

ADHD indicators modelling based on Dynamic Time Warping from RGBD data: A feasibility study

Antonio Hernández-Vela^{*†}, Miguel Reyes^{*†}, Laura Igual^{†*}, Josep Moya[‡], Verónica Violant[§] and Sergio Escalera^{†*}

^{*}Computer Vision Center, Edifici O, Campus UAB, 08193 Bellaterra, Spain

Email: {ahernandez,mreyes}@cvc.uab.cat

[†]Departament de Matemàtica Aplicada i Anàlisi

Universitat de Barcelona, Gran Via de les Corts Catalanes 585, 08007 Barcelona, Spain

Email: {ligual,sergio}@maia.ub.es

[‡]Observatori de Salut Mental de Catalunya, Parc Taulí 1, 08208 Sabadell, Spain

Email: jmoya@tauli.cat

[§]Campus Mundet, Edifici Llevant 2nd, Passeig de la Vall d'Hebron 171, 08035 Barcelona, Spain

Email: vviolant@ub.edu

Abstract—In this paper, we present a feasibility study for the automatic modelling of visual communicative indicators in video sequences of children with attention deficit hyperactivity disorder (ADHD). Our methodology is based on RGBD sequences (RGB + Depth), recorded with the recent Microsoft's Kinect sensor. More specifically, a feature vector is extracted from a previously fitted human skeleton model in the RGBD sequence, and compared to a training set using Dynamic Time Warping (DTW). Finally, some qualitative results are presented, concluding that the presented methodology is feasible for the modelling of ADHD visual indicators.

I. INTRODUCTION

Attention deficit disorder –with or without hyperactivity– is one of the main reasons of consultation in mental health centers for children and adolescents. The basic characteristics of ADHD are excessive and harmful levels of activity, inattention, and impulsiveness. Currently, the diagnosis is made through direct observation of patients for long periods of time, and it is often not feasible in practise; hence, we propose an automation in order to help doctors diagnose the disorder.

From the point of view of data acquisition, many methodologies treat images captured by visible-light cameras. Computer Vision is then used to detect, describe, and learn visual features of the human body [1], [2]. On the other hand, depth information is invariant to color, texture and lighting, making it easier to distinguish between the background and the foreground object. Nowadays, several works have been published related to this topic because of the emergence of inexpensive structured light technology, which is reliable and robust to capture depth information along with the corresponding synchronized RGB image. This technology has been developed by the PrimeSense [3] company and marketed by Microsoft XBox under the name of Kinect. Using this sensor, Shotton et al. [4] present one of the greatest advances in the extraction of the human body pose from depth images, representing the body as a skeletal form comprised by a set of joints. Moreover, some works have started a general behaviour study based on RGBD data [5].

Our proposal is to perform a temporal analysis of the al data given by [4], applied to some captured RGBD sequences of children with ADHD. More specifically, this temporal analysis is performed using DTW in order to segment and model visual communicative indicators which may be potentially useful for the diagnosis of ADHD.

The rest of the paper is organized as follows: Section 2 presents our methodology, section 3 shows some preliminary results we obtained, and finally section 4 concludes the paper.

II. METHODOLOGY

In this section, we first describe the feature vector extraction step based on the al model returned by the method of [4]. Secondly, we briefly introduce the DTW framework used for the temporal analysis, and finally, we apply the system to perform begin-end of gesture detection in large video data.

A. Feature Vector Extraction

The articulated human model is defined by the set of 15 reference points, and has the advantage of being highly deformable, and thus, able to fit to complex human poses. In order to normalize the descriptor, we use the neck joint of the model as the origin of coordinates (OC) and compute the rest of points relative to. Thus, the final feature vector \mathbf{V}_j at frame j that defines the human pose is described by 42 elements (14 joints \times three spatial coordinates),

$$\mathbf{V}_j = \{\{v_{j,x}^1, v_{j,y}^1, v_{j,z}^1\}, \dots, \{v_{j,x}^{14}, v_{j,y}^{14}, v_{j,z}^{14}\}\}$$

B. Dynamic Time Warping

The DTW algorithm [6] was defined to match temporal distortions between two models, finding an alignment warping path between the two time series $Q = \{q_1, \dots, q_n\}$ and $C = \{c_1, \dots, c_m\}$. In order to align these two sequences, a $M_{m \times n}$ matrix is designed, where the position (i, j) of the matrix contains the distance between c_i and q_j . The Euclidean distance is the most frequently applied. Then, a warping path,

$$W = \{w_1, \dots, w_T\}, \max(m, n) \leq T < m + n + 1$$

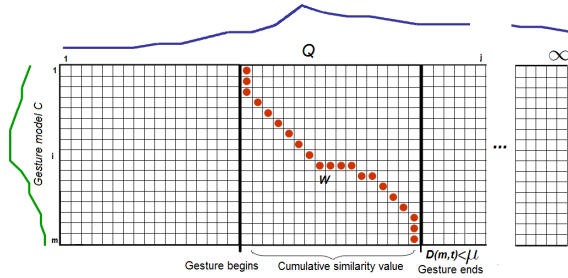


Figure 1. Begin-end of gesture recognition of a model C in an infinite sequence Q .

is defined as a set of "contiguous" matrix elements that defines a mapping between C and Q .

We are generally interested in the final warping path that satisfying these conditions minimizes the warping cost,

$$DTW(Q, C) = \min \left\{ \frac{1}{T} \sqrt{\sum_{t=1}^T w_t} \right\}, \quad (1)$$

where T compensates the different lengths of the warping paths. This path can be found very efficiently using dynamic programming to evaluate the following recurrence, which defines the cumulative distance $\gamma(i, j)$ as the distance $d(i, j)$ found in the current cell and the minimum of the cumulative distance of the adjacent elements,

$$\gamma(i, j) = d(i, j) + \min\{\gamma(i-1, j-1), \gamma(i-1, j), \gamma(i, j-1)\}. \quad (2)$$

Given the nature of our system to work in uncontrolled environments, we continuously review the stage for possible actions or gestures. In this case, our input feature vector Q is of "infinite" length, and may contain segments related to gesture C at any part.

C. Begin-end of gesture detection

In order to detect a begin-end of gesture $C = \{c_1, \dots, c_m\}$ in a possible infinite sequence $Q = \{q_1, \dots, q_\infty\}$, a $M_m \times \infty$ matrix is designed, where the position (i, j) of the matrix contains the distance between c_i and q_j , quantifying its value by the Euclidean distance, as commented before. Finally, our warping path is defined by $W = \{w_1, \dots, w_\infty\}$ as in the standard DTW approach. Our aim is focused on finding segments of Q sufficiently similar to the sequence C . The system considers that there is correspondence between the current block k in Q and a gesture if satisfying the following condition,

$$M(m, k) < \mu, k \in [1, \dots, \infty]$$

for a given cost threshold μ . This threshold value is estimated in advance for each of the categories of actions or gestures using leave-one-out cross-validation strategy.

Once a possible end of pattern of gesture or action is detected, the working path W can be found through backtracking of the minimum path from $M(m, k)$ to $M(0, z)$, being z the instant of time in Q where the gesture begins. Note that, in this case, $d(i, j)$ is the cost function which measures the difference among our descriptors V_i and V_j . An example of a begin-end

gesture recognition for a model and infinite sequence together with the working path estimation is shown in Figure 1.

III. PRELIMINARY RESULTS

In order to validate our proposal, we applied the method on five video sequences of one hour each one at 24 FPS, recorded with the Kinect device. Two different scenarios have been considered. In both of them there are three children between 8-11 years –half of them with ADHD diagnosis– in a classroom, but in the first one they are doing math exercises while in the second one they are playing videogames. Given this scenario, we defined some specific gestures and trained the system with them. Figure 2 shows a sequence where our method successfully detects the beginning and end for the gesture "lower head", which is one of the visual communicative indicators of ADHD related to the focus of attention and subject agitation.

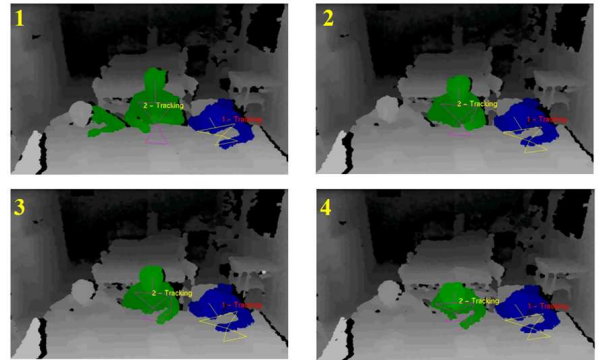


Figure 2. Sequence where the gesture "lower head" is found. Frames 1 and 4 show the beginning and end of the gesture, respectively. Frames 2 and 3 are intermediate frames.

IV. CONCLUSION

In this paper, we have presented a methodology for the detection of the beginning and end of gestures based human skeleton points described using RGBD representation from Kinect device. First results indicate that the presented methodology can successfully recognize a set of defined gestures related to ADHD indicators, showing the viability of the system to be considered for diagnostic support.

REFERENCES

- [1] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," *CVPR*, vol. 1, pp. 886–893, 2005.
- [2] E. N. Mortensen, H. Deng, and L. Shapiro, "A sift descriptor with global context," *CVPR*, vol. 1, pp. 184–190 vol. 1, 2005.
- [3] *Prime Sensor NITE 1.3 Algorithms notes*, PrimeSense Inc., 2010, last viewed 14-07-2011 13:19. [Online]. Available: <http://www.primesense.com>
- [4] J. Shotton, A. Fitzgibbon, M. Cook, T. Sharp, M. Finocchio, R. Moore, A. Kipman, and A. Blake, "Real-time human pose recognition in parts from single depth images," *CVPR*, 2011.
- [5] M. Reyes, G. Domnguez, and S. Escalera, "Feature weighting in dynamic time warping for gesture recognition in depth data," *HICV workshop, ICCV*, 2011.
- [6] M. Parizeau and R. Plamondon, "A comparative analysis of regional correlation, dynamic time warping, and skeletal tree matching for signature verification."