been proposed. In this part, we briefly review human action recognition approaches based on RGB-D videos. More details regarding the state-of-art can be found in [7]. Figure 1 presents an overview of the human motion descriptors proposed in the literature for RGB-cameras. In this way, We can essentially distinguish between two kinds of descriptors: depth-based descriptors and skeleton-based descriptors.

Inspired by earlier work, many studies had proposed the adaption of classical descriptors to depth videos. Following the work of Dollar et al. [6], Xia et al. [8] developed an algorithm which detects spatio-temporal interest points (STIP) on depth maps. The filter was called DSTIP. To describe the STIPs, they extended the cuboid descriptor, introduced by
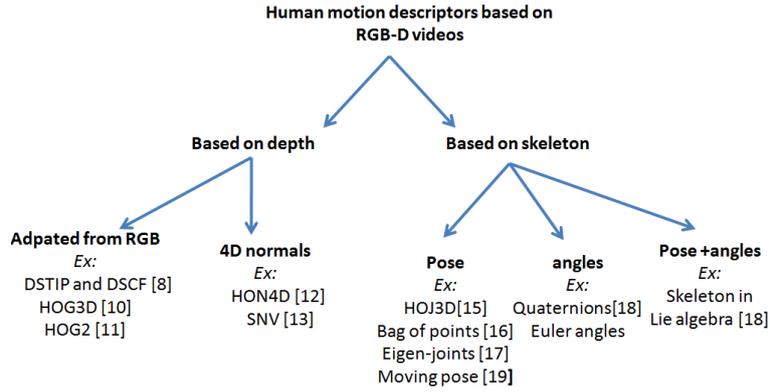
Fig. 1. An overview of the state-of-art

Dollar et al. [9], to a 4-Dimensional cuboid (Depth Cuboid Similarity Features-DCSF). Klaser et al. [10] proposed the 3D Histogram of Oriented Gradient (HOG3D) as an extension of the classical HOG. On the other hand, Ohn-Bar et al. [11] also extended the HOG to the HOG2, considering spatial and temporal histograms of gradients. Nevertheless, it has been shown in recent papers [12] that usual methods developed for RGB videos cannot work optimally in depth maps.

Thus, new approaches had been developed, considering the human body as a hyper-surface immersed in $\mathbb{R}^4$ and calculating the 4D normals to this hyper-surface. Oreifej et al. [12] proposed to build a 4-dimensional Histogram Oriented Normals (HON4D). The quantization step is performed using the Polychrons which represents the 4D extension of the 2D polygon. Using also 4D normals, Yang et al. [13] proposed to calculate Super Normal Vector which are obtained by calculating poly-normals around specific zones, and then, by learning a dictionary in order to encode poly-normals via *sparse coding*. This proposed method has shown better performance. It can be due to the local information gathered by the SNV in opposition to the HON4D where the neighboring information is lost.

The depth-based descriptors showed their high accuracy. Practically, they are more robust to occlusion and noise than joints skeleton descriptors. Despite of that, the joints are intuitively more discriminative. Many bio-mechanical studies had already used this kind of representation to model the human motion [14]. One of the first works using joints to describe human motion was proposed by Li et al. [16]. The underlying idea was to build an action graph using a 3D bag of points. Another earlier paper on action recognition using RGB-D videos represented the joints as a 3D Histogram Oriented joints (HOJ3D) [15]. Each posture is represented by an oriented Histogram and then a Hidden Markov Model (HMM) is built based on bag of postures. To reduce the orientation variability which biases the classification, many authors had chosen to use relative joints position. For example, Yand et al. [17] used Eigen-joints to describe human motion which include the temporal and spatial distance between joints.

Then, a Principal Components Analysis (PCA) is applied to the features to obtain compact and same-size vectors.

More efficient descriptors calculated from skeleton joints were recently introduced, inspired by bio-mechanical studies. Indeed, the human body can be represented as a sequence of local referential where each one is related to a segment. In [18], transformation matrices between segments expressed in $SE(3) \times SE(3) \times ... \times SE(3)$ are found, where $SE(3)$ is the special euclidean group. Each pose is considered as a point in $SE(3) \times SE(3) \times ... \times SE(3)$. The main idea is to interpolate these points in the Lie algebra $se(3) \times se(3) \times ... \times se(3)$ ,which is the tangent space of the Lie group $SE(3) \times SE(3) \times ... \times SE(3)$, to obtain trajectories and calculate the distance using Dynamic Time Warping (DTW). This type of approach is required here, because $SE(3) \times SE(3) \times ... \times SE(3)$ has not a space vector structure ,then, linear applications like interpolation cannot be applied. Zanfir et al. [19] also used a physical representation of the motion to find the appropriate features. The proposed features represent a concatenation of joints position, velocity and acceleration. An initial normalization is done to remove the anthropometric variability. An optimization step is necessary to find the weight of each term.

Inspired by the two last methods and motivated by finding a real-time and accurate descriptor, we propose to calculate the position, the velocity and the acceleration of joints and to interpolate them using a cubic spline interpolation in order to build a compact and same-size descriptor. In the literature, many descriptors based on depth and on skeleton had shown their accuracy and their robustness in the field of action and activity recognition. However, in our knowledge, none of the papers we found had reported the computational time required to build their descriptors, while this is a very important aspect, especially for Human Computer Interaction applications. In one hand, the depth-based descriptors are generally very slow to compute, due certainly to the high dimensionality of the initial data they use(depth videos). In the other hand, skeleton-based descriptors are often low-cost in terms of execution time. Sometimes optimization steps are added to make descriptors

more robust or to reduce their dimensions. But, the price to pay is that these additional steps participate to a considerable increase of computational time. It is important to notice that we do not take into account the skeleton capture latency which is almost instantaneous thanks to the developed libraries associated to the RGB-D cameras [2] [3]".

## III. PROPOSED APPROACH

In this section, we describe the proposed approach to calculate a descriptor for 3D human action recognition, designed to satisfy a trade off between accuracy and computational cost. The first sub-section explains the chosen pre-processing applied to the skeleton joints to palliate anthropometric variability, body orientation change and noise with fast processes. Then, in the second sub-section, the calculation of the first features which are position, velocity and acceleration of the joints is detailed. The last section describes the transition from the first-features to the final descriptors. The figure 2 presents an overview of the proposed approach.

### A. A pre-processing on skeleton joints

An action is made of a sequence of $N$ skeleton poses. Each pose itself is composed by $n$ joints with the knowledge of the 3D position $p_i(t)$ of each joint $i$ :

$$p_i(t) = [x_i(t), y_i(t), z_i(t)] \qquad (1)$$

with $t$ the index of the frame. The action represented by the skeleton sequence can be seen as a time series (2), with $t$ varying from 0 to $N$.

$$p(t) = [p_1(t), p_2(t), ..., p_n(t)] \qquad (2)$$

In most of the bio-mechanical studies, the origin is confused with the hip joint coordinates. The same way, we subtract the hip coordinates from joints position (3).

$$p_{\mathrm{norm}}(t) = [p_1(t) - p_{\mathrm{hip}}, p_2(t) - p_{\mathrm{hip}}, ..., p_n(t) - p_{\mathrm{hip}}] \qquad (3)$$

Care should to taken, owing to the fact that the length of the segments varies from a subject to another. For example, considering two subjects 1 and 2, with *S1(i)*, *S2(i)* respectively the size of their segments *i* and *S1(j)*, *S2(j)*, respectively the size of their segments *j*. Naturally, *S1(i)/S2(i)* differs from *S1(j)/S2(j)*.

To cope with this kind of variability, we apply a pose normalization inspired by the work of Zanfir et al.[19]. After obtaining the normalized pose, it still can happen that although the segments of different people have now the same "normalized" length, the body orientation can vary from a video to another. For this reason an alignment process is performed on the skeleton joints (both of these normalization and alignment processes are described a more detailed way below). Finally, small oscillations due to the instability of the skeleton can be observed in joints visualization. To reduce this resulting noise, we empirically apply a Gaussian filter on the normalized data with a mean equals to 3 and a standard deviation equals to 5 (values araised from experimentation).

*Data normalization*: The variation of the body morphology can affect the trajectories and consequently the classification step. Following the idea developed in [19], we normalize all the limbs. But instead of learning an average skeleton like in [19], and constraining all skeletons to have the same size as the learned model (knowing that this procedure will increase the execution time), we just normalize the limbs length in the Euclidean sense without imposing a new size. This process is done iteratively, beginning at the root (hip joint) and moving successively to neighbor segments. This succession allows the skeleton to keep its shape unchanged.

*Skeleton Alignment*: Even if we normalize the data, the joints positions can differ considerably between two instances because of the body orientation. To solve this problem, we align the data the following way: we consider that for each action the first skeleton of the sequence is in the rest state. To do that, we assume that we work in a specific scenario where the actions are already segmented and where each first skeleton is in the rest state. We then choose one of the first skeletons as a reference and we optimize the transformation matrix between the first pose in each sequence and the reference skeleton using a mean square approach. We apply the obtained transformation matrix to the rest of the sequence.

### B. Calculating the first-features

Following a bio-mechanic approach, we assume that the motion is characterized by the position, the velocity and the acceleration of the joints. Therefore, we use the joints position information to calculate the other first features composed by the velocity *V(t)* (4) and the acceleration *a(t)* (5).

$$V(t) = p_{\mathrm{norm}}(t+1) - p_{\mathrm{norm}}(t-1) \qquad (4)$$

$$a(t) = p_{\mathrm{norm}}(t+2) + p_{\mathrm{norm}}(t-2) - 2 \times p_{\mathrm{norm}}(t) \qquad (5)$$

We called them the first features because we will build our descriptors using these features but not in their original form. The problem is that the length of videos is not uniform so the vector size of the first features varies from an action sequence to another. For this reason, we choose to interpolate the first features and to sample the interpolated function to get a fix number of features dimension $k$. In our case, we choose *k=25*.In the following paragraph we give details of the proposed methodology.

### C. A Cubic Spline Interpolation approach

Considering that the human actions are continuous, we admit that the position, the velocity and the acceleration are also continuous functions varying over time. Knowing that a given action can be performed a more or less rapid way, we propose to interpolate the position, the velocity and the acceleration components of each joint *i* $x_i(t), y_i(t), z_i(t), Vx_i(t), Vy_i(t), Vz_i(t), ax_i(t), ay_i(t), az_i(t)$ using a cubic spline interpolation. The cubic spline interpolation is a
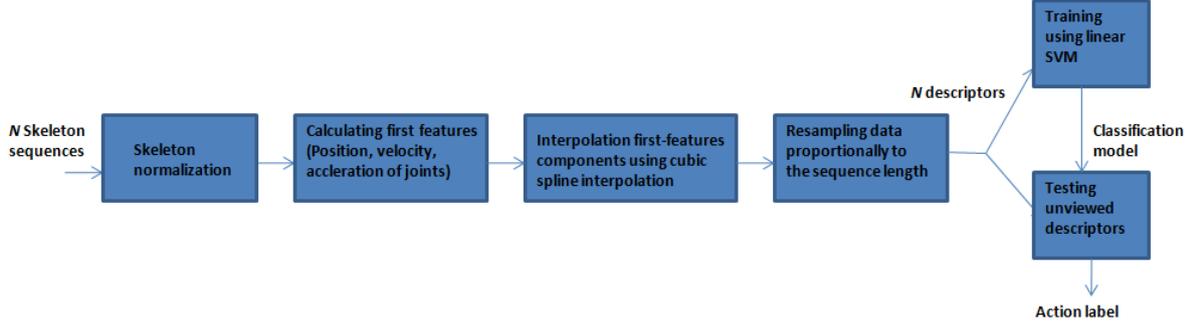
Fig. 2.   An overview of the proposed approach for human action recognition

well-known technique of interpolation which allows to connect discrete set of points coherently thanks to continuous functions and then supporting the use of a continuous variable to express the transition between the original discrete points. More precisely, this algorithm is a piecewise polynomial interpolation of the third degree (cubic). We concatenate all the interpolated function in $D_{spline,i}(t)$:

$$D_{\text{spline,i}}(t) = [spline(x_i(t)), spline(y_i(t)), spline(z_i(t)),$$
$$spline(Vx_i(t)), spline(Vy_i(t)), spline(Vz_i(t)), \quad (6)$$
$$spline(ax_i(t)), spline(ay_i(t)), spline(az_i(t))]$$

$$D_{\text{spline}}(t) = \cup_{i=1..n} D_{\text{spline,i}}(t) \quad (7)$$

After that, each term $D_{spline,i}$ is sampled with a step proportional to the video length, so that we always get descriptors having the same size equals to $n$, whatever the original length is. Such a process is mandatory, because the subsequent classification step we apply requires all the motion descriptors having the same length. The final descriptor $D_{final}$ is then obtained as expressed by equations (8) and (9).

$$D_{\text{final,i}} = \cup_{j=1..k} D_{\text{spline,i}}(j \times lengthvideo(i)/n) \quad (8)$$

$$D_{\text{final}}(t) = \cup_{i=1..n} D_{\text{final,i}}(t) \quad (9)$$

Because we are focusing on the improvement introduced by the descriptor we propose, actions classification is performed using an "out of the shelf" linear SVM classifier provided by the *SVMlibrary* in matlab [20].

## IV. EXPERIMENTS AND RESULTS

In this section, we evaluate the proposed skeleton-based descriptors using the challenging benchmark MSR-Action3D. MSR-Action 3D is a dataset captured by a depth-camera. This benchmark contains 20 human actions which are *high*

| AS1 | AS2 | AS3 |
|---|---|---|
| Horizontal arm wave | High arm wave | High throw |
| Hammer | Hand catch | Forward kick |
| Forward punch | Draw x | Side kick |
| High throw | Draw tick | Jogging |
| Hand clap | Draw circle | Tennis swing |
| Bend | Two arm wave | Tennis serve |
| Tennis serve | Forward kick | Golf swing pick up and thorw |
| Pick up and throw | Side boxing | |

Table 1. Subsets of MSR Action 3D dataset

*arm wave, horizontal arm wave, hammer, hand catch, forward punch, high throw, draw x, draw tick, draw circle, hand clap, two hand wave, side-boxing, bend, forward kick, side kick, jogging, tennis swing, tennis serve, golf swing, pickup and throw* . These actions are performed by 10 subjects two or three times for each. In total, MSR-Action 3D is composed by 567 actions. Like in [8], we use only 557 actions of them and exclude 10 very biased actions. The dataset provides the skeleton joints positions and the depth images sequences. These actions were chosen to interact with game consoles (Human Machine Interaction application). The challenging aspect of this dataset is partly due to the fact that it contains different actions that look very similar.

### A. Evaluation settings and parameters

To evaluate the recognition accuracy, we use the same cross tests subjects splitting as the one proposed in [16].In this case, the actions performed by subjects 1,3,5,7,9 are used to train the classifier, while the rest of the actions is used to test the classification model. Following the procedure used in [13], we divide the dataset into 3 subsets: AS1, AS2, AS3 according to table 1. The two first subsets groups similar actions, while AS3 groups complex actions together. The training and testing steps are done in each subset separately. Then, the global accuracy is calculated as the mean of the subsets accuracies. As explained earlier, accuracy alone is not a sufficient criterion

| Modality | descriptor | AS1 | AS2 | AS3 | Overall | Time execution per descriptor (s) | Ratio (accuracy/time execution) |
|---|---|---|---|---|---|---|---|
| Depth | HOG2 [10] | 90,47% | 84,82% | 98,20% | 91,16% | 6,44 | 14,15 |
| | HON4D [11] | 94,28% | 91,71% | 98,20% | 64,47% | 27,33 | 3,09 |
| | SNV [12] | 95,28% | 94,69% | 96,43% | **95,46%** | 146,57 | 0,65 |
| Skeleton | JP [18] | 82,86% | 68,75% | 83,73% | 78,44% | 0,58 | 135,44 |
| | RJP [18] | 81,90% | 71,43% | 88,29% | 80,53% | 2,15 | 37,44 |
| | Q [18] | 66,67% | 59,82% | 77,48% | 67,99% | 1,33 | 51,12 |
| | LARP [18] | 83,81% | 84,82% | 92,79% | **87,14%** | 17,61 | 4,94 |
| | **Ours (spline curves kinematics)** | 83,08% | 79,46% | 93,69% | 85,41% | **0,26** | **328,5** |

Table 2. Recognition performance and time execution per descriptor for different features on the MSR Action 3D dataset

when dealing with real time. That's why in this work, we also calculate the mean execution time per descriptor. To make it clear, a descriptor can have a very good ability to recognize actions but can be in the same time very slow to compute. This phenomenon is unsuitable for Human Machine Interaction applications. To include both information of rapidity and accuracy, we calculate a new indicator which represents the ratio of accuracy over the execution time per descriptor. To compare our proposed approach with the recent state-of-art methods and to overcome the problem of parameters evaluation variation, we collect available codes and run them in the same conditions with the same settings and parameters tuning. The tested descriptors can be divided into two groups. The first group represents the depth-based descriptors and includes HOG2[11], SNV [13] and HON4D[12], while the second group gathers skeleton-based descriptors using DTW to re-parametrize the features, available in the code proposed in [18]. This last group contains the following descriptors: absolute joints position (JP), Relative joints position (RJP), quaternion (Q) and Lie Algebra relative Positions (LARP).

*B. Results*

The accuracy and the mean time execution per descriptor (including prior and our proposed method) as well as their ratios are reported in Table 2. The results confirm the assumption formulated earlier, that is depth descriptors are generally more accurate but slower to compute. We can see clearly that our descriptor outperforms the other descriptors in terms of rapidity with an execution time per descriptor equals to 0.26s. Comparing with the skeleton representation, the LARP is the only descriptor that is more accurate than our, although the difference remains light. In contrast, the difference in time execution between the two descriptors cannot be neglected (17.41s against 0.26s). The ratio between accuracy and execution times allows to include both information of accuracy and rapidity
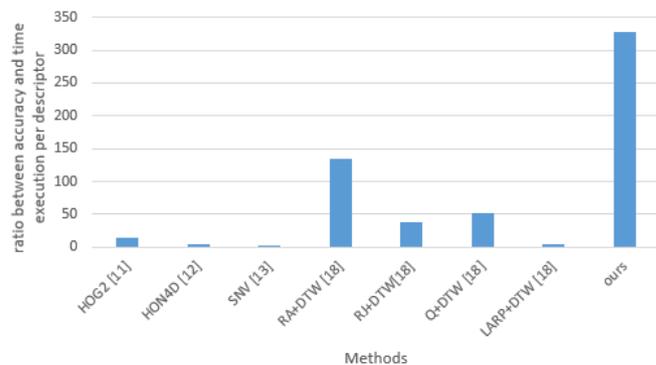


Fig. 3. The ratio between accuracy and time execution per descriptor on MSR Action 3D

and gives us a more discriminative way of comparison. Table 2 and figure 3 show that our descriptor is the more adapted for real-time actions recognition applications, because it exhibits the highest score (ratio between accuracy and execution time). Nevertheless, the proposed ratio is not always relevant to the global performance of the motion descriptor. Indeed, if the time execution is very low, even a classifier with 5 percent of accuracy can outperforms the state-of-art descriptors which is not conform to the reality. For this reason, we propose to use this criterion as a reference only when the descriptor exceeds a reasonable accuracy (superior to 80 percent). However, it is important to precise that the skeletons in MSR Action 3D are relatively stable and do not face to occlusions.

## V. CONCLUSION AND FUTURE WORK

In this paper, we proposed a real-time human motion descriptor based on skeleton provided by RGB-D cameras, such as Kinect. The proposed features represent a concatenation of the interpolated joints position, joints velocity and joints

acceleration. We began with a skeleton normalization to reduce the orientation and anthropometric variability. Then, joints velocity and acceleration were calculated using the joints position information. A 1D cubic spline interpolation was done on the first-features components. The recovered final descriptor was composed by sampled values on the obtained curves after interpolation. We showed the effectiveness of our method for real-time applications with respect to state-of-the-art approaches by introducing a criterion combining the information of execution time per descriptor and accuracy. In our work, we interpolated each 1D signal even if the action descriptor must be normally invariant to the 3D euclidean transformation which is not the case of features components. In this way, we are planning to test in a future work a 3D interpolation to overcome the limits of our method without increasing the calculation cost. On another side, a more detailed experimentation must be done to explain the role of each kinematic measure and to analyze the effect of pre-processes like normalization and skeleton alignment. The robustness of the proposed descriptor must be also tested with the use of different databases captured in a different context. For example, the kinematic spline curves (the proposed descriptor in this paper) can be sensitive to skeleton occlusions or to noise. In fact, the online case is unpredictable and our descriptor must be tested in different situations of skeleton capture. Finally, the execution time per descriptor gave us an idea about the behavior of each descriptor on a real-time application. However, this information could be improved by a more formal study of complexity, that will be a part of our future work.

## REFERENCES

[1] Poppe, R. A survey on vision-based human action recognition. *Image and vision computing*, 28(6), 976-990, 2010.
[2] Microsoft Kinect for Windows. *http://www.microsoft.com/en-us/kinectforwindows/a*
[3] Open Natural Interaction SDK. *http://openni.org/docs2/Reference/smpl user tracker.html*
[4] C. Cortes. and V. Vapnik. Support-vector networks. In *Machine learning*, 20(3), 273-297, 1995.
[5] M.E. Yuksel. A hybrid neuro-fuzzy filter for edge preserving restoration of images corrupted by impulse noise. In *IEEE Trans. on Image Processing*, 15(4): 928–936, 2006.
[6] P. Dollar., V. Rabaud., G. Cottrel. and S. Belongie. Monitoring animal behavior in the smart vivarium. In *Measuring Behavior*,2005.
[7] J.K. Aggarwal. and M.S. Ryoo. Human activity analysis: A review. In *ACM Computing Surveys*, 43(3):16:1-16:43,2011.
[8] L. Xia. and J. Aggarwal. Spatio-temporal depth cuboid similarity Feature for activity recognition using depth camera. In *CVPR*,2012.
[9] P. Dollar., V. Rabaud., G. Cottrel. and S. Belongie. Behavior recognition via sparse spatio-temporal features. In *PETS*,65-72,2005.
[10] A. Klaser., M. Marszelak. and C. Schmid. A spatio-temporal descriptor based on 3D-gradients. In *BMVC*,2010.
[11] E. Ohn-Bar. and M.M. Trivedi. Joint angles similarities and HOG2 for action recognition. In *CVPRW*, 465-470, 2013.
[12] O. Oreifej. and Z. Liu. Histogram of oriented 4D normals for activity recognition from depth sequences. In *CVPR*, 2013.
[13] X. Yang. and Y. Tian. Super normal vector for activity recognition using depth sequences. In *CVPR*,804-811, 2014.
[14] G. Johansson. Visual perception of biological motion and a model for its analysis. In *Perception and Psychophysics*,201-211, 1973.
[15] L. Xia., C.C. Chan and J.K. Aggarwal. View invariant human action recognition using histograms of 3d joints. In *CVPRW*,20-27, 2012.
[16] W. Li, Z. Zhang and Z. Liu. Action recognition based on a bag of 3d points. In *CVPRW*,9-14, 2010.
[17] X. Yang and Y. Tian. Eigen-joints based action recognition using naive-Bayes-Nearest-Neighbor. In *CVPRW*,14-19, 2012.
[18] R. Vemulapalli, F. Arrate. and R. Chellappa. Human action recognition by representing 3d skeletons as points in a .lie group In *CVPR*,588-595, 2014.
[19] M. Zanfir., M. Leordeanu. and C. Sminchisescu. The moving pose: An efficient 3d kinematics descriptor for low-latency action recognition and detection In *ICCV*,2752-2759, 2013.
[20] C. Chang. and C. Lin. Libsvm: a library for support vector machines In *TIST*,2(3):27, 2011.