

Gesture sequence recognition with one shot learned CRF/HMM hybrid model

Selma Belgacem^a, Clément Chatelain^a, Thierry Paquet^a

^a*LITIS EA 4108, University of Rouen, Saint-Etienne du Rouvray, France*

Abstract

In this paper, we propose a novel markovian hybrid system CRF/HMM for gesture recognition, and a novel motion description method called *gesture signature* for gesture characterisation. The gesture signature is computed using the optical flows in order to describe the location, velocity and orientation of the gesture global motion. We elaborated the proposed hybrid CRF/HMM model by combining the modeling ability of Hidden Markov Models and the discriminative ability of Conditional Random Fields. In the context of one-shot-learning, this model is applied to the recognition of gestures in videos. In this extreme case, the proposed framework achieves very interesting performance and remains independent from the moving object type, which suggest possible application to other motion-based recognition tasks.

Keywords: gesture recognition, one-shot-learning, hybrid system, hidden Markov model, conditional random field, gesture characterisation.

1. Introduction

2 Following the increasing demand for intuitive and simple human/computer
3 interaction, the gesture analysis and recognition research field has received
4 a lot of attention these last years. A gesture can be defined as a short
5 human body motion, achieved primarily with arms. In some particular
6 situations such as disability or constrained environment, the gesture is the
7 only human/machine communication channel. This study falls into gesture
8 characterization and recognition.

Email address: `belgacemselma@yahoo.fr` (Selma Belgacem)

9 The recognition of gesture sequences combines both segmentation and
10 classification. As stated by Sayre [1], segmentation and classification are two
11 tasks that must be performed simultaneously. The segmentation step has to
12 face the variability of the duration of gestures, while the classification step
13 has to face the variability of instances of a same gesture.

14 A video gesture can be represented in a simplified three-dimensional
15 space consisting of its two-dimensional projection and its variation through
16 time. The recognition system must be robust to recording environment
17 variations such as changes in brightness, backgrounds, colors, objects, signer
18 appearance (clothes, skin color, height, etc.).

19 Markov models, which are widely applied to the recognition and segmenta-
20 tion of sequential data, model the temporal dependencies in sequences. They
21 are based on the Markovian assumption that accounts for the short-term
22 dependencies only, omitting the long-term dependencies in the model.

23 Although introducing some simplification in the model, generative Markov
24 models such as Hidden Markov Models (HMM) [2] allow to introduce a
25 temporal structure between classes representing a high-level knowledge such
26 as a language model. The principle of the HMM is to model the observation
27 generation based on some hidden states. Each observation only depends
28 from the current hidden state (thus assuming observations to be conditionally
29 independent between each other) and each hidden state only depends from the
30 previous state (for an order 1 Markov model). Then, through the inference
31 phase, the most likely sequence of hidden states that describes the given
32 sequence of observations is determined using Viterbi algorithm [3]. On
33 the other hand, HMM's use Gaussian Mixtures (GM) to model the data
34 distribution. When training data are too few, modeling becomes poor and
35 inadequate with GMM, which is a major drawback of HMM's. However,
36 discriminative models, such as Conditional Random Fields (CRF) [4] which
37 are also Markov models, can remedy this problem. The CRF model was
38 proposed by Lafferty et al. [4] in 2001. It has some advantages that can address
39 HMM problems.

40 CRF's have been designed in order to model the decision process of
41 labelling a sequence. Therefore they account for the a posteriori probability
42 of a particular sequence of labels. Similar to HMM's, at each time step a label
43 depends on the the previous label (Markov assumption), but may depend on
44 the whole observation sequence making no requirement about the conditional
45 independence of the observation data. As opposed to HMM, CRF are not
46 able to model high level information such as a language model, or syntactical

47 rules. They are local classifiers in a sequential process. Thus, the high-level
48 knowledge must be introduced in post-processing as an additional step of
49 filtering in order to guaranty the structural labelling consistency. The HMM's
50 generative framework has this ability of coping with high level structuring
51 information.

52 Finally, if we compare the advantages and disadvantages of CRF and
53 HMM, we find a certain complementarity between the two models. Therefore,
54 in this work we propose to combine these two models in a hybrid framework,
55 allowing the integration of knowledge while being robust to different sources
56 of variability. We also propose to characterize gestures using an original
57 global description of shapes and motions in the video frames. This method
58 describes the location, the velocity and the direction of the motion, based on
59 the optical flow velocity information. In one-shot gesture learning context,
60 this system was tested using the " *Gesture Challenge 1-2*" dataset proposed
61 by ChaLearn 2011-2012 [5, 6]. We will show that the lack of training data is
62 another problem which can be solved by Markov models to a certain extent.

63 We will show mainly the principle of our hybrid model CRF/HMM and
64 explain how we adapt it to the one-shot learning context, in order to cope with
65 the lack of training data. We will describe also our gesture characterization
66 model and present the experimental protocol and the evaluation of our system
67 recognition results.

68 **2. Related works**

69 Human gesture analysis is an active research domain with a lot of applica-
70 tions. Among them, many studies have been devoted to gesture recognition,
71 especially the design of automatic systems for recognizing the sign language.
72 Such systems would allow deaf people to better communicate with machines
73 or with other humans.

74 For gesture sequences recognition, the use of global parallel HMM models
75 is common in the literature [7, 8, 9, 10, 11, 12]. For example Vogler et al. [7],
76 Agris et al. [8] and Ong et al. [9] designed a parallel HMM model for signed
77 sentences recognition. They distinguished gesture descriptors such as position,
78 orientation and distance to facilitate the learning process of the HMM and
79 optimize the use of these descriptors. This decomposition is manifested by
80 the generation of one HMM for each descriptor and for each sub-unit of the
81 model.

82 Another issue when dealing with real-world problem such as gesture
83 recognition is the lack of labeled examples.

84 Konecny [10] et al., Jackson [11] and Weiss [12] proposed a global
85 HMM model for gesture sequences recognition using single-instance learning
86 databases. The global model is a set of left-right interconnected HMM's
87 modeling each gesture. From each state of each HMM, it is possible to remain
88 in that state or to jump to a subsequent internal or external state. In the
89 model proposed by Jackson [11], each frame of the gesture video is represented
90 by a state. This model remain complex due to the large number of states
91 involved.

92 The idea of combining HMM with other classification scheme is not
93 new. Such hybrid framework is intended to introduce a better discrimina-
94 tion between classes, than pure generative models can do. One of the first
95 combination scheme was proposed in the 1990s by the integration of neural
96 networks to HMM's [13]. Such combination is prevalent in the literature in
97 various fields. This type of hybrid models was applied to speech recogni-
98 tion [14, 15, 16, 17, 18, 19], handwriting recognition [20, 21, 22, 23, 24, 25, 26]
99 and gesture recognition [27]. HMM models have also been combined with
100 SVM models for handwriting recognition [28] and with dynamic programming
101 methods for gesture recognition [29]. We noticed that the application of these
102 hybrid models to gesture recognition is recent and not much studied in the
103 literature.

104 To the best of our knowledge, the only work addressing CRF and HMM
105 combination is the work of Soullard et al. [30], based on the work of Gunawar-
106 dana et al. [31]. In this work, the authors constrain the learning step of a
107 hidden CRF by initialising it with the parameters of a pre-trained HMM.
108 This method ensures the convergence of the hidden CRF learning step and
109 shows the difficulty of learning convergence of such models.

110 The idea of our approach is different and is inspired from neuro-Markovian
111 approaches. The principle of these approaches is to replace the HMM data
112 model, consisting of a mixture of Gaussians, by a discriminative model that
113 classifies local observations. This model is traditionally composed of a neural
114 network which provides local posteriors associated to each local observation in
115 the sequence. In this work, we propose the use of a CRF in order to perform
116 this discriminative layer. The CRF layer will discriminate local observations
117 and provide local class posteriors to the HMM layer. These local posteriors
118 are then combined during the HMM decoding stage that integrates more
119 global information embedded in the HMM transition model (known as the

120 language model). According to the principle of our hybrid model, the HMM
 121 learning step and the CRF learning step are performed separately. Details of
 122 the new hybrid model we propose are presented in section 3.

123 3. Hybrid CRF/HMM model

124 3.1. Overview of the CRF/HMM model

125 In this section, we present our hybrid CRF/HMM system for gesture
 126 recognition. It combines the discriminative ability of CRF with the modeling
 127 ability of HMM. Combining the two models is performed in an easy and
 128 straightforward way derived from the literature. The discriminative CRF
 129 stage provides local class posterior probabilities that are fed to the HMM stage
 130 that account for more global constraints regarding the label sequence. Let us
 131 recall that a label is noted y_t , corresponding to a gesture segment which can
 132 span over multiple video frames. An observation is noted x_t , corresponding
 133 to a feature vector extracted from one frame. The feature vector is a real
 134 valued vector when using the first HMMs devoted to the frame labelling
 135 task (the dimension is the feature vector size, see experimental results). This
 136 feature vector is later quantified into multiple bins when used by the CRF
 137 (see section 3.3) for the gesture recognition task. Its size is defined in section
 138 4. The number of a sequence frames is noted T , it depends of the gesture size.
 139 $y_{1:T}$ and $x_{1:T}$ are respectively noted by Y , and X . X_d presents the quantified
 140 feature vector. Figure 1 shows the proposed hybrid system.

141 Following this model, the HMM probability $p(y_{1:T}, x_{1:T})$ (see Eq. 1)
 142 depends on the posteriors computed using the CRF.

$$p(y_{1:T}, x_{1:T}) = p(x_1|y_1)p(y_1) \prod_{t=2}^T p(x_t|y_t)p(y_t|y_{t-1}) \quad (1)$$

143 In the classic form of HMMs, $p(x_t|y_t)$ is a Gaussian mixture. In our new
 144 model, this distribution will be, in some way, replaced by the categorical
 145 distribution $p(y_t|x_t)$ computed by the CRFs. Indeed, $p(x_t|y_t)$ is a likelihood,
 146 while the CRF outputs posteriors $p(y_t|x_t)$. Therefore, $p(x_t|y_t)$ is computed
 147 from $p(y_t|x_t)$ using Bayes' rule :

$$p(x_t|y_t) = \frac{p(y_t|x_t)p(x_t)}{p(y_t)} \quad (2)$$

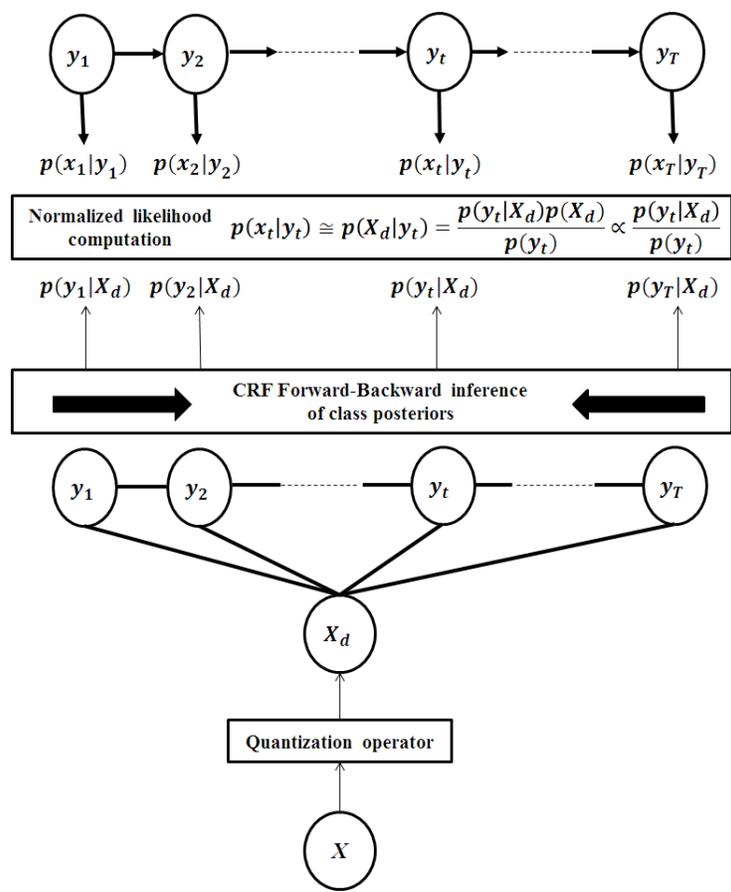


Figure 1: The graphical model CRF/HMM : the HMM joint probability $p(y_{1:T}, x_{1:T})$ for the observation sequence X and the state sequence $y_{1:T}$ is computed using CRF local class posterior probabilities $p(y_t|x_t)$

148 As every gesture class are considered to be equally likely, $p(y_t)$ is a constant
 149 $\forall t \in \mathbb{N}$. The aim of the decoding process is to find the state sequence $y_{1:T}$
 150 that maximises $p(y_{1:T}, x_{1:T})$. As the observation probability $p(x_t)$ is time
 151 independent, $p(x_t)$ is not involved in the maximization of $p(x_t|y_t)$. Hence, the
 152 maximization of $p(x_t|y_t)$ turns toward the maximization of $p(y_t|x_t)$.

153 Given that the CRF are able to take into account the whole observation
 154 sequence to compute the posteriors of each class, we assume that $p(y_t|x_t) =$
 155 $p(y_t|x_{1:T}) = p(y_t|X) \cong p(y_t|X_d)$.

156 This is computed within the CRF using the forward-backward algorithm
 157 [32], where the forward probability α_t and the backward probability β_t are
 158 computed using the following recurrences:

$$\alpha_t(i) = \sum_{j=1}^{N_s} \alpha_{t-1}(j) \psi_t(s_i, s_j, o_t) \quad (3)$$

$$\beta_t(i) = \sum_{j=1}^{N_s} \beta_{t+1}(j) \psi_{t+1}(s_i, s_j, o_t) \quad (4)$$

159 where

$$\psi_t(s_i, s_j, o_t) = \exp\left(\sum_{k=1}^K \lambda_k f_k(y_t = s_i, y_{t-1} = s_j, x_t = o_t)\right) \quad (5)$$

160 and s_i, s_j are hidden state that belong to \mathcal{S} , and o_t is an observation that
 161 belong to \mathcal{O} . Finally, following the forward-backward procedure, we have:

$$p(X_d) = \sum_{j=1}^{N_s} \alpha_T(j) = \sum_{j=1}^{N_s} \beta_1(j) = \sum_{j=1}^{N_s} \alpha_t(j) \beta_t(j) \quad (6)$$

$$p(y_t = s_i | X_d) = \frac{p(y_t = s_i, X_d)}{p(X_d)} = \frac{\alpha_t(i) \beta_t(i)}{\sum_{j=1}^{N_s} \alpha_t(j) \beta_t(j)} = \gamma_t(i) \quad (7)$$

162 3.2. Training the CRF/HMM model

163 We chose to achieve a separate training of the HMMs and the CRF. As
 164 first stage, HMMs are trained with the standard Baum Welch algorithm which
 165 means that the target function is the likelihood of the global gesture model.
 166 Transition probability Matrices are learned separately for each gesture class,

167 and gathered into a global model for decoding gesture sequences. This model
168 is described in section 3.5.

169 In the second stage, CRF are trained with the classic LBFGS algorithm.
170 As CRF do not benefit from an embedded training procedure like HMM, the
171 target function of this training phase is the local frame level classe (state)
172 posterior. Therefore local frame level labels are necessary. In this respect we
173 introduce a frame level labelling stage that consists in using the HMM model
174 of gesture trained on the dataset, in a forced alignment mode. The frame
175 labels produced serve as the objective target of the of the cost function for
176 training the CRF. During this second training phase the CRF learns a single
177 model for all gestures, considering as many classes in the model as there are
178 sub-gestures. The number of sub-gestures is equal to the number of states in
179 the HMM model of gesture.

180 The training chain is furthermore explained in figure 3.

181 *3.3. CRF/HMM adaptation to one-shot learning*

182 In this section, we focus on the learning of the recognition system using a
183 unique sample per class. These learning conditions are interesting since the
184 annotation efforts are extremely reduced in this case. Furthermore, using a
185 single sample per class allows to speed up the learning process.

186 The one-shot learning framework has been quite extensively used for
187 gesture analysis and recognition [10, 11, 12, 33, 6]. These system are generally
188 made of a standard recognition method that has been adapted to the one shot
189 learning framework. We now describe the adaptation of our models (HMM
190 and CRF) to one shot learning.

191 To model the feature space, the HMM relies on Gaussian mixtures esti-
192 mated on the learning database. When considering a very reduced number
193 of samples, the parameters of the Gaussian distribution $p(x_t|y_t)$ are very
194 difficult to estimate, especially the variance. Therefore, first we limited the
195 mixture to one Gaussian per gesture class. Second, the variance is computed
196 on every gesture class in order to increase the amount of data and improve the
197 estimation. Doing that, each gesture class has the same variance. Although
198 these two tricks are a limitation of the initial method, the experiments showed
199 the interest of such an adaptation.

200 In its initial form, the CRF method is mathematically able to deal with
201 either discrete or continuous features[34, 35]; however, since the CRF clas-
202 sification stage is derived from a logistic regression, it is more adapted to
203 discrete features than continuous [36]. This is even more true when the

204 number of samples is small. Indeed, in the context of one shot learning the
 205 loss of information induced by the discretization of continuous features may
 206 have a regularization effect when training the CRF with one single example.
 207 Feature quantization also allows to efficiently tune the parameters linked to
 208 each discrete feature value. Although quantization involves a loss of infor-
 209 mation, the integration of a large set of features allows to capture a global
 210 representation of the whole gesture. Therefore, we turned toward the use of a
 211 feature quantization procedure. Notice that some recent developments have
 212 introduced Hidden CRF models in order to cope with continuous features
 213 [37]. But such a framework would require more data than possible in the
 214 one-shot learning context.

215 The quantification is achieved using a uniform scalar quantifier that maps
 216 each continuous feature into N_q discrete features, according to the following
 217 equation:

$$\begin{aligned}
 Q : [-V_{\max}, V_{\max}] &\longrightarrow [-N_q, N_q] \\
 x &\longmapsto \frac{x \times N_q}{V_{\max}}
 \end{aligned} \tag{8}$$

218 We empirically tuned the value N_q in order to reach the best recognition
 219 performance using a validation procedure. We found that $N_q = 16$ was the
 220 best value.

221 3.4. Structure and parametrization of the CRF/HMM model

222 As for a standard HMM, the HMM of our hybrid structure is made of
 223 states describing each gesture. Although the gesture duration can be modelled
 224 through the state auto-transitions, it is known that a better modelization
 225 can be achieved by setting a variable number of states per gesture. We
 226 experimentally checked that this strategy outperforms the performance of
 227 the same system with a fixed number of states per gesture. The number of
 228 states of each gesture i is determined automatically depending on its frame
 229 length $\mathbf{f}_g(i)$. The theoretical number of frames per state, denoted \mathbf{f}_s , is one
 230 hyper-parameter of the system. We denote the number of states of a gesture
 231 model i ; $N_e(i) = \mathbf{f}_g(i)/\mathbf{f}_s$. As we already mentioned, we limit the data model
 232 to have only one Gaussian per state.

233 The CRF part of our hybrid model has a standard linear structure, as
 234 shown in figure 1. The CRF training leads to a single model that discriminates
 235 all the gestures of the dataset. As explained in the previous section, the
 236 CRF formulation allows to consider an observation window, including the

237 current observation and a neighbouring context to be determined. To adapt
 238 the system to the gesture duration variability, we chose a variable size \mathbf{f}_w of
 239 the observation window w_o (equation 9). \mathbf{f}_w is statistically estimated on the
 240 learning databases. In order to avoid overfitting the CRF, a regularization
 241 term has been empirically tuned to a value of 1.5.

$$\mathbf{f}_w(w_o(\mathbb{G})) = \min\left(\frac{3}{4} \min_{\mathbf{g} \in \mathbb{G}} \mathbf{f}_g(\mathbf{g}), \text{threshold}\right) \quad (9)$$

242 3.5. Decoding using the CRF/HMM model

243 The gesture sequence to recognize may contain an arbitrary number of
 244 gestures, in an arbitrary order. Therefore, the model should evenly switch
 245 between the gesture models. This can be modelled by gathering all the
 246 gesture model within a global sequence model, as shown in Figure 2. In this
 247 model, each line represents an isolated gesture, with a variable number of
 248 state. This global model allows to describe any arbitrary gesture sequence
 249 with equiprobable gesture transition probabilities.

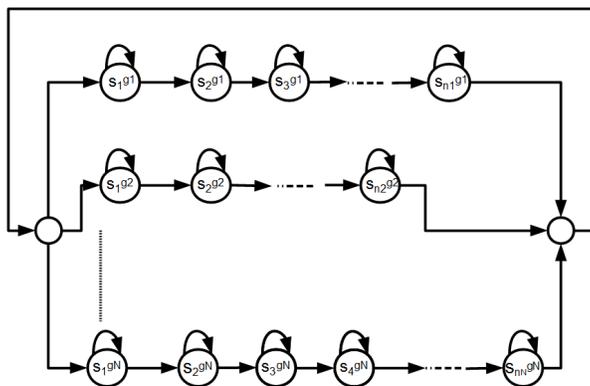


Figure 2: The recognition model of gesture sequences using HMM. $s_j^{g_i}$ represents the state j of the gesture i

250 3.6. CRF/HMM algorithm

251 Algorithm 1 summarizes the training and decoding process of our hybrid
 252 CRF/HMM model for gesture sequences recognition. Table 1 details algorithm
 253 functions and variables description. GSHOG characterisation function extract
 254 features from videos using Gesture Signature and HOG methods explained in

255 section 4. CRF/HMM procedure is furthermore explained by the diagram
 256 represented in figure 3.

```

Begin CRFHMMRecognition(Databases:videos):sequences
  for all Databases do
    for all GestClassVideos do
      | LFeatFile←GSHOG(GestClassVideo)
    end for
    BawmWelch(HMMPParams,LFeatFiles)
    UnifyHMMGaussVar(HMMPParams.GaussVar)
    LabGestClassVideos←Viterbi(HMMPParams,LabList,LFeatFile)
    QLFeatFiles←Quantify(LFeatFiles)
    LBFGS(CRFParams,QLfeatFiles,LabGestClassVideos)
  257 for all TestSeqVideos do
    | TFeatFile←GSHOG(TestSeqVideo)
  end for
  QTFeatFiles←Quantify(TFeatFiles)
  PosterioriProbas←ForwardBackward(CRFParams,TFeatFiles)
  GestSeqs←CRFHMMViterbi(HMMPParams.TransitionProbas,
    PosterioriProbas,LabList,SizeSeqs)
  end for
End
  
```

Algorithm 1: CRF/HMM learning and decoding algorithm

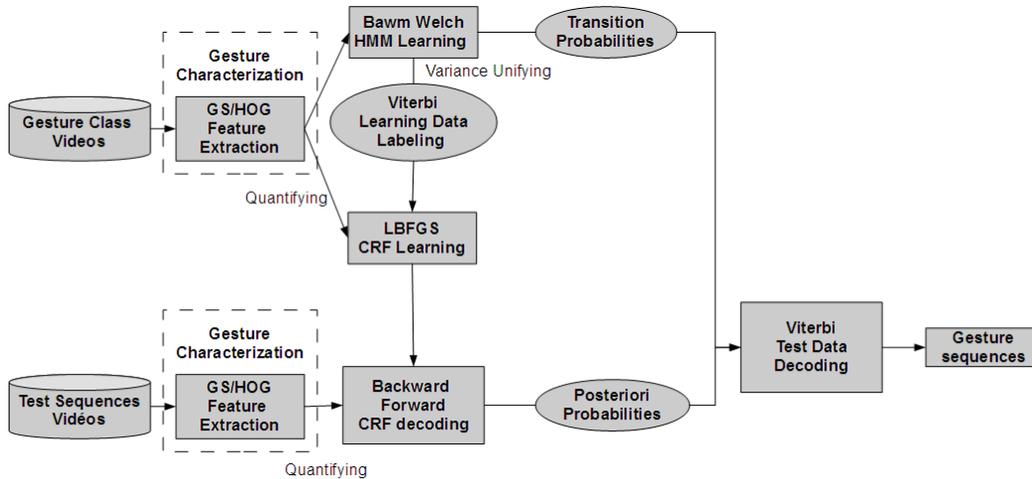


Figure 3: The CRF/HMM training and decoding process diagram

Function	Description
GSHOG	feature extraction (see section 4)
BawmWelch	HMM learning algorithm
UnifyHMMGaussVar	HMM adaptation
Viterbi	Learning videos labeling
Quantify	CRF adaptation
LBFGS	CRF learning method
ForwardBackward	posteriori probabilities computing
CRFHMMViterbi	Test sequences decoding
Variable	Description
Databases	learning and test databses
GestClassVideos	videos of gesture classes
LFeatFiles	feature files of gesture class videos
HMMPParams	HMM parameters
HMMPParams.GaussVar	HMM Gaussian variables
LabGestClassVideos	labeled video frames of gesture classes
LabList	gesture class labels
QLFeatFiles	quantified feature files of gesture class videos
CRFParams	CRF parameters
TestSeqVideos	videos of test sequences
TFeatFiles	feature files of test sequence videos
QTFeatFiles	quantified feature files of test sequence videos
PosterioriProbas	posteriori probabilities
HMMPParams.TransitionProbas	global transition probabilities
SizeSeqs	size of test sequences
GestSeqs	recognized gesture sequences

Table 1: Algorithm 1 Functions and variables description

258 4. Global gestures characterization

259 Gestures characterization requires velocity descriptors and shape descrip-
260 tors as well. Considering that signers can wear clothes in different colors
261 and have different skin colours, color descriptors are not included in our
262 characterization model.

263 In this section, we propose a second contribution presenting a model for
264 the gesture characterisation : a set of motion descriptors deduced from optical
265 flows velocities. We call this set of descriptors *Gesture Signature* (GS).

266 For a complete gesture characterization, we add global contour features
267 extracted with a classic shape descriptor; Histograms of Oriented Gradients
268 (HOG). To apply this descriptor, we resumed the implementation of Dalal et
269 al.[38]. 9 directions are used to quantify gradients inclination angles calculated
270 on the image. Such descriptors will account for shape descriptors.

271 4.1. Characterization with optical flows : *Gesture Signature*

272 Optical flows describe local velocities at the pixel level. They are known
273 for their robustness to brightness changes [39]. They are invariant to colors
274 and object distortion. Optical flows are able to describe simultaneously all
275 movements in the scene without any segmentation. Therefore, this method
276 seems adequate to simultaneously extract a maximum of information on body
277 motion, while being robust to variability of color, shape and brightness. In
278 what follows, we propose a feature vector whose components are combinations
279 of velocity values computed with optical flows.

280 Hand movements are usually located on the left and the right part of the
281 image, so it is advantageous to divide the image into two vertical sections.
282 Thus, the description of the movement is better localized and motions are
283 characterized in these two distinct regions.

284 Each part of the image is described by a gesture signature which consists
285 of 9 descriptors derived from positive and negative horizontal components
286 V_x^+ and V_x^- , and 9 descriptors derived from vertical components V_y^+ and V_y^- .
287 These components are derived from optical flows at each pixel of the image at
288 position p . Obviously, for each pixel p , two of these four values are null, one
289 pixel can have only one direction according to the x-axis and one direction
290 according to the y-axis.

291 For a given direction, these 9 descriptors consist of 4 movement *location*
292 descriptors, 2 movement *velocity* descriptors and 3 movement *orientation*
293 descriptors. Although these features are simple, they are complementary and

294 describe precisely the gesture changes since location, velocity and orientation
 295 are the main components of a gesture.

296 Table 2 shows the 18 features set (related variables are defined in table 3).
 297 The 8 horizontal and vertical location features are related to inertia center
 298 coordinates. They represent the vertical and horizontal positions of velocity
 299 centers with respect to the global movement of the considered portion of the
 300 image.

301 There are 4 features of movement velocity and strength. The first descrip-
 302 tor gives an energy information of the movement. It is inversely proportional
 303 to the quadratic mean of the moving pixels velocities. For normalization
 304 reasons, we use the inverse of this quadratic mean. The second descriptor
 305 gives information about the motion amplitude. It is the median of the moving
 306 pixels velocities. The median integrates information about the linear mo-
 307 mentum, where the mass is replaced in our case by the number of moving
 308 pixels. The median also reduces the noise effect. V_x^* and V_y^* components are
 309 the medians of a thresholded velocity vector which is computed with optical
 310 flows. Values of the threshold are given below.

$$S_{V_x} = \frac{\sum_{p=1}^{N_{px}^s} |V_x(p)|}{N_{px}^s} \quad (10)$$

$$S_{V_y} = \frac{\sum_{p=1}^{N_{px}^s} |V_y(p)|}{N_{px}^s} \quad (11)$$

311 The 6 movement orientation features are statistics on pixels moving in the
 312 same direction, positive or negative. The first two descriptors characterize
 313 the amount of pixels moving in the same direction. The third descriptor
 314 characterizes the dominant direction of the movement. Those three descriptors
 315 characterize the relationship or the symmetry between the two main movement
 316 groups whose orientations are opposite. Figure 4 shows the interest of these
 317 descriptors and illustrates the symmetry information. Thus, by analyzing
 318 the variation of these three descriptors, we can deduce the type of associated
 319 movement. Hence the importance and the complementarity of these three
 320 orientation descriptors.

321 5. Experimental protocol

322 In this section, we explain the experimental protocol : databases and
 323 evaluation methods

Table 2: The 8 movement **location** features, the 4 motion **velocity** features and the 6 movement **orientation** features of the *Gesture Signature* characterisation model.

	Descriptor	horizontally	vertically
Location	Average Abscissa of pixels moving in the Positive direction (AAP)	$\frac{1}{I_w} \times \frac{\sum_{p=1}^{N_{px}^+} v_x^+(p) x_p}{\sum_{p=1}^{N_{px}^+} v_x^+(p) }$	$\frac{1}{I_w} \times \frac{\sum_{p=1}^{N_{px}^+} v_y^+(p) x_p}{\sum_{p=1}^{N_{px}^+} v_y^+(p) }$
	Average Ordinate of pixels moving in the Positive direction (AOP)	$\frac{1}{I_h} \times \frac{\sum_{p=1}^{N_{px}^+} v_x^+(p) y_p}{\sum_{p=1}^{N_{px}^+} v_x^+(p) }$	$\frac{1}{I_h} \times \frac{\sum_{p=1}^{N_{px}^+} v_y^+(p) y_p}{\sum_{p=1}^{N_{px}^+} v_y^+(p) }$
	Average Abscissa of pixels moving in the Negative direction (AAN)	$\frac{1}{I_w} \times \frac{\sum_{p=1}^{N_{px}^-} v_x^-(p) x_p}{\sum_{p=1}^{N_{px}^-} v_x^-(p) }$	$\frac{1}{I_w} \times \frac{\sum_{p=1}^{N_{px}^-} v_y^-(p) x_p}{\sum_{p=1}^{N_{px}^-} v_y^-(p) }$
	Average Ordinate of pixels moving in the Negative direction (AON)	$\frac{1}{I_h} \times \frac{\sum_{p=1}^{N_{px}^-} v_x^-(p) y_p}{\sum_{p=1}^{N_{px}^-} v_x^-(p) }$	$\frac{1}{I_h} \times \frac{\sum_{p=1}^{N_{px}^-} v_y^-(p) y_p}{\sum_{p=1}^{N_{px}^-} v_y^-(p) }$
Velocity	Global Velocity Inverse (GVI)	$\sqrt{\frac{N_{px}}{\sum_{p=1}^{N_{px}} (v_x(p))^2}}$	$\sqrt{\frac{N_{px}}{\sum_{p=1}^{N_{px}} (v_y(p))^2}}$
	Maximum Velocities Median (MVM)	$\frac{1}{S_{V_X}} \times V_X^* $	$\frac{1}{S_{V_Y}} \times V_Y^* $
Orientation	Proportion of the Pixels moving in the Positive direction (PPP)	$PPP_X = \frac{N_{px}^+}{N_{px}}$	$PPP_Y = \frac{N_{px}^+}{N_{px}}$
	Proportion of the Pixels moving in the Negative direction (PPN)	$PPN_X = \frac{N_{px}^-}{N_{px}}$	$PPN_Y = \frac{N_{px}^-}{N_{px}}$
	Dominant Orientation (DO)	$DO_X = \frac{N_{px}^+ - N_{px}^-}{N_{px}}$	$DO_Y = \frac{N_{px}^+ - N_{px}^-}{N_{px}}$

Variable	Description
I_w	image width
I_h	image height
N_{px}	total pixel number
$N_{px}^{V_X^+}$	number of pixels moving in the positive horizontal direction
$N_{px}^{V_X^-}$	number of pixels moving in the negative horizontal direction
$N_{px}^{V_Y^+}$	number of pixels moving in the positive vertical direction
$N_{px}^{V_Y^-}$	number of pixels moving in the negative vertical direction
$V_X^+(p)$	positive horizontal velocity component of a pixel p
$V_X^-(p)$	negative horizontal velocity component of a pixel p
$V_Y^+(p)$	positive vertical velocity component of a pixel p
$V_Y^-(p)$	negative vertical velocity component of a pixel p
$V_X(p)$	horizontal velocity component of a pixel p
$V_Y(p)$	vertical velocity component of a pixel p
V_X^*	median of horizontal components (absolute value) of pixel velocities
V_Y^*	median of vertical components (absolute value) of pixel velocities
S_{V_X}, S_{V_Y}	velocity thresholds (see equations 10 and 11)
PPP_X	Proportion of the Pixels moving in the Positive horizontal direction
PPN_X	Proportion of the Pixels moving in the Negative horizontal direction
PPP_Y	Proportion of the Pixels moving in the Positive vertical direction
PPN_Y	Proportion of the Pixels moving in the Negative vertical direction

Table 3: Gesture Signature variables description

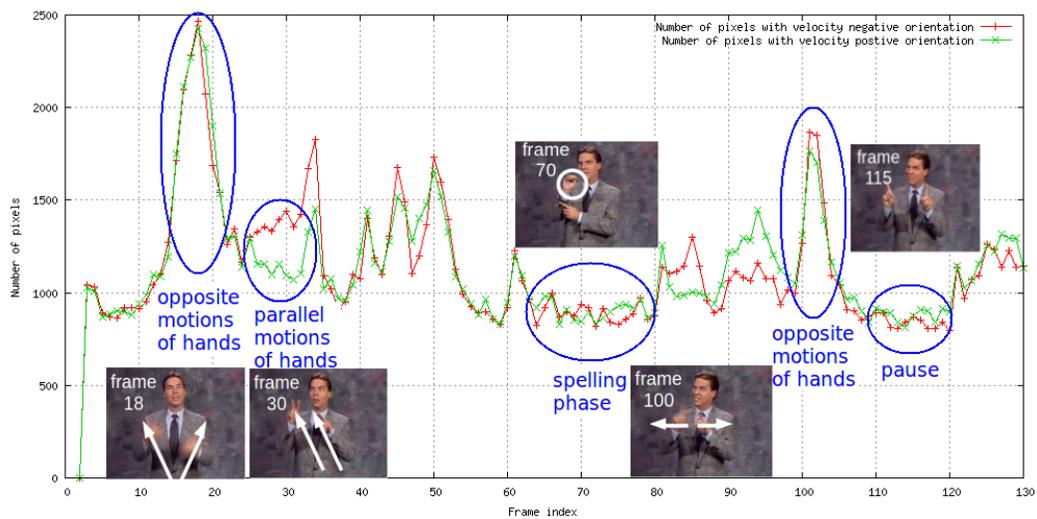


Figure 4: Evolution of the descriptors PPP_x (Proportion of the Pixels moving in the Positive horizontal direction) and PPN_x (Proportion of the Pixels moving in the Negative horizontal direction) in a video from SignStream database [40]. Two curves superimposed with a presence of a peak correspond to an opposite movement of the two hands. A strong difference between the two curves correspond to a parallel movement of both hands in the dominant direction. A stagnation of the two curves correspond to fixed hands (frame 70).

324 *5.1. Databases*

325 Our recognition system has been evaluated on public databases designed
326 for the ChaLearn 2011-2012 competition [5]. We did not participate to
327 this competition but we were able to compare our system to those of the
328 participants thanks to the evaluation platform proposed by the competition
329 organizers ¹. We detail the results of this evaluation in section 6.

330 ChaLearn databases are made of three types of resources: 480 system
331 development sub-databases named *devel*, 20 system validation sub-databases
332 named *valid* and 40 system final evaluation sub-databases named *final*. The
333 1-20 *final* sub-databases were tested in the first round of the competition and
334 21-40 *final* sub-databases were tested in the second round of the competition.
335 This final evaluation classifies participants in the ChaLearn competition.

336 Each of these sub-databases contains 47 pairs of videos. Each video pair
337 presents the same scene in two formats: RGB color format and depth format.
338 These videos are recorded using a Kinect (TM) camera. Videos of the same
339 sub-database share the same scenic features: same actor, same background,
340 same recording conditions, same theme and same gesture vocabulary. However,
341 these scenic characteristics vary from sub-database to another. 20 players
342 participated in the making of these databases, one actor per sub-database.
343 These databases present 30 vocabularies composed of 8-15 gestures belonging
344 to various themes such as video games, distance education, robot control,
345 sign language, etc.

346 Each sub-database includes two sets of video: a training set \mathbb{G} and a test
347 set \mathbb{S} . The training set \mathbb{G} consists of 10 videos. Each video contains a single
348 and isolated instance of a gesture: *one-shot learning* databases. The test set
349 \mathbb{S} consists of 40 videos. Each video includes a sequence of 1 to 5 successive
350 gestures separated by a common break point. Gestures organization in each
351 test sequences is random, there is no specific gestures grammar.

352 We summarize in the following subsection the various feature vectors used
353 for the tests.

354 *5.2. Feature vector variants*

355 Table 4 presents the different variants of the feature vector \vec{c} we used
356 in our experiments. We index each variant by its size $l(\vec{c})$. $l(\vec{v}(GS))$ is the
357 number of gesture signature features. $l(\vec{v}(HOG))$ is the number of HOG

¹<https://www.kaggle.com/c/GestureChallenge2>

358 features. Some variants of the feature vector \vec{c} are applied to two data formats
 359 (RGB image and depth image).

Table 4: Feature vector variants adopted in the experiments

total size $l(\vec{c})$	Descriptor			
	Gesture Signature GS		HOG	
	$l(\vec{c}(GS))$	description	$l(\vec{c}(HOG))$	description
52	16	no median, no image division	36	4 image blocks \times 9 gradient directions
54	18	no image division	36	4 image blocks \times 9 gradient directions
72	72	image division into 2 parts, 2 data formats	0	HOG not applied
180	36	image division into 2 parts	144	16 image blocks \times 9 gradient directions
360	72	image division into 2 parts, 2 data formats	288	16 image blocks \times 9 gradient directions, 2 data formats

360 5.3. Evaluation metric

361 The organizers of the ChaLearn competition defined a global evaluation
 362 metric on all test sequences based on the Levenshtein distance, also called
 363 edit distance [41]. This form of global error is denoted by \mathcal{L}_{ch} and given by
 364 equation 12.

$$\begin{aligned}
 \mathcal{L}_{ch} : \mathbb{D} &\longrightarrow \mathbb{R} & (12) \\
 \mathbb{S} &\longmapsto \frac{\sum_{\mathbf{s} \in \mathbb{S}} L(\mathcal{R}(\mathbf{s}), \mathcal{T}(\mathbf{s}))}{\sum_{\mathbf{s} \in \mathbb{S}} l(\mathcal{T}(\mathbf{s}))}
 \end{aligned}$$

365 where \mathbb{D} is the set of test databases, \mathbb{S} is the set of test sequences, \mathbf{s} is the
 366 sequence of gestures, $\mathcal{R}(\mathbf{s})$ is the system recognition result of sequence \mathbf{s} , \mathcal{T}
 367 is a function giving the ground truth sequence \mathbf{s} , $L(.,.)$ is the Levenshtein
 368 distance and $l(v)$ gives the size of a vector v .

369 We use the ChaLearn form of the error \mathcal{L}_{ch} to compare our recognition
 370 system to ChaLearn participants recognition systems. However, let us empha-
 371 size that \mathcal{L}_{ch} is slightly different from the classical Levenshtein distance (see

372 Equation 13), which is bounded and seems more generic. Thus, to present
 373 the main results of our various tests, we use the classic error form.

$$\begin{aligned} \mathcal{L} : \mathbb{D} &\longrightarrow [0, 1] \\ \mathbb{S} &\longmapsto \frac{1}{|\mathbb{S}|} \sum_{s \in \mathbb{S}} \frac{L(\mathcal{R}(s), \mathcal{T}(s))}{l(\mathcal{R}(s)) + l(\mathcal{T}(s))} \end{aligned} \quad (13)$$

374 6. Gesture recognition results

375 In this section we present the results of our system, using different variants.

376 We first compare the recognition results of the hybrid system CRF/HMM
 377 to the classic and adapted versions of HMM and CRF in subsection 6.1. Then,
 378 we present our rank compared to participants at the ChaLearn competition.
 379 Next, in subsection 6.2, we present some properties of the hybrid model
 380 CRF/HMM including its robustness with respect to the number of states and
 381 to the various feature vectors, and we conclude this section by demonstrating
 382 the advantage of the gesture signature model.

383 All recognition performance results of the hybrid system CRF/HMM
 384 presented in this section are obtained with tests performed with an adapted
 385 CRF/HMM as explained in section 3.3 unless otherwise stated. Adapted
 386 HMM and adapted CRF recognition systems cited in this section are also
 387 adapted as explained in section 3.3.

388 6.1. Evaluation of the CRF/HMM using the ChaLearn platform

389 We present in this subsection the recognition results of our best hybrid
 390 system CRF/HMM on the *valid* and *final* databases, as well as our ranking
 391 in the ChaLearn competition.

392 We first present a comparison of the performance of the main recognition
 393 systems that we studied (Table 5) on the *devel* databases. The feature vector
 394 is identical for all the systems ($l(\vec{c}) = 52$). The number of frames per state $\mathbf{f}_{\mathbf{g}}$
 395 is optimized for each system. $\mathbf{f}_{\mathbf{g}}(\mathbf{g})$ represents the size of the learned gesture,
 396 which means that every gesture is represented by a single class, subclasses
 397 that correspond to states in the case of HMM do not exist in the case of CRF.
 398 On the other hand, a post-processing step (see algorithm 2) is applied to the
 399 classic and adapted CRF in order to filter their recognition results. Without
 400 this step recognition error exceeds 0.5. Table 5 shows that the performance of
 401 the proposed hybrid system CRF/HMM clearly outperform the recognition
 402 performances of other systems. Indeed, for a data size equal to 750, we
 403 demonstrated with the statistical unilateral Student test that our hybrid

404 model CRF/HMM significantly outperforms classic and adapted HMM and
 405 CRF models with a confidence level detailed in table 5.

```

Begin CRFpostProcessing(RecognizedSeqs):FilteredSeqs
  for all RecognizedSeqs  $S_i$  do
    for all gestures  $G_j$  do
      Fix window size  $Sz(w_j) = \begin{cases} \frac{3}{4}Sz(G_j) & \text{if } 1 < \frac{Sz(S_i)}{Sz(G_j)} \\ Sz(S_i) & \text{if } 0 < \frac{Sz(S_i)}{Sz(G_j)} \leq 1 \end{cases}$ 
      for all shifted window positions (shifting step is  $Sz(w_j)$ ) do
        Search for the most occurent gesture  $G_m$  in the current window
        if ( $G_m = G_j$ ) then
          save  $G_m$  at the current window position (if conflict, keep the shortest gesture)
        end
      end for
    end for
    FilteredSeq  $\leftarrow$  all saved  $G_m$ 
  end for
End
  
```

Algorithm 2: Post-filtering and segmentation algorithm for CRF recognized sequences, where *FiltredSeqs* are filtered and segmented sequences extracted from *RecognizedSeqs* which are the recognized sequences with CRF Backward-Forward method

Table 5: a) The recognition results of various recognition systems based on HMM and CRF and tested on 20 *devel* databases (\mathbf{f}_s is optimized for each system, $l(\vec{c}) = 52$ and images are RGB). b) Unilateral student test results (confidence level): CRF/HMM compared to classic and adapted HMM and CRF models.

System	(a) Error : \mathcal{L}	(b) Student confidence level (%)
classic HMM	0.3615	99.95
adapted HMM	0.2354	85
classic CRF (continuous)	0.2930	99.95
adapted CRF (discrete)	0.2757	99.95
CRF/HMM (adapted)	0.2228	-

407 In order to rank our system in the ChaLearn 2011-2012 competition,

408 we tested the hybrid system on *valid* and *final* databases provided during
 409 the competition. Table 6 shows the hybrid system CRF/HMM recognition
 410 error values computed with both evaluation methods \mathcal{L} and \mathcal{L}_{ch} on *valid* and
 411 *final* databases. Table 6 presents the CRF/HMM system rank on both
 412 database categories using the \mathcal{L}_{ch} error. It appears that we ranked at the 7th
 413 position among 559 systems from 48 participants for both first and second
 414 rounds. The complete list with their score (the \mathcal{L}_{ch} error) is available on the
 415 Kaggle website for the first² and the second round³. We achieved this rank
 416 using only RGB format data.

Table 6: The recognition results of our best hybrid system CRF/HMM on 20 *valid* databases, 20 *final* 1-20 databases and 20 *final* 21-40 databases (each database category contains about 750 total sequences test in the order of 200 frames each, images are RGB and $l(\vec{c}) = 180$)

database category	Error		ranking
	\mathcal{L}	\mathcal{L}_{ch}	
<i>valid</i>	0.1772	0.3488	-
<i>final</i> 1-20 (1 st round)	0.1479	0.2964	7 th
<i>final</i> 21-40 (2 nd round)	0.1224	0.2523	7 th

417 6.2. Properties of the CRF/HMM system

418 6.2.1. Robustness to changes in the number of frames per state

419 Figure 5 shows the recognition error \mathcal{L} of the adapted HMM and the
 420 CRF/HMM systems with respect to the number of frames per state \mathbf{f}_s .
 421 Those systems were trained for each value of \mathbf{f}_s . One can observe that the
 422 CRF/HMM system outperforms the HMM, and that the CRF/HMM system
 423 provides extremely stable results, while the performance of the HMM is
 424 strongly variable. This is an interesting aspect of this system since it does
 425 not require a fine tuning of the hyper-parameter for reaching good results.

426 6.2.2. Robustness to changes in the gesture duration

427 The average number of frames per state has a direct impact on the
 428 CRF/HMM robustness to the gesture duration variation. With a large
 429 number of frames per state, the CRF/HMM system is able to handle the

²<https://www.kaggle.com/c/GestureChallenge/leaderboard>

³<https://www.kaggle.com/c/GestureChallenge2/leaderboard>

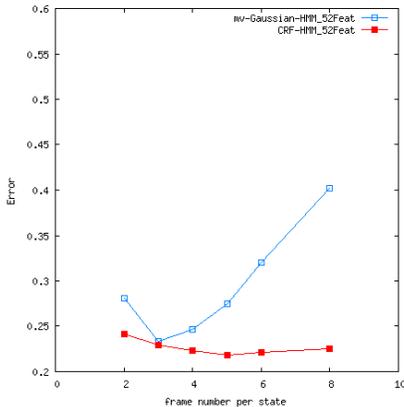


Figure 5: CRF/HMM and adapted HMM systems robustness to the variation of the number of frames per state.

Table 7: Recognition results with different variants of the feature vector on 20 *devel* databases (RGB image and depth image)

System	Error : \mathcal{L}	
	GS $l(\vec{c}) = 72$	(GS, HOG) $l(\vec{c}) = 360$
adapted HMM	0.2525	0.2425
CRF/HMM	0.2559	0.2255

430 temporal elasticity of a gesture. In other words, when a gesture expands
 431 or narrows through the number of frames in the test data, the CRF/HMM
 432 system is able to align the gesture model on the data and decode them. In
 433 addition, the CRF component is able to implicitly manage duration variation
 434 of gestures in a more straightforward way than HMMs do. It appears that
 435 the temporal elasticity of gestures is better captured by the simple structure
 436 of the hybrid model with a reduced number of states, compared to the totally
 437 connected structure of the HMM system alone as adopted by some participants
 438 of the ChaLearn competition [10, 11, 12].

439 6.2.3. Robustness to changes in the feature vector

440 Figures 6 present the variation of the error \mathcal{L} in terms of the number of
 441 frames per state f_s for two HMM systems (left) and for two CRF/HMM
 442 systems (right). Each pair of systems is evaluated with two different feature
 443 vectors. When the feature vector size decreases, CRF/HMM keep almost the
 444 same performance. In other words, a minimum of features is sufficient for
 445 CRF HMM, whereas for classic HMM, feature addition greatly increases the
 446 recognition performance. This recognition ability with a reduced number of
 447 features makes features extraction task easier and faster.

448 These three CRF/HMM robustness property prove that with a simple
 449 system, it is possible to reach high recognition performance thanks to CRF
 450 and HMM advantages combination and disadvantages compensation. We

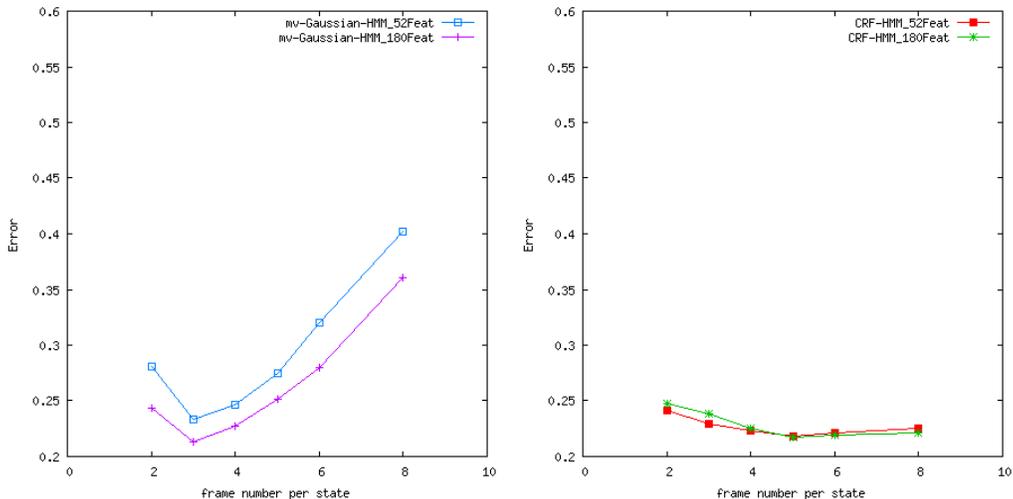


Figure 6: Adapted HMM (left) and CRF/HMM (right) robustness to the variation of the feature vector

451 can see the simplicity of the CRF/HMM system at three levels: (a) a simple
 452 model structure with a reduced number of state without jumps nor complete
 453 connection. (b) a reduced number of features. (c) a training data set reduced
 454 to an example by class.

455 6.2.4. Evaluation of the Gesture Signature

456 Table 7 shows the recognition results of two systems, adapted HMM and
 457 CRF/HMM, on *devel* databases applying three variants of the feature vector.
 458 The purpose of these tests is not to compare these two recognition systems
 459 but to validate the interest of the feature vector GS . According to table 7, we
 460 notice that the performance of recognition systems with the feature vector GS
 461 is very close to the performance of these recognition systems using a feature
 462 vector that combines GS features and the HOG features. Moreover, these
 463 error values are low and exhibit valuable recognition performance. Thus, the
 464 gesture signature GS can represent a complete characterization model.

465 Finally, these results and this study show that the CRF/HMM hybrid
 466 system is a system that has better performance than other classic systems
 467 (HMM and CRF), is robust to different variations and is interesting and
 468 practical in the real-world problem such as one shot learning.

469 7. Conclusion

470 In this article, an hybrid CRF/HMM system for gesture recognition has
471 been proposed. This HMM and CRF combination benefits from each model
472 advantages without undergoing its drawbacks. These systems have been
473 adapted to the one-shot learning context in order to suit to the real-world
474 constraints of small labelled datasets.

475 A new gesture characterization model has also been proposed, which is a
476 gesture signature that rely on optical flows. This model is able to describe
477 any dynamic scene using its motion, making it independent from the moving
478 object type.

479 We demonstrate that the CRF/HMM system are able to efficiently model
480 and manage spatio-temporal variations of sequential data and constitute a
481 robust recognition hybrid system that opens up new perspectives for sequential
482 Markov models. Among them, an interesting perspective concerns the gesture
483 detection task, called the gesture *spotting*, which consists on locating and
484 labelling specific gestures in videos. It can be applied in video retrieval and
485 indexing context. Our recognition model could be adapted to the spotting
486 task by representing false examples through an additional class to the gestures
487 vocabulary.

488 References

- 489 [1] K. M. Sayre, Machine recognition of handwritten words: A project report,
490 Pattern Recognition 5 (3) (1973) 213 – 228.
- 491 [2] L. R. Rabiner, A tutorial on hidden markov models and selected applica-
492 tions in speech recognition, Proceedings of the IEEE (1989) 257–286.
- 493 [3] A. Viterbi, Error bounds for convolutional codes and an asymptotically
494 optimum decoding algorithm, IEEE Transactions on Information Theory
495 13 (2) (1967) 260–269.
- 496 [4] J. Lafferty, A. McCallum, F. Pereira, Conditional random fields: Proba-
497 bilistic models for segmenting and labeling sequence data, ICML (2001)
498 282–289.
- 499 [5] I. Guyon, V. Athitsos, P. Jangyodsuk, B. Hamner, H. Escalante, Chalearn
500 gesture challenge: Design and first results., CVPR Workshops (2012)
501 1–6.

- 502 [6] Y. Yang, I. Saleemi, M. Shah, Discovering motion primitives for un-
503 supervised grouping and one-shot learning of human actions, gestures,
504 and expressions, *IEEE Transactions on Pattern Analysis and Machine*
505 *Intelligence* 35 (7) (2013) 1635–1648.
- 506 [7] C. Vogler, D. Metaxas, A framework for recognizing the simultaneous
507 aspects of american sign language, *Computer Vision and Image Under-*
508 *standing* 81 (2001) 358–384.
- 509 [8] U. von Agris, J. Zieren, U. Canzler, B. Bauer, K. K-F, Recent develop-
510 ments in visual sign language recognition, *Univers. Access Inf. Soc.* 6 (4)
511 (2008) 323–362.
- 512 [9] S. Ong, S. C. W. and Ranganath, Deciphering gestures with layered
513 meanings and signer adaptation, *IEEE International Conference on*
514 *Automatic Face and Gesture Recognition* (2004) 559.
- 515 [10] J. Konencny, M. Hagara, One-shot learning gesture recognition using
516 hog/hof features, *ICPR 2012 Workshop on Gesture Recognition*.
- 517 [11] E. Jackson, An hmm-based approach for gesture recognition using edge
518 features, *CVPR 2012 Workshop on Gesture Recognition*.
- 519 [12] D. Weiss, Hmm based one shot gesture recognition, *CVPR 2012 Work-*
520 *shop on Gesture Recognition*.
- 521 [13] E. Trentin, A survey of hybrid ANN/HMM models for automatic speech
522 recognition, *Neurocomputing* (2001) 91–126.
- 523 [14] N. Morgan, H. Bourlard, S. Renls, M. Cohen, H. Franco, Hybrid neural
524 network/hidden markov model systems for continuous speech recognition,
525 *IJPRAI* 7 (4).
- 526 [15] L. T. Niles, H. F. Silverman, Combining hidden Markov models and
527 neural network classifiers, in: *ICASSP, 1990*, pp. 417–420.
- 528 [16] J. Tebelskis, A. Waibel, B. Petek, O. Schmidbauer, Continuous speech
529 recognition by linked predictive neural networks, *NIPS* (1990) 199–205.
- 530 [17] F. Johansen, A comparison of hybrid hmm architectures using global
531 discriminative training., *ICSLP*.

- 532 [18] G. Rigoll, Maximum mutual information neural networks for hybrid
533 connectionist-hmm speech recognition systems., *IEEE Transactions on*
534 *Speech and Audio Processing* 2 (1) (1994) 175–184.
- 535 [19] G. Zavaliagos, S. Austin, J. Makhoul, R. Schwartz, A hybrid continuous
536 speech recognition system using segmental neural nets with hidden
537 markov models., *IJPRAI* 7 (4) (1993) 949–963.
- 538 [20] S. Thomas, C. Chatelain, L. Heutte, T. Paquet, Y. Kessentini, A deep
539 hmm model for multiple keywords spotting in handwritten documents,
540 *Pattern Analysis and Applications* 18 (4) (2015) 1003–1015.
- 541 [21] Y. Bengio, Y. LeCun, C. Nohl, C. Burges, Lerec: A nn/hmm hybrid for
542 on-line handwriting recognition, *Neural Comput.* 7 (6) (1995) 1289–1303.
- 543 [22] S. Knerr, E. Augustin, A neural network-hidden markov model hybrid
544 for cursive word recognition, *ICPR* 2 (1998) 1518–1520.
- 545 [23] S. Marukatat, T. Artieres, P. Gallinari, B. Dorizzi, Sentence recognition
546 through hybrid neuro-markovian modeling., in: *ICDAR, 2001*, pp. 731–.
- 547 [24] M. Gilloux, B. Lemarie, M. Leroux, A hybrid rbf network/hidden markov
548 model handwritten word recognition system., in: *ICDAR, 1995*, pp. 394–
549 397.
- 550 [25] M. Morita, R. Sabourin, F. Bortolozzi, C. Suen, Segmentation and
551 recognition of handwritten dates: an hmm-mlp hybrid approach., *IJDAR*
552 6 (4) (2003) 248–262.
- 553 [26] O. Matan, C. Burges, Y. Lecun, J. Denker, Multi-digit recognition using
554 a space displacement neural network, in: *Advances in Neural Information*
555 *Processing Systems* 4, 1992, pp. 488–495.
- 556 [27] A. Corradini, Real-time gesture recognition by means of hybrid recogniz-
557 ers., in: *Gesture Workshop, Vol. 2298*, Springer, 2001, pp. 34–46.
- 558 [28] A. Ganapathiraju, J. Hamaker, J. Picone, Hybrid svm/hmm architectures
559 for speech recognition., in: *INTERSPEECH, ISCA, 2000*, pp. 504–507.
- 560 [29] S. Rajko, G. Qian, A hybrid hmm/dpa adaptive gesture recognition
561 method., in: *ISVC, Vol. 3804*, 2005, pp. 227–234.

- 562 [30] Y. Soullard, T. Artieres, Hybrid hmm and hcrf model for sequence
563 classification, European Symposium on Artificial Neural Networks, Com-
564 putational Intelligence and Machine Learning.
- 565 [31] A. Gunawardana, M. Mahajan, A. Acero, J. Platt, Hidden conditional
566 random fields for phone classification., in: INTERSPEECH, ISCA, 2005,
567 pp. 1117–1120.
- 568 [32] S. Austin, R. Schwartz, P. Placeway, The forward-backward search
569 algorithm, IEEE ICASSP (1991) 697–700.
- 570 [33] D. Wu, F. Zhu, L. Shao, One shot learning gesture recognition from rgb
571 images., CVPR, IEEE (2012) 7–12.
- 572 [34] C. Sutton, A. McCallum, An introduction to conditional random fields,
573 Foundations and Trends in Machine Learning 4 (4) (2012) 267373.
- 574 [35] V. Radosavljevic, S. Vucetic, Z. Obradovic, Continuous conditional
575 random fields for regression in remote sensing, in: Proceedings of the
576 2010 Conference on ECAI 2010: 19th European Conference on Artificial
577 Intelligence, IOS Press, Amsterdam, The Netherlands, The Netherlands,
578 2010, pp. 809–814.
- 579 [36] D. Hebert, T. Paquet, S. Nicolas, Continuous crf with multi-scale quan-
580 tization feature functions application to structure extraction in old news-
581 paper, in: ICDAR, 2011, pp. 493–497.
- 582 [37] A. Quattoni, S. Wang, L. Morency, M. Collins, T. Darrell, Hidden
583 conditional random fields, Pattern Analysis and Machine Intelligence,
584 IEEE Transactions on 29 (10) (2007) 1848–1852.
- 585 [38] N. Dalal, B. Triggs, Histograms of oriented gradients for human detection,
586 in: International Conference on Computer Vision & Pattern Recognition,
587 Vol. 2, 2005, pp. 886–893.
- 588 [39] S. M. Bhandarkar, X. Luo, Integrated detection and tracking of multiple
589 faces using particle filtering and optical flow-based elastic matching,
590 CVIU 113 (6) (2009) 708–725.
- 591 [40] C. Neidle, S. Sclaroff, V. Athitsos, Signstream: A tool for linguistic
592 and computer vision research on visual-gestural language data, Behavior
593 Research Methods, Instruments, & Computers 33 (3) (2001) 311–320.

594 [41] V. Levenshtein, Binary Codes Capable of Correcting Deletions, Insertions
595 and Reversals, Soviet Physics Doklady 10 (1966) 707.