

# Recognition of Strings Using Nonstationary Markovian Models: An Application in ZIP Code Recognition

Djamel Bouchaffra, Venu Govindaraju, and Sargur N. Srihari

Center of Excellence for Document Analysis and Recognition (CEDAR)  
Department of Computer Science  
State University of New York at Buffalo  
Amherst, New York 14228-2567

## Abstract

*This paper presents Nonstationary Markovian Models and their application to recognition of strings of tokens, such as ZIP Codes in the US mailstream. Unlike traditional approaches where digits are simply recognized in isolation, the novelty of our approach lies in the manner in which recognitions scores along with domain specific knowledge about the frequency distribution of various combination of digits are all integrated into one unified model. The domain knowledge is derived from postal directory files. This data feeds into the models as n-grams statistics that are seamlessly integrated with recognition scores of digit images. We present the recognition accuracy (90%) achieved on a set of 20,000 ZIP Codes.*

## 1 Introduction

Researchers in all fields have shown the value of using contextual information to enhance the recognition of patterns. Speech recognition community was the first to model the contextual knowledge present in the English language [8]. However the signal and language information were used in two distinct but successive phases. It has been shown since that it is possible to integrate these bodies of information in a single phase [2, 3]. Besides, Shinghal and Toussaint [10] address the benefits of using transition probabilities that are dependent on both the word position and length in recognition of machine-printed text.

Let us consider all the ZIP Codes in the United States. If the source of all the ZIP Codes is Markovian, then once again, the identity of any numeric symbol at a given position will be confined to its neighborhood. This conjecture has not been verified either empirically or theoretically. Each numeral in a string

is created in isolation, independent of the neighboring numerals in the string. For that matter, it would not be unusual for a writer to write the numerals in the string in an arbitrary order (write the last 2 digits and then the first 3 to its left).

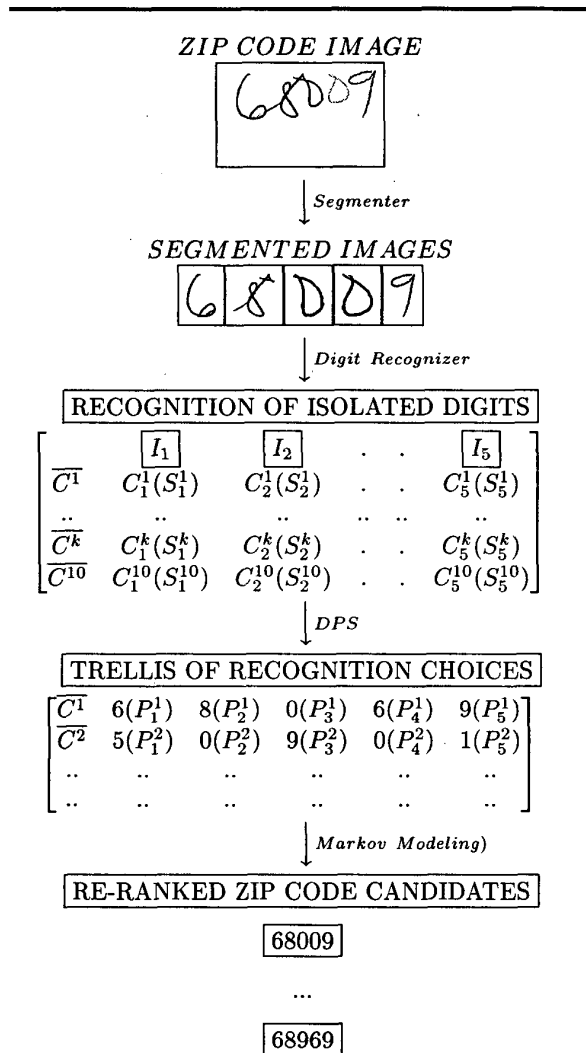
Figure 1 describes the flow-down of the various processes involved. Starting with an image of a ZIP Code the very first step is to isolate 5 digit images by a process of segmentation. Next is the task of recognizing each digit image by a digit recognizer. The recognizer returns a ranked list of classes for each digit along with its confidence associated with each choice. This paper assumes the results of this stage and hence all that precedes this stage is essentially a *black box*. We adopt a procedure called DPS (Deriving Probability from Scores), which we have described in detail elsewhere [2], that simply maps the confidence values to probability values.

This leads us to the subject matter of this paper: *Postprocessing* by Nonstationary Markov Modeling. The input is the trellis of class choices with associated probabilities and the output of the postprocessing is a ranked list of ZIP Code candidates.

## 2 Background

The recognition of handwritten numeric strings is usually easier than alphabetic words in a given address. There are fewer digits than alphabets (10 *versus* 26) and digits touch each other less often. Present ZIP Code recognition procedures routinely perform the task of locating the ZIP Code within the address, followed by segmenting and recognizing the ZIP Code digits.

What makes the task of finding 5 isolated digits non-trivial is the occasional occurrence of adjacent digits touching each other. We will assume in this paper that the digits are isolated by a segmenter [9, 11]



**Figure 1:** Starting with an image of a ZIP Code (numeric string of length  $L = 5$ ), segmentation and digit recognition provide  $k (=1 \dots 10)$  class choices ( $C_t^k$ ) for each isolated digit ( $t \in \{1 \dots 5\}$ ). Each class choice  $C_t^k$  is associated with a score  $S_t^k$ . A procedure called DPS converts the scores into probability values  $P_t^k$ . Postprocessing by Markov Modeling, which is the focus of this paper, takes as input the trellis of class choices with associated probabilities to output a ranked list of ZIP Code candidates.

with a correct segmentation rate of 92.9%.

We use a Gradient Structural Concavity (GSC) [5] recognizer to recognize the isolated digits. The accuracy of the recognizer on isolated digits is 96% on a particular set of images.

Let  $X_i$  represent the measurement made on an input pattern image  $I_i$ , where  $\bar{I} = [I_1, I_2, \dots, I_L]$  is a sequence of subimages (e.g., sequence of subimages of isolated digits). The classifier's  $k^{th}$  class choices are denoted by  $\bar{C}^k = [C_1^k, C_2^k, \dots, C_L^k]$  based on the measurement vector  $\bar{X}^k = [X_1^k, X_2^k, \dots, X_L^k]$  which corresponds to the subimages of  $\bar{I}$ .  $k$  ranges from 1 to 10 (the total number of classes) in the case of recognizing digits in a ZIP Code. When the pattern recognizer is treated as a blackbox, the measurement  $\bar{X}$  is taken as the vector of confidence scores  $\bar{S}$  returned by the recognizer rather than pattern feature vectors of individual classes.

Considering all the class choices, we obtain a trellis  $\mathcal{C} = \{\bar{C}_t^k\} (t \in \{1 \dots L\}, k \in \{1 \dots c\})$  represented in Table 1. The parameter  $k$  is the choice level of the recognizer and  $c$  is the total number of classes.  $C_t^k$  is  $k^{th}$  recognition choice in column  $t$  of the trellis. Each class choice is associated with a score  $S_t^k$  returned by the digit (symbol) recognizer.

We denote the set of all class paths through the trellis by  $\mathcal{P}$ . Let's  $\langle \bar{C}, \bar{X} \rangle$  be a particular class and its measurement path in  $\mathcal{P}$ . Our objective is to maximize  $P(\bar{C} | \bar{I})$ . We will drop  $\bar{X}$  from our notations henceforth, as it is implied by reference to  $\bar{C}$ . The task at hand is simply to maximize the conditional probability  $P(\bar{C} | \bar{S})$ . By Bayes' rule, this would be equivalent to maximizing the quantity in Equation 1 with respect to  $\bar{C}$ .

$$\max_{\bar{C}: \bar{C} \in \mathcal{P}} P(\bar{C} | \bar{S}) = \max_{\bar{C}: \bar{C} \in \mathcal{P}} \left[ \frac{P(\bar{S} | \bar{C}) \times P(\bar{C})}{P(\bar{S})} \right]. \quad (1)$$

Given that  $P(\bar{S})$  is independent of  $\bar{C}$ , the task at hand is one of maximizing  $P(\bar{S} | \bar{C}) \times P(\bar{C})$  over all possible paths  $\bar{C}$  in the trellis. For a vector of length  $L$ , this amounts to considering a total of  $10^L$  possibilities.

Typically, digit recognizers return just the top 2 or 3 recognition choices. Let us say the maximum is 3. This leads to  $3^5$  ZIP Code candidates that can be generated by the  $3 \times 5$  trellis.

Maximizing the expression  $P(\bar{S} | \bar{C}) \times P(\bar{C})$  over  $\bar{C}$  would be equivalent to maximizing any monotonically increasing function of the same expression. Hence, our objective could be re-stated as

$$\max_{\bar{C}: \bar{C} \in \mathcal{P}} [\log P(\bar{S} | \bar{C}) + \log P(\bar{C})]. \quad (2)$$

$$\max_{C_i} \left[ \sum_{t=1}^{t=L} (\log P(S_t | C_t)) + \left( \log P(C_1) + \sum_{t=2}^{t=L} \log P(C_t | C_1, C_2, \dots, C_{t-1}) \right) \right]. \quad (3)$$

	$I_1$	$I_2$	$I_3$	..	$I_L$
$\overline{C^1}$	$\langle C_1^1, S_1^1 \rangle$	$\langle C_2^1, S_2^1 \rangle$	$\langle C_3^1, S_3^1 \rangle$	..	$\langle C_L^1, S_L^1 \rangle$
$\overline{C^2}$	$\langle C_1^2, S_1^2 \rangle$	$\langle C_2^2, S_2^2 \rangle$	$\langle C_3^2, S_3^2 \rangle$	..	$\langle C_L^2, S_L^2 \rangle$
$\overline{C^3}$	$\langle C_1^3, S_1^3 \rangle$	$\langle C_2^3, S_2^3 \rangle$	$\langle C_3^3, S_3^3 \rangle$	..	$\langle C_L^3, S_L^3 \rangle$
..	..	..	..	..	..
$\overline{C^c}$	$\langle C_1^c, S_1^c \rangle$	$\langle C_2^c, S_2^c \rangle$	$\langle C_3^c, S_3^c \rangle$	..	$\langle C_L^c, S_L^c \rangle$

**Table 1:** 2D-Trellis representation of class candidates with their scores.

For sake of brevity, we will denote a path as  $(\overline{C}, \overline{S}) = (C_t, S_t)$  for  $(t = 1 \dots L)$  composed of a class  $C_t$  and its measurement (score) vector  $S_t$ . Using the assumption of class conditional independence of scores and the fact that each score  $S_i$  depends only on class  $C_i$ , the simplification of Equation 2 is presented in Equation 3.

Each of the two terms in the right hand side of Equation 3 bear diverse bodies of knowledge. The first term,  $P(S_t | C_t)$  is the probability of a measurement of a class at a particular node in the path  $\overline{C}$  given the class  $C_t$  at this node. Since each node  $(C_t, S_t)$  in a path of the trellis  $\mathcal{C}$  is a pair  $(C_t^k, S_t^k)$  where  $(k = 1, \dots, c)$  is the class choice and  $(t = 1, \dots, L)$  the position of this class in the path, the maximization problem of Equation 3 involves conditional probabilities  $P(S_t^k | C_t^k)$ . We assume that this conditional probability is independent of time (or position):  $(P(S_t^k | C_t^k) = P(S^k | C^k))$ . This latter quantity is the probability of a measurement returned by a digit recognizer given the class of the image. Computing this probability is non-trivial. Indeed, computing its Bayes derivative,  $P(C^k | S^k)$  is equally difficult. We have described a methodology of computing this probability elsewhere [2], in sufficient detail. For subsequent integration with the postprocessing module, it is crucial that the scores be converted to probability values ( $P_t^k$ ) so that the Bayesian framework can be adopted.

We use the nonstationary property of a Markov chain in the language part.

$$\log P(C_1) + \sum_{t=2}^{t=L} \log P(C_t | C_1, C_2, \dots, C_{t-1})$$

by highlighting the temporal aspect of the class distri-

bution. In fact, the position of a same class (or digit) is a discriminative factor in a ZIP Code ( $P(C_{t_i}^k) \neq P(C_{t_j}^k)$  if  $i \neq j$ ).

These facts allow the simplification of Equation 3 to Equation 4. Equation 4 shows that the language part is essentially a combination of  $n$ -gram probabilities, where  $1 \leq n \leq L$ .

### 3 Objective

Our objective can be simply stated as developing a postprocessing model that produces the maximally accurate recognition of a string of symbols given a *trellis* of symbols and their associated scores by a recognizer. In the case of ZIP Code images, a succinct description of the problem is illustrated by Figure 2.

### 4 Previous Work

[11] describes a “segment and then recognize” approach for the application of postal ZIP Code recognition. While there are applications of numeric string recognition other than ZIP Codes, such as recognizing courtesey amounts on bank checks [7], the ZIP codes have the distinguishing feature of having a fixed length. Much of the positional (temporal) aspect of our model depends on this very fact. Choosing the right model and of the right order has been discussed in the context of bank check recognition [7].

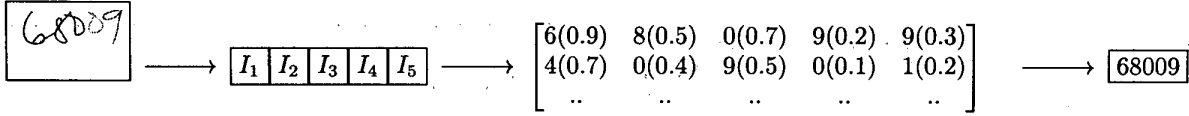
### 5 Modeling Contextual Knowledge

For the application of ZIP Code recognition, the appropriate model that represents the “structure” of the contextual knowledge must be determined. First, we experimentally test all possible models (in terms of degrees of Markovian processes) and all their combinations as an empirical way of determining the most suitable model.

#### 5.1 Empirical Evaluation of Models

If the numeric strings are all of a fixed length, ( $L$  in the discussion above) then the contextual information can be further qualified. By making the probabilities specific to their position within the string, one can better exploit the context. This adaptation is referred to in the literature as a *non-stationary* Markov Model [1]. The context that we are modeling is the probability of a particular class,  $\alpha$ , in a particular position,  $t$  in the numerical string, *given* the occurrence of a particular class (*say*,  $\beta$ ) in the *previous* position,  $t - 1$ :  $P(C_t = \alpha | C_{t-1} = \beta)$ .

$$\max_{C_t^{k_t}} \left[ \sum_{t=1}^{t=L} \log \frac{P(C_t^{k_t} | S^{k_t})}{P(C_t^{k_t})} + \log P(C_1^{k_1}) + \sum_{t=2}^{t=L} \log P(C_t^{k_t} | C_1^{k_1}, C_2^{k_2}, \dots, C_{t-1}^{k_{t-1}}) \right]. \quad (4)$$



**Figure 2:** I/O behavior of the postprocessing module. The input is a *trellis* of classifier choices associated with probability values. The output is the best ZIP Code candidate based on the probabilities of recognition as well as the Non-stationary Markovian modeling as described in this paper.

This notion of class conditional ( $P(C_t = \alpha | C_{t-1} = \beta)$ ) that we describe is similar to that of positional *n-grams* described in the literature [10].

- **Markovian Model of  $0^{\text{th}}$  order:** It represents the case where there is no contextual information shared between neighboring numerals in a string.

In the ZIP Code application, there are  $\binom{5}{1}$  possible position 1-grams. This is represented by 5 column vectors, one for each position ( $t$ ) of the ZIP Code. Each entry in the column vector,  $U_t[i]_{i \in \{0 \dots 9\}}$ , represents the *prior* probability  $P(C_t = \alpha)$ , the frequency of occurrence of digit  $\alpha$  in position  $t_{1 \dots 5}$ .

- **Markovian Model of  $1^{\text{st}}$  order:** Conditional probabilities of the first order are represented by 2-dimensional matrices  $B_{t_1, t_2}[i, j]$  where  $i, j \in \{0 \dots 9\}$ ;  $t_1, t_2 \in \{1 \dots 5\}$ ,  $t_1 \neq t_2$ . It is to be noted that  $t_1$  and  $t_2$  need not be adjacent, *i.e.*,  $t_2$  is not necessarily equal to  $t_1 \oplus 1$ <sup>1</sup>. Furthermore, we will distinguish between the ordering of  $t_1$  and  $t_2$ . Thus the entry  $B_{t_1, t_2}[i, j] \neq B_{t_2, t_1}[i, j]$ .

In the ZIP Code application there are  $2 \times \binom{5}{2}$  possible position pairs of  $t_1$  and  $t_2$ . This would result in 20 matrices of size  $10 \times 10$  each<sup>2</sup>.

- **Markovian Model of  $2^{\text{nd}}$  order:** Conditional probabilities of the second order are represented by 3-dimensional matrices  $T_{t_1, t_2, t_3}[i, j, k]$  where

<sup>1</sup>Modulo  $L$  addition (denoted henceforth by the operator  $\oplus$ ) will imply that the successive neighbor of the digit in the last ( $L^{\text{th}}$ ) position is the first digit ( $L \oplus 1 = 1$ ).

<sup>2</sup>Traditional Bigram tables would simply note the frequency of occurrence of each digit pair.

$$i, j, k \in \{0 \dots 9\}, \quad t_1, t_2, t_3 \in \{1 \dots 5\}, \\ t_1 \neq t_2 \neq t_3.$$

- **Markovian Model of  $3^{\text{rd}}$  order:** Conditional probabilities of the third order are represented by 4-dimensional matrices  $Q_{t_1, t_2, t_3, t_4}[i, j, k, l]$  where  $i, j, k, l \in \{0 \dots 9\}$ ,  $t_1, t_2, t_3, t_4 \in \{1 \dots 5\}$ ,  $t_1 \neq t_2 \neq t_3 \neq t_4$ .

- 5-grams information would correspond to a Markovian Model of the  $4^{\text{th}}$  order. In the ZIP Code application there is just one ( $\binom{5}{5}$ ) possible position 5-gram. The corresponding matrix

$$P_{t_1, t_2, t_3, t_4, t_5}[i, j, k, l] \quad \text{where } i, j, k, l \in \{0 \dots 9\}, \quad t_1, t_2, t_3, t_4, t_5 \in \{1 \dots 5\}, \quad t_1 \neq t_2 \neq t_3 \neq t_4 \neq t_5 \text{ would be a 5 dimensional matrix.}$$

## 5.2 Markov Random Fields

Relation 5 shows the assignment of a digit at a particular position *triggers* the assignment of a digit at any position to the right (front) of it, keeping in mind that the positions are in a circular list. Triggering allows us to model the correlation between digits at any two positions in the string. The information content of the non-contiguous combinations could be equally if not more informative. Unlike the case of text strings, numeric strings usually do not carry any special meaning in adjacent groups of digits. The concept of triggering that does not strictly follow the Markov chain is captured adequately by MRF [4].

If  $\zeta$  is the MRF over  $\Gamma$  describing the class distribution in a ZIP Code, then the condition upon  $\zeta$  is for every  $t=1, \dots, L$ .

$$P(\zeta_t = C_t | \zeta_u = C_u, t \neq u) = P(\zeta_t = C_t | \zeta_u = C_u, u \in \Gamma_t).$$

$$\begin{bmatrix} B_{t_1, t_2}[0, 0] & \dots & B_{t_1, t_2}[0, 9] \\ \vdots & \ddots & \vdots \\ B_{t_1, t_2}[9, 0] & \dots & B_{t_1, t_2}[9, 9] \end{bmatrix} = \begin{bmatrix} P(C_{t_1} = 0 | C_{t_2} = 0) & \dots & P(C_{t_1} = 0 | C_{t_2} = 9) \\ \vdots & \ddots & \vdots \\ P(C_{t_1} = 9 | C_{t_2} = 0) & \dots & P(C_{t_1} = 9 | C_{t_2} = 9) \end{bmatrix} \quad (5)$$

$$(t_1, t_2) \in \{1 \dots 5\}^2; t_1 \neq t_2$$

We use  $\Gamma_C$  (for class C) to define the neighborhood assigned to a MRF class distribution. The MRF for position shown by  $\boxtimes$  in the string is shown by  $\boxplus$ .  $\forall t \geq L, t = t \oplus L$ .

$$\begin{cases} \Gamma_1 = \{3\} \\ \begin{bmatrix} \boxtimes & \square & \boxplus & \square & \square \end{bmatrix} \\ \Gamma_2 = \{4\} \\ \begin{bmatrix} \square & \boxtimes & \square & \boxplus & \square \end{bmatrix} \\ \Gamma_t = \{t-2, t+2\} \quad 3 \leq t \leq L \\ \begin{bmatrix} \boxplus & \square & \boxtimes & \square & \boxplus \\ \boxplus & \boxplus & \square & \boxtimes & \square \\ \boxplus & \square & \boxplus & \square & \boxtimes \end{bmatrix} \end{cases} \quad (6)$$

### 5.3 Modeling the Structure of ZIP Codes

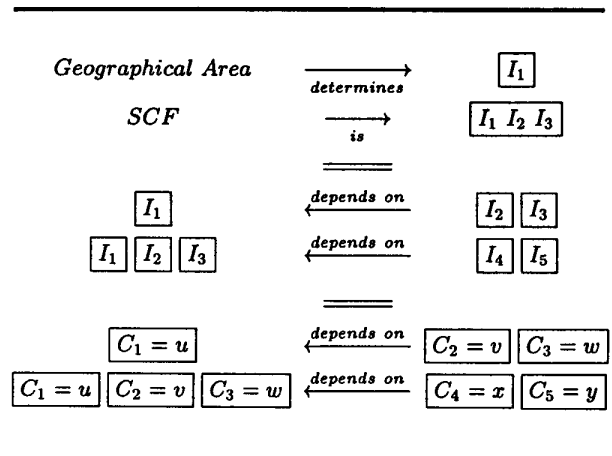
We have collected the n-gram tables (matrices) from the US Postal Delivery Point File. The ZIP Code consists of five digits. The first digit designates a broad geographical area of the United States, ranging from 0 for the Northeast to 9 for the West. The following two digits pinpoint population concentrations and reflect Sectional Center Facilities (SCF) that act as hubs in the transportation networks, (see Table 2). The final two digits designate small post offices. There are about 43,000 ZIP Codes in the United States that are valid. Only 43% of all possible 5 digit strings are valid.

An intuitive reading of the structure of the ZIP Code would suggest that the 2<sup>nd</sup> and 3<sup>rd</sup> digits of the ZIP Code are related to the 1<sup>st</sup> and the 4<sup>th</sup> and 5<sup>th</sup> digits depend on the first 3 digits (Relation 7).

Let us assume that the recognition choice of the ZIP Code under consideration is:  $[u v w x y]$ , i.e.,  $[I_1 \equiv u, I_2 \equiv v, I_3 \equiv w, I_4 \equiv x, I_5 \equiv y]$ . Then the objective is simply one of evaluating the probability of the string  $u v w x y$ . Equation 7 is what one would expect to be the probability of a sequence of 5 symbols. The probability of the ZIP Code  $[u v w x y]$  is given by Equation 8 using terms from all the 5 Markovian orders.

## 6 Experiments

We have experimented with various combinations of *triggers* and different orders of the Markov modeling.



**Table 2:** Structure of a 5 digit ZIP Code based on its design. The structure is enriched by the MRF modeling.

Not surprisingly, the best results were obtained for the modeling based on Equation 8 which is a derivative of Equation 7. This indeed is intuitively satisfying that the model which reflects our understanding of the correlation of digits in the ZIP code does perform the best. ZPP (ZIP Code Post Processing) software was tested on 4 different sets of ZIP Code images that were automatically extracted by the Handwritten Address Interpretation system [6]. GSC refers to the ZIP Code recognizer that is simply based on the recognition confidences of the GSC digit recognizer [5].

Table 3 shows a consistent 2% improvement on GSC. In fact, if the top 2 choice candidates are considered, the improvement margin is at 4%.

## 7 Summary

We have presented a nonstationary Markov model that merges images (signal) and context (language) in a fully Bayesian framework. The results obtained are promising.

## References

- [1] T.W. Anderson, L.A. Goodman, "Statistical inference about Markov chains", in *Ann. Math. Statist.*, 28, 89-111, 1957.

$$\begin{aligned}
P(C_1 = u, C_2 = v, C_3 = w, C_4 = x, C_5 = y) &= \\
&P(C_1 = u) \times \\
&P(C_2 = v, C_3 = w \mid C_1 = u) \times \\
&P(C_4 = x, C_5 = y \mid C_1 = u, C_2 = v, C_3 = w)
\end{aligned} \tag{7}$$

$$\begin{aligned}
&P(C_1 = u, C_2 = v, C_3 = w, C_4 = x, C_5 = y) \\
&= P(C_1 = u) \times P(C_2 = v, C_3 = w \mid C_1 = u) \times P(C_4 = x, C_5 = y \mid C_1 = u, C_2 = v, C_3 = w) \\
&= P(C_1 = u) \times P(C_2 = v \mid C_1 = u) \times P(C_3 = w \mid C_1 = u, C_2 = v) \times \\
&P(C_4 = x \mid C_1 = u, C_2 = v, C_3 = w) \times P(C_5 = y \mid C_1 = u, C_2 = v, C_3 = w, C_4 = x) \\
&= U_1[u] \times T_{3,2,1}[w, v, u] \times B_{2,1}[v, u] \times Q_{4,3,2,1}[x, w, v, u] \times P_{5,4,3,2,1}[y, x, w, v, u]
\end{aligned} \tag{8}$$

ZIP Code images	GSC	ZPP
1:	5,407	87.683%
2:	5,378	84.641%
3:	5,533	84.312%
4:	3,878	79.603%

**Table 3:** Results on 4 different sets consistently indicate the benefits of the Markov Modeling Postprocessing (ZPP) described in this paper compared to the method that does not make use of contextual information (GSC).

- [2] D. Bouchaffra, V. Govindaraju, and S. Srihari, "A methodology for determining probability of correctness of recognizer scores", in *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition, Santa Barbara, June 1998*.
- [3] D. Bouchaffra, E. Koontz, V. Kripasundar and R. K. Srihari, "Incorporating Diverse Information Sources in Handwriting Recognition PostProcessing", in *International Journal of Imaging Systems and Technology, vol. 7, issue 4, Published by John Wiley, June 1996*.
- [4] R.L. Dobrushin, "The description of random fields by means of conditional probabilities and conditions of its regularity", in *Theory Prob. Appl., vol. 13, pp. 197-224, 1968*.
- [5] J. Favata and G. Srikantan, "A Multiple Feature/Resolution Approach to Handprinted Digit and Character Recognition", in *International Journal of Imaging Systems and Technology, 7*, pp. 304-311, 1996.
- [6] V. Govindaraju, A. Shekhawat, and S. Srihari, "Determining Delivery Point Code in US mailstream", in *IWFHR III, 1993*.
- [7] G. Kim and V. Govindaraju, "Bankcheck recognition using cross validation of courtesy and legal amount fields", in *International Journal of Pattern Recognition and Artificial Intelligence 19(4)*, pp. 366-379, 1997.
- [8] L.R. Rabiner, "A tutorial on hidden Markov models and selected applications in speech recognition", in *Proceedings of the IEEE, 77(2)*, pp. 257-286, 1989.
- [9] Z. Shi and V. Govindaraju, "Segmentation and recognition of connected handwritten numeral strings", in *Pattern Recognition, 30*, pp. 1501-1504, 1997.
- [10] R. Singhal and G. Toussaint, "Experiments in text recognition with the modified Viterbi algorithm", in *IEEE Transaction on Pattern Analysis and Machine Intelligence, 1(2)*, 1979.
- [11] M. Sridhar and A. Badreldin, "Context-directed segmentation algorithm for handwritten numeral strings", in *Image and Vision Computing, 5(1)*, pp. 3-9, 1997.