

Document downloaded from:

<http://hdl.handle.net/10251/55613>

This paper must be cited as:

Ihab Alkhoury; Giménez Pastor, A.; Alfons Juan; Andrés Ferrer, J. (2015). Window repositioning for Printed Arabic Recognition. *Pattern Recognition Letters*. 51:86-93. doi:10.1016/j.patrec.2014.08.009.



The final publication is available at

<http://dx.doi.org/10.1016/j.patrec.2014.08.009>

Copyright Elsevier

Additional Information

Window repositioning for Printed Arabic Recognition

Ihab Khoury, Adrià Giménez, Alfons Juan and Jesús Andrés-Ferrer

Universitat Politècnica de València, Camí de Vera s/n, 46022 València, Spain

Abstract

Bernoulli HMMs are conventional HMMs in which the emission probabilities are modeled with Bernoulli mixtures. They have recently been applied, with good results, in off-line text recognition in many languages, in particular, Arabic. A key idea that has proven to be very effective in this application of Bernoulli HMMs is the use of a sliding window of adequate width for feature extraction. This idea has allowed us to obtain very competitive results in the recognition of both Arabic handwriting and printed text. Indeed, a system based on it ranked first at the ICDAR 2011 Arabic recognition competition on the Arabic Printed Text Image (APTI) database. More recently, this idea has been refined by using *repositioning* techniques for extracted windows, leading to further improvements in Arabic handwriting recognition. In the case of printed text, this refinement led to an improved system which ranked second at the ICDAR 2013 second competition on APTI, only at a marginal distance from the best system. In this work, we describe the development of this improved system. Following evaluation protocols similar to those of the competitions on APTI, exhaustive experiments are detailed from which

Email address: {ialkhoury, agimenez, ajuan, jandres}@dsic.upv.es (Ihab Khoury, Adrià Giménez, Alfons Juan and Jesús Andrés-Ferrer)

state-of-the-art results are obtained.

Keywords: Bernoulli HMMs, Printed Arabic Recognition, Sliding Window, Repositioning

1 Introduction

Hidden Markov Models (HMMs) are now widely used for off-line text recognition in many languages, in particular, languages with Arabic script (Dehghan et al., 2001; Günter and Bunke, 2004; Märgner and El Abed, 2007, 2009; Grosicki and El Abed, 2009). Following the conventional approach in speech recognition (Rabiner and Juang, 1993), HMMs at global (line or word) level are built from shared, *embedded*, HMMs at character (subword) level, which are usually simple in terms of number of states and topology. In the common case of real-valued feature vectors, state-conditional probability (density) functions are modeled as Gaussian mixtures since, as with finite mixture models in general, their complexity can be easily adjusted to the available training data by simply varying the number of components.

After decades of research in speech recognition, the use of certain real-valued speech features and embedded Gaussian (mixture) HMMs is a de-facto standard (Rabiner and Juang, 1993). However, in the case of text recognition there is no such standard. In fact, very different sets of features are in use today. In (Giménez and Juan, 2009) we proposed to by-pass feature extraction and directly feed columns of raw, binary pixels into *embedded Bernoulli (mixture) HMMs (BHMMs)*, that is, embedded HMMs in which the emission probabilities are modeled with Bernoulli mixtures. The basic idea is to ensure that no discriminative information is filtered out during feature

22 extraction, which in some sense is integrated into the recognition model.
23 In (Giménez et al., 2010), we improved our basic approach by using a sliding
24 window of adequate width to better capture image context at each horizontal
25 position of the text image. This improvement, to which we refer as *windowed*
26 *BHMMs*, achieved very competitive results on the well-known IfN/ENIT
27 database of Arabic town names (Märgner and El Abed, 2010). More recently,
28 very good results on the Arabic Printed Text Image (APTI) database were
29 also achieved using the same approach, which ranked first in the ICDAR
30 2011 Arabic recognition competition for printed Arabic text (Slimane et al.,
31 2011).

32 Although windowed BHMMs achieved good results on IfN/ENIT and
33 APTI, it was clear to us that text distortions are more difficult to model
34 with wide windows than with narrow (e.g. one-column) windows. In order
35 to circumvent this difficulty, we have considered new, adaptive window sam-
36 pling techniques, as opposed to the conventional, direct strategy in which
37 the sampling window center is applied at a constant height of the text image
38 and moved horizontally one pixel at a time. More precisely, these adaptive
39 techniques can be seen as an application of the direct strategy followed by a
40 *repositioning* step by which the sampling window is repositioned to align its
41 center to the center of gravity of the sampled image. This repositioning step
42 can be done horizontally, vertically or in both directions. Although vertical
43 repositioning is expected to have more influence on recognition results than
44 horizontal repositioning, we have studied both separately and in conjunction,
45 so as to confirm this expectation.

46 In (Giménez et al., 2014b), the repositioning techniques described above

47 are introduced and extensively tested on different databases for off-line hand-
 48 writing recognition. As expected, vertical repositioning provides excellent re-
 49 sults, not only on IfN/ENIT, but also on other well-known databases such as
 50 IAM words and RIMES. In the case of printed text, the use of repositioning
 51 techniques has allowed us to significantly improve our system at the ICDAR
 52 2011 first competition on APTI. Indeed, our improved system obtained much
 53 better results at the ICDAR 2013 second competition on APTI, in which it
 54 ranked second at a marginal distance from the first (Slimane et al., 2013). In
 55 this work, we describe the development of this improved system. Following
 56 evaluation protocols similar to those of the competitions on APTI, exhaustive
 57 experiments are described from which state-of-the-art results are obtained.

58 In what follows, we first review BHMMs (Sec. 2). Then, we describe
 59 the approach through which we are achieving the best results: windowed
 60 BHMMs with repositioning (Sec. 3) and its use for printed Arabic recognition
 61 by application of the Bayes decision rule (Sec. 4). In Sec. 5, we provide the
 62 results of a complete series of experiments on APTI as well as a comparison
 63 with results from other authors on this database. Finally, concluding remarks
 64 are given in Sec. 6.

65 2. Bernoulli HMMs

66 Let $O = (\mathbf{o}_1, \dots, \mathbf{o}_T)$ be a sequence of feature vectors. An HMM is a
 67 probability (density) function of the form:

$$P(O | \Theta) = \sum_{q_1, \dots, q_T} \prod_{t=0}^T a_{q_t q_{t+1}} \prod_{t=1}^T b_{q_t}(\mathbf{o}_t), \quad (1)$$

68 where the sum is over all possible *paths* (state sequences) q_0, \dots, q_{T+1} , such
69 that $q_0 = I$ (special *initial* or *start* state), $q_{T+1} = F$ (special *final* or *stop*
70 state), and $q_1, \dots, q_T \in \{1, \dots, M\}$, being M the number of regular (non-
71 special) states of the HMM. On the other hand, for any regular states i and j ,
72 a_{ij} denotes the *transition* probability from i to j , while b_j is the *observation*
73 probability (density) function at j .

74 A Bernoulli (mixture) HMM (BHMM) is an HMM in which the probabil-
75 ity of observing a binary feature vector \mathbf{o}_t , when $q_t = j$, follows a Bernoulli
76 mixture distribution for the state j

$$b_j(\mathbf{o}_t) = \sum_{k=1}^K \pi_{jk} \prod_{d=1}^D p_{jkd}^{o_{td}} (1 - p_{jkd})^{1-o_{td}}, \quad (2)$$

77 where o_{td} is the d -th bit of \mathbf{o}_t , π_{jk} is the prior of the k -th mixture component
78 in state j , and p_{jkd} is the probability that this component assigns to o_{td} to
79 be 1.

80 As discussed in the introduction, BHMMs at global (line or word) level
81 are built from shared, embedded BHMMs at character level. More precisely,
82 let C be the number of different characters (symbols) from which global
83 BHMMs are built, and assume that each character c is modeled with a dif-
84 ferent BHMM of parameter vector Θ_c . Let $\Theta = \{\Theta_1, \dots, \Theta_C\}$, and let
85 $O = (\mathbf{o}_1, \dots, \mathbf{o}_T)$ be a sequence of feature vectors generated from a sequence
86 of symbols $S = (s_1, \dots, s_L)$, with $L \leq T$. The probability of O can be calcu-
87 lated, using embedded HMMs for its symbols, as:

$$P(O | S, \Theta) = \sum_{i_1, \dots, i_{L+1}} \prod_{l=1}^L P(\mathbf{o}_{i_l}, \dots, \mathbf{o}_{i_{l+1}-1} | \Theta_{s_l}), \quad (3)$$

88 where the sum is carried out over all possible segmentations of O into L

89 segments, that is, all sequences of indices i_1, \dots, i_{L+1} such that

$$1 = i_1 < \dots < i_L < i_{L+1} = T + 1;$$

90 and $P(\mathbf{o}_{i_l}, \dots, \mathbf{o}_{i_{l+1}-1} \mid \Theta_{s_l})$ refers to the probability (density) of the l -th
91 segment, as given by (1) using the HMM associated with symbol s_l .

92 Maximum likelihood estimation (MLE) of BHMM parameters does not
93 differ significantly from the conventional Gaussian case, and it can be effi-
94 ciently performed using the well-known EM (Baum-Welch) re-estimation for-
95 mulae (Rabiner and Juang, 1993; Young et al., 1995). Please see (Giménez
96 et al., 2014b) for more details. Also as in the conventional Gaussian case,
97 BHMM parameters can be estimated by discriminative training (Giménez
98 et al., 2014a).

99 3. Windowed BHMMs with repositioning

100 Given a binary image normalized in height to H pixels, we may think of a
101 feature vector \mathbf{o}_t as its column at position t or, more generally, as a concate-
102 nation of columns in a window of W columns in width, centered at position
103 t . This generalization has no effect neither on the definition of BHMM nor
104 on its MLE, although it might be very helpful to better capture the image
105 context at each horizontal position of the image. As an example, the first
106 row in Fig. 1 shows a binary image of 4 columns and 5 rows, which is trans-
107 formed into a sequence of four 15-dimensional feature vectors by application
108 of a sliding window of width 3. For clarity, feature vectors are depicted as
109 3×5 subimages instead of 15-dimensional column vectors. Note that feature
110 vectors at positions 2 and 4 would be indistinguishable if, as in our previous
111 approach, they were extracted with no context ($W = 1$).

112 Although one-dimensional, “horizontal” HMMs for image modeling can
 113 properly capture non-linear horizontal image distortions, they are somewhat
 114 limited when dealing with vertical image distortions, and this limitation
 115 might be particularly strong in the case of feature vectors extracted with
 116 significant context. To overcome this limitation, we have considered three
 117 methods of window *repositioning* after window extraction: *vertical*, *horizon-*
 118 *tal*, and *both*. The basic idea is to first compute the center of mass of the
 119 extracted window, which is then repositioned (translated) to align its center
 120 to the center of mass. This is done in accordance with the chosen method,
 121 that is, horizontally, vertically, or in both directions. Obviously, the feature
 122 vector actually extracted is that obtained after repositioning. An example
 123 of feature extraction is shown in Fig. 1 in which the standard method (no
 124 repositioning) is compared with the three methods repositioning methods
 125 considered.

126 It is helpful to observe the effect of repositioning with real data. Fig. 2
 127 shows the sequence of feature vectors extracted from a real sample of the
 128 APTI database, with and without (both) repositioning. As expected, (ver-
 129 tical or both) repositioning has the effect of normalizing vertical image dis-
 130 tortions, especially translations.

131 4. Bernoulli HMMs for printed Arabic recognition

132 Given an observation O of unknown class, we use the Bayes decision
 133 rule to classify O into the class to which it belongs with highest (*posterior*)
 134 probability or, equivalently:

$$c^*(O) = \operatorname{argmax}_{c=1,\dots,C} \log P(c) + \log P(O | c) \quad (4)$$

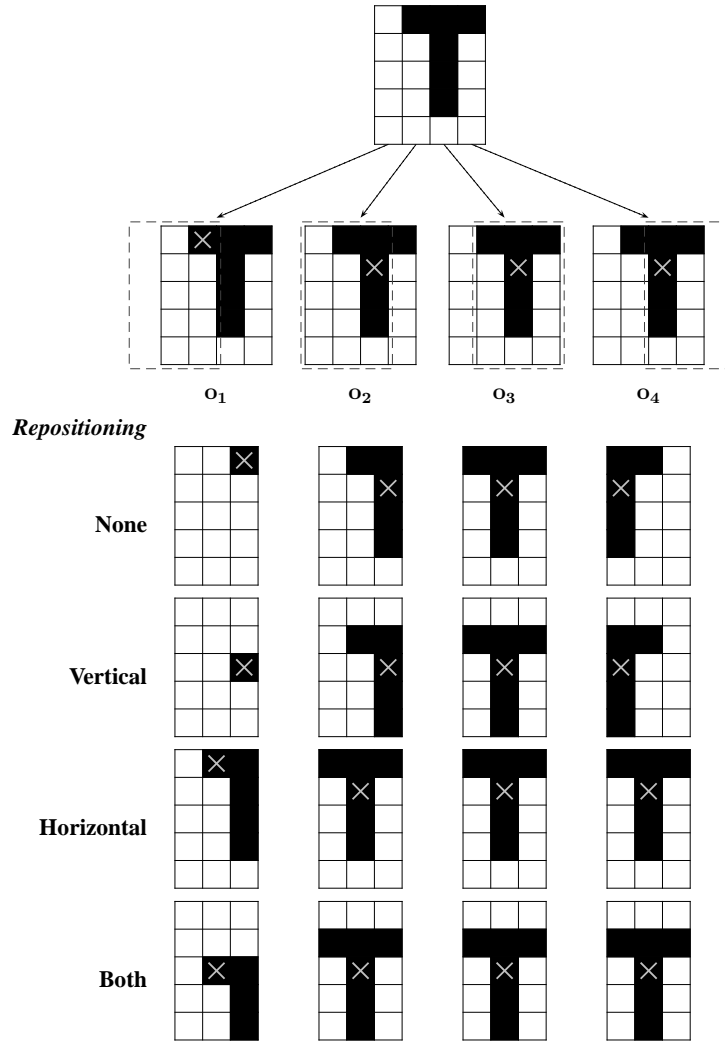


Figure 1: Example of transformation of a 4×5 binary image (top) into a sequence of 4 15-dimensional binary feature vectors $O = (o_1, o_2, o_3, o_4)$ using a window of width 3. After window extraction (illustrated under the original image), the standard method (no repositioning) is compared with the three repositioning methods considered: vertical, horizontal, and both directions. Mass centers of extracted windows are also indicated.

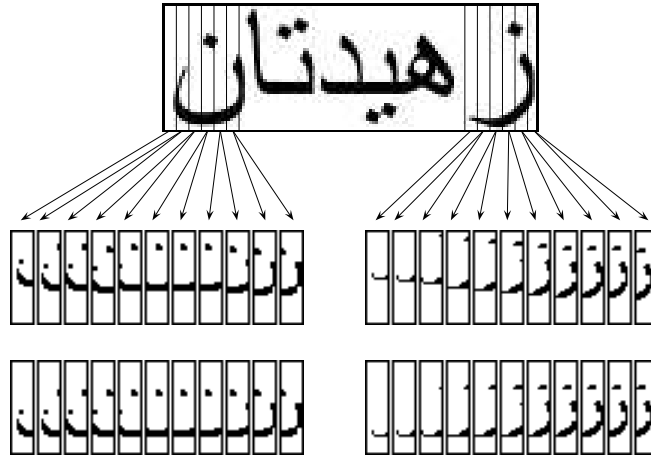


Figure 2: Original sample *Image_18_ArabicTransparent_5111* from set1 from APTI database (top) and its sequence of feature vectors produced with and without (both) repositioning (center and bottom, respectively).

135 where C is the total number of classes and, for each class $c = 1, \dots, C$,
 136 $P(c)$ is its *prior* probability and $P(O | c)$ is the class-conditional probability
 137 (density) for O to come from class c .

138 Class priors and class-conditional probability (density) functions are usu-
 139 ally estimated from a set of *training* observations. The conventional approach
 140 to estimate class priors is simply to compute their relative frequencies from
 141 the training set. However, the estimation of class-conditional probability
 142 (density) functions is more involved and depends on the type of representa-
 143 tion space for the observations. Usually, each class-conditional probability
 144 (density) function is modeled by an appropriate parametric function whose
 145 parameters are estimated by MLE from the training data. As an example,
 146 consider the problem of classifying images of *isolated printed Arabic charac-*
 147 *ters*. The number of classes is modest and it is not difficult to collect many

148 training examples for each class. Therefore, class priors can be accurately
 149 estimated by the conventional method. Also, if images are represented as
 150 sequences of feature vectors, each class-conditional probability function can
 151 be modeled by an independent BHMM (Eqs. (1) and (2)) with parameters
 152 estimated by MLE from training observations of its class (Giménez et al.,
 153 2014b).

154 The above approach for the estimation of class priors and class-conditional
 155 probability (density) functions is no longer applicable to classification prob-
 156 lems with large number of classes due to the lack of training data for each
 157 class. Consider, as we do in this work, the problem of classifying images
 158 of *printed Arabic words*. Collecting a number of training observations for
 159 each word will be really difficult if we are interested in recognizing a large
 160 number of different words. Indeed, it will be impossible if we are interested
 161 in building an *open-vocabulary* recognizer, that is, one even able to recognize
 162 words not “seen” (with no observations) in the training data. As with Arabic
 163 handwriting recognition in general, the usual approach in this case consists
 164 in using global (word) models defined in terms of local (subword) models.
 165 This is the approach followed in this work. Formally, given an observation
 166 O of an unknown word, we use Eq. (4) to decide to which word corresponds:

$$w^*(O) = \underset{w}{\operatorname{argmax}} \log P(S_w) + \log P(O | S_w, \Theta) \quad (5)$$

167 where, for each word w , S_w is its sequence of symbols (characters), $P(S_w)$
 168 is its prior probability and $P(O | S_w, \Theta)$ is the probability for O to be
 169 generated from a BHMM for w (Eq. 3). Word priors are modeled with n -
 170 gram language models at character level (Jelinek, 1997). Word-conditional
 171 probability functions are modeled by BHMMs built from shared, embedded

172 BHMMs at character level (Eq. (3)) with parameters trained by MLE.

173 Clearly, the direct way to measure the error of a word recognizer is to
174 count the (relative) number of misclassified observations in a collection of
175 *test* observations (i.e. samples held out during training). In what follows,
176 this is referred to as the Word Error Rate (WER). Apart from the WER, we
177 also use the Character Error Rate (CER), that is, the (relative) number of
178 misclassified characters. In practice, the CER can be considered equivalent
179 to the WER for comparison purposes.

180 5. Experiments

181 As indicated in the introduction, in this Section we provide the results of a
182 complete series of experiments on APTI as well as a comparison with results
183 from other authors on this database. APTI is briefly described in Section 5.1
184 together with its basic preprocessing for the experiments below. Then, two
185 experimental protocols are defined in Section 5.2, UPVPC1 and UPVPC2,
186 whose results are reported separately in Sections 5.3 and 5.4 respectively.
187 Finally, the idea of vertical repositioning is also tried on recent state-of-the-
188 art techniques based on neural networks in Section 5.5.

189 5.1. APTI database and preprocessing

190 The Arabic Printed Text Image (APTI) database is a collection of images
191 of Arabic Printed words. It was recently published by (Slimane et al., 2009)
192 for large-scale benchmarking of open-vocabulary, multi-font, multi-size and
193 multi-style text recognition systems in Arabic. It consists of 113284 different
194 single words, each one available in 10 different fonts, 10 different font sizes,
195 and also 4 different styles.

196 APTI is divided into six equilibrated sets (*set1, set2, . . . , set6*) to allow
197 for flexibility in the design of experimental protocols. Each set has different
198 words, but characters are equally distributed. The five first sets are available
199 for the scientific community. The sixth set is kept by the authors for future
200 evaluation of systems in blind mode.

201 For the experiments reported below, APTI was preprocessed by scaling
202 all images in the first five sets to a height of D pixels (for 10 different values
203 of D from 30 to 50) while keeping the aspect ratio. Scaled images were then
204 binarized by application of the Otsu’s method (Otsu, 1979).

205 5.2. Experimental protocols: UPVPC1 and UPVPC2

206 APTI was used first in the Arabic Recognition Competition of ICDAR
207 2011 (Slimane et al., 2011). Two experimental protocols were defined which
208 differ in the number of fonts used: APTIPC1 and APTIPC2. In APTIPC1,
209 only the Arabic Transparent font was used. In APTIPC2, however, five dif-
210 ferent fonts were used: Arabic Transparent (Trans), Andalus (Anda), Diwani
211 Letter (Diw), Simplified Arabic (Simp), and Traditional Arabic (Trad). In
212 both protocols, only the *Plain* font style was used, with sizes of 6, 8, 10,
213 12, 18 and 24 pixels. As indicated above, the first five sets were available
214 to participants for system training, while the sixth set was held-out by the
215 organizers for system comparison in blind mode.

216 In this paper, we could not use the training-test partition used at the
217 ICDAR 2011 competition because the sixth set is not publicly available.
218 Instead, we used the first four sets for training and the fifth set for testing.
219 More precisely, we defined two new protocols: UPVPC1 and UPVPC2. In
220 UPVPC1, 13000 images from the first four sets were randomly drawn (10000

221 for training and 3000 for testing). In UPVPC2, we used the whole first four
222 sets for training and the whole fifth set for testing. In particular, we used
223 2266500 images for training, and 566040 for testing.

224 5.3. Results using the UPVPC1 protocol

225 For (computational) simplicity, the UPVPC1 protocol was used in a first
226 series of experiments to study the effect on the CER of various key parame-
227 ters. We began with experiments for font size 6, which were then extended
228 to other font sizes. In particular, for each dimension D in $\{30, 32, \dots, 50\}$,
229 each sliding window width W in $\{1, 3, \dots, 11\}$, each number of states Q in
230 $\{4, 5, 6, 7, 8\}$ and each number of mixture components K in $\{1, 2, 4, \dots, 32\}$,
231 a BHMM-based word recognizer was trained from the training data of font
232 size 6 in the UPVPC1 protocol. For $K = 1$, BHMMs were initialized by
233 first segmenting training data with a “neutral” model, and then using the
234 resulting segments to perform a Viterbi initialization. Initialized BHMMs
235 were then trained with 4 EM iterations. For $K > 1$, BHMMs were initialized
236 by splitting the mixture components of the models trained with $K/2$ mixture
237 components per state. Again, after initialization, BHMMs were trained with
238 4 EM iterations. On the other hand, word priors were modeled with 5-gram
239 language models at character level.

240 The above training procedure led to a different recognizer for each combi-
241 nation of key parameter values (apart from the font size itself). Each of them
242 was of the form given by Eq. (5) though, as usual in (Arabic) text recognition,
243 a *Grammar Scale Factor (GSF)* was used to adjust the importance of class
244 priors with respect to word-conditional observation probabilities (i.e. the GSF
245 is a constant multiplier for log-priors). For each combination of parameter

246 values and each value of $GSF \in \{20, 30, 40, 50\}$, the corresponding recog-
247 nizer was assessed in terms of CER from the test data of font size 6 in the
248 UPVPC1 protocol.

249 Figure 3 shows the CER as a function of D (top left), K (top right), Q
250 (bottom left) and GSF (bottom right); for $W = 1, 3, 7$ and 11 (the curves
251 for $W = 5$ and 9 are similar and have been omitted for clarity). Each plotted
252 point shows the best CER obtained over all values tried for the parameters
253 not given. The best CER obtained is 3.4% for $D = 38, W = 7, Q = 7,$
254 $K = 32$ and $GSF = 50$. In the plot at the top left, it is shown for $D = 38$
255 and $W = 7$, as the minimum CER obtained for all values tried for Q, K and
256 GSF .

257 From the results in Fig. 3, it is clear that the use of windowed BHMMs
258 is of crucial importance. Indeed, the best CER obtained with no windows
259 ($W = 1$) is 6.6%; i.e. it nearly doubles the best CER with windows. Note
260 also that, as W , the number of mixtures components (K) has a strong effect
261 on the CER. The best error rates were obtained with the maximum value
262 of K tried (32). Therefore, this and larger values of K need to be tried in
263 further experiments with more training data. The dimension (D), number of
264 states (Q) and GSF are also key parameters to be adjusted, though Fig. 3
265 does not show wide fluctuations in CER for the ranges of values considered.

266 As discussed in (Dreuw et al., 2009), letters in Arabic script differ signif-
267 icantly in length, and thus it might not be appropriate to model all of them
268 using BHMMs of fixed number of states. With this idea in mind, an exper-
269 iment similar to that described above was carried out for $D = 38, W = 7,$
270 $K = 32, GSF = 50$ and variable number of states. To decide the number of

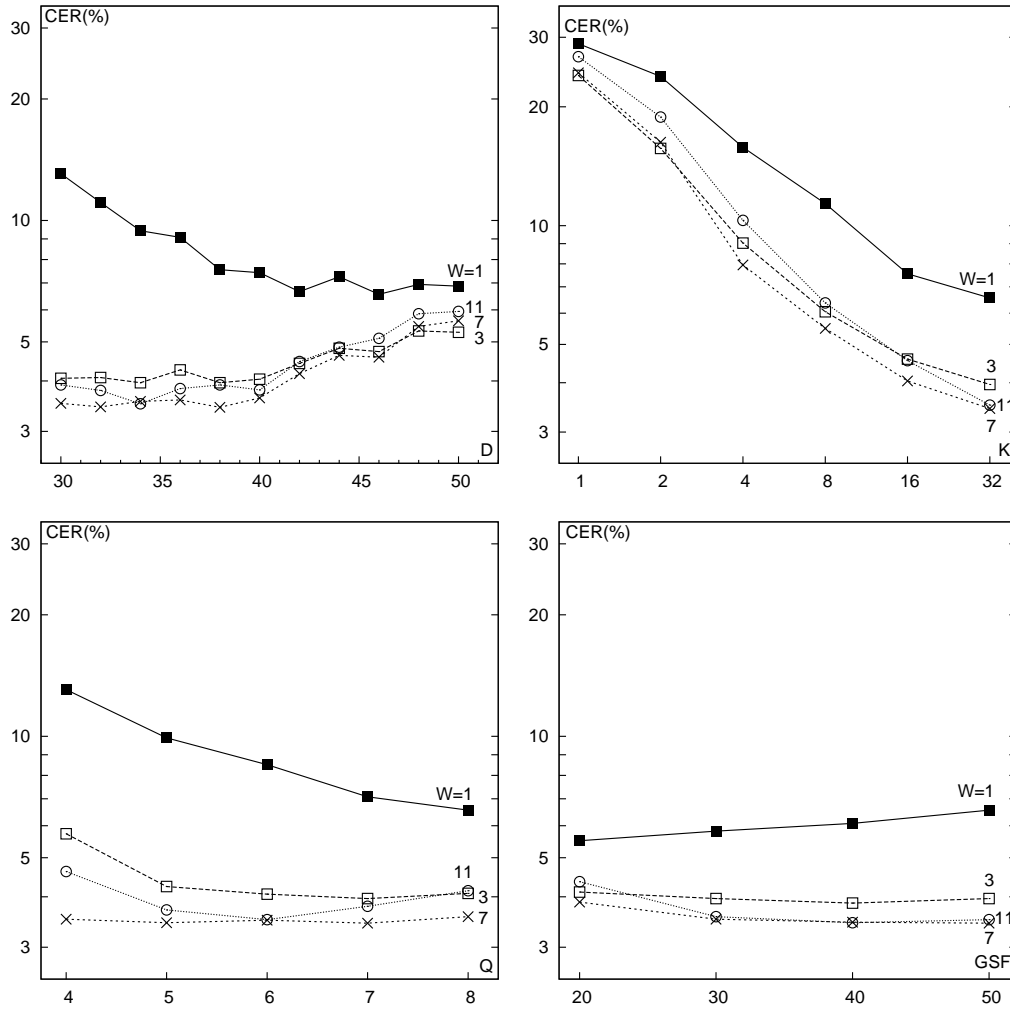


Figure 3: CER(%) as a function of the dimension D (top left), number of mixture components K (top right), number of states Q (bottom left) and GSF value (bottom right); for sliding window widths of $W = 1, 3, 7$ and 11 .

271 states for each character, we first Viterbi-segmented all training data using
 272 BHMMs of 7 states, and then computed the average length of the segments
 273 associated with each character. Given an average segment length for charac-
 274 ter c , \bar{T}_c , its number of states was set to $F \cdot \bar{T}_c$, where F is a *factor* measuring
 275 the average number of states that are required to emit a feature vector. Thus,
 276 its inverse, $\frac{1}{F}$, can be interpreted as a *state load*, that is, the average num-
 277 ber of feature vectors that are emitted in each state. For instance, $F = 0.2$
 278 means that only a fraction of 0.2 states is required to emit a feature vector
 279 or, alternatively, that $\frac{1}{0.2} = 5$ feature vectors are emitted on average in each
 280 state. We tried all values of F in $\{0.1, 0.2, \dots, 0.9\}$. The best result achieved
 281 is a CER of 3.2%, using $F = 0.5$, which is significantly better than the best
 282 result obtained above with fixed number of states (3.4%).

283 To complete our experiments with font size 6 data in the UPVPC1 pro-
 284 tocol, the best recognizer found above was also tested with the four reposi-
 285 tioning methods described in Sec. 3. As expected, the best CER, 1.1%, was
 286 obtained with *vertical* repositioning alone. Also as expected, it was similar
 287 to the CER achieved with repositioning in *both* directions (1.2%), and sig-
 288 nificantly better than those obtained with *horizontal* and *no* repositioning
 289 (3.2% for both).

290 The experiments described above in this Section were extended to all
 291 font sizes. More precisely, for each font size $S \in \{8, 10, 12, 18, 24\}$, each $D \in$
 292 $\{30, 32, \dots, 50\}$, $W \in \{1, 3, \dots, 11\}$, $Q \in \{5, 6, 7\}$ and $K \in \{1, 2, 4, \dots, 32\}$,
 293 a BHMM-based word recognizer was trained and tested, for each value of
 294 $GSF \in \{30, 40, 50\}$, as described above. Also as above, the best recognizer
 295 for each size was then tested with variable number of states ($F \in \{0.3,$

296 $\dots, 0.7\}$) and different repositioning techniques ($R = \{N, V, H, B\}$; where
 297 N =None, V =Vertical, H =Horizontal and B =Both vertical and horizontal).
 298 The results obtained were similar to those reported in Fig. 3 for font size
 299 6. More precisely, the best error rates were obtained with windows of width
 300 $W \in \{7, 9, 11\}$, $K = 32$ components, $GSF = \{40, 50\}$, variable number of
 301 states with $F \in \{0.4, 0.5, 0.6\}$, and vertical repositioning. For brevity, these
 302 error rates are not reported here in detail, as those in Fig. 3 for font size 6.
 303 Instead, only a summary of best error rates is reported in Table 1 (including
 304 font size 6 for completeness). Note that the best recognizer (combination of
 305 parameter values) for each font size is trained within the parameter ranges
 306 indicated above. Indeed, all recognizers trained within these ranges provide
 307 nearly identical error rates.

Table 1: Best recognizer (combination of parameter values) and its CER(%) for each size.

Size	D	W	R	F	K	GSF	CER(%)
6	38	7	V	0.5	32	50	1.1
8	40	7	V	0.6	32	40	0.6
10	44	9	V	0.5	32	40	0.6
12	40	9	V	0.5	32	40	0.4
18	40	9	V	0.5	32	40	0.5
24	42	11	V	0.4	32	40	0.8

308 To get some insight into the behavior of our windowed BHMMs, a real
 309 model for the character ﺍ is (partially) shown in Figure 4 (bottom) together
 310 with its Viterbi alignment with a real image of the character ﺍ , extracted from
 311 sample *Image_24-ArabicTransparent_562, set1* (top). Bernoulli prototypes

312 are represented as gray images where the gray level of each pixel measures the
 313 probability of its corresponding pixel to be black (white = 0 and black = 1).
 314 From these prototypes, it can be seen that the model works as expected, i.e.
 315 each state from right to left accounts for a different local part of ϕ , as if the
 316 sliding window was moving smoothly from right to left.

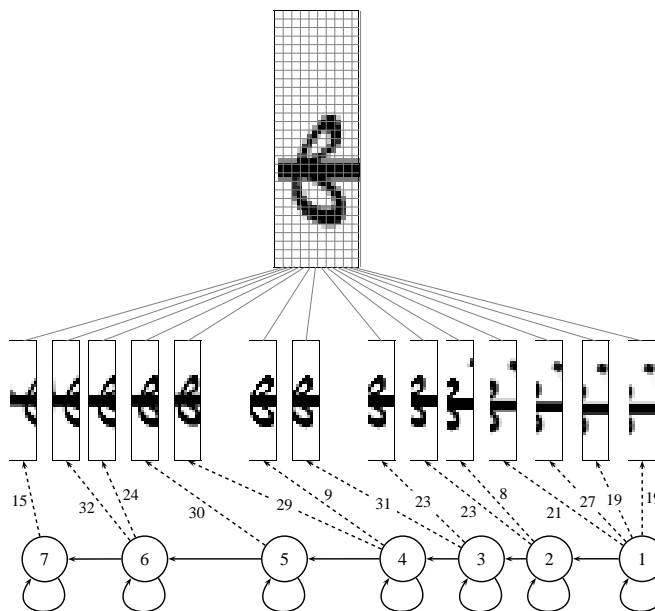


Figure 4: Real BHMM example for character ϕ and its Viterbi alignment with a real image
 of the character ϕ , extracted from sample *Image_24_ArabicTransparent_562* (top).

317 5.4. Results using the UPVPC2 protocol

318 The UPVPC1 protocol was used to study the effect on the CER of var-
 319 ious key parameters, variable number of states, and repositioning. Tak-
 320 ing into account the best results obtained with it, the UPVPC2 protocol
 321 was used in a new series of experiments to obtain results in conditions

322 similar to those used in the ICDAR 2011 Arabic Recognition Competi-
 323 tion (see Sec. 5.2). In particular, for each of the five font types consid-
 324 ered in UPVPC2, $T \in \{Trans, Anda, Diw, Simp, Trad\}$, and each font size
 325 $S \in \{6, 8, 10, 12, 18, 24\}$, a BHMM-based word recognizer was trained and
 326 tested from the data in UPVPC2 of font type T and size S . We used $D = 40$,
 327 $W = 9$, $R = V$, $F = 0.5$ (on a Viterbi segmentation produced by a recognizer
 328 trained with $Q = 7$, $K = 128$ and $GSF = 40$), $K = 128$ and $GSF = 40$.
 329 Except for the K , these parameter values are within the parameter ranges
 330 leading to the best error rates with the UPVPC1 protocol. However, in the
 331 case of K , we used 128 instead of 32. As discussed in Sec. 5.3, values of K
 332 larger than 32 had to be tried, especially with more training data as with
 333 the UPVPC2 protocol. Actually, we tried each $K \in \{1, 2, 4, \dots, 128\}$, though
 334 $K = 128$ provided the best error rates in all cases.

335 Table 2 shows CER results for each font type and size. The error rates
 336 labeled as 2013a in the Year column were obtained as described above. That
 337 is, each test sample was accompanied by its font type and size so as to select
 338 its appropriate recognizer. However, the error rates labeled as 2013b were
 339 obtained in a slightly different way, by only providing the font size of each
 340 test sample. In this case, given a test sample of size S , all the five font-
 341 dependent recognizers for size S were run in parallel and that producing the
 342 highest classification score (see Eq. (5)) was chosen to decide the recognized
 343 word. The error rates labeled as 2011 are the best results of the ICDAR 2011
 344 competition, which were also obtained by only providing the font size of each
 345 test sample.

346 A first conclusion that can be drawn from Table 2 is that the figures

Table 2: CER results for each font type and size (2013a="font type and size given"; 2013b="only font size given"; 2011="best results from the ICDAR 2011 competition").

Font/Size	Year	6	8	10	12	18	24	Mean
Andalus	2013a	0.9	0.2	0.1	0.1	0.0	0.0	0.2
	2013b	0.9	0.2	0.1	0.1	0.0	0.0	0.2
	2011	1.1	5.2	3.9	3.3	3.3	3.0	3.3
Arabic Transparent	2013a	0.6	0.1	0.1	0.0	0.0	0.1	0.2
	2013b	0.6	0.1	0.1	0.0	0.0	0.0	0.1
	2011	1.0	3.5	3.4	3.9	3.8	3.9	3.3
Simplified Arabic	2013a	0.5	0.1	0.1	0.0	0.0	0.0	0.1
	2013b	0.4	0.1	0.1	0.0	0.0	0.0	0.1
	2011	0.8	3.9	3.3	3.1	3.0	2.6	2.8
Traditional Arabic	2013a	6.4	1.3	0.5	0.3	0.2	0.2	1.5
	2013b	6.5	1.3	0.5	0.3	0.2	0.2	1.5
	2011	10.7	18.1	14.1	11.5	12.5	11.7	13.1
Diwani Letter	2013a	10.0	7.2	6.7	6.2	6.1	5.9	7.0
	2013b	10.0	7.2	6.7	6.2	6.1	5.9	7.0
	2011	9.1	24.2	16.6	10.9	5.1	7.4	12.2

347 labeled as 2013a and 2013b are virtually identical. Therefore, when font
348 size is known but font type is not, the procedure described above to obtain
349 the 2013b results seems absolutely reliable. Another important conclusion
350 from Table 2 is that the results of this work outperform by a large extent
351 those from the competition. Note that, on average, recognition of Andalus,
352 Arabic Transparent and Simplified Arabic is nearly perfect in terms of CER.
353 On the other hand, recognition of Traditional Arabic and Diwani Letter is
354 fairly good and comparatively much better than that of the ICDAR 2011
355 competition.

356 Apart from the above multi-font and mono-size recognition results, the
357 ICDAR 2011 competition also included mono-size results on only the Ara-
358 bic Transparent font. For this particular font, results were published for
359 both, competition participants (IPSAR and UPV) and organizers (DIVA-
360 REGIM). Also, more recent results have been published by (Awaida and
361 Khorsheed, 2012), and by (Dershowitz and Rosenberg, 2013). The most re-
362 cent results come from the ICDAR 2013 second competition on APTI, which
363 included three more participants than in its first edition: SID, THOCR and
364 Siemens (Slimane et al., 2013). All these results are shown in Table 3 in
365 terms of CER and WER. UPV-REC1, UPV-BHMM and UPV-2013 refer to
366 our system at, respectively, ICDAR 2011, ICDAR 2013 and this work. Note
367 that the results of UPV-BHMM and UPV-2013 are nearly identical and thus,
368 as expected, the UPVPC2 protocol provides a good approximation to the ex-
369 perimental conditions of the ICDAR competitions on APTI. These results
370 are much better than those of UPV-REC1 and only at a marginal distance
371 from the best system at the ICDAR 2013 second competition on APTI. They

372 are also much better than those reported in (Khoury et al., 2013), where an
 373 initial, preliminary part of the experiments and results described here can
 374 also be found.

Table 3: CER and WER results for the Arabic Transparent font in each size.

System	Year		6	8	10	12	18	24	Mean
IPSAR	2011	WER	94.3	26.7	25.0	16.9	22.9	22.5	34.7
		CER	40.6	5.8	4.9	3.1	4.3	3.2	10.3
UPV-REC1	2011	WER	5.5	2.6	3.3	7.5	15.4	15.6	8.3
		CER	1.0	0.4	0.6	1.3	3.1	4.0	1.7
DIVA-REGIM	2011	WER	13.1	4.1	4.3	6.1	2.1	1.1	5.1
		CER	2.0	0.8	0.7	1.2	0.3	0.3	0.9
Awaida et al.	2012	CER	-	-	-	-	-	-	3.4
Dershowitz et al.	2013	WER	72.4	21.1	10.2	6.0	1.0	1.5	18.7
		CER	31.8	5.6	2.5	2.4	0.2	0.4	7.2
UPV-BHMM	2013	WER	2.8	0.3	0.2	0.1	0.1	0.1	0.6
		CER	0.5	0.1	0.1	0.0	0.0	0.0	0.1
SID	2013	WER	5.7	3.8	1.8	1.2	3.4	2.6	3.1
		CER	0.3	0.0	0.0	0.1	0.0	0.0	0.1
THOCR	2013	WER	10.5	4.2	5.2	7.5	5.4	5.0	6.3
		CER	1.7	0.5	0.8	0.9	0.9	0.8	0.9
Siemens	2013	WER	0.1	0.1	0.0	0.1	0.0	0.0	0.1
		CER	0.0	0.0	0.0	0.0	0.0	0.0	0.0
UPV-2013	2013	WER	3.0	0.4	0.3	0.2	0.2	0.2	0.7
		CER	0.6	0.1	0.1	0.0	0.0	0.1	0.2

375 *5.5. New results using a DNN hybrid HMM system and vertical repositioning*

376 Previous experiments have shown that the results obtained by using BH-
377 MMs are improved by applying the vertical repositioning technique. In recent
378 work on handwritten recognition, vertical repositioning has also shown a sig-
379 nificant improvement when used with other models than Bernoulli HMMs. In
380 particular, in (Doetsch et al., 2012), a notable improvement was reported by
381 using a Long Short Term Memory recurrent neural network (LSTM-RNN)
382 tandem HMM and vertical repositioning on Arabic and French handwrit-
383 ing. This improvement is also observed in (Hamdani et al., 2014) where the
384 window repositioning is used as a preprocessing step.

385 In order to asses that the vertical repositioning is useful for printed Ara-
386 bic recognition with the current state-of-the-art techniques based on neural
387 networks, such as LSTM-RNN, we have carried out a new series of experi-
388 ments using the UPVPC2 protocol and a Deep Neural Network (DNN) hy-
389 brid HMM system (Dahl et al., 2012). This technique is similar to the Long
390 Short Term Memory (LSTM) technique applied in (Rashid et al., 2013). It
391 has been implemented in a recently released, open-source toolkit for auto-
392 matic speech recognition called TLK toolkit (The transLectures-UPV Team,
393 2013). On the basis of our experience on the application of TLK to speech
394 recognition tasks within the transLectures project, we decided to use it also
395 for the additional experiments discussed in this Section. The results of these
396 experiments, with and without vertical repositioning, are shown in Table 4.

397 As with the winner of ICDAR 2013 (Table 3), the results in Table 4 are
398 nearly perfect. Even though the error is nearly zero, vertical repositioning
399 still obtains slight improvements. In particular, for the more challenging

Table 4: CER and WER results for the Arabic Transparent font in each size.

System	Year		6	8	10	12	18	24	Mean
Vertical Rep.	2014	WER	0.16	0.13	0.12	0.12	0.13	0.15	0.14
		CER	0.03	0.03	0.02	0.02	0.03	0.03	0.03
Without Rep.	2014	WER	0.22	0.20	0.12	0.13	0.13	0.16	0.16
		CER	0.04	0.04	0.02	0.02	0.03	0.03	0.03

400 font sizes (6 and 8), a modest improvement is achieved when applying repo-
 401 sitioning. Specifically, for font size 6 results were 0.16% with repositioning
 402 and 0.22% without repositioning. (Note that, as we were using 19000 test
 403 samples approximately for each font size, a difference of 0.06% accounts for
 404 about 11 classification errors.) In a similar way, for font size 8, results were
 405 0.13 and 0.20 for repositioning and non-repositioning respectively.

406 6. Concluding remarks

407 Windowed Bernoulli HMMs with repositioning have been described and
 408 extensively tested for printed Arabic recognition on the Arabic Printed Text
 409 Image (APTI) database. A system based on these models, though with no
 410 repositioning, ranked first at the ICDAR 2011 Arabic recognition competi-
 411 tion for printed Arabic text, also based on the APTI database. Following
 412 evaluation protocols similar to those of the competition, this system has been
 413 largely improved by the use of repositioning and an exhaustive experimenta-
 414 tion to adjust various key parameters and model topology (variable number
 415 of states). Results comparatively much better than those of the competition
 416 have been reported on multi-font and mono-size recognition, with nearly per-
 417 fect performance for most fonts in terms of Character Error Rate. Indeed, a

418 second edition of the competition on APTI was recently held at the ICDAR
419 2013 and our improved system obtained results nearly identical to those re-
420 ported here. This second edition was harder than the first and our system
421 ranked second, though only at a marginal distance from the best.

422 For future work, we would be interested in carrying out a deep analysis
423 to compare repositioning with other, more complex, techniques for baseline
424 detection and correction.

425 **Acknowledgments**

426 The research leading to these results has received funding from the Eu-
427 ropean Union Seventh Framework Programme (FP7/2007-2013) under grant
428 agreement no 287755. Also supported by the Spanish Government (Plan E,
429 iTrans2 TIN2009-14511 and AECID 2011/2012 grant).

430 **References**

431 Awaida, S., Khorsheed, M., 2012. Developing discrete density Hidden Markov
432 Models for Arabic printed text recognition, in: Computational Intelligence
433 and Cybernetics (CyberneticsCom), 2012 IEEE International Conference
434 on, pp. 35–39.

435 Dahl, G.E., Member, S., Yu, D., Member, S., Deng, L., Acero, A., 2012.
436 Context-Dependent Pre-trained Deep Neural Networks for Large Vocab-
437 ulary Speech Recognition, in: IEEE Transactions on Audio, Speech, and
438 Language Processing.

- 439 Dehghan, M., Faez, K., Ahmadi, M., Shridhar, M., 2001. Handwritten Farsi
440 (Arabic) word recognition: a holistic approach using discrete HMM. Pat-
441 tern Recognition 34, 1057–1065.
- 442 Dershowitz, N., Rosenberg, A., 2013. Language, Culture, Computation:
443 Studies in Honor of Yaacov Choueka. Springer-Verlag. volume 8000 of *Lec-
444 ture Notes in Computer Science*. chapter Arabic Character Recognition.
- 445 Doetsch, P., Hamdani, M., Ney, H., Gimenez, A., Andres-Ferrer, J., Juan,
446 A., 2012. Comparison of Bernoulli and Gaussian HMMs using a vertical
447 repositioning technique for off-line handwriting recognition, in: Interna-
448 tional Conference on Frontiers in Handwriting Recognition, Bari, Italy.
449 pp. 3–7.
- 450 Dreuw, P., Heigold, G., Ney, H., 2009. Confidence-Based Discriminative
451 Training for Model Adaptation in Offline Arabic Handwriting Recognition,
452 in: ICDAR '09, Barcelona (Spain). pp. 596–600.
- 453 Giménez, A., Andrés-Ferrer, J., Juan, A., 2014a. Discriminative Bernoulli
454 HMMs for isolated handwritten word recognition. Pattern Recognition
455 Letters 35, 157 – 168. Frontiers in Handwriting Processing.
- 456 Giménez, A., Juan, A., 2009. Embedded Bernoulli Mixture HMMs for Hand-
457 written Word Recognition, in: ICDAR '09, Barcelona (Spain). pp. 896–900.
- 458 Giménez, A., Khoury, I., Andrés-Ferrer, J., Juan, A., 2014b. Handwriting
459 word recognition using windowed Bernoulli HMMs. Pattern Recognition
460 Letters 35, 149 – 156. Frontiers in Handwriting Processing.

- 461 Giménez, A., Khoury, I., Juan, A., 2010. Windowed Bernoulli Mixture
462 HMMs for Arabic Handwritten Word Recognition, in: ICFHR' 10, Kolkata
463 (India). pp. 533–538.
- 464 Grosicki, E., El Abed, H., 2009. ICDAR 2009 Handwriting Recognition
465 Competition, in: ICDAR '09, Barcelona (Spain). pp. 1398 – 1402.
- 466 Günter, S., Bunke, H., 2004. HMM-based handwritten word recognition: on
467 the optimization of the number of states, training iterations and Gaussian
468 components. *Pattern Recognition* 37, 2069–2079.
- 469 Hamdani, M., Doetsch, P., Kozielski, M., El-Desoky Mousa, A., Ney, H.,
470 2014. The RWTH Large Vocabulary Arabic Handwriting Recognition Sys-
471 tem, in: International Workshop on Document Analysis Systems, France.
- 472 Jelinek, F., 1997. *Statistical Methods for Speech Recognition*. MIT Press.
- 473 Khoury, I., Gimnez, A., Juan, A., Andrs-Ferrer, J., 2013. Arabic Printed
474 Word Recognition Using Windowed Bernoulli HMMs, in: Petrosino, A.
475 (Ed.), *Proc. of the 17th Int. Conf. on Image Analysis and Processing*
476 *ICIAP 2013*. Springer Berlin Heidelberg. volume 8156, pp. 330–339.
- 477 Märgner, V., El Abed, H., 2007. ICDAR 2007 - Arabic Handwriting Recog-
478 nition Competition, in: ICDAR '07, Curitiba (Brazil). pp. 1274–1278.
- 479 Märgner, V., El Abed, H., 2009. ICDAR 2009 Arabic Handwriting Recogni-
480 tion Competition, in: ICDAR '09, Barcelona (Spain). pp. 1383–1387.
- 481 Märgner, V., El Abed, H., 2010. ICFHR 2010 - Arabic Handwriting Recog-
482 nition Competition, in: ICFHR '10, Kolkata (India). pp. 709–714.

- 483 Otsu, N., 1979. A Threshold Selection Method from Gray-Level Histograms.
484 IEEE Trans. on Systems, Man and Cybernetics 9, 62–66.
- 485 Rabiner, L., Juang, B.H., 1993. Fundamentals of speech recognition.
486 Prentice-Hall.
- 487 Rashid, S.F., Schambach, M.P., Rottland, J., von der Nüll, S., 2013. Low
488 Resolution Arabic Recognition with Multidimensional Recurrent Neural
489 Networks, in: Proceedings of the 4th International Workshop on Multilin-
490 gual OCR, ACM, New York, NY, USA. pp. 6:1–6:5.
- 491 Slimane, F., Ingold, R., Kanoun, S., Alimi, A.M., Hennebert, J., 2009. A New
492 Arabic Printed Text Image Database and Evaluation Protocols, IEEE. pp.
493 946–950.
- 494 Slimane, F., Kanoun, S., Abed, H.E., Alimi, A.M., Ingold, R., Hennebert, J.,
495 2011. ICDAR 2011 - Arabic Recognition Competition: Multi-font Multi-
496 size Digitally Represented Text, IEEE. pp. 1449–1453.
- 497 Slimane, F., Kanoun, S., El Abed, H., Alimi, A.M., Ingold, R., Hennebert,
498 J., 2013. ICDAR 2013 Competition on Multi-font and Multi-size Digitally
499 Represented Arabic Text, CPS. pp. 1465–1469.
- 500 The transLectures-UPV Team, 2013. The translectures-upv toolkit (tlk).
501 <http://translectures.eu/tlk>. <http://www.translectures.eu/tlk/citing-tlk/>.
- 502 Young, S., et al., 1995. The HTK Book. Cambridge University Engineering
503 Department.