

Improving Offline Handwritten Text Recognition with Hybrid HMM/ANN Models

Salvador España-Boquera, Maria Jose Castro-Bleda,
Jorge Gorbe-Moya, and Francisco Zamora-Martinez

Abstract—This paper proposes the use of hybrid Hidden Markov Model (HMM)/Artificial Neural Network (ANN) models for recognizing unconstrained offline handwritten texts. The structural part of the optical models has been modeled with Markov chains, and a Multilayer Perceptron is used to estimate the emission probabilities. This paper also presents new techniques to remove slope and slant from handwritten text and to normalize the size of text images with supervised learning methods. Slope correction and size normalization are achieved by classifying local extrema of text contours with Multilayer Perceptrons. Slant is also removed in a nonuniform way by using Artificial Neural Networks. Experiments have been conducted on offline handwritten text lines from the IAM database, and the recognition rates achieved, in comparison to the ones reported in the literature, are among the best for the same task.

Index Terms—Handwriting recognition, offline handwriting, hybrid HMM/ANN, HMM, neural networks, multilayer perceptron, image normalization.

1 INTRODUCTION AND MOTIVATION

OFFLINE handwritten text recognition is one of the most active areas of research in computer science and it is inherently difficult because of the high variability of writing styles. High recognition rates are achieved in character recognition and isolated word recognition, but we are still far from achieving high-performance recognition systems for unconstrained offline handwritten texts [1], [2], [3], [4], [5], [6], [7].

Automatic handwriting recognition systems normally include several preprocessing steps to reduce variation in the handwritten texts as much as possible and, at the same time, to preserve information that is relevant for recognition. There is no general solution to preprocessing of offline handwritten text lines, but it typically relies on slope and slant correction and normalization of the size of the characters. With the slope correction, the handwritten word is horizontally rotated such that the lower baseline is aligned to the horizontal axis of the image. Slant is the clockwise angle between the vertical direction and the direction of the vertical text strokes. Slant correction transforms the word into an upright position. Ideally, the removal of slope and slant results in a word image that is

independent of these factors. Finally, size normalization tries to make the system invariant to the character size and to reduce the empty background areas caused by the ascenders and descenders of some letters.

This paper presents new techniques to remove the slope and the slant from handwritten text lines and to normalize the size of the text images by using Artificial Neural Networks (ANNs). Local extrema from a text image classified as belonging to the lower baseline by a Multilayer Perceptron (MLP) are used to accurately estimate the slope and the horizontal alignment. Slant is removed in a nonuniform way by also using ANNs. Finally, another MLP computes the reference lines of the slope and slant-corrected text in order to normalize its size.

Hidden Markov Models (HMMs) have been widely applied to offline handwriting recognition [8], [2], [4], [5], [6], [9], [10], [11] after their success in automatic speech recognition. The basic idea is that handwriting can be interpreted as a left-to-right sequence of ink signals which is analogous to the temporal sequence of acoustic signals in speech. The motivation for the work on the hybrid HMM/ANN models presented here originates from

- S. España-Boquera, M.J. Castro-Bleda, and J. Gorbe-Moya are with the Departamento de Sistemas Informáticos y Computación, Universitat Politècnica de València, Camino de Vera s/n, 46022 Valencia, Spain. E-mail: {sespana, mcastro, jgorbe}@dsic.upv.es.
- F. Zamora-Martinez is with the Departamento de Ciencias Físicas, Matemáticas y de la Computación, Universidad CEU-Cardenal Herrera, Alfara del Patriarca, Valencia, Spain, and the Departamento de Sistemas Informáticos y Computación, Universitat Politècnica de València, Camino de Vera s/n, 46022 Valencia, Spain. E-mail: fzamora@dsic.upv.es.

Manuscript received 9 Mar. 2009; revised 8 Jan. 2010; accepted 27 Apr. 2010; published online 9 Aug. 2010.

Recommended for acceptance by E. Saund.

For information on obtaining reprints of this article, please send e-mail to: tpami@computer.org, and reference IEEECS Log Number TPAMI-2009-03-0156.

Digital Object Identifier no. 10.1109/TPAMI.2010.141.

- a critical analysis of the state-of-the-art in offline handwritten text recognition [5], [10], [11],
- our earlier work on offline handwriting recognition using conventional HMMs [12],
- our own and others' experience in using hybrid HMM/ANN models for automatic speech recognition [13], [14], [15], [16], [17] and for online handwriting recognition [18], [19], [20], [21], [22], [23],
- and previous works which used hybrid HMM/ANN word models with remarkable success, although they were limited to digit recognition or small vocabulary tasks [24], [25], [26].

TABLE 1
MLPs for Preprocessing

Task	MLP	Input to MLP	Output to MLP	MLP topology
Image cleaning	Enhancer-MLP	Pixels in a window around current pixel	Restored value of current pixel	(11×11)-32-16-1
Slope removal	Slope-MLP	Pixels in a window around current pixel	Lower baseline/Not lower baseline pixel	(50×30)-64-16-2
Slant removal	Slant-MLP	Resized sheared image in a window around current column of pixels	Presence of slant angle	(40×40)-64-8-1
Size normalization	Normalize-MLP	Pixels in a window around current pixel	Ascender/Upper baseline/Lower baseline/Descender/None pixel	(50×30)-64-16-5

MLPs are used for regression (Enhancer-MLP) and for classification (Slope-MLP, Normalize-MLP, and Slant-MLP).

Hybrid HMM/ANN models compute the emission probabilities for the HMMs with a neural network instead of the commonly used Gaussian mixtures. This work is the first successful attempt, to the best of our knowledge, to use hybrid HMM/ANN models in unconstrained offline handwritten text recognition.

In many other works, artificial neural networks have been extensively applied to classify characters as part of isolated or continuous handwritten word recognizers [27], [28], [29], [30], [31].

In our experiments with hybrid HMM/ANN models, left-to-right Markov chains have been used to model graphemes, and a single neural network has been used to estimate the emission probabilities. The estimates of the posterior probabilities computed by the neural network are divided by the prior state probabilities, resulting in scaled likelihoods which are used as emission probabilities in the HMMs.

We have conducted experiments with the “large writer-independent text line recognition task” of the IAM database [32], [33] using our preprocessing and conventional HMMs as optical models. This baseline experiment achieves comparative performance with state-of-the-art systems (38.8 percent of word error rate). Next, experiments with our hybrid HMM/ANN system were performed and excellent results were achieved improving our baseline by a relative word error rate reduction of more than 42 percent (from a word error rate of 38.8 percent to 22.4 percent).

Section 2 introduces our approach to handwritten text preprocessing, showing illustrative examples. Section 3 describes the proposed hybrid HMM/ANN recognition system. The experimental setup is detailed in Section 4 and recognition results for the HMM and HMM/ANN systems with the IAM line task are shown in Section 5. Analysis of the experiments and comparison with other approaches are presented in Section 6. Finally, conclusions are drawn in Section 7.

2 PREPROCESSING AND FEATURE EXTRACTION

Handwritten image normalization from a scanned image includes several steps, which usually begin with image cleaning, page skew correction, and line detection [9]. A database of skew-corrected lines has been used in all the experiments [32]; thus page skew correction and line detection are skipped in this work. With the handwritten text line images, several preprocessing steps to reduce variations in writing style are usually performed: slope and slant removal and character size normalization. This paper

presents new techniques to remove the slope and the slant from handwritten text lines, and to normalize the size of the text images with ANNs. Table 1 outlines the key ideas of the MLPs which are used for preprocessing. These MLPs are described in more detail in Section 4.

2.1 Image Cleaning

Before any other preprocessing step, the text line scanned image is first cleaned and enhanced. Neural networks have been used in previous works for image restoration by learning the appropriate filters from examples [31]. Similarly, we have used a neural network filter to clean and enhance the handwritten text images by estimating the gray level of one pixel at a time [34]: An MLP (from now on called Enhancer-MLP) is fed with a square of pixels centered at the pixel to be cleaned, and the output is the restored value of the pixel. The entire image is cleaned by scanning all the pixels in this way. A scheme of this process is illustrated in Fig. 1. A real example of a cleaned image is shown in Fig. 2b.

2.2 Slope Removal

With the skew-corrected lines, most handwriting recognition systems require the detection of the different areas of the cursive script: the main body area (area between the upper baseline and the lower baseline), the ascenders, and the descenders (see Fig. 3 for an example). These areas can be detected by means of horizontal histogram projection [35], [36], [37] or also by obtaining the upper and lower contours of the image [38], [39]. Instead of relying on these geometric heuristics, our approach consists of automatically detecting a set of points from the image and

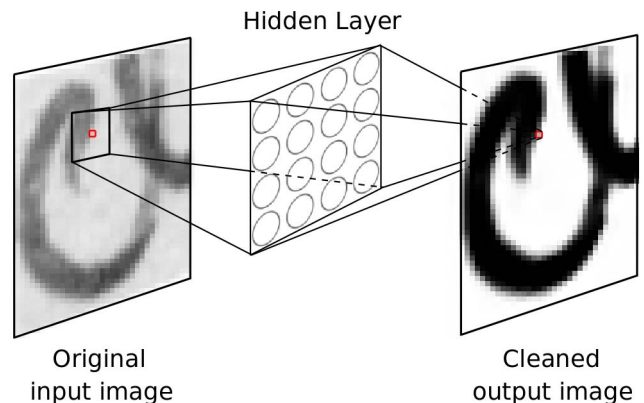


Fig. 1. Enhancer-MLP: An MLP to enhance and clean images. The entire input image is cleaned by scanning it with the neural network.

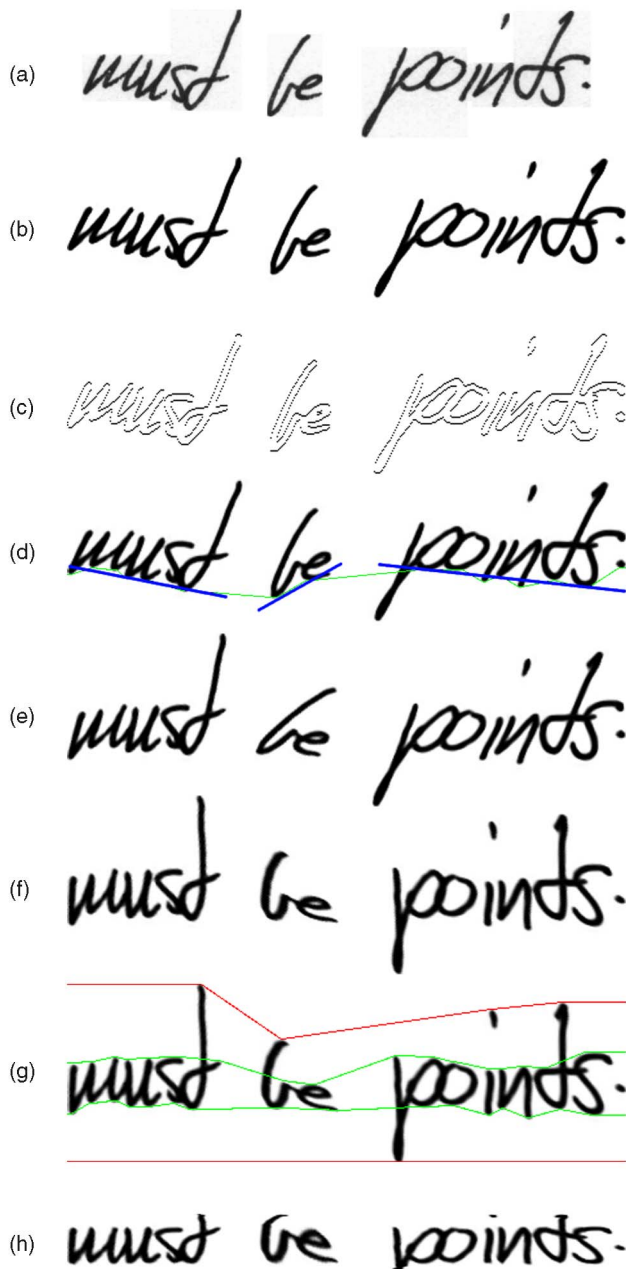


Fig. 2. Preprocessing example. (a) Original image from the IAM database. (b) Cleaned image. (c) Text contour extracted to obtain local extrema. (d) Local extrema classified by an MLP as the lower baseline to be used for slope correction. (e) Slope-corrected image. (f) Slant-corrected image. (g) Local extrema labeled by an MLP as belonging to the four reference lines to be used for size normalization. (h) Normalized final image.

classifying them by supervised machine learning techniques [40], [12]. The points to be classified are local vertical extrema of text contours which are used to determine the reference lines: line of ascenders, upper baseline, lower baseline, and line of descenders (see Fig. 3). These reference lines provide an efficient way to perform slope removal and size normalization.

In a first step, and once the text line image has been cleaned, the local extrema are obtained. First, a vertical

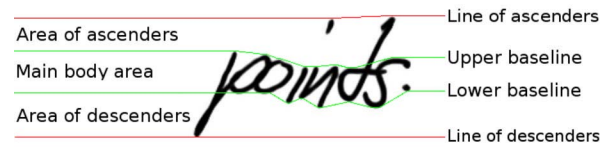


Fig. 3. Example of text line image with the different areas (ascenders, descenders, and main body areas) and the reference baselines (upper and the lower baselines, and the lines of ascenders and descenders) of the script.

contour of the image is extracted (see Fig. 2c). After that, the contour points are grouped into lines following a proximity criterion: Two pixels on adjacent columns are considered to belong to the same line when the difference between their vertical coordinates is less than 3. Finally, the maxima of the upper contours and the minima of the lower contours are computed.

Since the lower baseline suffices to correct the slope, an MLP that is trained to classify local extrema as belonging or not belonging to the lower baseline is used (this MLP will be called Slope-MLP). The input to this Slope-MLP is a window that is centered at the pixel to be classified.

Once the lower baseline points have been detected, the image is horizontally divided into segments in order to apply the slope correction to every segment. A vertical histogram projection is used to estimate the mean space width between ink regions, and this value is used as a threshold to split the image into segments. For each of these segments, the lower baseline is estimated by means of least-squares fitting of the lower baseline points. An example of the splitting process of the estimated lower baseline is shown in Fig. 2d. These lower baselines are used to correct the slope and the vertical relative positions of the segments. Fig. 2e illustrates an example of a slope-corrected image.

2.3 Slant Removal

After the slope correction, slant is removed by means of a two-step method. In the first step, a global slant angle is estimated and removed by performing a shear operation to the image for every integer angle between an interval (in this case, $[-50^\circ, 50^\circ]$) and choosing the sheared image whose vertical projection has the maximum variance [41]. In the second step, a novel nonuniform local slant correction method is applied: An MLP (from now on called Slant-MLP) is trained to estimate if a given column of the text line image is slanted. Using a sliding input window, the Slant-MLP is applied to every column of the image with some additional columns of context, for each integer angle in $[-50^\circ, 50^\circ]$. This procedure generates a matrix which contains the estimated score of correcting each column with each slant angle. Finally, a dynamic programming algorithm is applied over this matrix [42] to obtain the path with maximum score which also satisfies a smoothness criterion (the slant angle must not change more than $\pm 1^\circ$ per column). This sequence of angles conforms the input to a nonuniform shear that generates the final slant-corrected image. The total computational cost is linear with the size of the image and the number of shear angles considered. The whole slant removal process is illustrated in Fig. 4 and an example of a slant-corrected image is shown in Fig. 2f.

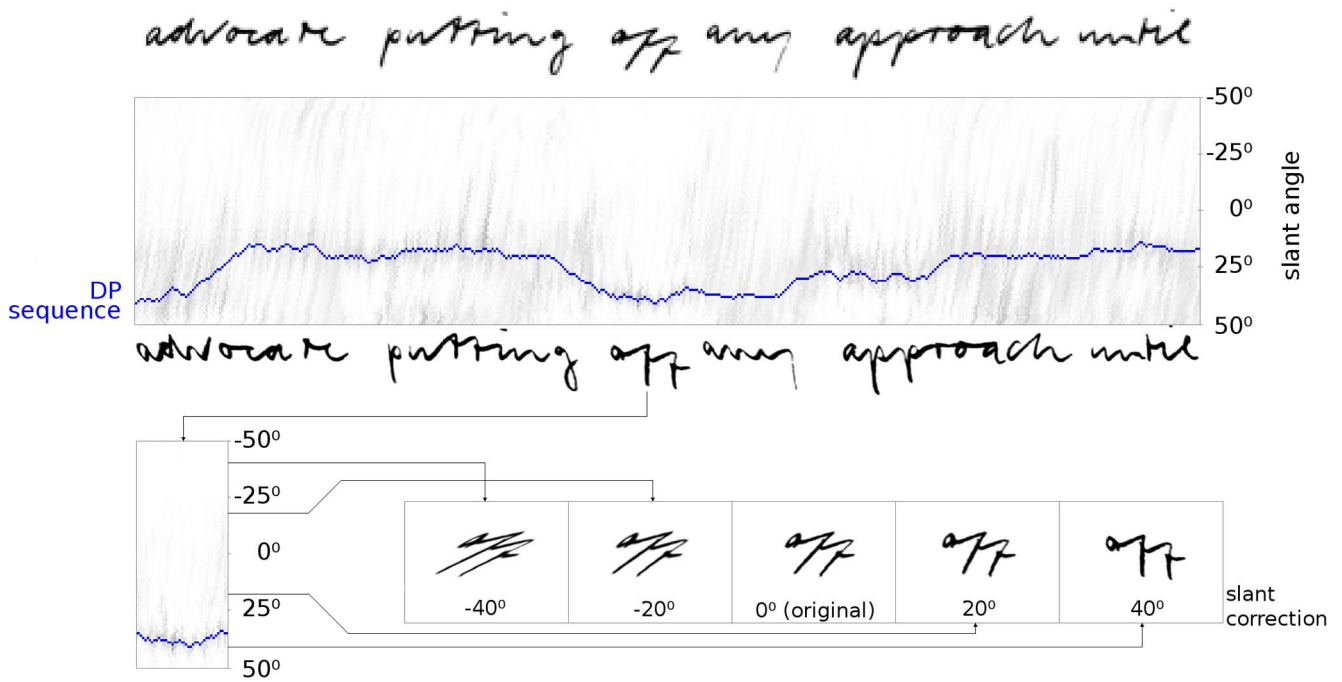


Fig. 4. An example of slant removal: the original text line image and the slant-corrected text line image. The Slant-MLP estimates a measure of the slant angle of each pixel column, shown as a gray level matrix. A dynamic programming (DP) algorithm obtains the optimal sequence of slant angles. Beneath the matrix is a detail of it for a segment of the text line image.

2.4 Size Normalization

When the image is slope and slant-corrected, the size of the text line is normalized in order to minimize the variations in size and position of its three zones (main body area, ascenders, and descenders). Furthermore, the normalized size of ascenders and descenders is reduced with respect to the body since they are not as informative (the presence or absence of ascenders and descenders is preserved, as well as the width, but not the actual height).

One approach to size normalization consists of detecting the reference baselines and normalizing the size according to them [18], [40], [12], [43]. Following this idea, our size normalization method detects and classifies the local extrema using the same method based on ANNs described in Section 2.2. This time, local extrema are classified into five classes (the four reference lines and points not belonging to any of these lines) by using another MLP (Normalize-MLP). Points belonging to the same class are used to obtain each reference line by linear interpolation (see Fig. 2g). These lines comprise the three zones to be normalized. The normalization process is performed for each column of the image by linearly scaling the three zones to a fixed height. Ascenders are reduced to 20 percent of the final image height and descenders are reduced to 10 percent. See Fig. 2h for an example of a normalized image.

2.5 Feature Extraction

After preprocessing, a feature extraction method is applied to capture the most relevant characteristics of the character to recognize. In our system, a handwritten text line image is converted into a sequence of fixed-dimension feature vectors. Following [10], features are extracted by applying a grid to the image and computing three values for each cell of the grid: the normalized gray level and the horizontal and vertical gray level derivatives. A grid of square cells

with 20 rows has been used, so every feature vector is composed of 60 values. An example of a graphical representation of the features obtained in this way is shown at the top of Fig. 5.

3 HYBRID HMM/ANN MODELING

For small vocabulary handwriting recognition tasks (for example, check amounts or postal addresses), it is possible to model words individually. But, for large vocabulary or even unconstrained tasks, the only feasible approach is to recognize individual graphemes and map them onto complete words belonging to a fixed vocabulary Ω . The same problem has to be addressed for automatic speech recognition, and HMMs have been accepted as the standard solution [44]. For offline handwritten text recognition, the image is converted into a sequence $X = (x_1 \dots x_m)$ of feature vectors and, under the statistical approach to pattern recognition [44], [45], the goal of general handwritten text recognition is to find the likeliest word sequence $W^* = (w_1 \dots w_n)$ maximizing the a posteriori probability:

$$W^* = \operatorname{argmax}_{W \in \Omega^+} P(W|X). \quad (1)$$

The application of the Bayes rule leads to a decomposition of $P(W|X)$ into the optical model $P(X|W)$ and the statistical language model $P(W)$. The problem can then be reformulated as:

$$W^* = \operatorname{argmax}_{W \in \Omega^+} P(X|W)P(W). \quad (2)$$

In state-of-the-art handwritten text recognition systems, $P(X|W)$ is usually estimated by an HMM-based recognizer and a word n -gram language model is usually used to

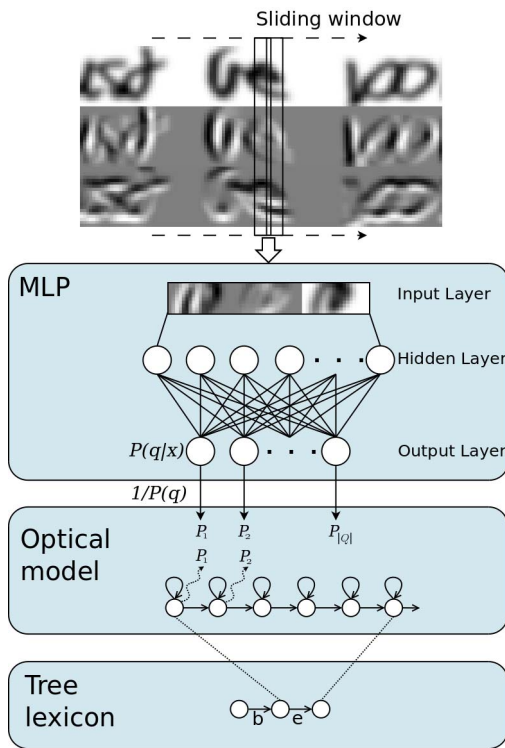


Fig. 5. A scheme of the proposed hybrid HMM/ANN recognition system. First, the image is preprocessed (see Fig. 2) and the resulting feature vector, plus a left and right context, is processed by an MLP. The $|Q|$ outputs of the MLP (after dividing by the prior state probabilities) are used as emission probabilities in the HMMs.

approximate $P(W)$. Typically, each grapheme is modeled by a left-to-right HMM and the number of states is chosen globally or individually for each grapheme. Gaussian mixtures are used to model the output distributions in each state q given the feature vector x , $P(x|q)$. The Baum-Welch algorithm is used for training the HMMs, whereas the Viterbi algorithm is used for recognition.

3.1 The Hybrid HMM/ANN Approach

In the hidden Markov modeling approach, the emission probability density $P(x|q)$ must be estimated for each state q of the Markov chains, that is, the probability of the observed feature vector x given the hypothesized state q of the model. In the proposed hybrid HMM/ANN approach, the emission probabilities are provided with a neural network since ANNs can be trained to estimate probabilities that are related to these emission probabilities. In particular, an MLP can be trained to approximate the a posteriori probabilities of states, $P(q|x)$, if each MLP output unit is associated with a specific state of the model and if it is trained as a classifier [14], [46]. In order to obtain such distribution for every state q from the set Q of Markov chain states, the softmax activation function has been chosen at the output layer:

$$f(y_q) = \frac{\exp(y_q)}{\sum_{i \in Q} \exp(y_i)}, \quad (3)$$

where y_q is the q th output value before applying the softmax function. This activation function enables the

estimation of valid probability values, i.e., to lie between zero and one and to sum to one.

The a posteriori probability estimates from the MLP outputs, $P(q|x)$, can be converted to emission probabilities $P(x|q)$ by applying Bayes rule:

$$P(x|q) = \frac{P(q|x)P(x)}{P(q)}. \quad (4)$$

The class priors $P(q)$ can be estimated from the relative frequencies of each state from the information produced by a forced Viterbi alignment of the training data. Thus, the scaled likelihoods $P(x|q)/P(q)$ can be used as emission probabilities in the proposed system since, during recognition, the scaling factor $P(x)$ is a constant for all classes [14]. This allows MLPs to be integrated into hybrid structural-connectionist models via a statistical framework.

The advantages of this approach are the discriminate training criterion (all MLP parameters are updated in response to every input feature vector) and the fact that it is no longer necessary to assume an a priori distribution of the data. Furthermore, if left and right contexts are used at the input of the MLP, important contextual information can be incorporated into the probability estimation process.

Another strength of this approach is that computing emission probabilities with hybrid HMM/ANN models is usually faster than conventional HMMs with Gaussian emissions since it only requires a forward pass of the MLP for all states of the Markov chains.

On the other hand, one of the weaknesses of this hybrid HMM/ANN approach is that every feature vector must be labeled to train the MLP. However, this is not a serious drawback since these labels can be generated by running a previously trained handwriting recognition system in a forced alignment mode in order to initialize these labels.

In our experiments, we modeled graphemes with left-to-right Markov chains and a single neural network with one output unit for each state of the Markov chains was used to estimate the distribution probabilities. An MLP with sigmoid hidden units and softmax output units is used. The estimates of the posterior probabilities computed by the MLP are divided by the prior state probabilities resulting in scaled likelihoods which are used as emission probabilities in the HMMs. A scheme of the proposed hybrid HMM/ANN recognition system is shown in Fig. 5.

3.2 Training the HMM/ANN Models

The training of the MLP is discriminant at the state level of Markov chains since each output is optimized during training by samples of its own class as well as by samples of the other classes. Training of the whole hybrid HMM/ANN system is done by an iterative Expectation-Maximization algorithm, where the training of the supervised ANN and either Baum-Welch or Viterbi alignment of the training corpus are alternated. We have opted for Viterbi alignment and the training procedure proceeds as follows:

1. Assign an initial labeling of desired MLP outputs to every feature vector of the training and validation data sets. This labeling can be computed by dividing the image into equal parts or by using a previously

trained handwritten recognition system in a forced alignment mode.

2. Assign an initial nonzero value to transition probabilities of the Markov chains.
3. Train the supervised ANN with the training pairs, using the mean-square error (MSE) on the validation data set as the stopping criterion.
4. Use the partially trained hybrid ANN/HMM models to perform a forced Viterbi alignment of the training data. This Viterbi procedure uses the class priors estimated from the relative frequencies of each class in the training data. This Viterbi alignment is used both for obtaining a new segmentation or labeling of the training and validation sets and also for counting the number of times each HMM transition has been used. These counts are used to reestimate the transition probabilities.
5. Go to step 3 until convergence, that is, until the difference between two consecutive iterations is below a threshold.

4 EXPERIMENTAL SETUP

4.1 The IAM Database

All experiments reported in this paper are conducted on handwritten text lines from the IAM database [32]. The version 3.0 of this database includes over 1,500 scanned forms of handwritten text from more than 650 different writers, for a total of more than 13,000 fully transcribed handwritten lines, without restrictions on the writing style or the writing instrument used. The sentences have been extracted from the Lancaster-Oslo/Bergen (LOB) text corpus [47].

A writer-independent text line recognition task has been considered. The subset of the IAM database used in this work consists of 6,161 training lines (from 283 writers), 920 validation lines (56 writers), and 2,781 test lines (161 writers). All of these data sets are disjoint, and no writer has contributed to more than one set. These partitions are the same as those used in several works by Bunke et al. [33], [48], [49].

A total of 87,967 instances of 11,320 distinct words occur in the union of the training, validation, and test sets. Lexicon is modeled with 78 characters: 26 lowercase letters, 26 uppercase letters, 10 digits, 14 punctuation marks, the space, and a character for garbage symbols.

4.2 MLP for Image Cleaning

As described in Section 2.1, an MLP has been used for image cleaning by learning the appropriate filter from examples.

4.2.1 Training Data

Original noisy images from the IAM database and the same images that were cleaned by hand formed the training pairs. Additionally, artificially noised images (created by following the ideas presented in [34]) were also used as training data.

4.2.2 Enhancer-MLP

In this case, the MLP is used for regression: The input is a fixed-sized moving window of 11×11 pixels centered at the pixel to be cleaned, and the output is the restored value of the current pixel (see Fig. 1). The Enhancer-MLP has two hidden layers of 32 and 16 sigmoid units and one output linear unit. Training was performed using the stochastic

version of the backpropagation algorithm with momentum term [46], using the mean-square error (MSE) function. The last column of Table 1 shows the topology of the MLPs which are used for preprocessing.

4.3 MLPs for Slope Removal and Size Normalization

As pointed out in Sections 2.2 and 2.4, two MLPs to classify local extrema as belonging to one of the reference lines (lower line, upper line, line of descenders, or line of ascenders) are needed as part of the slope removal and image size normalization processes.

4.3.1 Training Data

We needed supervised training patterns to train MLPs to classify interest points as belonging to the reference lines. A subset of 1,000 images from the IAM training set has been used. Local extrema of the 1,000 images were semi-automatically labeled using an active learning approach: First, a horizontal projection algorithm was used to classify the points belonging to each reference line of a subset of the 1,000 images; second, the subset of images was manually corrected using a graphical tool designed for this purpose [12]; third, these images were used to train an MLP to classify interest points. With this "pretrained" MLP, interest points of the 1,000 images were automatically classified, and, afterward, all of the images were manually supervised. At the end of this process, we had a training set composed of the interest points of the 1,000 images: 800 lines were used as training data and the remaining 200 lines were used as validation data.

4.3.2 Slope-MLP

The Slope-MLP was trained to classify local extrema as belonging to or not belonging to the lower baseline. The Slope-MLP input is a moving window around the current pixel, being the choice of an appropriate window size, a trade-off between context and input size. To partially overcome this problem, we have opted to use a fisheye distortion centered at the pixel to classify (see Fig. 6 for an example) [12]. The fisheye distortion maintains a very accurate resolution near the center and, at the same time, has a much smaller size than using the original image. In this way, a detailed image near the interest point and a coarse representation of the relative position of the surrounding text is obtained: The input to this Slope-MLP is a window of 500×250 pixels centered at the point to be classified, downsampled to 50×30 values using the fisheye distortion. Two output units with a softmax activation function were used to determine whether or not the current pixel belonged to the lower baseline. After doing a scanning of topologies, two hidden layers of 64 and 16 sigmoid units were used.

4.3.3 Normalize-MLP

Size normalization was achieved by using a second MLP, which classifies the local extrema into five classes (the four reference lines and points not belonging to any of these lines). The input to this Normalize-MLP is the same as the Slope-MLP input and the output corresponded to five output units with softmax activation function. We used two hidden layers of 64 and 128 sigmoid units.



Fig. 6. Fisheye lens example (from up-to-down): original image of 500×250 pixels centered at the pixel to be classified; the same image with a fisheye lens distortion; the image downsampled to 50×30 values used by the neural network classifier.

Both MLPs, Slope-MLP and Normalize-MLP, were trained using the stochastic version of the backpropagation algorithm with momentum term and the cross-entropy error function.

4.4 MLP for Slant Removal

As described in Section 2.3, part of the process of slant removal needs an MLP to determine whether or not an image has slant.

4.4.1 Training Data

The same set of 1,000 images was manually slant-corrected in a nonuniform way by using a graphical tool. The user specifies a series of slant angles which are interpolated for every image column. This information is used to train the Slant-MLP. As before, 200 images were used for validation.

4.4.2 Slant-MLP

Each image is sheared for different integer angles from -50° to $+50^\circ$ and resized to 40 pixels height, preserving the aspect ratio. The input to this Slant-MLP is a square of 40×40 pixels centered at the column to be evaluated, and the output is a measure of the local slant presence (shown as gray levels in Fig. 4). After doing a parameter and topology scanning, two hidden layers of 64 and 8 units were used. Training was performed using the stochastic version of the backpropagation algorithm with momentum term, using the mean-square error function.

5 EXPERIMENTS

5.1 Dictionary and Language Model

A word bigram language model was trained with three different text corpora: the LOB corpus [47] (excluding those sentences that contain lines from the test set of the IAM database), the Brown corpus [50], and the Wellington corpus

[51]. In order to cope with the fact that lines are fragments of sentences, we have randomly broken each sentence from the corpus into fragments to resemble lines. All of this text is supplemented with the training lines from the IAM database. Then, the final training material is comprised of:

- Sentences: 51,560 LOB sentences (2,134 sentences which contained IAM test lines were eliminated), 51,763 Brown sentences, and 20,592 Wellington sentences.
- Fragments of sentences to resemble lines: More than 400,000 lines randomly obtained from the above set of sentences.
- Lines: Finally, the 6,161 IAM training lines were also added.

The bigram language model used in the recognition systems was generated, using the SRI Language Modeling Toolkit [52] with the modified Kneser-Ney back-off discounting.

To achieve unconstrained handwriting recognition, an open dictionary, composed of the 20,000 most frequently occurring (case insensitive) words in the training material, was used to test our recognition systems.

5.2 Measuring Recognition Performance

The recognition performance was measured in terms of the Word Error Rate (WER), which is computed by comparing the output of the recognizer with the reference transcription. WER is defined as the number of word errors (insertions, substitutions, and deletions) summed over the whole test set and divided by the total number of words in the transcriptions of the reference set:

$$\text{WER} = 100 \times \frac{\text{insertions} + \text{substitutions} + \text{deletions}}{\text{total number of words}}. \quad (5)$$

A null WER is only reached if the recognizer output matches the reference transcription exactly.

The Character Error Rate (CER) was also measured for the final test experiments. CER is defined as expression (5), but with characters instead of words.

In order to properly compare different systems, it is highly desirable to provide not only the value of the WER (or CER) but also a confidence interval for it. In [53], the author proposes a method for computing these intervals without simulations. Following his work, we have computed the confidence interval for every experiment with the IAM validation and test sets, which are composed of 920 and 2,781 lines, respectively. In every experiment, the computed intervals correspond to a 95 percent confidence level.

5.3 Baseline Experiments: HMMs

Baseline recognition HMM experiments were conducted, using continuous density HMMs with diagonal covariance matrices of 64 Gaussians in each state, and with a left-to-right topology without skips. The 78 optical models were trained and tested with the HTK toolkit [54].

The validation set of the IAM database was used to optimize the number of states of the optical HMMs and the integration of the statistical language model. A bigram language model and an open dictionary were used as

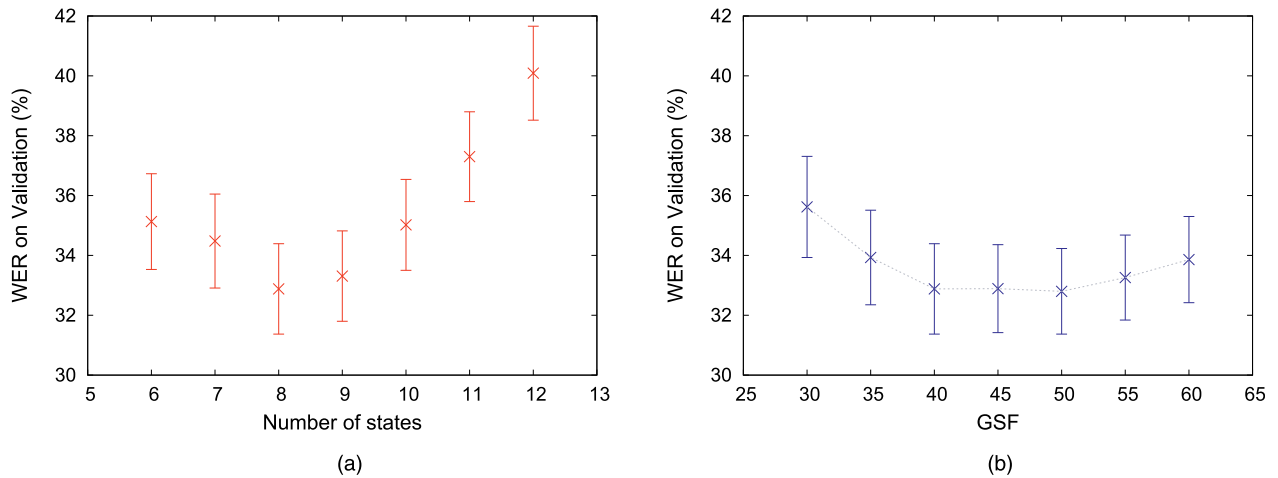


Fig. 7. Tuning HMMs: Word Error Rate of the HMMs on the validation set (a) varying the number of states and (b) for different Grammar Scale Factors. WER is given with a 95 percent confidence interval.

explained in Section 5.1.¹ Table 2 summarizes the experiments varying the number of states of the HMMs, and Fig. 7 plots the word error rate with a 95 percent confidence level interval. From these figures, it can be observed that several configurations achieved equivalent statistical performance, the 8-state HMMs being the best topology.

To compensate for scale differences between the likelihood values $P(X|W)$ from the HMMs and the probabilities $P(W)$ provided by the language model in (2), a grammar scale factor is used to weight the influence of the bigram language model against the optical model. The grammar scale factor used for these experiments was fixed to 40. Afterward, the grammar scale factor was optimized by systematically testing values from 20 to 60 on the 8-state HMMs. Table 3 shows the performance of this experiment on the validation set, using bigrams. It can be observed that performance is almost identical between a grammar scale factor of 40 and 50, with the lowest word error rate of 32.8 percent being achieved with a grammar scale factor of 50. All of these results are also plotted in Fig. 7 with a 95 percent confidence level interval.

5.4 Experiments with Hybrid HMM/ANN Models

Hybrid HMM/ANN models, with a different number of states and different topologies and parameters of MLP, were tested. In all cases, the MLP input consisted of nine consecutive feature vectors (the central feature vector and a context of four vectors at each side). The softmax outputs (after being divided by the prior state probabilities) were used as emission probabilities of the states of the 78 optical models. Thus, we trained fully connected MLPs of 540 input units (the 60-dimensional nine feature vectors). The number of output units is determined by the total number of states of the 78 optical models (from 78×6 output units for 6-state HMMs to 78×9 output units for 9-state HMMs) since each output unit of the MLP is related to one state of the HMMs. The number of hidden units was determined empirically by measuring the MSE on the validation set. Other parameters, such as the learning rate and the momentum term, were also empirically tuned with the validation data.

1. Those LOB sentences which contained IAM validation lines were also excluded to estimate the bigram language model for the tuning experiments.

Training was performed using stochastic backpropagation with momentum and the mean-square error function. In order to monitor the generalization performance during learning and to stop training when there was no longer an improvement, the validation set was used. More than five million training patterns (corresponding to the 6,161 training lines) and close to 800,000 validation patterns (from 920 lines) composed the training and validation data sets, respectively. Due to the large time requirements to train the MLPs, we used a resampling algorithm: Only 300,000 training patterns and 200,000 validation patterns were used in each training epoch. These subsets were randomly selected in each run.

The emission probabilities are obtained by dividing the a posteriori probability estimates from the MLP outputs by the class priors. The a priori probabilities of the states are estimated from the relative frequencies of each state, which are computed from the segmentation given by a forced Viterbi alignment of the training data.

TABLE 2
Tuning the Number of States of the HMMs

Model	WER on Validation (%)
6-state HMMs	35.1 \pm 1.6
7-state HMMs	34.5 \pm 1.6
8-state HMMs	32.9 \pm 1.5
9-state HMMs	33.3 \pm 1.5
10-state HMMs	35.0 \pm 1.5
11-state HMMs	37.3 \pm 1.5
12-state HMMs	40.0 \pm 1.6

Word Error Rate of the HMMs on the validation set.

TABLE 3
Tuning the GSF

Model	GSF	WER on Validation (%)
8-state HMMs	30	35.6 \pm 1.7
8-state HMMs	35	33.9 \pm 1.6
8-state HMMs	40	32.9 \pm 1.5
8-state HMMs	45	32.9 \pm 1.5
8-state HMMs	50	32.8 \pm 1.4
8-state HMMs	55	33.3 \pm 1.4
8-state HMMs	60	33.8 \pm 1.4

Word Error Rate of the best HMMs on the validation set for different grammar scale factors.

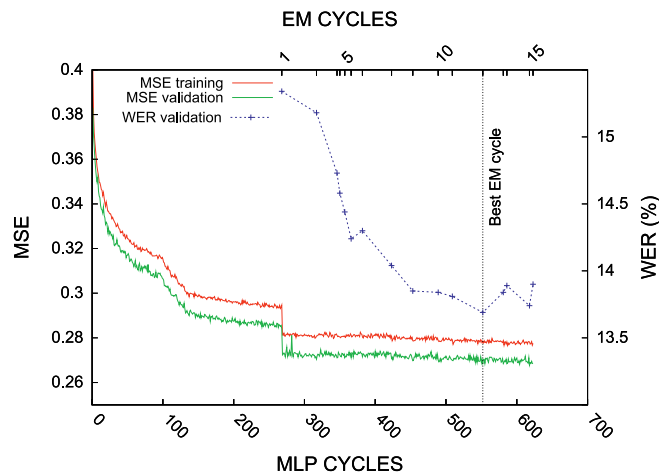


Fig. 8. Evolution of the EM training algorithm for the 7-state HMMs and an MLP of two hidden layers of 192 and 128 units each. The mean-square error (MSE) is shown for the training and validation data, along with the WER on validation of each iteration of the EM algorithm.

The initial segmentation at state level that is required to train an MLP at the first step of the EM algorithm (see Section 3.2) was generated by running a “pretrained” hybrid HMM/ANN handwriting recognition system in a forced alignment mode. The word error rate on the validation set, using a closed dictionary from the IAM task composed of 10,353 words and a bigram language model estimated with the LOB corpus, was used as the stopping criterion. Around 10 iterations of the EM training algorithm were enough for all the configurations. Fig. 8 illustrates a typical evolution of the EM training, showing the evolution of the mean-square error on the training and validation set. The evolution of the WER on validation is also shown in the graphic. The hybrid HMM/ANN models were trained and tested with the APRIL toolkit [55], which was developed for neural networks and pattern recognition tasks in our research group.

Table 4 shows the performance achieved for the different configurations of the hybrid HMM/ANN systems on the validation data set, using the bigram language model and the open dictionary, as explained in Section 5.1. The lowest WER was obtained by using 7-state Markov chains and an MLP topology of two hidden layers of 192 and 128 units, respectively. Further experiments were conducted with the optimized configuration. Table 5 shows the performance of this hybrid HMM/ANN system on the validation data set, using the bigram language model with different grammar scale factors. The grammar scale factor has been optimized by systematically testing values from 6 to 16 on the validation text lines. Small changes in word error rate are observed, and the best performance, 19.0 percent, was

TABLE 4
Tuning the Topology of the Hybrid HMM/ANN Models

Model	WER on Validation (%)
6-state HMMs, MLP 192-128	19.5 ± 1.3
7-state HMMs, MLP 192-128	19.0 ± 1.2
8-state HMMs, MLP 384-128	19.1 ± 1.2
9-state HMMs, MLP 384-128	19.5 ± 1.2

Word Error Rate of the hybrid HMM/ANN models on the validation set.

TABLE 5
Tuning the GSF

Model	GSF	WER on Validation (%)
7-state HMMs, MLP 192-128	6	21.3 ± 1.4
7-state HMMs, MLP 192-128	8	19.6 ± 1.3
7-state HMMs, MLP 192-128	10	19.0 ± 1.2
7-state HMMs, MLP 192-128	12	19.0 ± 1.2
7-state HMMs, MLP 192-128	14	19.3 ± 1.2
7-state HMMs, MLP 192-128	16	20.1 ± 1.1

Word Error Rate of the best hybrid HMM/ANN models on the validation set for different grammar scale factors.

achieved with a grammar scale factor of 10 or 12. All of these word error rate results are plotted in Fig. 9 with a 95 percent confidence level interval.

6 DISCUSSION AND COMPARISON

This section describes the performance of the optimized systems on the test set, and a comparison is made with the best published results. Experiments to study the influence of the dictionary on the recognizer were also carried out.

6.1 HMM versus Hybrid HMM/ANN Models

Table 6 shows the error rate of the recognized test lines of the IAM task using the best HMM and the best hybrid HMM/ANN systems. We tested each system with the open dictionary of 20,000 words and a bigram language model, as explained in Section 5.1. For each recognition experiment, two performance figures were obtained: the word error rate and the character error rate. No parameters were optimized on the test set.

Our baseline experiment achieved comparative performances with state-of-the-art HMM systems: a WER of 38.8 percent ±1.0 with a 95 percent confidence interval and a CER of 18.6 percent in the interval (18.0, 19.2). Our final hybrid HMM/ANN system achieved excellent results: a WER of 22.4 percent in the interval (21.6, 23.2) and a CER of 9.8 percent in the interval (9.4, 10.2). The hybrid system outperforms our baseline in 16 points in WER, which represents a relative error rate reduction of 42 percent. Similarly, the character error rate improved nearly 9 points, which represents a relative percentage of improvement that is greater than 47 percent.

Besides the word and character error rates, another measure to consider when evaluating a recognition engine is the decoding time, since a high value might diminish the usability of the recognition system in practical applications. Unfortunately, this time is not reported by most authors. Our hybrid HMM/ANN prototype required an average time of 0.76 second per word for preprocessing and 0.65 second per word on decoding [56]. This CPU time was measured in a single core of an Intel Core 2 Quad CPU Q6600 @ 2.40 GHz using DDR2-800 MHz memory. These times could be reduced in the production stage of the recognition engine, and the latency could also be reduced by using several cores since many steps are highly parallelizable.

6.2 Influence of the Dictionary

A series of experiments to study the influence of the dictionary size were carried out. Open dictionaries with between 10,000 and 30,000 words were generated by taking

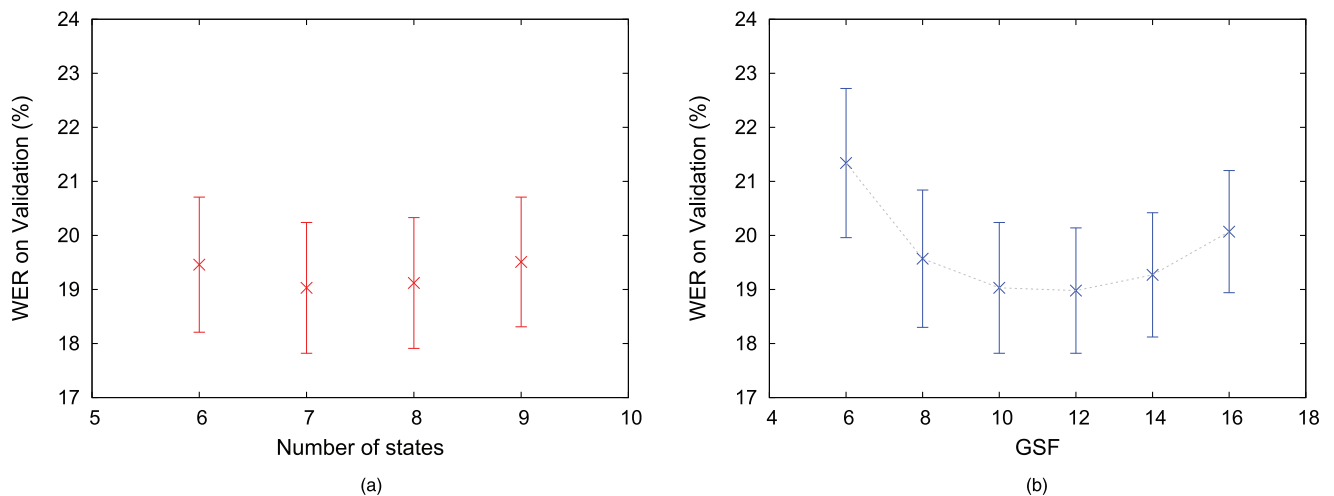


Fig. 9. Tuning HMM/ANN models: Word Error Rate of the HMM/ANN models on the validation set (a) varying the topology and (b) for different Grammar Scale Factors. WER is given with a 95 percent confidence interval.

the N most frequently occurring words in the training material for the language model.

Table 7 shows the word error rate and the character error rate of the test set when the size of the open dictionary was increased. The second column of this table shows the test set coverage, and the last two columns show the word and the character error rate of the test set, using a bigram language model estimated for each lexicon. To give an idea of the meaning of coverage in the IAM line task, consider, for example, that with the open dictionary of 20,000 words, a coverage of 78.86 percent was achieved, that is, 21.14 percent of the words in the test set were out-of-vocabulary words. However, if we measure the out-of-vocabulary running words, this figure falls to 4.99 percent. (The out-of-vocabulary running words were 1,268, that is, 1,268 running words from the 25,424 running words in the test set were not in the lexicon.) As expected, performance increased with lexicon size and test coverage. Fig. 10 plots this experiment against test set coverage.

Another experiment was also carried out to study the influence of using closed dictionaries. Two closed dictionaries were generated: one containing only the 4,953 words in

the IAM test line set and another one padding the first one up to 20,000 words (the most frequently occurring words in the training material for the language model). The influence of using the closed dictionaries is shown in Table 8, using a 95 percent confidence interval for WER. In this case, a bigram language model estimated for each lexicon was also used. Not surprisingly, the closed dictionary containing only the test set words achieved the best performance. The score with the 20,000 word closed dictionary was still better than those reached with open dictionaries.

6.3 Comparison with Other Systems

Comparisons to other recognition systems in the literature are difficult due to the lack of availability of common databases. With regard to the publications using the IAM database, a more detailed comparison is possible. In a very recent work, Graves et al. [33] presented a novel handwriting recognition system based on recurrent neural networks which achieved the best published recognition rates to date, a WER of 25.9 percent with bigrams. In order to compare both systems, we contacted the authors to exactly reproduce the same experimental conditions. They

TABLE 6
Testing the Systems

Best model	Results of Test (%)	
	WER	CER
8-state HMMs	38.8 \pm 1.0	18.6 \pm 0.6
7-state HMMs, MLP 192-128	22.4 \pm 0.8	9.8 \pm 0.4

Error Rate of the HMMs and the hybrid HMM/ANN models on the test set.

TABLE 7
Influence of the Dictionary Size (with Open Dictionaries)

Dictionary size	Coverage (%)	Results of Test (%)	
		WER	CER
10,000	66.57	26.9 \pm 0.9	11.7 \pm 0.5
15,000	74.28	24.2 \pm 0.8	10.5 \pm 0.4
20,000	78.86	22.4 \pm 0.8	9.8 \pm 0.4
25,000	81.97	21.9 \pm 0.8	9.4 \pm 0.4
30,000	84.39	21.2 \pm 0.8	9.1 \pm 0.4

Word Error Rate of hybrid HMM/ANN models on the test set.

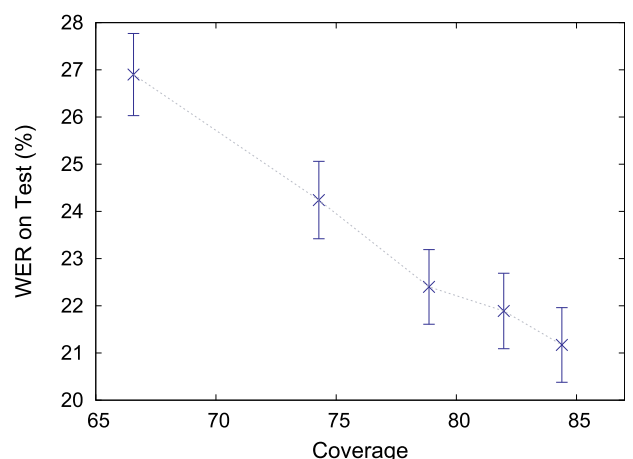


Fig. 10. HMM/ANN performance with open dictionaries plotted against test set coverage. WER on test is given with a 95 percent confidence interval.

TABLE 8
Influence of Using Closed Dictionaries

Dictionary size	Results of Test (%)	
	WER	CER
4,953	15.5 \pm 0.7	6.9 \pm 0.4
20,000	16.8 \pm 0.7	7.5 \pm 0.4

Word Error Rate of hybrid HMM/ANN models on the test set.

provided us with their dictionary and language model, and we retrained our optical models, using the same grapheme set as them and normalizing each input feature to zero mean and unit variance. Also, case differences in words were taken into account during recognition and in the WER computation. Surprisingly, we obtained the very same 25.9 percent of WER for the test set, despite the fact that recurrent neural networks and HMM/ANN are very different approaches.

7 CONCLUSIONS

In this paper, we have presented a hybrid HMM/ANN system for recognizing unconstrained offline handwritten text lines. The key features of the recognition system are the novel approach to preprocessing and recognition, which are both based on ANNs. The preprocessing is based on using MLPs:

- to clean and enhance the images,
- to automatically classify local extrema in order to correct the slope and to normalize the size of the text lines images, and
- to perform a nonuniform slant correction.

The recognition is based on hybrid optical HMM/ANN models, where an MLP is used to estimate the emission probabilities.

The main property of ANNs which is useful for preprocessing tasks is their ability to learn complex nonlinear input-output relationships from examples. Used for regression, an MLP can learn the appropriate filter from examples. We have exploited this property to clean and enhance the text images. Used for classification, MLPs can be used to determine the membership of interest points from the image to the reference lines, which is useful for slope correction and size normalization, and to locally detect slant in a text image. This preprocessing behaved favorably when compared to other preprocessing techniques. We tested our HMM and HMM/ANN systems, performing the same experiments that those presented here, but by using more classical techniques to correct slope, slant, and size normalization [10]. We obtained a 54.3 percent and 29.8 percent test WER, respectively, which represent a percentual decrease of 29 percent and 25 percent when compared to the test results from Table 6.

The proposed hybrid HMM/ANN recognition system outperformed our baseline experiment, which is a state-of-the-art HMM-based system that includes our preprocessing. The novel hybrid HMM/ANN approach obtained an impressive 42 percent relative improvement in WER over our baseline. We compared our system with the recurrent neural network approach presented in [33] under the same experimental conditions, and we obtained the same results.

Our next goal is to upgrade our recognition engine by using ensembles of MLPs [46], [16], by combining several

recognizers [49], [57], and by using deep connectionist architectures [58], [59]. The first very basic idea is to use several MLPs rather than just a single one to solve a given pattern classification or regression task [46], [16]. This idea can be directly applied to the optical hybrid HMM/ANN models, using an ensemble of MLPs to estimate the emission probabilities of the Markov chains as well as using ensembles of MLPs in every preprocessing step. Another idea is to combine several individual recognition systems (based on HMMs, HMM/ANN models, or recurrent ANNs) and specialized classifiers [49], [57]. Finally, as pointed out in [58], [59], using deep learning methods would lead us to better trained ANNs, which could improve every step of our recognition engine.

ACKNOWLEDGMENTS

The authors acknowledge the valuable help provided by Moisés Pastor, Juan Miguel Vilar, Alex Graves, and Marcus Liwicki. Thanks are also due to the reviewers and the Editor-in-Chief for their many valuable comments and suggestions. This work has been partially supported by the Spanish Ministerio de Educación y Ciencia (TIN2006-12767) and by the BPFI 06/250 Scholarship from the Conselleria d'Empresa, Universitat i Ciencia, Generalitat Valenciana.

REFERENCES

- [1] T. Steinherz, E. Rivlin, and N. Intrator, "Offline Cursive Script Word Recognition—A Survey," *Int'l J. Document Analysis and Recognition*, vol. 2, no. 2, pp. 90-110, 1999.
- [2] R. Plamondon and S.N. Srihari, "On-Line and Off-Line Handwriting Recognition: A Comprehensive Survey," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 22, no. 1, pp. 63-84, Jan. 2000.
- [3] N. Arica and F. Yarman-Vural, "An Overview of Character Recognition Focused on Off-Line Handwriting," *IEEE Trans. Systems, Man, and Cybernetics, Part C: Applications and Rev.*, vol. 31, no. 2, pp. 216-233, May 2001.
- [4] A. Vinciarelli, "A Survey on Off-Line Cursive Word Recognition," *Pattern Recognition*, vol. 35, no. 7, pp. 1433-1446, 2002.
- [5] H. Bunke, "Recognition of Cursive Roman Handwriting—Past, Present, and Future," *Proc. Seventh Int'l Conf. Document Analysis and Recognition*, vol. 1, pp. 448-459, Aug. 2003.
- [6] A. Koerich, R. Sabourin, and C. Suen, "Large Vocabulary Off-Line Handwriting Recognition: A Survey," *Pattern Analysis and Applications*, vol. 6, no. 2, pp. 97-121, 2003.
- [7] H. Fujisawa, "Forty Years of Research in Character and Document Recognition—An Industrial Perspective," *Pattern Recognition*, vol. 41, no. 8, pp. 2435-2446, 2008.
- [8] A. El-Yacoubi, M. Gilloux, R. Sabourin, and C.Y. Suen, "An HMM-Based Approach for Off-Line Unconstrained Handwritten Word Modeling and Recognition," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 21, no. 8, pp. 752-760, Aug. 1999.
- [9] U.-V. Marti and H. Bunke, "Using a Statistical Language Model to Improve the Performance of an HMM-Based Cursive Handwriting Recognition Systems," *Int'l J. Pattern Recognition and Artificial Intelligence*, vol. 15, no. 1, pp. 65-90, 2001.
- [10] A.H. Toselli, A. Juan, J. González, I. Salvador, E. Vidal, F. Casacuberta, D. Keysers, and H. Ney, "Integrated Handwriting Recognition and Interpretation Using Finite-State Models," *Int'l J. Pattern Recognition and Artificial Intelligence*, vol. 18, no. 4, pp. 519-539, 2004.
- [11] A. Vinciarelli, S. Bengio, and H. Bunke, "Offline Recognition of Unconstrained Handwritten Texts Using HMMs and Statistical Language Models," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 26, no. 6, pp. 709-720, June 2004.
- [12] J. Gorbe-Moya, S. España-Boquera, F. Zamora-Martínez, and M.J. Castro-Bleda, "Handwritten Text Normalization by Using Local Extrema Classification," *Proc. Eighth Int'l Workshop Pattern Recognition in Information Systems*, pp. 164-172, 2008.

- [13] Y. Bengio, "A Connectionist Approach to Speech Recognition," *Int'l J. Pattern Recognition and Artificial Intelligence*, vol. 7, no. 4, pp. 647-667, 1993.
- [14] H. Boullard and N. Morgan, *Connectionist Speech Recognition—A Hybrid Approach*. Kluwer Academic, 1994.
- [15] M.J. Castro and F. Casacuberta, "Hybrid Connectionist-Structural Acoustical Modeling in the ATROS System," *Proc. Sixth European Conf. Speech Comm. and Technology*, vol. 3, pp. 1299-1302, 1999.
- [16] M.J. Castro and F. Casacuberta, "Committees of MLPs for Acoustic Modeling," *Proc. Fifth Iberian Symp. Pattern Recognition*, pp. 797-807, 2000.
- [17] R. Gemellovo, F. Mana, and D. Albesano, "Hybrid HMM/Neural Network Based Speech Recognition in Loquendo ASR," http://www.loquendo.com/en/brochure/Speech_Recognition_ASR.pdf, 2008.
- [18] Y. Bengio, Y. LeCun, C. Nohl, and C. Burges, "LeRec: A NN/HMM Hybrid for On-Line Handwriting Recognition," *Neural Computation*, vol. 7, no. 6, pp. 1289-1303, 1995.
- [19] M. Schenkel, I. Guyon, and D. Henderson, "On-Line Cursive Script Recognition Using Time Delay Neural Networks and Hidden Markov Models," *Machine Vision and Applications*, vol. 8, no. 4, pp. 215-223, 1995.
- [20] S. Jaeger, S. Manke, and A. Waibel, "Npen++: An On-Line Handwriting Recognition System," *Proc. Seventh Int'l Workshop Frontiers in Handwriting Recognition*, pp. 249-260, 2000.
- [21] S. Marukatat, T. Artieres, B. Dorizzi, and P. Gallinari, "Sentence Recognition through Hybrid Neuro-Markovian Modelling," *Proc. Int'l Conf. Document Analysis and Recognition*, pp. 731-735, 2001.
- [22] É. Caillault and C. Viard-Gaudin, "Mixed Discriminant Training of Hybrid ANN/HMM Systems for Online Handwritten Word Recognition," *Int'l J. Pattern Recognition and Artificial Intelligence*, vol. 21, no. 1, pp. 117-134, 2007.
- [23] A. Graves, S. Fernandez, M. Liwicki, H. Bunke, and J. Schmidhuber, "Unconstrained Online Handwriting Recognition with Recurrent Neural Networks," *Advances in Neural Information Processing Systems*, vol. 20, pp. 577-584, MIT Press, 2008.
- [24] S. Knerr and E. Augustin, "A Neural Network-Hidden Markov Model Hybrid for Cursive Word Recognition," *Proc. 14th Int'l Conf. Pattern Recognition*, vol. 2, pp. 1518-1520, 1998.
- [25] A.W. Senior and A.J. Robinson, "An Off-Line Cursive Handwritten Recognition System," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 20, no. 3, pp. 309-321, Mar. 1998.
- [26] J.H. Kim, K.K. Kim, and C.Y. Suen, "An HMM-MLP Hybrid Model for Cursive Script Recognition," *Pattern Analysis and Applications*, vol. 3, pp. 314-324, 2000.
- [27] C. Burges, O. Matan, Y. LeCun, J. Denker, L. Jackel, C. Stenard, C. Nohl, and J. Ben, "Shortest Path Segmentation: A Method for Training a Neural Network to Recognize Character Strings," *Proc. Int'l Joint Conf. Neural Networks*, vol. 3, pp. 165-172, 1992.
- [28] C. Burges, J. Ben, J. Denker, Y. LeCun, and R. Nohl, "Off-Line Recognition of Handwritten Postal Words Using Neural Networks," *Int'l J. Pattern Recognition and Artificial Intelligence*, vol. 7, no. 4, pp. 689-704, 1993.
- [29] A. Koerich, Y. Leydier, R. Sabourin, and C. Suen, "A Hybrid Large Vocabulary Handwritten Word Recognition System Using Neural Networks with Hidden Markov Models," *Proc. Eighth Int'l Workshop Frontiers in Handwriting Recognition*, pp. 99-104, 2002.
- [30] Y. Tay, M. Khalid, R. Yusof, and C. Viard-Gaudin, "Offline Cursive Handwriting Recognition System Based on Hybrid Markov Model and Neural Networks," *Proc. IEEE Int'l Symp. Computational Intelligence in Robotics and Automation*, pp. 1190-1195, July 2003.
- [31] S. Marinai, M. Gori, and G. Soda, "Artificial Neural Networks for Document Analysis and Recognition," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 27, no. 1, pp. 23-35, Jan. 2005.
- [32] U.-V. Marti and H. Bunke, "The IAM-Database: An English Sentence Database for Offline Handwriting Recognition," *Int'l J. Document Analysis and Recognition*, vol. 5, no. 1, pp. 39-46, 2002.
- [33] A. Graves, M. Liwicki, S. Fernández, R. Bertolami, H. Bunke, and J. Schmidhuber, "A Novel Connectionist System for Unconstrained Handwriting Recognition," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 31, no. 5, pp. 855-868, May 2009.
- [34] J.L. Hidalgo, S. España, M.J. Castro, and J.A. Pérez, "Enhancement and Cleaning of Handwritten Data by Using Neural Networks," *Proc. Second Iberian Conf. Pattern Recognition and Image Analysis*, pp. 376-383, 2005.
- [35] D.J. Burr, "A Normalizing Transform for Cursive Script Recognition," *Proc. Sixth Int'l Conf. Pattern Recognition*, pp. 1027-1030, 1982.
- [36] R.M. Bozinovic and S.N. Srihari, "Off-Line Cursive Script Word Recognition," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 11, no. 1, pp. 68-83, Jan. 1989.
- [37] A. Vinciarelli and J. Luettin, "A New Normalization Technique for Cursive Handwritten Words," *Pattern Recognition Letters*, vol. 22, no. 9, pp. 1043-1050, 2001.
- [38] K.Y. Wong, R.G. Casey, and F.M. Wahl, "Document Analysis System," *IBM J. Research and Development*, vol. 26, no. 6, pp. 647-655, 1982.
- [39] V. Romero, M. Pastor, A.H. Toselli, and E. Vidal, "Improving Handwritten Off-Line Text Slant Correction," *Proc. Sixth IASTED Int'l Conf. Visualization, Imaging, and Image Processing*, pp. 389-394, 2006.
- [40] P. Simard, D. Steinkraus, and M. Agrawala, "Ink Normalization and Beautification," *Proc. Eighth Int'l Conf. Document Analysis and Recognition*, pp. 1182-1187, 2005.
- [41] M. Pastor, A. Toselli, and E. Vidal, "Projection Profile Based Algorithm for Slant Removal," *Proc. Int'l Conf. Image Analysis and Recognition*, pp. 183-190, 2004.
- [42] S. Uchida, E. Taira, and H. Sakoe, "Nonuniform Slant Correction Using Dynamic Programming," *Proc. Sixth Int'l Conf. Document Analysis and Recognition*, vol. 1, pp. 434-438, 2001.
- [43] J. Schenk, J. Lenz, and G. Rigoll, "On-Line Recognition of Handwritten Whiteboard Notes: A Novel Approach for Script Line Identification and Normalization," *Proc. 11th Int'l Workshop Frontiers in Handwriting Recognition*, pp. 540-543, 2008.
- [44] L. Rabiner and B.H. Huang, *Fundamentals of Speech Recognition*. Prentice-Hall, 1993.
- [45] F. Jelinek, *Statistical Methods for Speech Recognition*. MIT Press, 1997.
- [46] C.M. Bishop, *Neural Networks for Pattern Recognition*. Oxford Univ. Press, 1995.
- [47] S. Johansson, E. Atwell, R. Garside, and G. Leech, *The Tagged LOB Corpus: User's Manual*. Norwegian Computing Centre for the Humanities, 1986.
- [48] R. Bertolami and H. Bunke, "Ensemble Methods to Improve the Performance of an English Handwritten Text Line Recognizer," *Proc. Conf. Arabic and Chinese Handwriting Recognition*, pp. 265-277, 2008.
- [49] R. Bertolami and H. Bunke, "Hidden Markov Models-Based Ensemble Methods for Offline Handwritten Text Line Recognition," *Pattern Recognition*, vol. 41, no. 11, pp. 3452-3460, 2008.
- [50] W. Francis and H. Kucera, "Brown Corpus Manual, Manual of Information to Accompany a Standard Corpus of Present-Day Edited American English," technical report, Dept. of Linguistics, Brown Univ., 1979.
- [51] L. Bauer, "Manual of Information to Accompany the Wellington Corpus of Written New Zealand English," technical report, Dept. of Linguistics, Victoria Univ., 1993.
- [52] A. Stolcke, "SRILM: An Extensible Language Modeling Toolkit," *Proc. Int'l Conf. Spoken Language Processing*, pp. 901-904, 2002.
- [53] J.M. Vilar, "Efficient Computation of Confidence Intervals for Word Error Rates," *Proc. IEEE Int'l Conf. Acoustics, Speech, and Signal Processing*, pp. 5101-5104, 2008.
- [54] S.J. Young, P.C. Woodland, and W.J. Byrne, "HTK: Hidden Markov Model Toolkit V1.5," technical report, Cambridge Univ. Eng. Dept. Speech Group and Entropic Research Laboratories, Inc., 1993.
- [55] S. España-Boquera, F. Zamora-Martínez, M.J. Castro-Bleda, and J. Gorbe-Moya, "Efficient BP Algorithms for General Feedforward Neural Networks," *Proc. Second Int'l Work-Confer. Interplay between Natural and Artificial Computation, Part I, Bio-Inspired Modeling of Cognitive Tasks*, pp. 327-336, 2007.
- [56] S. España-Boquera, M.J. Castro-Bleda, F. Zamora-Martínez, and J. Gorbe-Moya, "Efficient Viterbi Algorithms for Lexical Tree Based Models," *Proc. Int'l Conf. Advances in Non-Linear Speech Processing*, pp. 179-187, 2007.
- [57] F. Zamora-Martínez, M.J. Castro-Bleda, S. España-Boquera, and J. Gorbe-Moya, "Improving Isolated Handwritten Word Recognition Using a Specialized Classifier for Short Words," *Current Topics in Artificial Intelligence*, pp. 61-70, 2010.
- [58] Y. Bengio, P. Lamblin, D. Popovici, and H. Larochelle, "Greedy Layer-Wise Training of Deep Networks," *Proc. Neural Information Processing Systems Conf.*, pp. 153-160, 2006.

[59] Y. Bengio, "Learning Deep Architectures for AI," *Foundations and Trends in Machine Learning*, vol. 2, no. 1, 2009.



Salvador España-Boquera received the Licenciado degree in computer science in 1998 from the Universitat Politècnica de València, Spain, and the Licenciado degree in mathematics in 2004 (best student award) from the Universitat de València, Spain. He joined the Departamento de Sistemas Informáticos y Computación of the Universitat Politècnica de València in 1999, where he is currently teaching and finishing the PhD degree. His current research interests include handwriting and speech recognition, image processing, sequence learning, language technologies, and algorithmics.



Maria Jose Castro-Bleda received the PhD degree in computer science from the Universitat Politècnica de València, Spain, in 1998. Since 1993, she has been with the Departamento de Sistemas Informáticos y Computación at the Universitat Politècnica de València, Spain, where she is currently an associate professor. Her main research interests include machine learning, speech and handwritten text recognition, and language technologies.



interests include document analysis, handwritten text normalization techniques, and handwriting recognition.

Jorge Gorbe-Moya received the MSc degree in artificial intelligence, pattern recognition, and digital imaging from the Universitat Politècnica de València, Spain, in 2008. The topic of his master's thesis was the application of neural networks to handwritten text preprocessing and recognition. He is currently working as a researcher in the Departamento de Sistemas Informáticos y Computación at the Universitat Politècnica de València.



Computación. His current research interests include the integration of neural network language models for handwritten text recognition and machine translation.

Francisco Zamora-Martinez received the MSc degree in artificial intelligence, pattern recognition, and digital imaging from the Universitat Politècnica de València, Spain, in 2008. The topic of his master's thesis was language modeling based on neural networks. He is now teaching at the Universidad CEU-Cardenal Herrera in Valencia and finishing the PhD degree at the Universitat Politècnica de València in the Departamento de Sistemas Informáticos y

► For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/publications/dlib.