



Writer independent on-line handwriting recognition using an HMM approach

Jianying Hu*, Sok Gek Lim, Michael K. Brown

Bell Laboratories, Lucent Technologies, 700 Mountain Avenue, Murray-Hill, NJ 07974, USA

Received 26 August 1998; accepted 12 January 1999

Abstract

In this paper we describe a Hidden Markov Model (HMM) based writer independent handwriting recognition system. A combination of signal normalization preprocessing and the use of invariant features makes the system robust with respect to variability among different writers as well as different writing environments and ink collection mechanisms. A combination of point oriented and stroke oriented features yields improved accuracy. Language modeling constrains the hypothesis space to manageable levels in most cases. In addition a two-pass *N*-best approach is taken for large vocabularies. We report experimental results for both character and word recognition on several UNIPEN datasets, which are standard datasets of English text collected from around the world. © 1999 Pattern Recognition Society. Published by Elsevier Science Ltd. All rights reserved.

Keywords: Handwriting recognition; Hidden Markov models; Invariant features; Segmental features; *N*-best decoding; UNIPEN

1. Introduction

Hidden Markov Models (HMMs) have been used with great success for stochastic modeling of speech for more than a decade now [1]. More recently they have also been applied to on-line handwriting recognition [2–5] with varying degrees of success. Some earliest works (i.e., Ref. [3]) consisted essentially of direct application of speech recognition systems to handwriting recognition with only substitution of some simple point oriented features for the speech features.

In the last few years we have been experimenting with HMM-based methods for on-line handwriting recognition. For some characteristics, the results have been better than those achieved with speech recognition. For example, writer independence is achieved for handwriting recognition by a combination of preprocessing to remove much of the variation in handwriting due to varying personal styles and writing influences, and by

feature invariance to reduce sensitivity to the remaining variations. Both of these techniques are more mature than corresponding methods for speech recognition. Indeed, except for pre-emphasis there is almost no processing before feature extraction in most speech recognition systems.

Probably a major reason for greater success in preprocessing handwriting is due to the visual nature of the results, where uniformity can be easily seen even by untrained observers, while speech preprocessing may yield a normalized voice that would require great expertise to judge uniformity. Invariance in handwriting features is also more readily obtained by taking advantage on known geometrical invariance relations.

In this paper we first give a complete overview of an on-line handwriting recognition system which we have developed over the last few years. The main novel aspects of this system compared to other HMM-based systems [2–4,6,7] include: signal preprocessing with a new word normalization method; application of invariant features; combination of high-level and low-level features; stochastic modeling with HMMs that incorporates a dynamically evolving language model; and *N*-best decoding combined with two-stage delayed stroke modeling for

*Corresponding author. Tel.: + 908-582-5660; fax: + 908-582-7308.

E-mail address: jianhu@bell-labs.com (J. Hu)

large vocabulary word recognition. We then present experimental results for both character and word recognition for the newly released UNIPEN data [8,9], making these results useful for future comparison to other systems.

The next section gives a brief overview of preprocessing methods. A new method using the Hough transform is used; details are referenced. In Section 3 HMMs for handwriting recognition are described. The Viterbi decoding method is used for efficiency. A four-stage segmental means training method is discussed. Section 4 covers handwriting features of two types: point based or local features and stroke based or high-level, regional features. We do not use any global features. Global properties of the handwriting are addressed only in preprocessing. Delayed strokes and efficiency issues are discussed in Section 5. Experimental results are presented in Section 6 and we conclude in Section 7.

2. Preprocessing

Preprocessing of on-line handwriting can be classified into two types: noise reduction and normalization [10]. Noise reduction attempts to reduce imperfections caused mainly by hardware limits of electronic tablets, through operations such as smoothing, wild point reduction, hook removal, etc. We employ spline filtering for smoothing after the standard wild point reduction and dehooking procedures. A spline kernel is convolved with the input data points which results in new sample points from a local approximate cubic spline fitting of the original data points [5]. Cusps (points of sharp directional change) are detected beforehand using a dynamic detection algorithm [5] and treated as boundary points during smoothing, so that the important information contained in cusps is preserved.

Normalization refers to the reduction of geometric variance in the handwriting data that is typically due to differences in writing style and rendering environment. Examples of normalization include scaling to a standard size, rotation of the text baseline, deskewing of slanted text, etc. The key to high performance writer-dependent recognition is to either perform preprocessing well enough that the writer-dependent variations are removed or to select features that are invariant with respect to writer-dependent variations. Since neither strategy is perfect by itself, some aspects of both are applied in our work. Different normalization procedures are used for words and isolated characters.

2.1. Word normalization

Four *boundary lines* are defined for handwritten words: the base line (joining the bottom of small lower-case letters such as “a”); the core line (joining the top of small lower-case letters); the ascender line (joining the top of

letters with ascenders such as “l”), and the descender line (joining the bottom of letters with descenders such as “g”). For each input word these boundary lines are detected (the base and core lines are assumed to be present in all words, while the ascender and descender lines are not), then the word is rotated such that the base line is horizontal, and scaled such that the *core height* (the distance between the base and core lines) equals a predefined value. Deskewing is then applied to correct the overall slant of the word.

The two main aspects of word normalization, size and orientation normalization, are closely related: a good estimate of core height relies on accurate estimation of the orientation. Both tasks become very difficult for unconstrained handwritten words, where the boundary lines for those words are often not well defined, as shown in Fig. 1. Various techniques have been developed in the past to tackle this difficult problem. The main approaches include the histogram-based methods [10–13], linear-regression-based methods [7], and model-based methods [14].

We have developed a new approach based on the Hough transform [15]. First, local maximum and minimum points along the y coordinate are extracted, then a modified Hough transform is applied to extract parallel lines corresponding to the boundary lines. In order to handle the problems caused by the large variation in natural handwriting coupled with the sparse nature of extremum points in each word sample, one-dimensional Gaussian smoothing with adjustable variance is applied in the Hough space, followed by parameter refinement using linear regression. Compared with previous techniques, the Hough transform based method has the advantage that it simultaneously provides the optimal estimates of both the orientation and the core height of an input word. Fig. 2 illustrates the estimated base line and core line of the word sample shown in Fig. 1.



Fig. 1. A handwritten word whose boundary lines are not well defined.

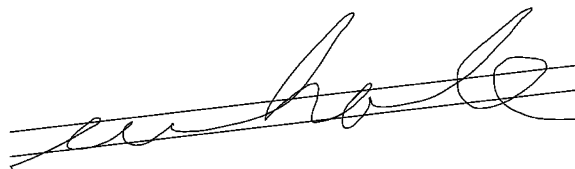


Fig. 2. A handwritten word and the estimated base line and core line.

The skew angle estimation algorithm is similar to that of Brocklehurst [16]. First points of maximum and minimum y coordinates are identified and downward strokes (strokes going from a maximum point to a minimum point and longer than a threshold) are isolated. The straight piece of each downward stroke is then obtained by truncating the stroke from both ends until the angles formed by the two halves of the stroke differ by less than 10° . The estimated skew angle θ is the average of the skew angles of all the straight pieces, each weighted by the length of the piece. After the skew angle estimation, each point in the script is then corrected by replacing x by $x' = x - y \tan \theta$. This deskewing algorithm is simple and fast, and proves to be very effective.

2.2. Character normalization

When the input consists of isolated characters, the above normalization procedures can no longer be applied reliably. In this case we simply normalize the bounding box of the character to a standard size, and no orientation or skew correction is attempted.

2.3. Resampling

At the end of preprocessing, an equa-arc-length resampling procedure is applied to remove variations in writing speed. The sampling distance is proportional to the core height in the case of word recognition, and the height of the bounding box in the case of character recognition. Local features such as tangent slope angle are computed at each sample point. Each handwriting sample is then represented as a time-ordered sequence of observations in the form of feature vectors: $\mathbf{O} = [\mathbf{o}_1, \mathbf{o}_2, \dots, \mathbf{o}_T]$. The details of the features used in our system will be explained in Section 4.

3. The recognition framework

The general recognition framework is composed of Hidden Markov Models (HMMs), representing strokes and characters, embedded in a grammar network representing the vocabulary. The main characteristic of the system is that segmentation and recognition of handwritten words are carried out simultaneously in an integrated process, which provides an ideal mechanism to handle cursive and mixed handwriting.

3.1. Model descriptions

A discrete hidden Markov model with N states and M distinct observation symbols v_1, v_2, \dots, v_M is described by the state transition probability matrix $\mathbf{A} = [a_{ij}]_{N,N}$, state conditional observation probabilities: $b_j(k) = \Pr(\mathbf{o}_t = v_k | s_t = j)$, $k = 1, \dots, M$, where \mathbf{o}_t is the observa-

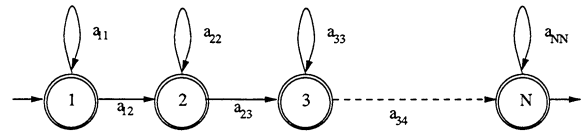


Fig. 3. A left-to-right HMM.

tion and s_t is the active state at time t , and initial state distribution $\pi_i = \Pr(s_1 = i)$.

The subcharacter and character models we have adopted are so called left-to-right HMMs without state skipping ($a_{ij} = 0$ for $j > i + 1$; $\pi_i = 0$ for $i > 1$), as shown in Fig. 3. We have selected this relatively simple topology because it has been shown to be successful in speech recognition, and there has not been sufficient proof that more complex topologies would necessarily lead to better recognition performance. Furthermore, in unconstrained handwriting, skipping of segments seem to happen most often for ligatures, which in our system is handled by treating ligatures as special “connecting” characters and allowing them to be bypassed in the language model.

Given an HMM as described above, the probability of an observation sequence $\mathbf{O} = [\mathbf{o}_1, \mathbf{o}_2, \dots, \mathbf{o}_T]$ along a particular state sequence $\mathbf{q} = [s_1, s_2, \dots, s_T]$ is:

$$P(\mathbf{O}, \mathbf{q} | \lambda) = P(\mathbf{O} | \mathbf{q}, \lambda)P(\mathbf{q} | \lambda) = \prod_{i=1}^T b_{s_i}(\mathbf{o}_i) \prod_{i=1}^{T-1} a_{s_i s_{i+1}},$$

where λ is the set of model parameters. The overall probability of this observation sequence is

$$P(\mathbf{O} | \lambda) = \sum_{\text{all } \mathbf{q}} P(\mathbf{O} | \mathbf{q}, \lambda)P(\mathbf{q} | \lambda).$$

For left-to-right models, this sum is usually dominated by the *optimal path probability*:

$$\mathbf{q}_{\max} = \operatorname{argmax}_{\mathbf{q}} P(\mathbf{O} | \mathbf{q}, \lambda)P(\mathbf{q} | \lambda).$$

The last probability can be efficiently computed using the Viterbi Algorithm. More detailed description of HMMs can be found in Refs. [1,17].

3.2. Model definitions

The basic model units used in our system are subcharacter models called *nebulous stroke models*. A character model is the concatenation of several such stroke models as specified in a *character lexicon*. The advantage of using subcharacter models is that sharing among characters can be easily implemented, simply by referring to the same stroke models in the character lexicon.

A stroke could be any segment of a handwritten script. We do not attempt to impose rigidly defined, presegmented strokes as is the case in Ref. [2]. Instead, we make the system learn indefinite stroke models through training, hence the term *nebulous*. No manual segmentation is

involved in training; the stroke models are trained first on isolated characters and later on whole word samples, but never on isolated or segmented strokes (which would be impossible to obtain since there is no specific definition of a stroke).

Using shared stroke models among characters has two advantages. The first is that it results in a reduced model set. The second one is an extension to the first: since shared models also share training samples, the models can be adequately trained with fewer training samples. On the other hand, excessive sharing can also result in smearing of the model space and thus compromise the recognition performance. Our experiments based on the UNIPEN [8,9] data show that when there is large amount of training data, better results are obtained with smaller amount of sharing – the best performance was obtained with sharing only among different classes of the same character which are composed of the same strokes but rendered in different orders (e.g., the “t” with the cross bar drawn before the vertical hook and the “t” with the cross bar drawn after the vertical hook). If the model size becomes an issue (e.g., in applications run on thin clients such as PDA devices), various clustering techniques can be applied to cluster and merge similar stroke models to achieve more aggressive sharing.

Currently each stroke is modeled by one HMM state. The number of classes for each character and the number of strokes in each class model are selected by hand first and then adjusted empirically. As expected, our experiments show that assigning more states to character classes with more complicated shape (e.g., “g”) than those with simpler shape (e.g., “i”) leads to better recognition results than assigning same number of states to all classes. Currently the number of states per character ranges from 1 to 8. How to automate the process of choosing the number of classes as well as number of states per class is a difficult problem and is a topic of future research.

Ligatures are attached to characters only during training. At recognition time they are treated as special, one stroke “connecting” characters inserted between “core” characters and can be skipped with no penalty. This treatment insures that our system can handle mixed style handwriting as opposed to pure cursive only.

The handling of delayed strokes is explained in detail in Section 5.

3.3. Decoding with grammar constraints

The left-to-right character models as described above are embedded in a grammar network, which can represent a variety of grammatical constraints, e.g., word dictionaries, statistical character N -gram models and context free grammars. In the following we will explain the concepts in the case of a dictionary. Other forms of grammar constraints can be implemented similarly.

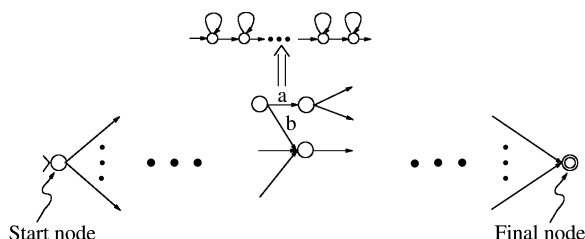


Fig. 4. A grammar network implementing a word dictionary.

Fig. 4 gives the simplified diagram of the grammar network representing a word dictionary. Here each arc represents a character (in the actually network a group of parallel arcs represent all character classes corresponding to the same character). Each path from the start node to the final node corresponds to a word in the dictionary. At recognition time, each arc is replaced by the corresponding HMM.

Given an observation sequence \mathbf{O} , the task of recognition is to find the word W in the given dictionary which maximizes the posterior probability $P(W|\mathbf{O}) = P(\mathbf{O}|W)P(W)/P(\mathbf{O})$. Since $P(\mathbf{O})$ is the same for all words, and assume that all words in the dictionary occurs with equal probability, the problem is reduced to maximizing: $P(\mathbf{O}|W) = \sum_{q \in Q_W} P(\mathbf{O}|q)P(q)$, where Q_W is the set of state sequences that correspond to the same word W . As mentioned before, when left-to-right models are used $P(\mathbf{O}|W)$ is often dominated by the probability along the optimal path, thus the problem can be further reduced to maximizing the term: $P(\mathbf{O}|W, \mathbf{q}_{\max}) = P(\mathbf{O}|\mathbf{q}_{\max} \in Q_W)P(\mathbf{q}_{\max} \in Q_W)$. The solution to the last problem can be found efficiently via Viterbi search, as explained in more detail later.

The embedding of character HMMs with grammar constraints are implemented using an *automatic evolutionary grammar interpretation system* (AEGIS). The three primary components of the evolutionary grammar (EG) are a grammar arc table, a null arc (unlabeled grammar arc) table, and a recursive transition network (RTN) table which records the grammar rules of expanding the non-terminals [18]. Labeled grammar arcs can have either a non-terminal label representing another RTN subnetwork or an HMM label representing an HMM defined in another table called the lexicon, which describes the structure of HMM for each character class. The EG initially contains a degenerate grammar consisting of only a start node, an end node and a single arc with a non-terminal label. During the evolution process, upon encountering a non-terminal arc, the arc is first replaced with the subnetwork it represents and is examined again. This process continues until all non-terminal references on the earliest arcs are eliminated. If a resulting label references an HMM, the appropriate model structure is built as indicated by the lexicon. Once all leading HMM references are built, HMM score integration proceeds.

As a score emerges from an HMM and needs to be propagated further in the network, additional evolution of the EG may occur. In this way, only those regions of the grammar touched by the HMM search are expanded. Beam search methods [19] can be used to limit the amount of grammar expansion.

For any given dictionary a grammar compiler [20] is used to convert a list of word specifications into an optimized network with shared prefixes and suffixes. In the case of isolated character recognition, the grammar network is simply composed of parallel arcs leading from the start node to the final node, each representing one character in the alphabet.

The Viterbi algorithm is used to search for the most likely state sequence corresponding to the given observation sequence and to give the accumulated likelihood score along this best path [1]. Suppose that for any state i , $q_i(t)$ denotes the selected state sequence (hypothesis) leading to i at sample point t , and $\delta_i(t)$ denotes the accumulated log-likelihood score of that hypothesis. O_t represents the observation at sample point t , and $\Delta_i(O_t)$ represents the log-likelihood score of O_t in state i . Since each character model is a left-to-right HMM with no state skipping, updating the hypothesis and likelihood scores within each character model is straightforward. Score propagation through grammar nodes is a little more complicated. Suppose g is a grammar node and $p(g)$ and $s(g)$ denote the sets of preceding and succeeding character classes corresponding to the incoming and outgoing arcs respectively. For each character class l , $m(l)$ denotes the HMM used to model the class; $h(l)$ denotes the initial state of the model; and $f(l)$ denotes the final state of the model. At each sample point t during the Viterbi search, the maximum of all the accumulated scores at the final states of the preceding character models, also called *incoming scores*, is found and propagated to the initial state of each of the succeeding models, along with the corresponding state sequence. The operation is carried out as follows:

$$k = \operatorname{argmax}_{l \in p(g)} \delta_{f(l)}(t-1); \quad (1)$$

and for each state $j = h(l)$; $l \in s(g)$:

$$q_j(t) = \begin{cases} q_{f(k)}(t-1), j & \text{if } \delta_{f(k)}(t-1) > \delta_j(t-1), \\ q_j(t-1), j & \text{Otherwise,} \end{cases} \quad (2)$$

$$\delta_j(t) = \begin{cases} \delta_{f(k)}(t-1) + \Delta_j(O_t) & \text{if } \delta_{f(k)}(t-1) > \delta_j(t-1), \\ \delta_j(t-1) + \Delta_j(O_t) & \text{Otherwise.} \end{cases} \quad (3)$$

3.4. Model training

Models are trained using the well known iterative segmental training method based on Viterbi decoding [1]. Given a set of training samples, the HMM for each sample is instantiated by concatenating the HMMs for

the appropriate character classes, ligatures and delayed strokes, which are in turn instantiated by concatenating the composing stroke models specified in the character lexicon. The training procedure then is carried out through iterations of segmentation of training samples by Viterbi algorithm using the current model parameters, followed by parameter re-estimation using the means along the path. The iterative procedure stops when the difference between the likelihood scores of the current iteration and those of the previous one is smaller than a present threshold.

Training is conducted first on isolated character samples and then, whole word samples. No manual segmentation is involved throughout the whole process. We define four different stages of training based on the training samples used and the granularity of training labels.

1. *Initial character training* is carried out on isolated character samples including ligatures and delayed strokes. Each sample is labeled by the specific character class, thus the corresponding HMM instantiated consists of a linear sequence of states. The initial model parameters are computed using equal-arc-length segmentations of all samples. This stage essentially serves as a model initializer.
2. *Lattice character training* is also carried out on isolated character samples, however, here each sample is labeled only by the character truth as opposed to the specific class. The HMM instantiated for each sample is no longer a linear model, but rather a lattice containing all classes used to model the same character. Starting from parameters obtained from initial character training, the Viterbi algorithm is used to “pick” the best class for each sample and update model parameters accordingly. Since this stage requires much less supervision than the first, it can easily accommodate large amount of training data and thus capture variations among many different writers.
3. *Linear word training* is carried out on whole word samples. Similar to stage 1, each sample is labeled not only by the corresponding word, but also by the exact character class sequence representing the particular style of that sample, which is then converted to a unique stroke sequence according to the character lexicon. This highly constrained training can be used to provide a reliable starting point for whole word training. However, the labeling process is tedious and error prone.
4. *Lattice word training* is carried out on whole word samples. Similar to stage 2, each sample is labeled by the word truth only. Each word is represented by a lattice, or finite state network, that includes all possible character class sequences that can be used to model the word. The sequence that best matches the sample is chosen by the decoding algorithm and the

resulting segmentation is used for parameter re-estimation. Again like stage 2, this process can be applied to a large number of training samples, and this is where most characteristics of cursive or mixed handwriting are captured.

4. Features for on-line handwriting

After preprocessing and resampling, seven features are computed at each sample point. Two of these features are conventional ones widely used in HMM based on-line handwriting recognizers: tangent slope angle and normalized vertical coordinate. The others are more complex features which are classified into two categories: two invariant features and three high-level features (with local representations). On top of these features computed at each sample point, the recognition system is further augmented with a character level segmental feature designed to capture whole character level shape characteristics. These features are explained in detail in the following.

4.1. Invariant features

The concept of invariant features arises frequently in various machine vision tasks. Depending on the specific task, the geometric transformation ranges from simple rigid plane motion to general affine transformation, to prespective mapping, etc. [21]. In the case of handwriting recognition, the transformation of interest is *similitude* transformation, which is a combination of translation, rotation and scaling. Although the variance caused by similitude transformation can be somewhat reduced through the normalization process in preprocessing as discussed in Section 2, it is well known that normalization is never perfect, especially for unconstrained cursive or mixed handwriting. Therefore, including invariant features can greatly improve the robustness of the system. We introduce two features, *normalized curvature* and *ratio of tangents*, which are invariant under arbitrary similitude transformation. Similar features have appeared perviously in the literature for planar shape recognition under partial occlusion [21].

A similitude transformation of the Euclidean plane $\mathbb{R}^2 \rightarrow \mathbb{R}^2$ is defined by $\mathbf{w} = c\mathbf{U}\mathbf{r} + \mathbf{v}$, where c is a positive scalar,

$$\mathbf{U} = \begin{bmatrix} \cos \omega & -\sin \omega \\ \sin \omega & \cos \omega \end{bmatrix}, \quad \mathbf{v} = [v_x, v_y]^T,$$

representing a transformation that includes scaling by c , rotation by angle ω and translation by \mathbf{v} . We regard two curves as equivalent if they can be obtained from each other through a similitude transformation. Invariant features are features that have the same value at corresponding points on different equivalent curves.

Suppose that a smooth planar curve $\mathbf{P}(t) = (x(t), y(t))$ is mapped into $\tilde{\mathbf{P}}(\tilde{t}) = (\tilde{x}(\tilde{t}), \tilde{y}(\tilde{t}))$ by a reparameterization $t(\tilde{t})$ and a similitude transformation, i.e.

$$\tilde{\mathbf{P}}(\tilde{t}) = c\mathbf{U}\mathbf{P}(t(\tilde{t})) + \mathbf{v}. \quad (4)$$

Without loss of generality, assume that both curves are parameterized by arc length (natural parameter), i.e. $t = s$ and $\tilde{t} = \tilde{s}$. Obviously, $d\tilde{s} = c ds$. It can be shown [21] that curvature (the reciprocal of radius) at the corresponding points of the two curves is scaled by $1/c$, i.e. $\tilde{\kappa}(\tilde{s}) = \kappa((\tilde{s} - \tilde{s}_0)/c)/c$. It follows that

$$\frac{\tilde{\kappa}'(\tilde{s})}{(\tilde{\kappa}(\tilde{s}))^2} = \frac{\kappa'((\tilde{s} - \tilde{s}_0)/c)}{(\kappa((\tilde{s} - \tilde{s}_0)/c))^2}, \quad (5)$$

where $\tilde{\kappa}' = d\tilde{\kappa}/d\tilde{s}$ and $\kappa' = d\kappa/ds$, thus eliminating the scale factor from the value of the ratio. Eq. (5) defines an invariant feature which we call *normalized curvature*.

The computation of the normalized curvature defined above involves derivative estimation of up to the third order. Another set of invariants that require lower orders of derivatives can be obtained by using the invariance of distance ratios between corresponding points. Consider again the two equivalent curves $\mathbf{P}(t)$ and $\tilde{\mathbf{P}}(\tilde{t})$ defined above. Suppose P_1 and P_2 are two points on $\mathbf{P}(t)$ whose tangent slope angles differ by θ ; \tilde{P}_1 and \tilde{P}_2 are two points on $\tilde{\mathbf{P}}(\tilde{t})$ with the same tangent slope angle difference. P and \tilde{P} are the intersections of the two tangents on $\mathbf{P}(t)$ and $\tilde{\mathbf{P}}(\tilde{t})$, respectively (Fig. 5). Since angles and hence turns of the curve are invariant under the similitude transformation, it can be shown that if point \tilde{P}_1 corresponds to point P_1 , then points \tilde{P}_2 and \tilde{P} correspond to points P_2 and P respectively [21]. It follows from Eq. (4) that

$$|\tilde{P}\tilde{P}_2|/|\tilde{P}_1\tilde{P}| = |PP_2|/|P_1P|. \quad (6)$$

Eq. (6) defines another invariant feature which we call *ratio of tangents*. In order to use ratio of tangents as an invariant feature in handwriting recognition, a fixed angle difference $\theta = \theta_0$ has to be used for all sample points in all scripts. Since in real applications we normally have only scattered sample points instead of continuous script and in general, we cannot find two sample points whose slope angle difference is equal to θ_0 , a

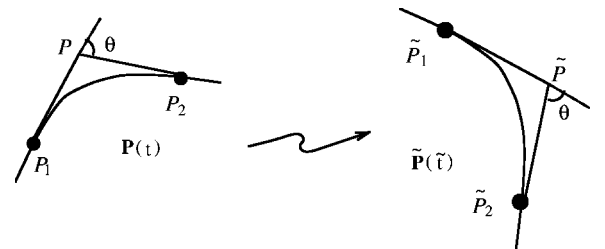


Fig. 5. Ratio of tangents.

carefully designed interpolation algorithm is used to compute this feature [22]. Obviously the choice of θ_0 greatly affects the tangent ratio values. If θ_0 is too small, the feature tends to be too sensitive to noise. On the other hand, if θ_0 is too large, the feature becomes too global, missing important local shape characteristics. This value is chosen heuristically and $\theta_0 = 10^\circ$ is used in the current system.

To evaluate accurately the invariant features described above, high-quality derivative estimates up to the third order have to be computed from the sample points. We used spline filters of up to the 5th degree for this purpose and obtained satisfying results [22].

4.2. High-level features

The trade-off between high-level, long-range features and low level, local features is common among many pattern recognition problems: the former are usually more powerful but less robust, while the latter is less informative but more reliable. In on-line handwriting recognition, some attempts have been made on approaches based on high-level features such as loops, crossings and cusps [2,23], where each handwritten word is converted into a sequence of primitives based on extraction of high-level features. These methods do not work well on unconstrained handwriting, because the results of high-level feature extraction tend to be highly erroneous due to the large shape variations in natural cursive handwriting, especially among different writers.

More recently, many powerful systems have been developed using local features such as slope angle and curvature, sampled at equa-arc-length points along an input script [3,4,6,7,24]. Our system falls into this category. In such approaches, each sample point along the script is uniformly represented by a local feature vector. The sequence of time ordered feature vectors is then processed through a network of statistical models (HMMs [3,4,5] or TDNN [6,7]) and dynamic programming techniques are applied to find the best alignment between the input sequence and the models, thus providing the recognition result. These approaches, which we call *point oriented* methods, are much more robust because no premature shape classification is attempted before the recognition of the whole word; in other words, segmentation and recognition are integrated into a single process. On the other hand, they suffer from the loss of information carried by high-level features. While features formed by temporally nearby points such as cusps are represented implicitly to a certain extent through the slope angle feature, those formed by points that are spatially or temporally far apart, such as crossings and loops, are not presented to the recognition system at all. It is clear that to further improve such systems, high-level long-range features need to be incorporated.

Several approaches have been proposed before to incorporate certain long-range features into a point oriented system. Manke et al. [25] proposed using local bitmap images called context bitmaps to model temporally long-range but spatially short-range phenomena such as crossings. However, this method cannot model features such as loops, which are long-range both spatially and temporally.

We have developed a new method for combining high-level long-range features and local features. First, high-level features such as crossings, loops and cusps are extracted. Then, a *localization* procedure is applied to *spread* these high-level features over the neighboring sample points, resulting in local representations of nearby high-level features. These features are then combined with the usual local features at each sample point.

We extract three high-level features commonly used in handwriting recognition: cusps, crossings and loops. Examples of these features are shown in Fig. 6. Details of the extraction algorithm can be found in Ref. [26].

The difficult problem in using these features in an HMM-based system is how to represent them in the feature vector assigned to each sample point. The loop feature is relatively straightforward – we can simply assign a binary value to each sample point indicating whether or not the point is on a loop. However, the cusp and crossing features are not so trivial. If we apply the same scheme and define a feature with value 1 at a cusp



cusps



crossings



loops

Fig. 6. Examples of high-level features.

(crossing) and 0 at all other points, the corresponding feature sequence will be composed of consecutive 0's with occasional occurrences of single 1's. These occasional 1's will most likely not make any difference in recognition because they are statistically insignificant.

Our solution to this problem is to “spread” the influence of the crossing or cusp feature to its neighboring points. Instead of using the feature itself, we use a derived feature which measures the distance along the script from each sample point to the nearest cusp or crossing point. The value of this derived feature is 0 at the labeled high-level feature point and positive for all nearby points, where a larger value indicates that the point is further away from the nearest labeled sample point. If at any point there is no labeled high-level feature point within a predefined maximum distance w , then the value w is assigned for that point. The two features derived using this scheme are called *cusp distance* and *crossing distance*, respectively.

Although high-level features are not robust features themselves, our experiments show that by combining their local representations with the “true” local features, more information is represented by the feature vectors while the robustness of the system is still preserved, resulting in improved recognition performance [26].

In our discrete HMM system, the seven features computed at each sample point are treated as being independent from each other. During training, a separate probability distribution is estimated for each feature at each model state. At recognition time, a combined log-likelihood score of a particular feature vector being generated by a certain state is calculated as the weighted sum of the individual feature log-likelihood scores.

4.3. Segmental features

The purpose of introducing character level segmental features, which are features computed over the whole segment representing a character, is to capture shape characteristics of the whole character. This is desirable because, after all, characters are the basic shape units used to construct a word. However it is also difficult to implement in a system based on integrated segmentation and recognition. As we stressed before, in such a system no segmentation is carried out before recognition. Therefore, segmental features cannot be computed before recognition begins, since character boundaries are unknown.

Our solution to this problem is a method called *interleaved segmental matching*, where a modified Viterbi algorithm is applied to generate partial segmentation hypotheses based on the local features (including local representations of high-level features), which are then augmented with segmental features computed on the hypothesized segments. The resulted system is called an *augmented HMM system*.

In Section 3.3 we explained how hypotheses are propagated through a grammar node g in a commonly used implementation (Eqs. (1)–(3)). In order to incorporate character shape information into the search, we augment the incoming scores with character matching scores, computed using global character shape models. To be more specific, let $\alpha_l(t_1, t_2)$ be the likelihood score of the segment from sample point t_1 to t_2 on the observation sequence being matched to character class l , the augmented incoming scores are defined as $\tilde{\delta}_{f(l)}(t-1) = \delta_{f(l)}(t-1) + \alpha_l(t_l, t-1)$, where $t_l = t - d_{m(l)}(t-1)$ and $d_{m(l)}(t-1)$ is the number of sample points assigned to character model $m(l)$ up to sample point $t-1$ (character duration). Using these augmented scores, Eqs. (1)–(3) are modified simply by replacing δ with $\tilde{\delta}$.

It should be pointed out that in such an augmented HMM system, the state sequence resulting from Viterbi search is no longer guaranteed to be the optimal sequence, because now the accumulated score of a sequence leading to state i at sample point t not only depends on the previous state, but also on *how* the previous state was reached (history reflected in the character duration). This dependence violates the basic condition for the Viterbi algorithm to yield optimal solution. However, our experiments showed that the gain obtained by incorporating segmental features by far outweighs the loss in the optimality of the algorithm [27].

The segmental matching score $\alpha(t_1, t_2)$ is computed using a correlation based metric inspired by the metrics developed by Sinden and Wilfong [28]. Given a segment a with coordinate sequence $\eta_a = \langle (x_1, y_1), (x_2, y_2), \dots, (x_n, y_n) \rangle$, the *instance vector* of a is defined as: $\mathbf{V}_a = (\bar{x}_1, \bar{y}_1, \bar{x}_2, \bar{y}_2, \dots, \bar{x}_n, \bar{y}_n)$, where $\bar{x}_i = x_i - x_a$ and $\bar{y}_i = y_i - y_a$ for $1 \leq i \leq n$, and (x_a, y_a) is the centroid of a . The *normalized instance vector* of a , $\mathbf{u}_a = \mathbf{V}_a / |\mathbf{V}_a|$, is a translation and scale independent representation of segment a in R^{2n} . Through a resampling procedure, a sample segment of arbitrary length can be mapped to a vector in R^{2N} where N is a predetermined number. The difference between any two sample segments a and b is defined as: $D(a, b) = 0.5(1 - \mathbf{u}_a \cdot \mathbf{u}_b)$, whose value ranges from 0 (when a and b are identical) to 1. The segmental matching score $\alpha(t_1, t_2)$ is then defined by $\alpha_l(t_1, t_2) = -w_\alpha D(a_l, a_{l,t_1,t_2})$, where a_l is the *model segment* for character class l , a_{l,t_1,t_2} is the segment from sample point t_1 to t_2 on the input sample sequence and w_α is a weight factor.

In order to compute the above segment matching score, a single model segment needs to be derived for each character class. Let $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_M$ be the normalized instance vectors of a set of prototypes for the character class l (which can be easily obtained as side products of segmental training). A single model segment representing this class is represented by vector \mathbf{w} which minimizes the sum of distances from individual prototypes. It can be easily shown that $\mathbf{w} = \tilde{\mathbf{u}} / |\tilde{\mathbf{u}}|$, where $\tilde{\mathbf{u}} = \sum_{i=1}^M \mathbf{u}_i$.

An alternative method to incorporate character level shape information is character matching in post-processing. In this method, an N -best decoding algorithm is applied with the local features to obtain the N most likely candidates. This stage is followed by a post-processing step, where character level matching is applied on character segments of the top candidates. The character matching score is then combined with the likelihood score to obtain a weighted sum score for each candidate and the candidate with the highest weighted sum score is chosen as the final result. Our experiments showed that this approach gave inferior recognition performance than the interleaved segmental matching method [29]. The reason is that in the interleaved segmental method local and segmental shape features are considered simultaneously during decoding to yield a solution that is favorable in both aspects. On the other hand, in the post-processing approach local features alone are considered during decoding and as a result the segmentation hypotheses are dictated by local features. During post-processing, only one particular segmentation for each candidate word is considered for character matching, although a different segmentation could result in a better character matching score. The premature limitation of hypothesis space reduces the chance of the correct word being chosen.

5. A two-stage approach for large vocabulary word recognition

When recognition has to be carried out on a large vocabulary, the size of the grammar network increases as well as the recognition time. This poses a serious problem especially when explicit delayed stroke modeling is used, as is the case in our system.

Delayed strokes refer to strokes such as the cross in “t” or “x” and the dot in “i” or “j”, which are sometimes drawn last in a handwritten word, separated in time sequence from the main body of the character. In most on-line handwriting recognition systems delayed strokes are first detected in preprocessing and then either discarded or used in postprocessing. There are two drawbacks to these approaches. First, the information contained in delayed strokes is wasted or inadequately utilized because stroke segmentation cannot be influenced by this data. Second, it is often difficult to detect delayed strokes reliably during preprocessing.

We have developed an approach where delayed strokes are treated as special characters in the alphabet. A word with delayed strokes is given alternative spellings to accommodate different sequences with delayed strokes drawn in different orders. During recognition delayed strokes are considered as inherent parts of a script just like normal characters and contribute directly to the scoring of the hypotheses. This is the technique which we

refer to as *exact delayed stroke modeling*. Good recognition results have been obtained using this approach in a small vocabulary task [5].

However exact delayed stroke modeling method has the disadvantage that it could dramatically increase the hypothesis space. Since each possible position of delayed stroke causes a new branch in the grammar network, potentially, the total number of paths (hypotheses) increases exponentially with the number of delayed strokes. Although the actual increase of the search space can be controlled to a certain extent through a beam search mechanism, the grammar network could still become unmanageable for a large vocabulary task. Furthermore, exact delayed stroke modeling can be achieved only with a full grammar representing a limited vocabulary. Since statistical grammars such as N -gram grammars used to represent unlimited vocabularies inevitably contain loops, exact delayed stroke modeling becomes infeasible when these grammars are applied.

We use a two-stage approach to solve this problem. First an N -best decoding algorithm is applied using *simplified delayed stroke modeling* to narrow down the choices of words. Then, in the detail matching stage, a grammar network is constructed covering the top N candidates using exact delayed stroke modeling, and a best-path search is applied to this reduced network to find the final optimal candidate. In simplified delayed stroke modeling, a character with a delayed stroke is given two groups of parallel sub-paths, one group containing the patterns corresponding to the complete character written in sequence, with the delayed stroke immediately preceding or following the main body of the character; the other group containing only patterns corresponding to the main body of the character, representing the scenario when the delayed stroke is not rendered until the end of the word. A shared delayed stroke module is attached at the end of the grammar network which allows zero, one or multiple occurrences of all types of delayed strokes. To demonstrate the idea, partial diagrams of two different grammar networks, both containing the word “ate”, are given in Fig. 7. Fig. 7a shows exact delayed stroke modeling, and Fig. 7b shows simplified delayed stroke modeling. Labels “dot”, “crs” and “sls” refer to the dot in “i” or “j”, the cross in “t” and the cross in “x”, respectively. Each labeled arc actually represents a group of parallel arcs corresponding to different models for different ways of drawing the symbol, however for the sake of clarity only one representative is shown in the examples. As can be seen, simplified delayed stroke modeling is inexact, e.g., not every path in the network corresponds to a “valid” rendering of a word; however, it is much more efficient and avoids the problem of exponential growth of hypothesis space.

Since the interleaved segmental method requires more computation than conventional Viterbi decoding, for

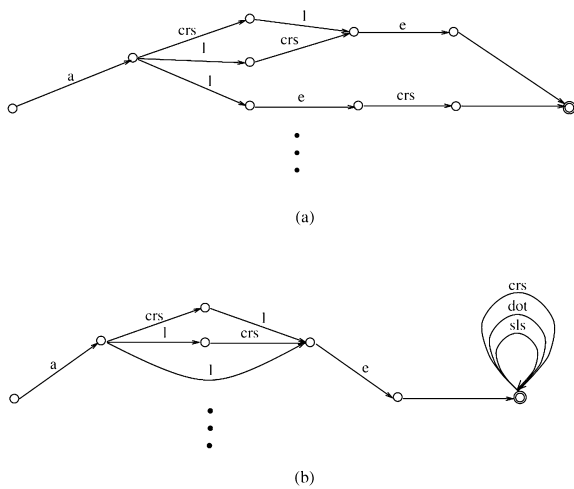


Fig. 7. Example grammar networks.

a large vocabulary task we propose to apply it only at the detailed matching stage as well.

N-best decoding has become a standard technique in speech recognition and, more recently, handwriting recognition. Many algorithms have been developed over the years. Most of these algorithms can be classified into three categories – those that produce the *N*-best candidates in one forward, frame-synchronous, Viterbi-like search [30–32], those that are based on *A**, or stack decoding [33,34], and the tree-trellis based approach [35]. We have developed a new algorithm based on the tree-trellis approach [35]. The difference is that instead of always retrieving the top *N* best candidates, where *N* is a predefined value, our algorithm retrieves a variable number of candidates depending on the distribution of the top scores, thus achieving higher efficiency [29].

6. Experiments on UNIPEN data

Experiments have been carried out using the UNIPEN data. UNIPEN is an international project for on-line handwriting data exchange and recognizer benchmarks which was initiated at the 12th ICPR in 1994 [8]. Several standard training and test sets of the UNIPEN data have just been made available, but we have not yet seen any report on experimental results using those datasets. The UNIPEN data is organized into 11 categories. We have used 4 of them in our experiments: 1a (isolated digits), 1b (isolated upper-case letters), 1c (isolated lower-case letters) and 6 (isolated cursive or mixed-style words in mixed cases).

6.1. Character recognition

Character recognition experiments were carried out using the 1a–c categories of UNIPEN training data

Table 1
Character recognition results

	Training samples	Test samples	Recognition rate (%)	
			Top 1	Top 5
1a (digits)	10,098	2684	96.8	100.0
1b (upper-case)	13,610	7351	93.6	99.4
1c (lower-case)	31,956	12,498	85.9	98.4

release train_r01_v06 [9]. Some cleaning was carried out to remove samples that were mislabeled (which constitute about 4% of the total data). The cleaned data sets were then divided into 2 parts: around 2/3 of the samples in each category were used for training and the remaining 1/3 used for testing. The partitions were made such that there are no shared writers between training and test sets.

For character recognition, models were trained using training stage 1 (initial character training) followed by stage 2 (lattice character training), as described in Section 3.4. Initial character training was carried out using 10 samples per character class, and lattice character training was then carried out on all character training samples. The results are tabulated in Table 1. The whole recognition process (including preprocessing) takes about 0.3 second per character on a 180 MHz SGI station.

6.2. Word recognition

Word recognition experiments were carried out using category 6 of the latest UNIPEN training and development test data releases train_r01_v07 and devtest_r01_v02. No cleaning was carried out for either training or testing.

Models were trained using training stage 1 (initial character training), 2 (lattice character training) and 4 (lattice word training), as described in Section 3.4. Stage 3 (linear word training) was skipped at this point because we did not have the time and resources for the massive manual labeling of word samples. Stages 1 and 2 were carried out in the same way as in character recognition experiments using the same datasets. The models were then further trained using around ~ 45,000 word samples from ~ 320 writers randomly drawn from train_r01_v07.

Three different word test sets were composed, of vocabulary sizes 500, 1000 and 2000 respectively. All test sets contain only samples from writers not used in training. In order to examine the effect of vocabulary and sample set sizes on the recognition as well as reflecting the sample distribution of the complete devtest_r01_v02 set, the following constraints were imposed in composing the test sets. First, the test sets contain samples from the

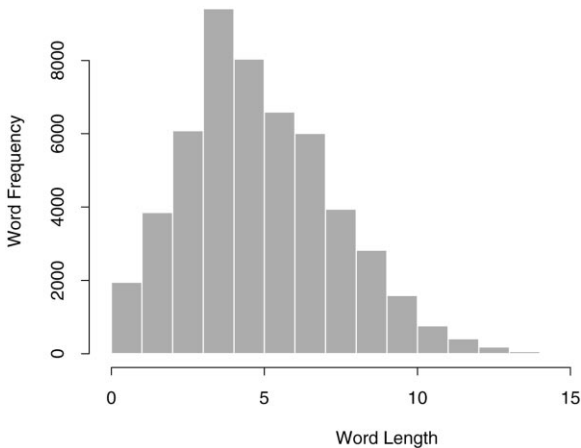


Fig. 8. Word length histogram for UNIPEN data set devtest_r01_v02.

Table 2
Distribution of UNIPEN data set devtest_r01_v02 in terms of cases

Exclusively lower-case	Mixed case	Exclusively upper-case
82.5%	12.5%	5.0%

same set of writers. Second, each vocabulary is a strict subset of the vocabulary of the next larger size. Third, the samples in each set were chosen such that the word length (in terms of number of characters) histogram of each set bears similar profile to that of the word length histogram of the complete devtest_r01_v02 set. Fig. 8 shows the word length histogram of the complete devtest_r01_v02 set. Fourth, the percentages of word samples written in lower-case only, mixed-case and upper-case only for each test set are kept close to those of the complete devtest_r01_v02 set, which are given in Table 2. Finally, the size of each test set was made to be larger than 1500 samples to ensure that the recognition results are statistically meaningful. In the end, samples from 100 writers (out of the 120 writers in devtest_r01_v02 that are not in train_r01_v07) were used to compose the three test sets satisfying these constraints.

Recognition was carried out in 2 stages: N -best decoding with simplified delayed stroke modeling and no segmental matching, followed by re-scoring using exact delayed stroke modeling and segmental matching. The results are summarized in Table 3.

A close examination of the errors reveals that the recognizer does a much better job on exclusively lower-case samples than on samples containing upper-case letters. For example, on the 1000 vocabulary test set, the top 1 recognition rate on exclusively lower-case samples is 91.5%, while that on samples containing upper-case

Table 3
Word recognition results on small to medium-size vocabularies

Vocabulary size	Test samples	Recognition rates (%)	
		Top 1	Top 10
500	1685	91.8	98.4
1000	2447	90.5	97.6
2000	3267	87.2	96.5

letters is only 84.1%. We believe there are two reasons for this discrepancy. First, lower-case letter models are better trained since the majority (75%) of training word samples are in lower-case only, and another 18% of the samples contain only one upper case letter at the beginning. Although there are enough upper-case samples for character training (stages 1 and 2), there are not adequate instances of upper-case letters in lattice training (Stage 4), thus the upper-case letter models are not well trained in the context of words. One possible way to alleviate this problem is to artificially increase the instances of upper-case letters in the training set. Second, the current size and orientation normalization algorithm tends to perform poorly on words containing upper case letters. This is because one basic assumption in the boundary line detection algorithm is that the base and core lines are present and dominant in all words, while the ascender and descender lines may or may not be present (see Section 2.1). This assumption is often violated in words containing upper-case letters, especially those written in upper-case only (in which case the core line is not present at all and instead the ascender line is always present). How to solve this problem and improve the robustness of the size and orientation normalization algorithm remains a topic of future research.

6.3. Large vocabulary experiments

Preliminary experiments have been carried out on large vocabularies of up to 20,000 words. Speed is a major concern with large vocabulary recognition when full dictionaries are used. We tested two different methods to reduce the computational load. Since vocabulary size mainly affects the first stage N -best decoding, both methods modify this stage only. The rescoring stage remains the same as described in the previous section. The idea behind both method is to apply somewhat “coarser” processing for N -best decoding to reduce the computational load without degrading the performance too much. This seems feasible because for the first stage processing we only need to make sure that the correct answer is among the top N , as opposed to being precisely at the top.

In the first method, “coarse” HMM models are defined where the number of states used to model each character

class is reduced to no more than 2. These models are trained in the same way as the original models and then used in the first stage *N*-best decoding procedure. The rational is that by reducing the number of states per character model the network size is reduced and so is the amount of computation. Unfortunately our experiments showed that this is not quite the case. Although the total network size is reduced by roughly half, *N*-best decoding takes even longer to complete while giving worse recognition results. Careful analysis lead to the following explanation. Since beam search mechanism is used in our system to prune the network during decoding, the computational load is determined by the number of *active* states as opposed to the total number of states in the network. When only 2 states are used per model, the models provide rather poor discrimination among different characters. As a result, many more competing paths have very similar accumulated scores even after decoding has proceeded far along, causing more paths and thus more HMMs to remain active until the very end. Thus the computational load, which depends on the number of active states, is increased despite the smaller size of the network.

In the second method, down-sampling is applied in feature extraction to reduce the number of observations to be processed for each sample word. In our experiments the average number of sample points per character in *N*-best decoding was reduced to ~ 9 from ~ 14 , which was the sampling rate used for rescoring. This turned out to be much more effective than the first method – the decoding time was reduced by roughly one third with only a small degradation in performance. This method was tested using three dictionaries, of sizes 5000, 10,000 and 20,000, respectively. For all three taks the test set is the same as the one used for the 1000 word vocabulary test, and each dictionary is composed by adding randomly chosen words from the standard UNIX dictionary to the original 1000 words. Therefore, unlike the previous group of experiments, here not all words in the dictionaries have instances in the test set. However, these tests still provide valid indications of the systems performance for large vocabulary tasks. The recognition results are summarized in Table 4.

To get a rough idea of the system’s performance as a function of vocabulary size, we plot all the recognition results in Fig. 9. Bear in mind that this is only a rough plot since the test set varies for different vocabulary sizes (as specified in Tables 3 and 4).

The plot shows that performance degradation for top one candidate is much more severe than that for the top ten. This indicates that the second stage rescoring still has much room for improvement, possibly through more careful modeling, better features, and more sophisticated training techniques such as discriminative training [36,37]. Speed is still a problem for large vocabulary tasks – it currently takes about two minutes on average

Table 4
Word recognition results on large vocabularies

Vocabulary size	Test samples	Recognition rates (%)	
		Top 1	Top 10
5000	2447	83.2	94.8
10,000	2447	79.8	92.9
20,000	2447	76.3	91.0

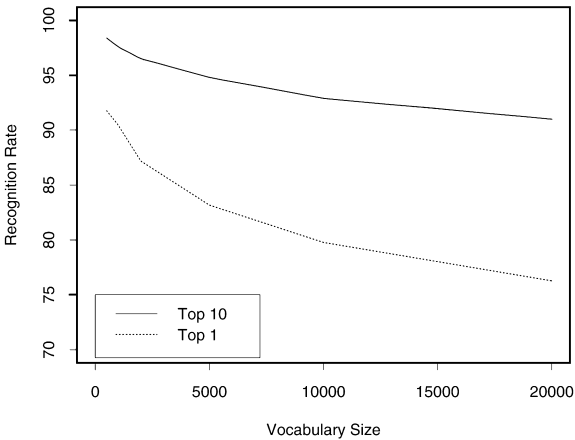


Fig. 9. Word recognition rates at different vocabulary sizes.

to recognize a word with the 20,000 word dictionary on an 180 MHz SGI station. Obviously down-sampling can only be applied to a limited extent. Other techniques, such as smarter pruning, will have to be investigated. Another possibility is to use statistical language models such as variable-length *N*-grams [38,39], which require much smaller networks than complete large dictionaries, for *N* best decoding.

7. Conclusions

These are early results on the UNIPEN data. We think the results are good, but have not seen any published UNIPEN results from other researchers yet, so we look forward to seeing comparable reports.

Error rates for the character recognition test results reported in Table 1 are accurate to within less than 1% with 95% confidence. Test 1a has 10 classes and tests 1b and c are 26-class problems. We think the error rate for test 1c is higher because there is more shape confusion among lower-case letters compared to upper-case ones.

The word recognition test sets are of modest size for small vocabularies due to the nature of UNIPEN data. The database typically contains only one or a very small number of examples of each word. Examination of the

test data indicates that it is a challenging test because the writing style and legibility varies widely among the 100 different writers. Furthermore, the test sets were taken directly from devtest_v01_r02 with no cleaning. A quick examination shows that they indeed contain quite a few mislabeled or misspelled samples. These results are comparable to those published for other testing data but of course there is no way to make a direct comparison due to probable differences in the quality and other characteristics of the data.

Future work will be directed primarily toward refinement of the character models and training methods to improve the recognition performance. Improvements in speed are also needed. Toward this end we are considering a major overhaul of the HMM decoder and automatic grammar reconstruction for the N -best two-stage processing method. We would also like to investigate more sophisticated pruning methods as well as the use of statistical language models for the first stage N -best decoding. Our goal will be real-time performance on very large vocabularies.

8. Summary

Hidden Markov Models (HMMs) have been used with great success for stochastic modeling of speech for more than a decade now. More recently they have also been applied to on-line handwriting recognition with varying degrees of success. In this paper we describe an on-line handwriting recognition system which we have developed over the last few years. The main novel aspects of this system compared to other HMM-based systems include: signal preprocessing with a new word normalization method; application of invariant features; combination of high-level and low-level features; stochastic modeling with HMMs that incorporates a dynamically evolving language model; and N -best decoding combined with two-stage delayed stroke modeling for large vocabulary word recognition.

We present experimental results for both character and word recognition for the newly released UNIPEN data. Character recognition experiments were carried out using dataset train_r01_v06. Word recognition experiments were carried out using datasets train_r01_v07 and devtest_r01_v02, with results on dictionary sizes varying from 500 words to 20,000 words. The use of the standard UNIPEN datasets makes these results useful for future comparison to other systems.

References

- [1] L.R. Rabiner, B.H. Juang, Fundamentals of Speech Recognition, Prentice-Hall, Englewood Cliffs, NJ, 1993.
- [2] S. Bercu, G. Lorette, On-line handwritten word recognition: an approach based on hidden markov models, in: Proceedings of the Third International Workshop on Frontiers in Hand-writing Recognition, Buffalo, USA, May 1993, pp. 385–390.
- [3] J. Makhoul, T. Starner, R. Schartz, G. Chou, On-line cursive handwriting recognition using speech recognition methods, in: Proceedings of IEEE ICASSP'94, Adelaide, Australia, April 1994, pp. v125–v128.
- [4] K.S. Nathan, H.S.M. Beigi, J. Subrahmonia, G.J. Clary, H. Maruyama. Real-time on-line unconstrained handwriting recognition using statistical methods, in: Proceedings of IEEE ICASSP'95, Detroit, USA, June 1995, pp. 2619–2622.
- [5] J. Hu, M.K. Brown, W. Turin, Handwriting recognition with hidden Markov models and grammatical constraints, in: Proceedings of the Fourth IWFHR, Taipei, Taiwan, December 1994, pp. 195–205.
- [6] S. Manke, U. Bodenhausen, Npen + + : a writer independent, large vocabulary on-line cursive handwriting recognition system, in: Proceedings of the Third ICDAR, Montreal, Canada, August 1995, pp. 403–408.
- [7] M. Schenkel, I. Guyon, D. Henderson, On-line cursive script recognition using time delay neural networks and hidden Markov models, Machine Vision and Applications, Special Issue on Cursive Script Recognition, vol. 8, 1995.
- [8] I. Guyon, L. Schomaker, R. Plamondon, M. Liberman, S. Janet, Unipen project of on-line data exchange and recognizer benchmarks, in: Proceedings of the 12th ICPR, Jerusalem, Israel, October 1994, pp. 29–33.
- [9] L. Schomaker, The unipen project, in <http://unipen.nici.kun.nl> 1998.
- [10] W. Guerfali, R. Plamondon, Normalizing and restoring on-line handwriting, Pattern Recognition 26 (3) (1993) 419–431.
- [11] D.J. Burr, A normalizing transform for cursive script recognition, in: Proceedings of the Sixth ICPR, vol. 2, Munich, October 1982, pp. 1027–1030.
- [12] M.K. Brown, S. Ganapathy, Preprocessing techniques for cursive script recognition, Pattern Recognition 16 (5) (1983) 447–458.
- [13] H.S.M. Beigi, K. Nathan, G.J. Clary, J. Subrahmonia, Size normalization in on-line unconstrained handwriting recognition, in: Proceedings of ICASSP'94, Adelaide, Australia, April 1994, pp. 169–172.
- [14] Y. Bengio, Y. LeCun, Word normalization for on-line handwritten word recognition, in: Proceedings of 12th ICPR, vol. 2, Jerusalem, October 1994, pp. 409–413.
- [15] Amy S. Rosenthal, J. Hu, M.K. Brown, Size and orientation normalization of on-line handwriting using Hough transform, in: Proceedings of ICASSP'97, Munich, Germany, April 1997.
- [16] E.R. Brocklehurst, P.D. Kenward, Preprocessing for cursive script recognition, NPL Report DITC 132/88, November 1988.
- [17] L.R. Rabiner, A tutorial on hidden markov models and selected applications in speech recognition, Proc. IEEE 77 (2) (1989) 257–286.
- [18] M.K. Brown, S.C. Glinski, Stochastic context-free language modeling with evolutionary grammars, ICSLP'94 vol. 2, September 1994, 779–782.

- [19] B.T. Lowerre, D.R. Reddy, The HARPY speech understanding system, in: W.A. Lea (Ed.), Trends in Speech Recognition, Prentice-Hall, Englewood Cliffs, NJ, 1980, pp. 340–360 Chapter 15.
- [20] M.K. Brown, J.G. Wilpon, A grammar compiler for connected speech recognition, IEEE Trans. Signal Process. 39 (1) (1991) 17–28.
- [21] A.M. Bruckstein, R.J. Holt, A.N. Netravali, T.J. Richardson, Invariant signatures for planar shape recognition under partial occlusion, CVGIP: Image Understanding 58 (1993) 49–65.
- [22] J. Hu, M.K. Brown, W. Turin, Invariant features for HMM based handwriting recognition, in: Proceedings of ICIAP'95, Sanremo, Italy, September 1995, pp. 588–593.
- [23] S.A. Guberman, V.V. Rozentsveig, Algorithm for the recognition of handwritten text, Automat. Remote Control 37 (5) (1976) 751–757.
- [24] J. Hu, M.K. Brown, W. Turin, HMM based on-line handwriting recognition, IEEE PAMI 18 (10) (1996) 1039–1045.
- [25] S. Manke, M. Finke, A. Waibel, Combining bitmaps with dynamic writing information for on-line handwriting recognition, in: Proceedings of the 12th ICPR, Jerusalem, October 1994, pp. 596–598.
- [26] J. Hu, A.S. Rosenthal, M.K. Brown, Combining high-level features with sequential local features for on-line handwriting recognition, in: Proceedings of ICIAP'97, Florence, Italy, September 1997, pp. 647–654.
- [27] J. Hu, M.K. Brown, W. Turin, Use of segmental features in HMM based handwriting recognition, in: Proceedings of IEEE SMC'95, Vancouver, Canada, October 1995, pp. 2778–2782.
- [28] F. Sinden, G. Wilfong, L. Ruedisueli, On-line recognition of handwritten symbols, IEEE Trans. PAMI 18 (9) (1996).
- [29] J. Hu, M.K. Brown, On-line handwriting recognition with constrained n -best decoding, in: Proceedings of 13th ICPR, vol. C, Vienna, Austria, August 1996, pp. 23–27.
- [30] C.-H. Lee, L.R. Rabiner, A frame-synchronous network search algorithm for connected word recognition, IEEE Trans. ASS 37 (11) (1989) 1649–1658.
- [31] V. Steinbiss, Sentence-hypotheses generation in a continuous-speech recognition system, in: Proceedings of EURO-SPEECH'89 1989, pp. 51–54.
- [32] R. Schwartz, Y.-L. Chow, The N -best algorithm: an efficient and exact procedure for finding the N most likely sentence hypotheses, in: Proceedings of ICASSP'90, 1990, pp. 81–84.
- [33] F. Jelinek, L.R. Bahl, R.L. Mercer, Design of a linguistic statistical decoder for the recognition of continuous speech, IEEE Trans. Inform. Theory IT-21 (3) (1975) 250–256.
- [34] D. Sturtemant, A stack decoder for continuous speech recognition, in: Proceedings of Speech and Natural Language Workshop, October 1989, pp. 193–198.
- [35] F.K. Soong, E.-F. Huang, A tree-trellis based fast search for finding the N best sentence hypotheses in continuous speech recognition, in: Proceedings of ICASSP'91, Toronto, Canada, May 1991, 705–708.
- [36] W. Chou, B.H. Juang, C.H. Lee, Segmental gpd training of HMM based speech recognizers, in: Proceedings of ICASSP'92, San Francisco, CA, March 1992.
- [37] B.H. Juang, P.C. Chang, Discriminative training of dynamic programming based speech recognizers, IEEE Trans. Speech and Audio Process. SAP-1 (2) (1993) 135–143.
- [38] I. Guyon, F. Pereira, Design of a linguistic postprocessor using variable memory length Markov models, in: Proceedings of the Third ICDAR, Montreal, Canada, August 1995, pp. 454–457.
- [39] J. Hu, W. Turin, M.K. Brown, Language modeling using stochastic automata with variable length contexts, Comput. Speech Language 11 (1) (1997) 1–16.

About the Author—JIANYING HU studied Electrical Engineering at Tsinghua University in Beijing, China, from 1984 to 1988, and received her Ph.D. in Computer Science from SUNY Stony Brook in 1993. Since then she has been a member of technical staff at Bell Laboratories, Murray Hill, formerly of AT&T, now of Lucent Technologies.

Dr. Hu has worked extensively on curve indexing, on-line handwriting recognition, language modeling and document processing. Her current research interests include document understanding, handwriting analysis, information retrieval and content based image retrieval. She is a member of the IEEE computer society.

About the Author—SOK GEK LIM received the B.Sc. in Computer Science and Ph.D. in Electrical Engineering from The University of Western Australia in 1991 and 1997, respectively.

She worked with NEC in Singapore in cooperation with Kent Ridge Digital Lab (Formerly know as the Institute of Systems Science (ISS)) and the Supercomputing Centre at the National University of Singapore, from 1991 to 1992. She was a research officer with the centre for intelligent Information Processing Systems (CIIPS) in the Department of Electrical and Electronic Engineering at The University of Western Australia, from 1992 to 1996. She was awarded a post-doctoral fellowship from the University of West Florida in 1997.

Dr. Lim is currently with Lucent Technologies Bell Laboratories, Murray, NJ. Her research interests are image analysis and pattern recognition, <http://ciips.ee.uwa.edu.au/~gek>.

About the Author—MICHAEL K. BROWN is a Member of Technical Staff at Bell Labs and Senior Member of the IEEE. He has over 50 publications and more than a dozen patents.

From 1973 to 1976 Dr. Brown was with the Burroughs Corporation in Detroit (now Unisys Corp.) developing of ink jet printing systems. This work resulted in the University of Michigan's first Master's Thesis, "Control Systems for Ink Jet Printers" that describes advances in high speed electrostatically controlled jet printing techniques. From 1976 to 1980 he was self-employed as an Engineering Consultant with Burrough's Corp. working on unconstrained handwritten character recognition (Courtesy Amount Reading of bank cheques) while pursuing his Ph.D. degree at the University of Michigan. His published Ph.D. thesis addressing on-line unconstrained "Cursive Script Word Recognition" describes new techniques in feature extraction, pattern recognition and writer adaptation.

In 1980 Dr. Brown joined the Speech Processing Group at Bell Telephone Laboratories, Murray Hill, where he developed speech recognition algorithms, designed VLSI ASIC's and designed speech processing systems. In 1983 he got transferred to the Robotics Principles Research Department where he pursued interests in sensory perception, sensor integration, distributed computing, natural language (speech) understanding, dialogue management, task planning, dynamic simulation and electro-mechanical control. Since 1990, when the department was renamed Interactive Systems Research, Dr. Brown has concentrated more on human-machine and machine-environment interaction, also has resumed his interest in unconstrained on-line handwriting recognition. In 1996 he joined the Multimedia Communications Research Laboratory to continue work on natural language understanding systems and pursue new interests in multimedia pattern analysis.