

Transcript mapping for handwritten Arabic documents

Liana M. Lorigo Venu Govindaraju
Center for Unified Biometrics and Sensors
University at Buffalo, Amherst, NY 14228
govind@buffalo.edu

ABSTRACT

Handwriting recognition research requires large databases of word images each of which is labeled with the word it contains. Full images scanned in, however, usually contain sentences or paragraphs of writing. The creation of labeled databases of images of isolated words is usually tedious, requiring a person to drag a rectangle around each word in the full image and type in the label. Transcript mapping is the automatic alignment of words in a text file with word locations in the full image. It can ease the creation of databases for research. We propose the first transcript mapping method for handwritten Arabic documents. Our approach is based on Dynamic Time Warping (DTW) and offers two primary algorithmic contributions. First is an extension to DTW that uses true distances when mapping multiple entries from one series to a single entry in the second series. Second is a method to concurrently map elements of a partially aligned third series within the main alignment. Preliminary results are provided.

Keywords: Handwriting recognition, alignment, Arabic script

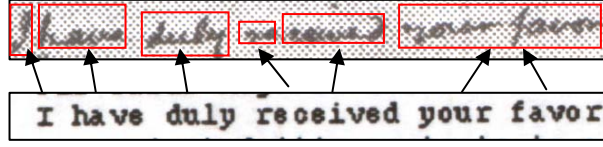
1. INTRODUCTION

We consider the task of aligning a text file with an image of handwritten text. Transcript mapping enables one to automatically generate databases of word images with text labels for use in recognition research. This avoids the usual tedious process of manually dragging a box around each word in an image of a full page and keying in the annotations separately. Transcript mapping is also called “alignment”. While it is much easier than full recognition, it requires some recognition capability. For free handwriting in Arabic, the state of the art does not yet include full recognition but is ideal for this work since this work builds on partial recognition work and facilitates research.

Figure 1a shows sample input to a transcript mapping system. Such input includes an image of text and a corresponding text file. The first step in transcript mapping is the detection of word locations in the image of handwriting. This step is called “word segmentation” or “segmentation”. If this step achieves perfect results, the problem would be trivial since the correct alignment would be one-to-one. However, segmentation is never perfect due to variation in handwriting. Figure 1b shows an English example with locations indicated by bounding boxes. An English example is shown for ease of presentation only, as the method of this paper operates on the Arabic script exclusively. “Received” has been split into two boxes, and “your” and “favor” have been combined. Splitting of a word is called “oversegmentation” and merging multiple words is called “undersegmentation”. The arrows indicate the desired mapping.

وأدويةهم
ونظائهم، ثم قدم العربي مستعمراً حاملاً معه ثقافة وطباً جلياً،
فأضاح الكثير على المفاهيم الطبية السائدة منذ الصور المدهية.
ولقد تم تسجيل ما حمله المسلمون من عرب وفرنس في كتاب
للشاعر الشرق أوسطي (ها ي باو) الذي عاش في الصين قبل

وأدويتهم وغذائهم، ثم قدم الغرب مستعمراً حاملاً معه ثقافته وطبه الحديث، فأضاف الكثير إلى المفاهيم الطبية السائدة منذ العصور القديمة. ولقد تم تسجيل ما حملة المسلمون من عرب وفرنس في كتاب للشاعر الشرق أوسطي (هاي يوا) الذي عاش في الصين قبل



(b)

Figure 1: (a) Image of Arabic handwriting and its transcript. (b) Illustration on an English example. This example is from a Thomas Jefferson letter used in [1], but the boxes shown do not reflect that method.

The Arabic alphabet has 28 letters. Arabic text is written from right to left, and adjacent letters in a word are joined except when the first is one of the six highlighted in Figure 2. In such cases, there is a space before the next letter. A subword is a connected sequence of letters. Letters are connected along the baseline (also called the lower baseline). Each has two to four shapes corresponding to start, middle, or end of a subword, or alone (Figure 3a). The first letter is “alif”, leftmost in Figure 2. Alif is frequent and easily recognizable when alone so plays a role in our method.

Our method uses the subword as the basic unit instead of the word as used in the previous English work and in Figure 1b. In handwritten text, spaces within Arabic words often have the same width as spaces that indicate word breaks, so spaces cannot be used as an indicator for word endings. They are, however, a good estimate of subword endings. The presence of dots and other markings complicates the task for Arabic compared with English. While dots such as in “i” and “j” can usually be discarded as noise in English handwriting recognition systems, many Arabic letters differ only by the number and position of dots. These dots and additional markings or diacritics (Figure 3b) must be assigned to a letter. In ideal penmanship, each would appear above or below the corresponding letter body, but in practice they are often shifted so their alignment with letter bodies is an additional challenge.

ي و ه ن م ل ك ق ف غ ع ظ ط ض ص ش س ز ر ذ د خ ح ج ث ت ب ا

Figure 2: The Arabic alphabet (shown left to right for display only). The highlighted letters cannot be connected to a following letter.



Figure 3: (a) The isolated, initial, medial, and final shapes of the letter "ayn". (b) Markings in handwriting.

2. PREVIOUS WORK

Little work has been done on the alignment task for handwritten images. In 2002, Tomai *et al.* [1] proposed a recognition-based method that was tested on the image of which a portion is shown in Figure 1b. In it, the word segmentation method produced approximately 15 hypotheses per line. A subset of words from the transcript is assigned as a lexicon for each line, and a lexicon is built for each word image. The average size of the lexicon for a line is 13 words. This small size enables recognition on degraded images that could not be recognized by a standard recognizer.

For our purposes, however, the drawback is the requirement of a handwriting recognizer. Recognition systems are not yet available for free Arabic handwriting ([2]). Moreover, a transcript mapping system could assist research toward such systems for all Arabic handwriting, not only for degraded documents.

Kornfield *et al.* [3] presented a method that was also applied to historical handwritten English documents but which does not require a recognizer. It is based on Dynamic Time Warping (DTW). DTW aligns two time series using dynamic programming, and here those series are the image locations from the segmentation step and the text words. Let b_i indicate location (or bounding box) i and w_j word j . Let $d(i,j)$ be some measure of the distance between b_i and w_j . Let $\Delta(i,j)$ be the cost of the best assignment up to b_i and w_j , starting from the beginning of the text. Initialize the base case, so $\Delta(0,0) = d(0,0)$. Then $\Delta(i,j)$ is computed as the minimum of the three previous paths: either ending with b_{i-1} at w_j , b_i and w_{j-1} , or b_{i-1} and w_{j-1} , plus $d(i,j)$. After computing the full recursion, the alignment of a full page is recovered by tracing back the best path. Their method was tested on 70 pages of letters authored by George Washington and achieved 74.5% accuracy when aligning single lines of text (that is, when locations of line breaks were given in the transcript) and 60.5% accuracy when aligning full pages at a time.

A 2006 method by Rothfeder et al. [4] used the same features and test set as [3] but a Hidden Markov model instead of DTW for alignment. It achieved 72.8% accuracy when aligning full pages at a time, an improvement over the previous method. A drawback of both methods, however, is that comparisons are limited to those between individual bounding boxes (feature vectors) and individual transcript words. For example, to model the alignment of “received” in Figure 1, the DTW method approximates the distance between the pair of boxes showing “re” and “ceived” and the transcript word “received” by the sum of the distances between the transcript word and each individual box. This is a poor estimate since we expect the distances using the partial words to be high, and their sum even higher, but the true distance that would use the concatenation of the boxes would be small. The HMM method similarly uses the probabilities of observing each individual feature vector (box) given a transcript word. At no time is a new feature vector that represents the concatenation of two boxes used. Likewise at no time is the concatenation of two transcript words used to handle the case of undersegmentation as in “your favor”.

We found no prior work in transcript mapping for handwritten Arabic documents. The differentiating aspect most relevant is the large number of dots and diacritics in Arabic. The assignment of these markings to subwords must be determined. While it is less variable than the task of assigning subimages to text words since an ambiguity usually involves only the subimages to the left and right of the marking, it may be impossible to detect beforehand. Our challenge was to develop a transcript mapping system that does not require a full recognizer and which operates on Arabic text. Our first experiments involved DTW with a small number of features. They achieved low accuracy due to the poor distance approximations computed for merged or split words. We present a variation on DTW to overcome this problem and a “mini-DTW” step within the larger alignment to address the assignment of dots and diacritics.

3. APPROACH

3.1 True Distance DTW

Cases of undersegmentation or oversegmentation are the cause of the difficulty in this task. However, they are exactly where the distance estimates computed in a standard DTW algorithm are poorest since the algorithm sums individual distances that would likely be high, yet the true distance should be low. A better approach would compute the distance between multiple boxes and one transcript unit or one box and multiple units simultaneously.

We developed a variation on DTW that we call True Distance DTW. True Distance DTW (TD) stores two Δ matrices instead of one. $\Delta_{1-1}(i,j)$ is the total cost of the alignment up to box i and w_j , with b_i and w_j as a one-to-one match. $\Delta(i,j)$ is the total cost of the alignment up to b_i and w_j with no constraint, the same as in standard DTW. Δ_{1-1} is computed as $\Delta(i-1,j-1) + d(i,j)$. In this implementation we allow at most two subimages per text subword and at most two subwords per subimage. Distances are computed by a function that allows two subwords ($j-1$ and j) and one subimage or two subimages ($i-1$ and i) and one subword. $\Delta(i,j)$ is the minimum of these and $\Delta_{1-1}(i,j)$:

$$\Delta(i,j) = \min \begin{cases} \Delta(i-1,j-2) + d(i,j-1..j) \\ \Delta(i-2,j-1) + d(i-1..i,j) \\ \Delta_{1-1}(i,j) \end{cases}$$

The choice at which the minimum was attained is stored for final alignment. We call these values the direction grid. After $\Delta(i,j)$, $\Delta_{1-1}(i,j)$, and the direction grid are computed for the entire transcript and all subimages in the full image, the system backtracks through the direction grid to generate the alignment. The backtracking starts at the direction value corresponding to the last subimage and the last transcript word. Three constraints are added. First, transcript subwords of a single letter with no marks or with only a single dot cannot map to two subimages. It is presumed that the segmentation step would not split such a subword. Second, subimages from different lines cannot map to the same subword. Third, the backtracking routine forces a shift to the previous subword and the previous subimage after a merge to prevent three-to-one matches that would be inconsistent with the TD calculations.

A sample TD run for the image in Figure 4 is shown in Figure 5. Columns indicate subwords and rows indicate image locations. The bold entries indicate the best path. Note that the first two subwords (rightmost) are merged. The method detected this as shown by “1 0” in the bounding box in Figure 4. Subword 8 is split, as also detected.

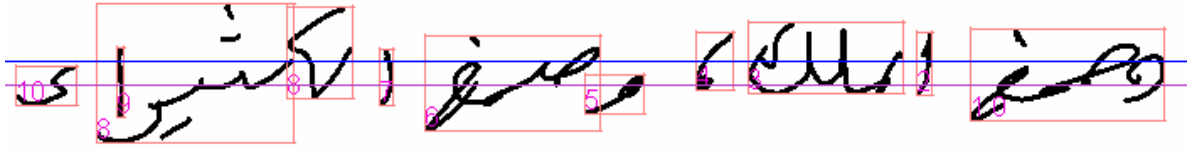


Figure 4: Image corresponding to TD illustration. Numbers in boxes indicate the alignment result. They correspond to the indices of the transcript subwords as show at the top of the chart in Figure 5.

	0	1	2	3	4	5	6	7	8	9	10
	و	صغ	أ	ملك	،	و	صغ	أ	أكثر	أ	ء
وهو	2.25 2.25	4 - ←	5 - *	26 - *	35 - *	39.25 - *	69 - *	80 - *	119 - *	136 - *	136 - *
أ	14.25 - *	12.5 12.5 \	4 4 \	12 12 \	27 27 \	38.25 38.25 \	49.5 49.5 \	58.5 69 ←	85 85 \	98 119 ←	130 138 ←
ملك	48.25 - *	47.5 47.5 \	29.5 29.5 \	10 10 \	7 28 ←	29.25 53.25 ←	66 71.5 ←	62.5 66.5 ←	56.5 70.5 ←	65.5 102 ←	101 123 ←
،	61.25 - *	57.5 57.5 \	48.5 48.5 \	15 35.5 ↑	10 10 \	9.25 9.25 \	22 38.5 ←	37.5 67 ←	66.5 66.5 \	57.5 57.5 \	56.5 66.5 ←
و	64.25 - *	56.5 68.5 up	60.5 60.5 \	36.5 60.5 up	17 17 \	10.25 10.25 \	16.5 16.5 \	17.5 25 ←	33 43.5 ←	48.5 69.5 ←	58.5 58.5 \
صغ	99.25 - *	63.5 65.5 ↑	64.5 68.5 ←--	60.5 69.5 UP	47.5 47.5 \	25.25 25.25 \	11.5 11.5 \	10.5 28.5 ←	28.5 28.5 \	31.5 45 ←	44 60.5 ←
أ	123.25 - *	80.5 108.5 ↑	64.5 64.5 \	72.5 72.5 \	62.5 62.5 \	49.75 49.75 \	26.5 34.5 ↑	12.5 12.5 \	14.5 14.5 \	27.5 29.5 ←	34.5 34.5 \

ك	163.25 - *	125.5 135.5 ↑	82.5 82.5 \	67.5 67.5 \	70.5 73.5 <--	67.75 67.75 \	51.5 62 ↑	28.5 28.5 \	15.5 15.5 \	16.5 16.5 \	15.5 31.5 <--
ش	178.25 - *	134.5 170.5 ↑	130.5 130.5 \	68.5 94.5 ↑	71.5 71.5 \	70.75 72.75 <--	61 75 ↑	56.5 56.5 \	12.5 32.5 ↑	20.5 20.5 \	19.5 19.5 \
ا	212.25 - *	173.5 187.5 ↑	135.5 135.5 \	97.5 138.5 ↑	70.5 70.5 \	73.75 73.75 \	78 80 ↑	62 62 \	33.5 60.5 ↑	13.5 13.5 \	23.5 23.5 \
و	225.25 - *	197.5 219.5 ↑	176.5 176.5 \	147.5 147.5 \	99.5 99.5 \	70.75 70.75 \	81 81 \	81 81 \	68 68 \	36.5 36.5 \	14.5 14.5 \

Figure 5: Alignment example. In each cell, the first number is Δ and the second is Δ_{1-1} . The direction for backtracking is also shown. “\” indicates the diagonal direction. “*” indicates an invalid correspondence.

3.2 Features

We use the term “mark” to mean a dot group or any other marking that is placed above or below a letter. TD does not distinguish between types of marks except according to whether they appear above or below the letter they modify.

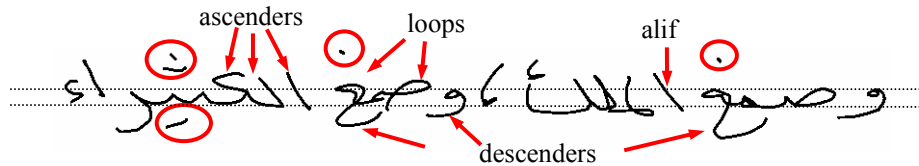


Figure 6: Features used to represent subimages. Marks are circled and other features are indicated. Upper and lower baselines are indicated.

We use six features building on [5]. Each is computed per subimage and generated automatically per transcript subword. Let $f=[f_1 \dots f_6]$ and $g=[g_1 \dots g_6]$ be feature vectors for a subimage and a subword respectively. The first five dimensions are number of loops, number of ascenders, number of descenders, number of marks above, and number of marks below. The sixth is binary and indicates whether or not the subimage or subword is an alif (Figure 6). The calculation of f_6 uses empirically determined ratios between height and width and between height and stroke width. Failing those tests it compares the subimage with a template of a diagonal line of the corresponding size. The system stores values for $g_1 \dots g_5$ for all letter shapes and punctuation marks. The values were set manually according to letter definitions and observations and can be learned automatically in the future. When a value can vary between 0 and 1, 0.5 is used. For example the loop in the letter waaw (subword 5 in Figure 4) is often omitted.

The distance used is the L2 distance between f and g , plus a penalty related to misalignment of an isolated alif or a period. The penalty is large if f indicates an alif and g does not. This has the effect of using alif as a “stop word” which means that the alignment is anchored by alifs since they are easier to align than other subwords. The penalty is also large when attempting to fit alif and another box with one transcript word. This is a heuristic that prevents an error observed in initial experiments. Last, it is large when aligning a wide subimage with a period. Large penalties have a fixed value of 10. The penalty is zero in all other cases. Since the segmentation algorithm may merge a period with a subword, the system allows one to match as a dot with the previous subword or as its own subword. Parentheses are allowed to pass the alif-related tests without penalty since their shapes are close. Punctuation marks are likely to be easier to detect than letters, and results below suggest that future work would benefit from additional analysis of them. Let f indicate the features for b_i and g the features for w_j . Then

$$d(i,j) = \sum_{k=1}^6 (f_k - g_k)^2 + \text{penalty}$$

3.3 Mini-DTW for Dot Groups and Markings

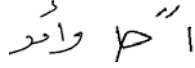


Figure 7: Correspondence of marks with letter bodies cannot be determined without at least partial recognition.

In tests of the system as described up to this point, the most common error was in assigning marks to subwords. Our initial method based on horizontal distances to subwords allowed incorrect assignments when marks were offset (Figure 7). It is difficult and in some cases impossible to determine the correct assignments only from locations. Since we have a transcript, however, we can use it to assign the marks. While one could frame the problem as three-dimensional DTW where the three series are the transcript as above, the subimages (as above except without marks) and the list of marks, this would entail unnecessary complexity since (assuming the line of text is known, which we do) a mark has at most two options: the subwords to its left and right.

Instead, we added a small DTW step whenever an ambiguous mark is encountered. It computes two rows of the grid twice, once with each assignment. For each column (subword in the transcript) the assignment with the smaller $\Delta(i,j)$ is selected. We store the choice for use in backtracking. When computing $d(i-1..i,j)$ for the next j , the system needs to know which assignment to use for b_{i-1} . The ambiguity can only be between two neighboring subimages, say b_i and b_{i-1} . The mark is initially assigned to b_i and an assignment to b_{i-1} is regarded as the alternate assignment.

3.4 Algorithm

The input is a binary image of multiple lines of text and an ASCII file containing the transcript. Line separation is achieved by projecting the image onto the vertical axis and regarding rows whose projection is zero as possible breaks. A more robust line finding algorithm can be used when needed. For each line, projections are used to compute upper and lower baselines (lines in Figure 4) and stroke width s , marks are located, and feature vectors are calculated. Marks from all lines are appended into one list, and features vectors from all lines are appended into another list V . The transcript is converted to an array L of lists of indices. Each list represents a subword and each index a letter shape. TD is run on V and L (Figure 8). Marks are reassigned if needed according to the mini-DTW step. The output of TD is the direction grid and another grid saying whether marks need to be shifted per cell.

In more detail, standard two-pass labeling of connected components (CCs) is performed, and very small CCs are discarded as noise. Each CC is initialized as a mark if its size is less than a threshold $t (=4s^2)$. This label is modified based on the lower baseline. If the CC is larger than t and its vertical range does not overlap with a small range above and below the baseline, it is classified as a mark. If its size is between $t/2$ and t , its vertical range overlaps with the baseline range and it extends above that range or if it has an alif shape, it is classified as non-mark. Upper and lower marks are separately merged into groups according to proximity (threshold $3s$), to treat multiple dots of the same letter as one mark in TD. Marks are associated with non-marks as in Section 3.3 and subimages are extracted. Ascenders and descenders are computed from the number of crossings at heights above/below the upper/lower baselines.

Initialization:

$\Delta(0,0) = \Delta_{1-1}(0,0) = \text{Distance}(0,0)$;
 $\Delta(0,j) = \text{Distance}(0,0..j)$;
 $\Delta(i,0) = \text{Distance}(0..i,0)$;
 $\Delta_{1-1}(i>0,0) = \Delta_{1-1}(0,j>0) = \text{Direction}(i>1,0) = \text{Direction}(0,j>1) = \text{INVALID}$;
 $\text{Direction}(0,1) = (i,j-1)$;
 $\text{Direction}(1,0) = (i-1,j)$ if able to split b_i

Body:

```

for i=1:n
{
    SetRow( $\Delta$ ,  $\Delta_{1-1}$ , Direction, i);
    if need to try alternate mark assignment for  $b_i$  and  $b_{i-1}$ 
    {
        Move mark from  $b_i$  to  $b_{i-1}$  and store alternate feature vector
        SetRow( $\Delta'$ ,  $\Delta_{1-1}'$ , Direction', i-1)
        SetRow( $\Delta'$ ,  $\Delta_{1-1}'$ , Direction', i)
        Move mark back from  $b_{i-1}$  to  $b_i$ 
        for j=1:m
            if  $\Delta'(i,j) < \Delta(i,j)$ 
                ( $\Delta$ ,  $\Delta_{1-1}$ , Direction) (i-1:i,m) = ( $\Delta'$ ,  $\Delta_{1-1}'$ , Direction') (i-1:i,m)
            }
        }
    }
}
End: Backtrack through Direction grid and force a move to (i-1,j-1) cell after a previous merge or split.
SetRow( $\Delta$ ,  $\Delta_{1-1}$ , Direction, i):
for j=1:m
{
     $\Delta_{1-1}(i,j) = \Delta_{1-1}(i-1, j-1) + \text{Distance}(i,j)$ ;
    if ( Ok to split  $w_j$  && ( $b_i$  and  $b_{i-1}$  on same line))
    {
        Using alternate mark assignment for  $b_{i-1}$  if indicated for j, set:
        tmp =  $\Delta_{1-1}(i-1,m) - \text{Distance}(i-1,j) + \text{Distance}(i-1..i,j)$ ;
    }
    ( $\Delta$ , Direction) (i,j) = (min, argmin) ( $\Delta_{1-1}(i,j)$ ,  $\Delta_{1-1}(i,j) - \text{Distance}(i, j-1) + \text{Distance}(i..j-1..j)$ , tmp);
}
}

```

Figure 8: Pseudo code for True Distance DTW, with m the number of subwords and n the number of subimages.

3.5 Limitation

Our method allows at most two boxes per subword or vice versa. Arabic subwords are shorter on average than English words, lessening the amount of oversegmentation expected. Also, the connected nature of Arabic allows one to use a simple segmentation method with relatively few segmentation errors. Thus, the two-to-one restriction is less significant than it would be for English text. It has not caused errors on this small test set and can be relaxed in future work when more flexibility is needed. By comparison, Rothfeder *et al.* allow at most three words per box [4]. They do not limit the number of boxes per word, but the restriction in probability calculation to single boxes with whole transcript words could prevent matches in oversegmentation since the subimages may not match the full word well.

4. RESULTS

We provide preliminary results on five images of the same text by different writers. One is shown in Figure 1a. The text contains 95 subwords and punctuation marks.

Table 1 shows the number of splits and merges detected, the number of alignment errors and the number of mislabeled subwords. In these examples, each error involved two subwords. The average number of alignment errors per image was 2.4. Two images had errors at the period in the text. Three had errors at the right parenthesis. These can be prevented with a shape matching method. Other errors include one due to poor baseline detection caused by nonlinear writing and one due to an unusually large diacritic. Two occur at the same subword and two at the same isolated letter (baa, second from left in Figure 2), and they will be investigated further.

Table 1: Performance on five images.

	Image 1	Image 2	Image 3	Image 4	Image 5
Number of true merges detected	0	2	4	2	2
Number of true splits detected	2	2	2	1	6
Number of alignment errors	1	2	4	3	2
Number of subwords mislabeled	2	4	8	6	4

Table 2: Mark reassignment. *Neither used mini-DTW: line separation error, shift too large to trigger method.

	Image 1	Image 2	Image 3	Image 4	Image 5
Number of marks reassigned correctly	1	5	2	3	1
Remaining mark assignment errors	1	0	2*	0	2

Table 2 shows results of the mini-DTW method for reassigning ambiguous marks. Most are correctly reassigned. Besides enabling word localization precise enough to include correct marks, this method enables correct scores to be propagated through the grid for a better overall alignment.

Note that handwriting styles varied significantly. The amount of difference in appearance between the text in Figure 1a and that in Figure 4 is approximately the amount observed across the images in the test set. The results presented are from preliminary tests, thus the small size of the test set. Future work will evaluate the method on large image sets in accordance with its intended domain. Future work will also measure robustness to variability regarding presence or absence of diacritical marks. However, when the method is used to create a database of annotated images, we expect the transcript to match the images in whether or not marks are used, so such variability would be limited.

Acknowledgments

The test dataset was collected by Faisal Farooq. Dr. Lorigo was supported by a DCI Postdoctoral Fellowship.

REFERENCES

- [1] C. Tomai, B. Zhang, and V. Govindaraju, "Transcript Mapping for Historic Handwritten Document Images," in *Proc. International Workshop on Frontiers in Handwriting Recognition*. Niagara-on-the-Lake, Ontario, Canada, 2002, pp. 413-418.
- [2] L. Lorigo and V. Govindaraju, "Off-line Arabic Handwriting Recognition: A Survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 28, pp. 712-724, 2006.
- [3] E. M. Kornfield, R. Manmatha, and J. Allan, "Text Alignment with Handwritten Documents," in *Proc. Document Image Analysis for Libraries*. Palo Alto, California, USA, 2004, pp. 195-211.
- [4] J. L. Rothfeder, R. Manmatha, and T. M. Rath, "Aligning Transcripts to Automatically Segmented Handwritten Manuscripts," in *Proc. Document Analysis Systems*, 2006, pp. 84-92.
- [5] L. Lorigo and V. Govindaraju, "Segmentation and Pre-Recognition of Arabic Handwriting," in *Proc. International Conference on Document Analysis and Recognition*. Seoul, Korea, 2005, pp. 605-609.