

Transcript Mapping for Historic Handwritten Document Images

Catalin I. Tomai, Bin Zhang and Venu Govindaraju
CEDAR

UB Commons, 520 Lee Entrance, Suite 202,
Amherst, NY, 14228-2567

E-mail: {catalin,binzhang,govind}@cedar.buffalo.edu

Abstract

There is a large number of scanned historical documents that need to be indexed for archival and retrieval purposes. A visual word spotting scheme that would serve these purposes is a challenging task even when the transcription of the document image is available. We propose a framework for mapping each word in the transcript to the associated word image in the document. Coarse word mapping based on document constraints is used for lexicon reduction. Then, word mappings are refined using word recognition results by a dynamic programming algorithm that finds the best match while satisfying the constraints.

1 Introduction and Previous Work

Historical documents are a valuable resource for scholars and their indexing for archival and retrieval purposes is highly-desired. This indexing problem can be treated differently depending on several factors: the documents were written by one author or by multiple authors, the availability of the text transcript of the document, the degree of noisiness of the document, etc. The processing in the multiple-authors case is much harder because of the high in-class word variability. The scanned documents present noise introduced by the photocopying and scanning processes together with underlines, overlapping lines and words, etc.

In case transcripts of the documents are available, the current systems index the image documents at the document level, that is, given a certain query word, a document image or a set of document images that contain that word are returned. What we want is to design a retrieval system that returns the exact document image word or line that corresponds to the particular query word. Since in most cases a one to one mapping between the lines of text in the transcript and the lines of the document image doesn't exist, the **word image - transcript word** mapping task is not evi-

dent. In our work we assume the transcript is one (long) line of text.

The proposed system's goal is to locate (spot) words in noisy historical documents written by multiple authors for which a transcription is present.

Other authors that have addressed this problem have decided against the use of OCR, for considerations of inadequacy. Since OCR systems depend on accurate word segmentation and recognition their usage was deemed inappropriate for this type of documents. Keaton and Goodman ([1]) developed an alternative strategy based on learning a set of keyword signatures of particular words of interest. There is no page segmentation step, instead a cross-correlation of the document image with a set of keyword prototypes which have been extracted from a training set of documents is executed.

Manmatha and others ([2]) deal with a single-author problem by matching word images with each other to create equivalence classes. Each equivalence class consists of multiple instances of the same word. They use a word segmentation step that extracts the bounding boxes of the word images by a sequence of window operations of smoothing and thresholding.

While we agree that a perfect line and word segmentation is impossible for historical documents, we believe that by using algorithms that process several word segmentation hypotheses we can satisfactorily solve the segmentation problem.

Unconstrained word recognition is a difficult task. For specific domains local or global constraints like mail address format or check layout or properties of specific fields (e.g. postal codes composed of digits only) are used to reduce the lexicons and consequently the recognition errors. For our problem, the transcripts availability allows us to reduce the number of word candidates for word recognition, based on some constraints that we are going to define later.

The question is: How to optimally utilize the information from a document image and its corresponding transcript to get the best mapping results?

A typical input to our system is presented in Figure 1, which consists of a letter written by Thomas Jefferson in 1787 and its corresponding transcript.

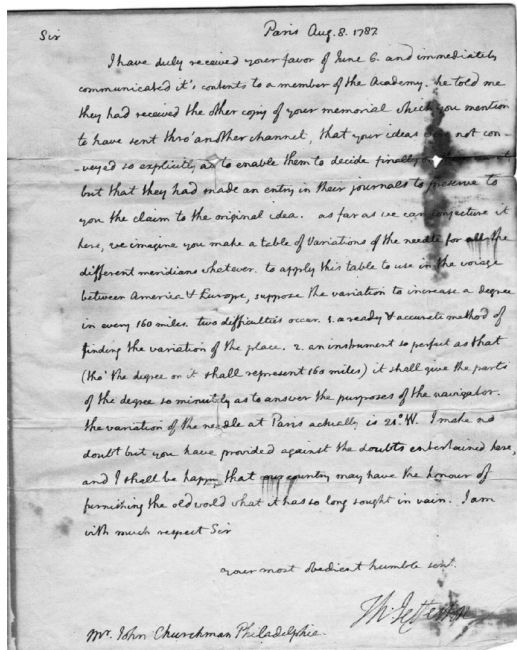


Figure 1. Image of a 1787 Thomas Jefferson letter and its corresponding transcript

The paper is organized as follows: Section 2 succinctly discusses the line and word separation modules. The word mapping problem is formalized in Section 3. Section 4 describes the proposed algorithm. The experimental results are presented and analyzed in Section 5. Conclusions are included in Section 6.

2 Line and word separation

2.1 Line separation

The goal of this module is to correctly divide the handwritten text into lines so that each line can be further divided into words (see [3] and [4] for algorithm details).

Variability in the baseline position, line skew, character-size and inter-line distance make this a very difficult task for historical documents.

The process does not always return the expected result, sometimes words from different lines are grouped together. The proposed system is robust with respect to these problems.

2.2 Word Separation

Word Separation is the problem of segmenting a line into words. We assume that inter-word spacing is greater than inter-character spacing. Punctuation information together with inter-word gaps is used for word separation.

For the task at hand a correct word separation is important since we expect a poor performance of the word recognizer on the word images of the historical documents. We generate multiple word separation hypotheses for each line.

In [2] the word separation is done without generating multiple hypotheses. Also, in [1] the authors avoid line separation and word separation while identifying candidate locations by cross-correlating the document with a set of keyword prototypes which have been extracted from a set of documents.

We believe that it is essential to generate multiple word separation hypotheses for a certain line to be more confident that the correct word separation configuration is included as one of the hypotheses or, more probably, composed from words from different hypotheses. Otherwise, we may miss the right configuration which would negatively impact the later stages of the matching process.

The algorithm used essentially ranks the gaps between centers of adjacent components, that is, the distance between the components' convex hulls ([4] for more details). Then, the hypotheses for choosing n words (where n can take values between n_{min} and n_{max}) from the given line are ranked and returned.

3 Problem Description

The result of the line and word segmentation of the handwritten document image I_D is a set of hypotheses. For each line we obtain several hypotheses, each one containing a set of word images - the results of a different word separation. The set of hypotheses for a line is matched against a subset of words from the transcript. Using the word recognition results we assign to each word image of a hypothesis a transcript word with a certain confidence. The set of word images from the hypotheses set, together with a sequence of transcript words are given to a Dynamic Programming algorithm that will return the highest scoring sequence of word images that correspond to the sequence transcript words.

Before we formalize the problem, we give the following definitions:

A **line hypothesis** W for a line image I_l consists of a list of ordered disconnected word images (see Figure 2 for an example):

$$W = (h_1, h_2, \dots, h_{n_l})$$



Figure 2. Two line hypotheses for the second line of the image

A **line hypothesis set**, \overline{W} , contains multiple hypotheses for a line image I_l :

$$\overline{W} = (W_1, W_2, \dots, W_{m_l})$$

A **document hypothesis set**, Ψ , contains the hypotheses from all line sub-images of a document image I_D :

$$\Psi = (\overline{W}_1, \overline{W}_2, \dots, \overline{W}_k)$$

A **truth transcript**, Γ , for a document image I_D is an ordered list of words, expressed as:

$$\Gamma = (\varpi_1, \varpi_2, \dots, \varpi_p)$$

A **mapping** \mathfrak{S} of a transcript Γ and it's corresponding image I_D for which a document hypotheses set Ψ was obtained, is defined as:

$$\mathfrak{S} = ((\varpi_1, h_1), (\varpi_2, h_2), \dots, (\varpi_p, h_p))$$

where $\Gamma = (\varpi_1, \varpi_2, \dots, \varpi_p)$, and $h_i \in \Psi, \forall 1 \leq i \leq p$.

After the line hypotheses generation, each word image h has assigned the following values: the line index ($h.ln$), a word sequence number ($h.wn$); the bounding box coordinates ($h.left, h.right, h.top, h.bottom$) specifying its position in the original image. In the end, each word image h will be associated a word transcript $h.v$ and a confidence value $h.f$, as returned by the word recognition.

Given a set of line hypotheses Ψ for a handwritten document image I and it's truth transcript Γ , the goal is to find an ordered word list $\psi \subset \Psi$ which best matches the transcript Γ , i.e.,

$$\begin{aligned} \psi = \{ & (h_1, h_2, \dots, h_p) \mid \min(\sum_{j=1}^p \frac{d(h_j.v, \varpi_j)}{h_j.f}), \\ & \forall 1 \leq j \leq p, h_j \in \Psi \text{ and } \varpi_j \in \Gamma, \text{ and } \forall r < s \\ & \text{and } h_r.ln = h_s.ln, h_r.left < h_s.left, \forall r < s, \\ & h_r.ln \leq h_s.ln \} \end{aligned}$$

In the formula above, $d(\cdot, \cdot)$ represents the distance between two words, while the second and third constraints assure that the order of words in ψ conforms to the word/line layout in the document image I_D .

4 Algorithm Design for Word Recognition/Mapping

An algorithm called Mixture of Word Recognition and Word Mapping (MiWRM), was designed to solve the problem described above.

4.1 MiWRM Algorithm

The diagram of MiWRM is shown in Figure 3

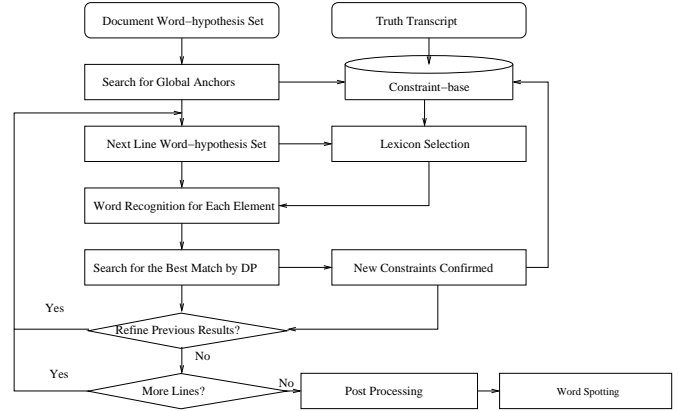


Figure 3. Algorithm Diagram of Word Recognition/Mapping

To each line we assign a subset of words from the transcript as the line lexicon. MiWRM builds an lexicon for each word image in Ψ . The goal is to have the right transcript word as one of the entries in the lexicon attached to the word image. The input to the word recognizer is the word image and the previously computed lexicon and the output is the ranked list of lexicon words. Therefore, for each line's hypothesis set $\overline{W} \subset \Psi$, any word $h \in \overline{W}$ will be associated with the word recognition result, which consists of a list of pairs (character string, confidence) sorted in descending order by confidence. A dynamic programming (DP) algorithm is used to find an ordered word list $\varphi, \forall h \in \varphi, h \in \overline{W}$ which best matches the transcript corresponding to that particular line. Word images in each line will be associated to their corresponding entries in φ and have attached a confidence value that reflects the degree of recognition certainty. We obtain a list $\mathfrak{S}_l = (\varpi_{k1}, h_{k1}), (\varpi_{k2}, h_{k2}), \dots, (\varpi_{kn}, h_{kn})$ for a certain line l . Each entry in the list corresponds to a unique lexicon word ($k1, k2, \dots, kn$) and the entries follow the word order in the transcript.

We are looking for a sequence of entries in the above list that contain a sequence of consecutive words images that were recognized with high confidence. Such a sequence is called an **anchor**.

Once a line is successfully processed (we have found at least one anchor), we enforce a set of new constraints on the contents of the previous and following lines. The new constraints set will be included in the constraints-base. If we have a high confidence in the current line processing results, some previous processing results will be refined, based on the updated constraint database.

The position of the lexicon word of the last entry of the last anchor is used for the computation of the lexicon for the following line.

4.2 Constraint-base

Besides the anchors we use other information that would further help shape the matching process. These constraints are mainly by-products of the recognition process. Examples are: number of components for a line (used to impose a lower limit on the size of the next line lexicon), average character size, average gap size, estimated line width. They are stored in a constraint-base (**CB**). CB is initialized before the recognition/mapping process and will be updated dynamically within the process once new constraints have been confirmed.

4.3 Coarse and Fine Word Mapping

Given a set of word hypotheses Ψ and the corresponding transcript Γ , *coarse word mapping* consists of finding the set of transcript words from Γ for each word image $h \in \Psi$ of a certain line hypothesis. In this scope we use the position of anchors, the constraint information and the word image position in the document.

The word recognizer used ([5]) takes as inputs a word image h and a lexicon, and returns a list of lexicon entries ranked in descending order of their recognition confidence.

Given a hypothesis set $\overline{W} \subset \Psi$ of a line image I_l , for which each word $h \in \overline{W}$ has been tagged with the word recognition result, the goal of *fine word mapping* is to find an ordered word list φ , $\forall h \in \varphi, h \in \overline{W}$, that best matches the transcript ℓ corresponding to I_l . The transcript ℓ of I_l is obtained in the coarse word mapping stage.

A dynamic programming algorithm, Longest Common Subsequence (LCS), is designed to find the common sub word-sequence (CSWS) of $W_i \subset \overline{W}$ and ℓ , $\forall 1 \leq i \leq m_l$. Here, for each W_i , only the top entry with the highest confidence is considered. Moreover, if the confidence of the top entry in W_i is lower than a threshold, W_i will be ignored (e.g. small word parts, noise).

In this way, each word w in the transcript ℓ will be associated with several (or none for some cases) word images (hypotheses), and the word image h with the highest confidence is chosen as a mapping to w . Therefore, we get a line mapping \mathfrak{S}_l of ℓ and \overline{W} .

Finally, we need to examine the correctness of the mapping \mathfrak{S}_l . Given

$$\mathfrak{S}_l = ((\varpi_1, h_1), (\varpi_2, h_2), \dots, (\varpi_p, h_p))$$

a legal mapping should satisfy the following condition:

$$\forall 1 \leq r < s \leq p, h_r.left \leq h_s.left$$

We use \mathfrak{S}_l to impose new constraints on the following lines, i.e., we can narrow down the searching range of coarse word mapping for the next line hypotheses' sets by starting from the word next to the last element in \mathfrak{S}_l . \mathfrak{S}_l is included in the constraint-base (CB).

4.4 Post Processing

The goal of this step is to finalize the process by determining the final positions of the word images corresponding to the transcript words.

Until now we have partly mapped the transcript Γ to the image I_D . The anchors are scattered inside the mapping \mathfrak{S} , any word between two anchors being in a dangling (unconfirmed) state. Let's consider two consecutive anchors in \mathfrak{S} , (ϖ_{a1}, h_{a1}) and (ϖ_{a2}, h_{a2}) .

$$\mathfrak{S} = (\dots(\varpi_{a1}, h_{a1}), (\varpi_{c1}, h_{c1}) \dots (\varpi_{ck}, h_{ck}), (\varpi_{a2}, h_{a2}) \dots)$$

We assume that the two anchors belong to the same line and the line bounding box is defined by the coordinates of the upper-left corner and bottom-right corner: (Lleft, Ltop, Lright, Lbottom). A rough mapping for any pair (ϖ_{ci}, h_{ci}) , $1 \leq i \leq k$ located between the two anchors is given by the following coordinates: $h_{ci}.left = h_{a1}.right + 1$, $h_{ci}.top = Ltop$, $h_{ci}.right = h_{a2}.left - 1$, $h_{ci}.bottom = Lbottom$. Every word in the transcript will be assigned an exact or rough mapping. Word images inside the anchors are not always assigned an exact mapping, because of the recognizer's weak performance on the noisy image. Figure 5 displays the bounding boxes for the mapped word images of two lines from the binarized example image. The entire set of mappings constitutes the document's mapping.

After building up the **document mapping** \mathfrak{S} , word spotting is straightforward. As defined before, $\mathfrak{S} = ((\varpi_1, h_1), (\varpi_2, h_2), \dots, (\varpi_p, h_p))$, so, given a keyword x , we just need to compare x with ϖ_i , $\forall 1 \leq i \leq p$, the word image h_y is located in the document I_D for any matched transcript word ϖ_y .

5 Experimental Results and Analysis

We have evaluated the performance of the system on the image from Figure 1. To accurately measure the algorithm performance, we have built a truth database that stores the

exact bounding-box information for the word images corresponding to the words in the transcript. From the total of 249 words, 217 words are included in the database while 32 words are excluded because of their extreme noisiness.

A mapping (ϖ_i, h_i) is evaluated as correct when the bounding-box of the word hypothesis h_i contains the bounding-box of ϖ_i in the truth database.

First, the original image is binarized using the Quadratic Integral Ratio (QIR) algorithm([6]). Then, the binarized image is divided into 23 lines. For each line, line hypotheses are computed. Finally, by applying the MiWRM algorithm we map each word in the transcript to a word image contained in some line hypothesis.

The lexicon size for each line is shown in Figure 4 (a). The average lexicon size per line is 13 words.

Figure 4 (b) displays the number of different line hypotheses. There is a total of 2039 word images in all the line hypotheses generated. The lexicon size for each word image in a line hypothesis is presented in Figure 4 (c). For each word we use a subset of the line lexicon of average size 6 words. The reason for using a sub-lexicon is the word recognizer's poor performance on a larger lexicon (given the poor quality of historic document images).

Before the post processing step of MiWRM, anchors that contain a total of 69 words were generated. After post processing, 180 words (82.95%) out of the total of 217 words are mapped (17 exactly mapped and 163 roughly mapped). This performance shows the effectiveness of the proposed algorithm.

From our knowledge this is the first attempt to address the word spotting problem with an available transcript, therefore making it impossible to compare our results with others. Our immediate goal is to improve on the present performance to obtain a better matching performance and to test the validity of our approach on a larger set of images.

One of the drawbacks of MiWRM is the large number of word images produced by the line hypotheses. However, using a cache mechanism, recognition is not repeated for the same word images.

6 Conclusions

Word Recognition/Mapping (WRM) is the key component of a system seeking to index historical handwritten documents. In this work we formalize the WRM problem and design an algorithm (MiWRM) to solve it. In MiWRM, word recognition and word mapping work in tandem, the lexicon size for word recognition is reduced by coarse word mapping (lexicon selection) based on the document constraints, and fine (exact) word mappings is done based on the word recognition results together with the constraints. The high accuracy of the mapping for a historic document image is an indication of the effectiveness of the proposed

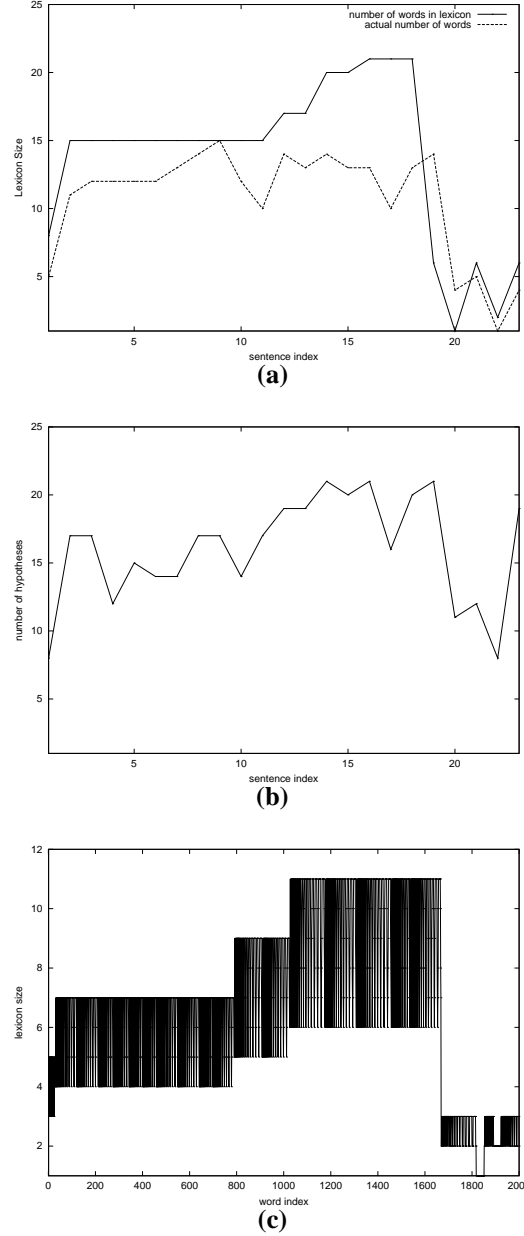


Figure 4. (a) Lexicon Size for Each Line (b) Number of Word-break Hypotheses in Each Line (c) Lexicon size for each word image

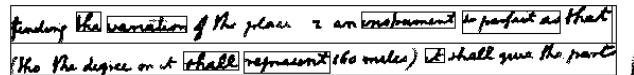


Figure 5. Rough and Exact mapping for two lines from the binarized image

system. Given the poor quality of historic document images, there is much room for improvement, such as seeking a better binarization algorithm, training the word recognizer for the type of characters found in these documents, more efficient mapping algorithms, etc.

7 Acknowledgments

The authors would like to thank Dr. Graham Leedham for helping them with the QIR binarization code

References

- [1] P. Keaton, H. Greenspan, and R. Goodman. Keyword spotting for cursive document retrieval. In *Proceedings of the Workshop on Document Image Analysis - DIA'97*, 1997.
- [2] R. Manmatha, Chengfeng Han, and E.M. Riseman. Word spotting: A new approach to indexing handwriting. In *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition 96, San Francisco*,, pages 631–637, June 1996.
- [3] G. Seni and E. Cohen. External word segmentation of off-line handwritten text lines. *PR*, 27(1):41–52, January 1994.
- [4] U. Mahadevan and R. Nagabushnam. Gap metrics for word separation in handwritten lines. In *ICDAR*, pages 124–127, 1995.
- [5] G. Kim and V. Govindaraju. A lexicon driven approach to handwritten word recognition for real-time applications. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(4), April 1997.
- [6] Y. Solihin and C.G. Leedham. Integral ratio: A new class of global thresholding techniques for handwriting images. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 21(8):761–768, 1999.