

Probability-based Dynamic Time Warping and Bag-of-Visual-and-Depth-Words for Human Gesture Recognition in RGB-D

Antonio Hernández-Vela^{a,b}, Miguel Ángel Bautista^{a,b}, Xavier Perez-Sala^{b,c,d}, Víctor Ponce-López^{a,b,e}, Sergio Escalera^{a,b}, Xavier Baró^{b,e}, Oriol Pujol^{a,b}, Cecilio Angulo^d

^a*Dept. MAIA, Universitat de Barcelona, Gran Via 585, 08007 Barcelona, Spain.*

^b*Computer Vision Center, Campus UAB, Edifici O, 08193 Bellaterra, Barcelona, Spain.*

^c*Fundació Privada Sant Antoni Abat, Rambla de l'Exposició, 59-69, 08800 Vilanova i la Geltrú, Spain.*

^d*UPC - BarcelonaTECH, Av. Víctor Balaguer 1, 08800 Vilanova i la Geltrú, Spain.*

^e*EIMT/IN3, Universitat Oberta de Catalunya, Rbla. del Poblenou 156, 08018 Barcelona, Spain.*

Abstract

We present a methodology to address the problem of human gesture segmentation and recognition in video and depth image sequences. A Bag-of-Visual-and-Depth-Words (BoVDW) model is introduced as an extension of the Bag-of-Visual-Words (BoVW) model. State-of-the-art RGB and depth features, including a newly proposed depth descriptor, are analysed and combined in a late fusion form. The method is integrated in a Human Gesture Recognition pipeline, together with a novel Probability-based Dynamic Time Warping (PDTW) algorithm which is used to perform prior segmentation of idle gestures. The proposed DTW variant uses samples of the same gesture category to build a Gaussian Mixture Model driven probabilistic model of that gesture class. Results of the whole Human Gesture Recognition pipeline in a public data set show better performance in comparison to both standard

BoVW model and DTW approach.

Keywords: RGB-D, Bag-of-Words, Dynamic Time Warping, Human
Gesture Recognition

1. Introduction

Nowadays, human gesture recognition is one of the most challenging tasks in computer vision. Current methodologies have shown preliminary results on very simple scenarios, but they are still far from human performance. Due to the large number of potential applications involving human gesture recognition in fields like surveillance [1], sign language recognition [2], or clinical assistance [3] among others, there is a large and active research community devoted to deal with this problem. Independently of the application field, the usual human gesture recognition pipeline is mainly formed by two steps: *gesture representation* and *gesture classification*.

Regarding the gesture representation step, literature shows a variety of methods that have obtained successful results. Commonly applied in image retrieval or image classification scenarios, *Bag-of-Visual-Words* (BoVW) is one of the most used approaches. This methodology is an evolution of *Bag-of-Words* (BoW) [4] representation, used in document analysis, where each document is represented using the frequency of appearance of each word in a dictionary. In the image domain, these words become visual elements of a certain visual vocabulary. First, each image is decomposed into a large set of patches, either using some type of spatial sampling (grids, sliding window, etc.) or detecting points with relevant properties (corners, salient regions, etc.). Each patch is then described obtaining a numeric descriptor. A set of

22 V representative visual words are selected by means of a clustering process
23 over the descriptors. Once the visual vocabulary is defined, each new image
24 can be represented by a global histogram containing the frequencies of visual
25 words. Finally, this histogram can be used as input for any classification
26 technique (i.e. k -Nearest Neighbor or SVM) [5, 6]. In addition, extensions
27 of BoW from still images to image sequences have been recently proposed in
28 the context of human action recognition, defining Spatio-Temporal-Visual-
29 Words (STVW) [7].

30 The release of the Microsoft KinectTM sensor in late 2010 has allowed
31 an easy and inexpensive access to almost synchronized range imaging with
32 standard video data. Those data combine both sources into what is com-
33 monly named RGB-D images (RGB plus Depth). This data fusion has re-
34 duced the burden of the first steps in many pipelines devoted to image or
35 object segmentation, and opened new questions such as how these data can
36 be effectively described and fused. Motivated by the information provided
37 by depth maps, several 3-D descriptors have been recently developed [8, 9]
38 (most of them based on codifying the distribution of normal vectors among
39 regions in the 3D space), as well as their fusion with RGB data [10] and
40 learning approaches for object recognition [11]. This depth information has
41 been particularly exploited for gesture recognition and human body segmen-
42 tation and tracking. While some works focus on just the hand regions for
43 performing gesture recognition [12, 13, 14, 15, 16, 17], in [18] Shotton intro-
44 duced one of the greatest advances in the extraction of the human body pose
45 using RGB-D, which is provided as part of the KinectTM human recognition
46 framework. The method is based on inferring pixel label probabilities through

47 Random Forest from learned offsets of depth features. Then, mean shift is
48 applied to estimate human joints and representing the body in skeletal form.
49 Hernández-Vela et al. [19] extended Shotton’s work applying Graph-cuts to
50 the pixel label probabilities obtained through Random Forest, in order to
51 compute consistent segmentations in the spatio-temporal domain. Girshick
52 and Shotton [20] proposed later a different approach in which they directly
53 regress the positions of the body joints, without the need of an intermediate
54 pixel-wise body limb classification as in [18]. The extraction of body pose in-
55 formation opens the door to one of the most challenging problems nowadays,
56 i.e. human gesture recognition.

57 In the gesture classification step there exists a wide number of methods
58 based on dynamic programming algorithms for both alignment and clustering
59 of temporal series [21]. Other probabilistic methods such as Hidden Markov
60 Models (HMM) or Conditional Random Fields (CRF) have been commonly
61 used in the literature [2]. Nevertheless, one of the most common methods for
62 Human Gesture Recognition is Dynamic Time Warping (DTW) [22], since it
63 offers a simple yet effective temporal alignment between sequences of differ-
64 ent lengths. However, the application of such methods to gesture detection in
65 complex scenarios becomes a hard task due to the high variability of the envi-
66 ronmental conditions among different domains. Some common problems are:
67 wide range of human pose configurations, influence of background, continu-
68 ity of human movements, spontaneity of human actions, speed, appearance
69 of unexpected objects, illumination changes, partial occlusions, or different
70 points of view, just to mention a few. These effects can cause dramatic
71 changes in the description of a certain gesture, generating a great intra-class

72 variability. In this sense, since usual DTW is applied between a sequence
73 and a single pattern, it fails when taking into account such variability.

74 The problem of gesture recognition in which an idle or reference ges-
75 ture is performed between gestures is addressed in this paper. In order to
76 solve this problem, we introduce a continuous human gesture recognition
77 pipeline based on: First, a new feature representation by means of a Bag-
78 of-Visual-and-Depth-Words (BoVDW) approach that takes profit of multi-
79 modal RGB-D data to tackle the gesture representation step. The BoVDW
80 is empowered by the combination of both RGB images and a new depth
81 descriptor which takes into account the distribution of normal vectors with
82 respect to the camera position, as well as the rotation with respect to the
83 roll axis of the camera. Next, we propose the definition of an extension of
84 DTW method to a probability-based framework in order to perform temporal
85 gesture segmentation. In order to evaluate the presented approach, we com-
86 pare the performances achieved with state-of-the-art RGB and depth feature
87 descriptors separately, and combine them in a late fusion form. All these
88 experiments are performed in the proposed framework using the public data
89 set provided by the ChaLearn Gesture Challenge¹. Results of the proposed
90 BoVDW method show better performance using late fusion in comparison to
91 early fusion and standard BoVW model. Moreover, our BoVDW approach
92 outperforms the baseline methodology provided by the ChaLearn Gesture
93 Recognition Challenge 2012. In the same way, the results obtained with the
94 proposed PDTW outperform the ones from the classical DTW approach.

¹<http://gesture.chalearn.org/>

95 The BoVDW model for gesture recognition is introduced in Section 2, as
96 well as the PDTW. Experimental results and their analysis are presented in
97 Section 3. Finally, Section 4 concludes the paper.

98 **2. BoVDW and Probability-based DTW for Human Gesture Recog-** 99 **niton**

100 As pointed out in the Introduction, we address the problem of gesture
101 recognition, with the constraint that an idle or reference gesture is performed
102 between gestures. The main reason for such constraint is that in many real-
103 world settings there always exists an idle gesture between movements rather
104 than a continuous flux of gestures. Some examples are sports like tennis,
105 swordplay, boxing, martial arts, or choreographic sports. However, the exist-
106 tence of an idle gesture is not only related to sports, some other daily tasks
107 like cooking or dancing contain idle gestures in certain situations. Moreover,
108 the proposed system can be extended to be applied to other gesture recogni-
109 tion domains without the need of modelling idle gestures, but any other kind
110 of gesture categories.

111 In this sense, our approach consists of two steps: *a temporal gesture*
112 *segmentation* step (the detection of the idle gesture), and *the gesture clas-*
113 *sification* step. The former one aims to provide a temporal segmentation
114 of gestures. To perform such temporal segmentation, a novel probabilistic-
115 based DTW models the variability of the idle gesture by learning a GMM
116 on the features of the idle gesture category. Once the gestures have been
117 segmented, the latter step is gesture classification. Segmented gestures are
118 represented and classified by means of a BoVDW method, which integrates

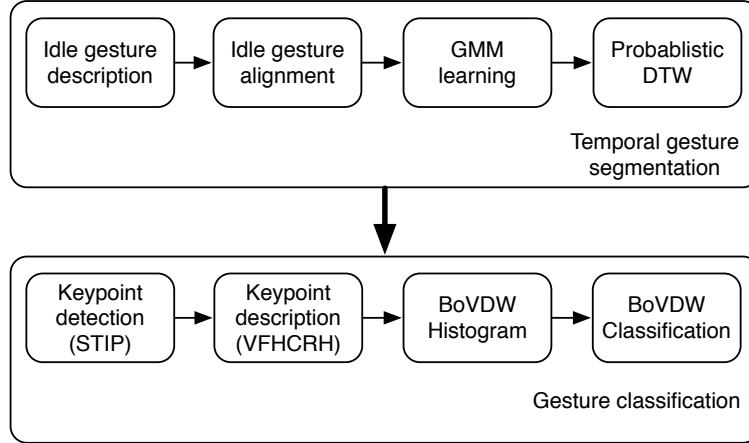


Figure 1: General pipeline of the proposed approach.

119 in a late fusion form the information of both RGB and Depth images.

120 The global pipeline of the approach is depicted in Figure 1. The proposal
 121 is divided in two blocks, the temporal gesture segmentation step and the
 122 gesture classification step, which are detailed in next sections.

123 2.1. Gesture Segmentation: Probability-based DTW

124 The original DTW is introduced in this section, as well as its common
 125 extension to detect a certain sequence given an indefinite data stream. In
 126 the following subsections, DTW is extended in order to align patterns taking
 127 into account the probability density function (PDF) of each element of the
 128 sequence by means of a Gaussian Mixture Model (GMM). A flowchart of the
 129 whole methodology is shown in Figure 2.

130 2.1.1. Dynamic Time Warping

131 The original DTW algorithm was defined to match temporal distortions
 132 between two models, finding an alignment/warping path between two time

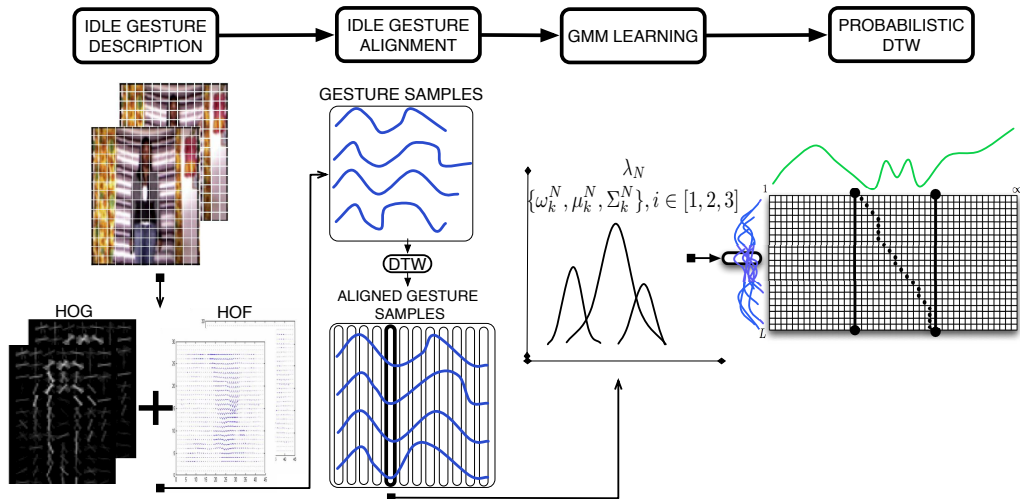


Figure 2: Flowchart of the Probabilistic DTW gesture segmentation methodology.

133 series: an input model $Q = \{q_1, \dots, q_n\}$ and a certain sequence $C = \{c_1, \dots, c_m\}$.
 134 In our particular case, the time series Q and C are video sequences, where
 135 each q_j and c_i will be feature vectors describing the j -th and i -th frame
 136 respectively. In this sense, Q will be an input video sequence and C will be
 137 the gesture we are aiming to detect. Generally, in order to align these two
 138 sequences, a $M_{m \times n}$ matrix is designed, where position (i, j) of the matrix
 139 contains the alignment cost between c_i and q_j . Then, a warping path of
 140 length τ is defined as a set of contiguous matrix elements, defining a mapping
 141 between C and Q : $W = \{w_1, \dots, w_\tau\}$, where w_i indexes a position in the cost
 142 matrix M . This warping path is typically subject to several constraints,
 143 *Boundary conditions:* $w_1 = (1, 1)$ and $w_\tau = (m, n)$.
 144 *Continuity and monotonicity:* Given $w_{\tau'-1} = (a', b')$, $w_\tau = (a, b)$, then
 145 $a - a' \leq 1$ and $b - b' \leq 1$. This condition forces the points in the cost matrix
 146 with the warping path W to be monotonically spaced in time.

147 Interest is focused on the final warping path that, satisfying these condi-
 148 tions, minimizes the warping cost,

$$DTW(M) = \min_W \left\{ \frac{M(w_\tau)}{\tau} \right\}, \quad (1)$$

149 where τ compensates the different lengths of the warping paths at each time
 150 t . This path can be found very efficiently using dynamic programming. The
 151 cost at a certain position $M(i, j)$ can be found as the composition of the
 152 Euclidean distance $d(i, j)$ between the feature vectors c_i and q_j of the two
 153 time series, and the minimum cost of the adjacent elements of the cost matrix
 154 up to that position, as,

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

155 However, given the streaming nature of our problem, the input video
 156 sequence Q has no definite length (it may be an infinite video sequence) and
 157 may contain several occurrences of the gesture sequence C . In this sense,
 158 the system considers that there is correspondence between the current block
 159 k in Q and the gesture when the following condition is satisfied, $M(m, k) <$
 160 $\theta, k \in [1, \dots, \infty]$ for a given cost threshold θ . At this point, if $M(m, k) < \theta$ k
 161 is consider a possible end of a gesture sequence C .

162 Once detected a possible end of the gesture sequence, the warping path W
 163 can be found through backtracking the minimum cost path from $M(m, k)$ to
 164 $M(0, g)$, being g the instant of time in Q where the detected gesture begins.
 165 Note that $d(i, j)$ is the cost function which measures the difference among
 166 descriptors c_i and q_j , which in standard DTW is defined as the euclidean
 167 distance between c_i and q_j . An example of a begin-end gesture recognition

168 together with the warping path estimation is shown in Figure 2 (last 2 steps:
 169 GMM learning and Probabilistic DTW).

170 2.1.2. Handling variance with Probability-based DTW

171 Consider a training set of N sequences, $S = \{S_1, S_2, \dots, S_N\}$, that is, N
 172 gesture samples belonging to the same gesture category. Then, each sequence
 173 $S_g = \{s_1^g, \dots, s_{L_g}^g\}$, (each gesture sample) is composed by a feature vector ²
 174 for each frame t , denoted as s_t^g , where L_g is the length in frames of sequence
 175 S_g . In order to avoid temporal deformations of the gesture samples in S ,
 176 all sequences are aligned with the median length sequence using the classical
 177 DTW with Euclidean distance. Let us assume that sequences are ordered
 178 according to their length, so that $L_{g-1} \leq L_g \leq L_{g+1}, \forall g \in [2, \dots, N-1]$, then,
 179 the median length sequence is $\bar{S} = S_{\lceil \frac{N}{2} \rceil}$.

180 It is worth noting that this alignment step by using DTW has no relation
 181 to the actual gesture recognition, as it is consider a pre-processing step to ob-
 182 tain a set of gesture samples with few temporal deformations and a matching
 183 length.

184 Finally, after this alignment process, all sequences have length $L_{\lceil \frac{N}{2} \rceil}$. The
 185 set of warped sequences is defined as $\tilde{S} = \{\tilde{S}_1, \tilde{S}_2, \dots, \tilde{S}_N\}$ (See Figure 3(b)).
 186 Once all samples are aligned, the N feature vectors corresponding to each
 187 sequence element at a certain frame t , denoted as $\tilde{F}_t = \{f_t^1, f_t^2, \dots, f_t^N\}$, are
 188 modelled by means of a G -component Gaussian Mixture Model (GMM)
 189 $\lambda_t = \{\alpha_k^t, \mu_k^t, \Sigma_k^t\}$, $k = 1, \dots, G$, where α_k^t is the mixing value, and μ_k^t and
 190 Σ_k^t are the parameters of each of the G Gaussian models in the mixture. As

²HOG/HOF descriptors in our particular case, see Sec. 3.2.1 for further details.

191 a result, each one of the GMMs that model each \tilde{F}_t is defined as follows,

$$p(\tilde{F}_t) = \sum_{k=1}^G \alpha_k^t \cdot e^{-\frac{1}{2}(x-\mu_k^t)^T \cdot (\Sigma_k^t)^{-1} \cdot (x-\mu_k^t)}. \quad (3)$$

192 The resulting model is composed by the set of GMMs that model each
 193 set \tilde{F}_t among all warped sequences of a certain gesture class. An example of
 194 the process is shown in Figure 3.

195 2.1.3. Distance measures

196 In the classical DTW, a pattern and a sequence are aligned using a dis-
 197 tance metric, such as the Euclidean distance. However, since our gesture
 198 samples are modelled by means of probabilistic models, in order to use the
 199 principles of DTW, the distance must be redefined. In thise sense, a soft-
 200 distance based on the probability of a point x belonging to each one of the
 201 G components in the GMM is consider, i.e. the posterior probability of x is
 202 obtained according to Eq. (3). Therefore, since $\sum_{k=1}^G \alpha_k^t = 1$, the probability
 203 of a element $q_j \in Q$ belonging to the whole GMM λ_t can be computed as,

$$P(q_j, \lambda_t) = \sum_{k=1}^G \alpha_k^t \cdot P(q_j)_k, \quad (4)$$

204

$$P(q_j)_k = e^{-\frac{1}{2}(q_j-\mu_k^t)^T \cdot (\Sigma_k^t)^{-1} \cdot (q_j-\mu_k^t)}, \quad (5)$$

205 which is the sum of the weighted probability of each component. Never-
 206 theless, an additional step is required since the standard DTW algorithm
 207 is conceived for distances instead of similarity measures. In this sense, a
 208 soft-distance based measure of the probability is used, which is defined as,

$$D(q_j, \lambda_t) = \exp^{-P(q_j, \lambda_t)}. \quad (6)$$

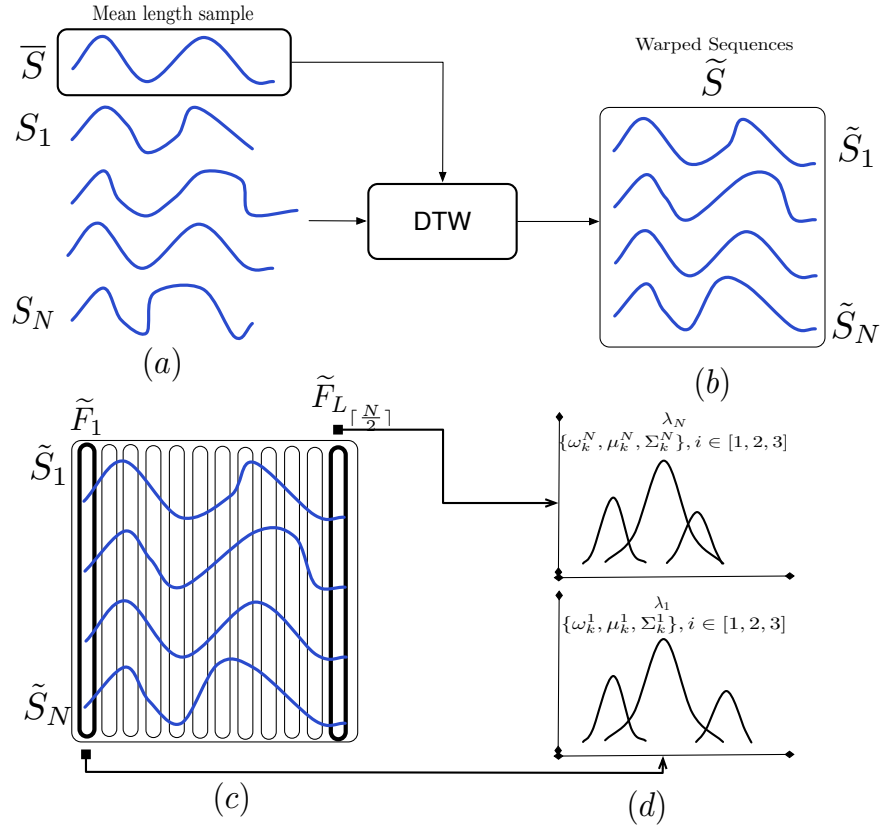


Figure 3: (a) Different sequences of a certain gesture category and the median length sequence. (b) Alignment of all sequences with the median length sequence by means of Euclidean DTW. (c) Warped sequences set \tilde{S} from which each set of t -th elements among all sequences are modelled. (d) Gaussian Mixture Model learning with 3 components.

209 In conclusion, possible temporal deformations of different samples of the
 210 same gesture category are taken into account by aligning the set of N gesture
 211 samples with the median length sequence. In addition, by modelling with
 212 a GMM each set of feature vectors which compose the resulting warped
 213 sequences, we obtain a methodology for gesture detection that is able to deal
 214 with multiple deformations in gestures both temporal (which are modelled

215 by the DTW alignment), or descriptive (which are learned by the GMM
 216 modelling). The algorithm that summarizes the use of the probability-based
 217 DTW to detect start-end of gesture categories is shown in Table 1. Figure 6
 218 illustrates the application of the algorithm in a toy problem.

Table 1: Probability-based DTW algorithm.

<p>Input: A set of GMM models $\lambda = \{\lambda_1, \dots, \lambda_m\}$ corresponding to a gesture category, a threshold value μ, and the streaming sequence $Q = \{q_1, \dots, q_\infty\}$. Cost matrix $M_{m \times \infty}$ is defined, where $\mathcal{N}(x), x = (i, t)$ is the set of three upper-left neighbor locations of x in M.</p> <p>Output: Warping path W of the detected gesture, if any.</p> <pre> // Initialization for $i = 1 : m$ do for $j = 1 : \infty$ do $M(i, j) = \infty$ end end for $j = 1 : \infty$ do $M(0, j) = 0$ end for $j = 0 : \infty$ do for $i = 1 : m$ do $x = (i, j)$ $M(x) = D(q_j, \lambda_i) + \min_{x' \in \mathcal{N}(x)} M(x')$ end if $M(m, j) < \mu$ then $W = \{\operatorname{argmin}_{x' \in \mathcal{N}(x)} M(x')\}$ return end end </pre>
--

219 *2.2. Gesture Representation: BoVDW*

220 In this section, the BoVDW approach for Human Gesture Representation
 221 is introduced. Figure 4 contains a conceptual scheme of the approach. In

222 this figure, it is shown that the information from RGB and Depth images
 223 is merged, while circles representing the spatio-temporal interest points are
 described by means of the proposed novel VFHCRH descriptor.

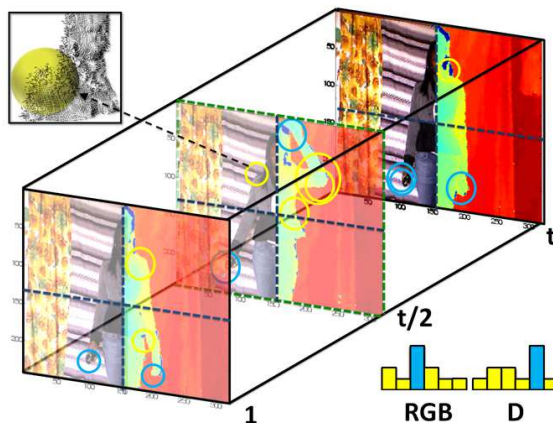


Figure 4: BoVDW approach in a Human Gesture Recognition scenario. Interest points in RGB and depth images are depicted as circles. Circles indicate the assignment to a visual word in the shown histogram – computed over one spatio-temporal bin. Limits of the bins from the spatio-temporal pyramids decomposition are represented by dashed lines in blue and green, respectively. A detailed view of the normals of the depth image is shown in the upper-left corner.

224

225 2.2.1. Keypoint detection

226 The first step of BoW-based models consists of selecting a set of points
 227 in the image/video with relevant properties. In order to reduce the amount
 228 of points in a dense spatio-temporal sampling, the Spatio-Temporal Interest
 229 Point (STIP) detector [23] is used, which is an extension of the well-known
 230 Harris detector in the temporal dimension. The STIP detector firstly com-
 231 puts the second-moment 3×3 matrix η of first order spatial and temporal

232 derivatives. Finally, the detector searches regions in the image with signif-
 233 icant eigenvalues $\lambda_1, \lambda_2, \lambda_3$ of η , combining the determinant and the trace of
 234 η ,

$$H = |\eta| - K \cdot T_r(\eta)^3, \quad (7)$$

235 where $|\cdot|$ corresponds to the determinant, $T_r(\cdot)$ computes the trace, and K
 236 stands for a relative importance constant factor. As multi-modal RGB-D
 237 data is employed, the STIP detector is applied separately on the RGB and
 238 Depth volumes, so two sets of interest points S_{RGB} and S_D are obtained.

239 *2.2.2. Keypoint description*

240 In this step, the interest points detected in the previous step should be
 241 described. On one hand, state-of-the-art RGB descriptors are computed for
 242 S_{RGB} , including Histogram of Gradients (HOG) [24], Histogram of Optical
 243 Flow (HOF), and their concatenation HOG/HOF [25]. On the other hand,
 244 a new descriptor VFHCRH (Viewpoint Feature Histogram Camera Roll His-
 245 togram) is introduced for S_D , as detailed below.

246 *2.2.3. VFHCRH*

247 The recently proposed Point Feature Histogram (PFH) and Fast Point
 248 Feature Histogram (FPFH) descriptors [8] represent each instance in the 3-
 249 D cloud of points with a histogram encoding the distribution of the mean
 250 curvature around it. Both PFH and FPFH provide $\mathcal{P}6$ DOF (Degrees of
 251 Freedom) pose invariant histograms, being \mathcal{P} the number of points in the
 252 cloud. Following their principles, Viewpoint Feature Histogram (VFH)[9]
 253 describes each cloud of points with one descriptor of 308 bins, variant to
 254 object rotation around pitch and yaw axis. However, VFH is invariant to

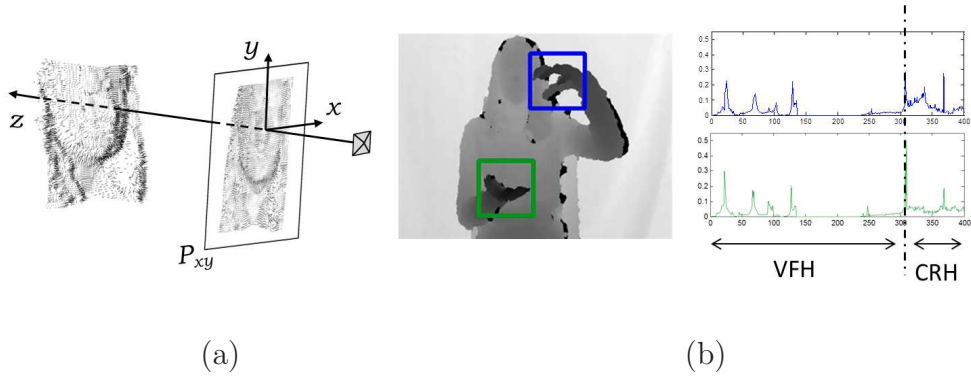


Figure 5: (a) Point cloud of a face and the projection of its normal vectors onto the plane P_{xy} , orthogonal to the viewing axis z . (b) VFHCRH descriptor: Concatenation of VFH and CRH histograms resulting in 400 total bins

255 rotation about the roll axis of the camera. In contrast, Clustered Viewpoint
 256 Feature Histogram (CVFH) [26] describes each cloud of points using a dif-
 257 ferent number of descriptors r , where r is the number of stable regions found
 258 on the cloud. Each stable region is described using a non-normalized VFH
 259 histogram and a Camera's Roll Histogram (CRH), and the final object de-
 260 scription includes all region descriptors. CRH is computed by projecting the
 261 normal of the point cloud $\tau^{(i)}$ of the i -th point $\rho^{(i)}$ onto a plane P_{xy} that is
 262 orthogonal to the viewing axis z , the vector between the camera center and
 263 the centroid of the cloud, under orthographic projection,

$$\tau_{xy}^{(i)} = \|\tau^{(i)}\| \cdot \sin(\phi), \quad (8)$$

264 where ϕ is the angle between the normal $\tau^{(i)}$ and the viewing axis. Finally,
 265 the histogram encodes the frequencies of the projected angle ψ between $\tau_{xy}^{(i)}$
 266 and y -axis, the vertical vector of the camera plane (see Fig. 5(a)).

267 In order to avoid descriptors of arbitrary lengths for different point clouds,

268 the whole cloud is described using VFH. In addition, a 92 bins CRH is
 269 computed for encoding $6DOF$ information. The concatenation of both his-
 270 tograms results in the proposed VFHCRH descriptor of 400 bins shown in
 271 Figure 5(b). Note how the first 308 bins of the concatenated feature vector
 272 correspond to the VFH, that encode the normals of the point cloud. Fi-
 273 nally, the remaining bins corresponding to the CRH descriptor, encode the
 274 information of the relative orientation of the point cloud to the camera.

275 2.2.4. BoVDW histogram

276 Once all the detected points have been described, the vocabulary of V
 277 visual/depth words is designed by applying a clustering method over all the
 278 descriptors. Hence, the clustering method – k -means in our case– defines
 279 the words from which a query video sequence will be represented, shaped
 280 like a histogram h that counts the occurrences of each word. Additionally,
 281 in order to introduce geometrical and temporal information, spatio-temporal
 282 pyramids are applied. Basically, spatio-temporal pyramids consist of dividing
 283 the video volume in b_u , b_v , and b_p bins along the u , v , and p dimensions of the
 284 volume, respectively. Then, $b_u \times b_v \times b_p$ separate histograms are computed
 285 with the points lying in each one of these bins, and they are concatenated
 286 jointly with the general histogram computed using all points.

287 These histograms define the model for a certain class of the problem –in
 288 our case, a certain gesture. Since multi-modal data is considered, different
 289 vocabularies are defined for the RGB-based descriptors and the depth-based
 290 ones, and the corresponding histograms, h^{RGB} and h^D , are obtained. Finally,
 291 the information given by the different modalities is merged in the next and
 292 final classification step, hence using *late fusion*.

293 *2.2.5. BoVDW-based classification*

294 The final step of the BoVDW approach consists of predicting the class
295 of the query video. For that, any kind of multi-class supervised learning
296 technique could be used. In our case, a simple k -Nearest Neighbour classifi-
297 cation is used, computing the complementary of the histogram intersection
298 as a distance,

$$d^F = 1 - \sum_i \min(h_{model}^F(i), h_{query}^F(i)), \quad (9)$$

299 where $F \in \{RGB, D\}$. Finally, in order to merge the histograms h^{RGB}
300 and h^D , the distances d^{RGB} and d^D are computed separately, as well as the
301 weighted sum,

$$d_{hist} = (1 - \beta)d^{RGB} + \beta d^D, \quad (10)$$

302 to perform late fusion, where β is a weighting factor.

303 **3. Experiments and Results**

304 To better understand the experiments, firstly the data, methods, and
305 evaluation measurements are discussed.

306 *3.1. Data*

307 Data source used is the ChaLearn [27] data set, provided by the CVPR2011
308 Workshop’s challenge on Human Gesture Recognition. The data set consists
309 of 50,000 gestures each one portraying a single user in front of a fixed cam-
310 era. The images are captured by the Kinect device providing both RGB and
311 depth images. A subset of the whole data set has been considered, formed
312 by 20 development batches with a manually tagged gesture segmentation,
313 which is used to obtain the idle gestures. Each batch includes 100 recorded

314 gestures grouped in sequences of 1 to 5 gestures performed by the same
315 user. The gestures from each batch are drawn from a different lexicon of 8
316 to 15 unique gestures and just one training sample per gesture is provided.
317 These lexicons are categorized in nine classes, including: (1) body language
318 gestures (scratching your head, crossing your arms, etc.), (2) gesticulations
319 performed to accompany speech, (3) illustrators (like Italian gestures), (4)
320 emblems (like Indian Mudras), (5) signs (from sign languages for the deaf),
321 (6) signals (diving signals, marshalling signals to guide machinery or vehicle,
322 etc.), (7) actions (like drinking or writing), (8) pantomimes (gestures made
323 to mimic actions), and (9) dance postures.

324 For each sequence, the actor performs an idle gesture between each gesture
325 to classify. These idle gestures are used to provide the temporal segmentation
326 (further details are shown in the next section). For this data set, background
327 subtraction was performed based on depth maps, and a 10×10 grid approach
328 was defined to extract HOG+HOF feature descriptors per cell, which are
329 finally concatenated in a full image (posture) descriptor. Using this data set,
330 the recognition of the idle gesture pattern will be tested, using 100 samples
331 of the pattern in a ten-fold validation procedure.

332 *3.2. Methods and Evaluation*

333 The experiments are presented in two different sections. The first section
334 considers the temporal segmentation experiment while the second section
335 aims the gesture classification experiments.

336 *3.2.1. Temporal Segmentation Experiments*

337 In order to provide with quantitative measures of the temporal segmenta-
338 tion procedure, we first describe the subset of the data used and the feature
339 extraction.

340 **• Data and Feature extraction**

341 For the temporal segmentation experiments we used the 20 development
342 batches provided by the organization of the challenge. These batches contain
343 a manual labelling of gesture start and end points. Each batch includes 100
344 recorded gestures, grouped in sequences of 1 to 5 gestures performed by the
345 same user. For each sequence the actor performs an idle gesture between
346 each gesture of the gestures drawn from lexicons. Finally, this means that
347 we have a set of approximately 1800 idle gestures.

348 Each video sequence of each batch was described using a 20×20 grid
349 approach. For each patch in the grid we obtain a 208 feature vector consisting
350 of HOG (128 dimensions) and HOF (80 dimensions) descriptors which are
351 finally concatenated in a full image (posture descriptor). Due to the huge
352 dimensionality of the descriptor of a single frame (83200 dimensions), we
353 utilized a Random Projection to reduce dimensionality to 150 dimensions.

354 **• Experimental Settings**

355 For both of the DTW approaches the cost-threshold value θ is estimated
356 in advance using ten-fold cross-validation strategy on the set of 1800 idle
357 gesture samples. This involves using 180 idle gestures as the validation data,
358 and the remaining observations as the training data. This is repeated such
359 that each observation in the sample is used once as the validation data.
360 Finally, the threshold value θ chosen is the one associated with the largest

361 overlapping performance. For the probabilistic DTW approach, each GMM
 362 was fit with 4 components. The value of G was obtained using a ten-fold
 363 cross-validation procedure on the set of 1800 idle gestures as well. In this
 364 sense, the cross-validation procedure for the probability-based DTW is a
 365 double loop (optimizing on the number of GMM components G , and then, on
 366 the cost-threshold θ). In the HMM case, we used the Baum-Welch algorithm
 367 for training, and 3 states were experimentally set for the idle gesture, using
 368 a vocabulary of 60 symbols computed using K-means over the training data
 369 features. Final recognition is performed with temporal sliding windows of
 370 different wide sizes, based on the idle gesture samples length variability.

371 **• Methods, Measurements and Results**

372 Our probability-based DTW approach using the proposed distance D shown
 373 in Eq. (6) is compared to the usual DTW algorithm and the Hidden Markov
 374 Model approach. The evaluation measurements presented are *overlapping*
 375 and *accuracy* of the recognition for the idle gesture, considering that a gesture
 376 is correctly detected if overlapping in the idle gesture sub-sequence is greater
 377 than 60% (the standard overlapping value, computed as the intersection over
 378 the union between the temporal bounds in the ground truth, and the ones
 379 computed by our method). The accuracy is computed frame-wise as

$$Acc = \frac{TruePositives + TrueNegatives}{TruePositives + TrueNegatives + FalsePositives + FalseNegatives}. \quad (11)$$

380

381 The results of our proposal, HMM and the classical DTW algorithm are shown
 382 in Table 2. It can be seen how the proposed probability-based DTW outperforms
 383 the usual DTW and HMM algorithms in both experiments. Moreover, confidence
 384 intervals of DTW and HMM do not intersect with the probability-based DTW in

Table 2: *Overlapping and accuracy* results.

	Overlap.	Acc.
Probability-based DTW	0.3908 ± 0.0211	0.6781 ± 0.0239
Euclidean DTW	0.3003 ± 0.0302	0.6043 ± 0.0321
HMM	0.2851 ± 0.0432	0.5328 ± 0.0519

385 any case. From this results it can be concluded that performing dynamic program-
 386 ming increases the generalization capability of the HMM approach, as well as a
 387 model defined by a set of GMMs outperforms the classical DTW on RGB-Depth
 388 data without increasing the computational complexity of the method. Figure 6
 389 shows qualitative results from two sample video sequences.

390 3.2.2. *BoVDW Classification Experiments*

391 In all the experiments shown in this section, the vocabulary size was set to
 392 $N = 200$ words for both RGB and depth cases. For the spatio-temporal pyramids,
 393 the volume was divided in $2 \times 2 \times 2$ bins (resulting in a final histogram of 1800
 394 bins). Since the nature of our application problem is one-shot learning (only one
 395 training sample is available for each class), a simple Nearest Neighbor classification
 396 is employed. Finally, for the late fusion, the weight $\beta = 0.8$ was empirically set, by
 397 testing the performance of our method in a small subset of development batches
 398 from the dataset. We observed that when increasing β , starting from $\beta = 0$, the
 399 performance keeps increasing in a linear fashion, until the value $\beta = 0.45$. From
 400 $\beta = 0.45$ to $\beta = 0.8$ the performance keeps improving more slightly, and finally,
 401 from $\beta = 0.8$ to $\beta = 1$ the performance drops again.

402 For the evaluation of the methods, in the context of Human Gesture Recogni-
 403 tion, the Levenshtein distance or edit distance was considered. This edit distance
 404 between two strings is defined as the minimum number of operations (insertions,

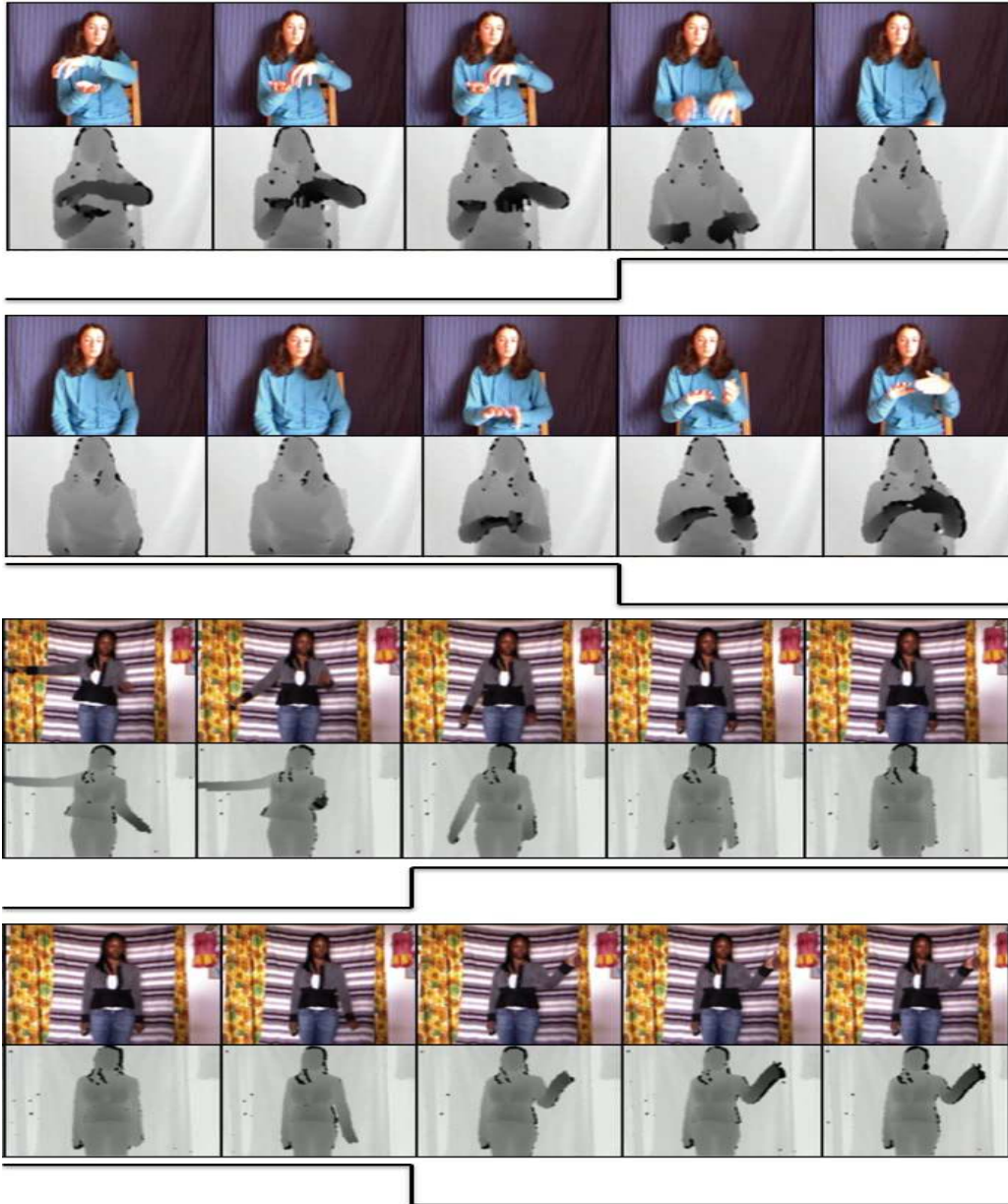


Figure 6: Examples of idle gesture detection on the Chalearn data set using the probability-based DTW approach. The line below each pair of depth and RGB images represents the detection of a idle gesture (step up: beginning of idle gesture, step down: end)

405 substitutions or deletions) needed to transform one string into the other. In our
406 case, strings contain gesture labels detected in a video sequence. For all the com-
407 parison, the mean Levenshtein distance (MLD) was computed over all sequences
408 and batches.

409 Table 3 shows a comparison between different state-of-the-art RGB and depth
410 descriptors (including our proposed VFHCRH), using our BoVDW approach. More-
411 over, we compare our BoVDW framework with the baseline methodology provided
412 by the ChaLearn 2012 Gesture Recognition challenge. This baseline first computes
413 differences of contiguous frames, which encode movement information. After that,
414 these difference images are divided into cells forming a grid, each one containing
415 the sum of movement information among it. These 2D grids are then transformed
416 then into vectors, one for each difference image. Moreover, the model for a gesture
417 is computed via Principal Component Analysis (PCA), using all the vectors be-
418 longing to that gesture. The eigenvectors are just computed and stored, so when a
419 new sequence arrives, its movement signature first is computed, and then projected
420 and reconstructed using the different PCA models from each gesture. Finally, the
421 classification is performed by choosing the gesture class with lower reconstruction
422 error. This baseline obtains a MLD of 0.5096. Table 4 shows the results in all the
423 20 development batches separately.

424 When using our BoVDW approach, in the case of RGB descriptors, HOF alone
425 performs the worst. In contrast, the early concatenation of HOF to HOG descrip-
426 tor outperforms the simple HOG. Thus, HOF contributes adding discriminative
427 information to HOG. In a similar way, looking at the depth descriptors, it can be
428 seen how the concatenation of the CRH to the VFH descriptor clearly improves
429 the performance compared to the simpler VFH. When using late fusion in order
430 to merge information from the best RGB and depth descriptors (HOGHOF and
431 VFHCRH, respectively), a value of 0.2714 for MLD is achieved. Figure 7 shows

Table 3: Mean Levenshtein distance for RGB and depth descriptors.

RGB desc.	MLD	Depth desc.	MLD
HOG	0.3452	VFH	0.4021
HOF	0.4144	VFHCRH	0.3064
HOGHOF	0.3314		

432 the confusion matrices of the gesture recognition results with this late fusion con-
433 figuration. In general, the confusion matrices follow an almost diagonal shape,
434 indicating that the majority of the gestures are well classified. However, the re-
435 sults of batches 3, 16, 18, 19 are significantly worse, possibly due to the static
436 characteristics of the gestures in these batches. Furthermore, late fusion was also
437 applied in a 3-fold way, merging HOG, HOF, and VFHCRH descriptors separately.
438 In this case the weight β was assigned to HOG and VFHCRH descriptors (and
439 $1 - \beta$ to HOF), improving the MLD to 0.2662. From this result it can be concluded
440 that HOGHOF late fusion performs better than HOGHOF early fusion.

441 4. Conclusion

442 In this paper, the BoVDW approach for Human Gesture Recognition has been
443 presented using multi-modal RGB-D images. A new depth descriptor VFHCRH
444 has been proposed, which outperforms VFH. Moreover, the effect of the late fu-
445 sion has been analysed for the combination of RGB and depth descriptors in the
446 BoVDW, obtaining better performance in comparison to early fusion. In addition,
447 a probabilistic-based DTW has been proposed to asses the temporal segmentation
448 of gestures, where different samples of the same gesture category are used to build
449 a Gaussian-based probabilistic model of the gesture in which possible deformations

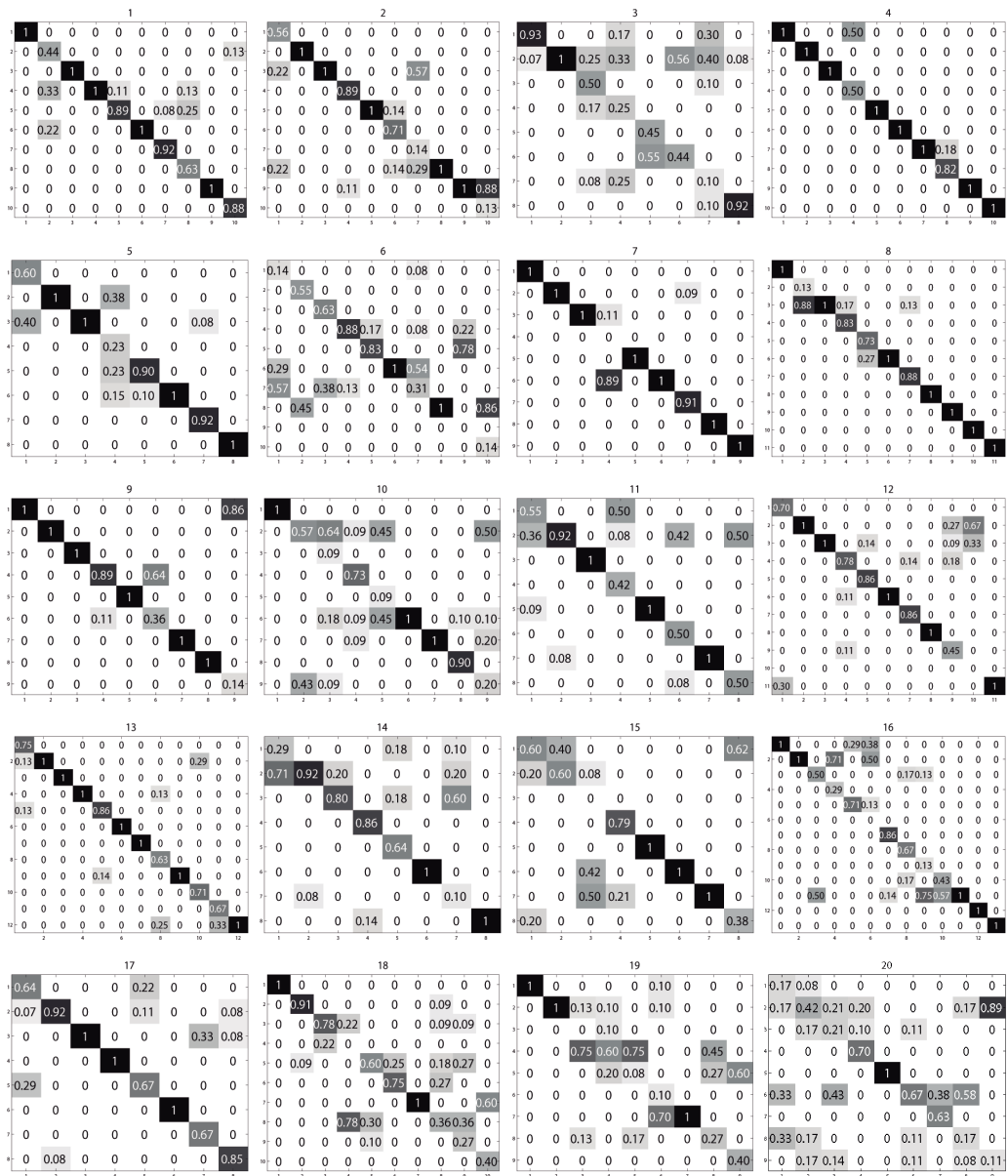


Figure 7: Confusion matrices for gesture recognition in each one of the 20 development batches.

Table 4: Mean Levenshtein Distance of the best RGB and depth descriptors separately, as well as the 2-fold and 3-fold late fusion of them. Results obtained by the baseline from the ChaLearn challenge are also shown. Rows 1 to 20 represent the different batches.

	HOGHOF	VFHCRH	2-fold L.F.	3-fold L.F.	Baseline
1	0.19	0.17	0.12	0.20	0.42
2	0.24	0.30	0.24	0.26	0.57
3	0.76	0.39	0.40	0.49	0.78
4	0.14	0.08	0.08	0.11	0.32
5	0.08	0.33	0.17	0.17	0.25
6	0.41	0.47	0.44	0.34	0.54
7	0.10	0.18	0.11	0.13	0.64
8	0.12	0.26	0.14	0.08	0.40
9	0.11	0.18	0.15	0.13	0.30
10	0.57	0.40	0.39	0.46	0.79
11	0.47	0.36	0.27	0.34	0.54
12	0.37	0.20	0.21	0.17	0.42
13	0.16	0.14	0.10	0.09	0.34
14	0.41	0.34	0.30	0.30	0.69
15	0.38	0.28	0.34	0.28	0.54
16	0.22	0.41	0.34	0.29	0.42
17	0.38	0.16	0.15	0.17	0.55
18	0.38	0.43	0.40	0.38	0.53
19	0.67	0.50	0.50	0.44	0.61
20	0.46	0.57	0.56	0.48	0.52

450 are implicitly encoded. In addition, to embed these models into the DTW frame-
451 work, a soft-distance based on the posterior probability of the GMM was defined.
452 In conclusion, a novel methodology for gesture detection has been presented, which
453 is able to deal with multiple deformations in data.

454 **Acknowledgements**

455 This work has been partially supported by the “Comissionat per a Universitats
456 i Recerca del Departament d’Innovació, Universitats i Empresa de la Generalitat de
457 Catalunya” and the following projects: IMSERSO Ministerio de Sanidad 2011 Ref.
458 MEDIMINDER, RECERCAIXA 2011 Ref. REMEDI, TIN2012-38187-C03-02 and
459 CONSOLIDER INGENIO CSD 2007-00018. The work of Antonio is supported by
460 an FPU fellowship from the Spanish government. The work of Víctor is supported
461 by the 2013FI-B01037 fellowship. The work of Miguel Ángel is supported by an
462 FI fellowship from SUR-DEC of Generalitat de Catalunya and FSE.

- 463 [1] A. Hampapur, L. Brown, J. Connell, A. Ekin, N. Haas, M. Lu, H. Merkl,
464 S. Pankanti, Smart video surveillance: exploring the concept of multiscale
465 spatiotemporal tracking, *Signal Processing Magazine, IEEE* 22 (2) (2005)
466 38–51.
- 467 [2] T. Starner, A. Pentland, Real-time american sign language recognition from
468 video using hidden markov models, in: *Computer Vision, 1995. Proceedings.,
469 International Symposium on, 1995*, pp. 265 –270.
- 470 [3] A. Pentland, Socially aware computation and communication, *Computer* 38
471 (2005) 33–40.
- 472 [4] D. D. Lewis, Naive (bayes) at forty: The independence assumption in infor-
473 mation retrieval, in: *ECML, Springer Verlag, 1998*, pp. 4–15.
- 474 [5] G. Csurka, C. Dance, L. Fan, J. Willamowski, C. Bray, Visual categorization
475 with bags of keypoints, in: *ECCV, 2004*, pp. 1–22.

- 476 [6] M. Mirza-Mohammadi, S. Escalera, P. Radeva, Contextual-guided bag-of-
477 visual-words model for multi-class object categorization, in: CAIP, 2009, pp.
478 748–756.
- 479 [7] J. C. Niebles, H. Wang, L. Fei-Fei, Unsupervised learning of human action
480 categories using spatial-temporal words, IJCV 79 (3) (2008) 299–318.
- 481 [8] R. Bogdan, N. Blodow, M. Beetz, Fast point feature histograms (fpfh) for 3d
482 registration, in: ICRA, 2009, pp. 3212–3217.
- 483 [9] R. Rusu, G. Bradski, R. Thibaux, J. Hsu, Fast 3d recognition and pose using
484 the viewpoint feature histogram, in: Intelligent Robots and Systems (IROS),
485 2010 IEEE/RSJ International Conference on, 2010, pp. 2155 –2162.
- 486 [10] K. Lai, L. Bo, X. Ren, D. Fox, Sparse distance learning for object recognition
487 combining rgb and depth information, in: ICRA, 2011, pp. 4007 –4013.
- 488 [11] L. Bo, X. Ren, D. Fox, Depth kernel descriptors for object recognition, in:
489 IROS, 2011, pp. 821–826.
- 490 [12] T. Wan, Y. Wang, J. Li, Hand gesture recognition system using depth data,
491 in: Consumer Electronics, Communications and Networks (CECNet), 2012
492 2nd International Conference on, 2012, pp. 1063–1066.
- 493 [13] Y. Li, Hand gesture recognition using kinect, in: Software Engineering and
494 Service Science (ICSESS), 2012 IEEE 3rd International Conference on, 2012,
495 pp. 196–199.
- 496 [14] F. Pedersoli, N. Adami, S. Benini, R. Leonardi, Xkin -: extendable hand pose
497 and gesture recognition library for kinect, in: ACM Multimedia, 2012, pp.
498 1465–1468.

- 499 [15] K. K. Biswas, S. Basu, Gesture recognition using microsoft kinect, in: Au-
500 tomation, Robotics and Applications (ICARA), 2011 5th International Con-
501 ference on, 2011, pp. 100–103.
- 502 [16] P. Doliotis, A. Stefan, C. McMurrough, D. Eckhard, V. Athitsos, Comparing
503 gesture recognition accuracy using color and depth information, in: Proceed-
504 ings of the 4th International Conference on PErvasive Technologies Related
505 to Assistive Environments, PETRA '11, 2011, pp. 20:1–20:7.
- 506 [17] C. Keskin, F. Kıraç, Y. E. Kara, L. Akarun, Real time hand pose estimation
507 using depth sensors, in: Consumer Depth Cameras for Computer Vision,
508 Springer, 2013, pp. 119–137.
- 509 [18] J. Shotton, A. Fitzgibbon, M. Cook, T. Sharp, M. Finocchio, R. Moore,
510 A. Kipman, A. Blake, Real-time human pose recognition in parts from single
511 depth images, in: CVPR, 2011, pp. 1297 –1304.
- 512 [19] A. Hernández-Vela, N. Zlateva, A. Marinov, M. Reyes, P. Radeva, D. Di-
513 mov, S. Escalera, Human limb segmentation in depth maps based on spatio-
514 temporal graph-cuts optimization, *Journal of Ambient Intelligence and Smart*
515 *Environments* 4 (6) (2012) 535–546.
- 516 [20] R. Girshick, J. Shotton, P. Kohli, A. Criminisi, A. Fitzgibbon, Efficient re-
517 gression of general-activity human poses from depth images, in: ICCV, 2011,
518 pp. 415 –422.
- 519 [21] F. Zhou, F. De la Torre, J. Hodgins, Hierarchical aligned cluster analysis for
520 temporal clustering of human motion, *IEEE TPAMI* (2012) 1.
- 521 [22] M. Reyes, G. Dominguez, S. Escalera, Feature weighting in dynamic time
522 warping for gesture recognition in depth data, in: ICCV, 2011.

- 523 [23] I. Laptev, On space-time interest points, *Int. J. Comput. Vision* 64 (2-3)
524 (2005) 107–123.
- 525 [24] N. Dalal, B. Triggs, Histograms of oriented gradients for human detection,
526 *CVPR* 1 (2005) 886–893.
- 527 [25] I. Laptev, M. Marszalek, C. Schmid, B. Rozenfeld, Learning realistic human
528 actions from movies, in: *CVPR*, 2008, pp. 1–8.
- 529 [26] A. Aldoma, M. Vincze, N. Blodow, D. Gossow, S. Gedikli, R. Rusu, G. Brad-
530 ski, Cad-model recognition and 6dof pose estimation using 3d cues, in: *Com-
531 puter Vision Workshops (ICCV Workshops)*, 2011 IEEE International Con-
532 ference on, 2011, pp. 585–592.
- 533 [27] Chalearn gesture dataset, california (2011).
534 URL <http://gesture.chalearn.org/data>