

Incorporating Contextual Character Geometry in Word Recognition

Hanhong Xue and Venu Govindaraju
CEDAR, SUNY at Buffalo, Buffalo, NY 14260
{hxue,govind}@cedar.buffalo.edu

Abstract

Contextual character geometry is the geometric information available only when a character presents in the context of a word. Such information includes a character's location and relative size in the entire word image, forming a bounding box of the character. The differences between the geometry of an image segment and the expected geometry of a candidate character are considered as additional features to refine the recognition of individual characters. A typical word recognizer based on over-segmentation and segment-combination is used to illustrate the use of these new features and experimental results have shown significant improvement of recognition accuracy, especially on large lexicons.

1 Introduction

Handwritten word recognition can be broadly divided into two categories [4]: (i) analytical approaches that identify individual characters in the word, and (ii) holistic approaches that recognize the entire word without breaking it down to characters. The analytical approaches build word recognition on character recognition and do not require training samples for all possible words. Therefore, they are capable of solving a wider range of problems especially when the lexicon size is large or no lexicon presents.

A general purpose word recognizer [2] has been designed on top of character recognition and employed in real-time applications such as postal address recognition [6] and bankcheck reading [1]. Its architecture is pictured in Figure 1 and outlined as follows.

1. *Preprocessing.* This step includes noise removal, slant correction and smoothing.
2. *Over-segmentation.* The image is over-segmented so that a character can have at most 4 segments.
3. *Segment-combination.* Nearby segments (up to 4) are combined to make a super-segment and the super-

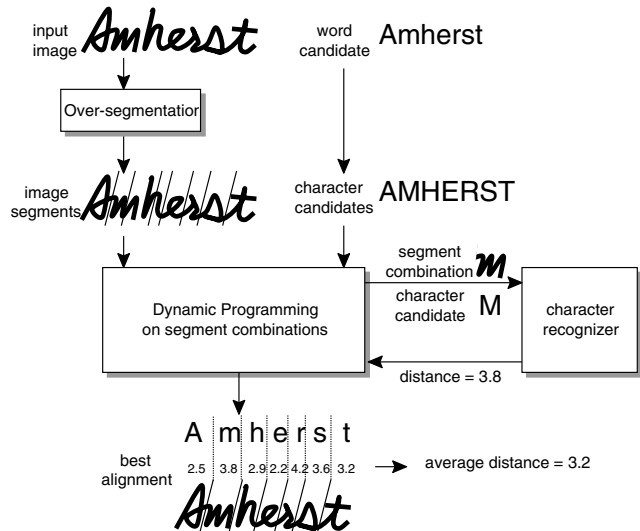


Figure 1. Architecture of a typical word recognizer built on top of a character recognizer

segments is divided into 3x3 sub-images. Aspect ratio, stroke ratio and features based on the distribution of eight directional slopes in the nine sub-images are extracted, resulting in a vector of 74 features for each super-segment.

4. *Character alignment.* A dynamic programming procedure is employed to find the best alignment of characters in the word. An underlying character recognizer is called to match super-segments against character candidates, returning the square distances between them.

As can be seen, the major advantage of this architecture is that any character recognizer can be easily plugged in to replace the current one. However, since it is purely based on a character recognizer which does not consider the context, i.e. the word, in which the character is matched, some unexpected alignment can be generated especially when the segmentation result is not optimal. For example, when the

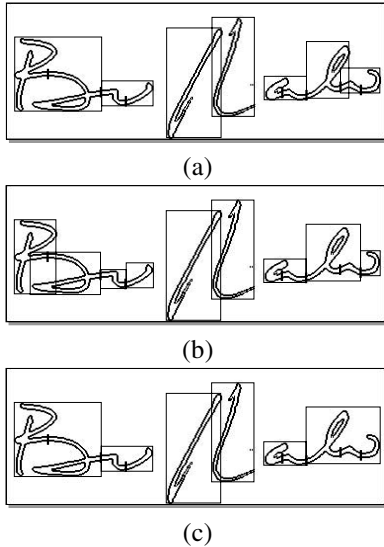


Figure 2. The best alignments for three word candidates: (a) "Buffalo", (b) "Fairfield" and (c) "Bville".

recognizer tries to recognize the image in Figure 2 given three word candidates, "Buffalo", "Fairfield" and "Bville", it ranks "Bville" before "Fairfield" and "Fairfield" before "Buffalo" by confusing "B" with "Fa", "l" with "i" and "lo" with "e". Though the sub-optimal segmentation, such as the "B" segmented in Figure 2, already accounts for some part of the mistake, one much more important reason is the lack of contextual character geometry in character recognition.

Contextual character geometry is the geometric information available only when a character presents in the context of a word. This information is important to word recognition because people usually write consistently, producing characters in similar sizes and in approximately even alignment. In character recognition, the recognizer is given a character image alone, without any information about its location and its relative size in the word image from which it is extracted. While in word recognition, such extra information is available and can be utilized by the underlying character recognizer to refine its judgment.

Therefore, this paper intends to incorporate contextual character geometry in word recognition and to illustrate the process by using the above-described word recognizer.

2 Contextual character geometry

Contextual character geometry is basically obtained from character bounding boxes which can be easily extracted by tracing the contours of super-segments. The absolute position and size of a character are not helpful since

the image size is variable in width and height, so we consider relative position and relative size instead.

2.1 Relative character width

We already know that a character "m" generally takes more space than a character "i" when they both present in the same word and perhaps the "m" is twice as much in width as the "i". Relative character widths are intended to capture such relation in width between characters and they are obtained by a special training procedure outlined as follows.

1. Assume equal expected relative width ($=1$) for all characters.
2. For each training word image,
 - (a) the recognizer is given an one-entry lexicon containing *only* the truth and produces the best alignment of characters for the truth;
 - (b) the unit width (in pixels) of the word image is calculated as image width (also in pixels but with spaces between characters removed) divided by the sum of expected relative widths of characters in the word;
 - (c) the relative width of a character is calculated as the character width (in pixels) divided by the unit width of the word.
3. Compute the expected relative width of each character by averaging the relative widths of all its occurrences.

This procedure is based on the Estimation and Maximization (EM) algorithm and Step 2 and Step 3 can be repeated until there is no significant change to the estimation. Table 1 gives expected relative character widths obtained after training. As can be seen, uppercase characters are generally wider than lowercase characters and "m" and "w" are much wider than "i" and "l". This conforms to the visual judgment of relative character widths by human eyes.

2.2 Relative character height

Relative character height is defined as the number of writing zones that a character is spanning vertically and it depends on the detection of reference lines. An algorithm of using linear regression on foreground vertical extrema [5] is applied to the detection of reference lines and, afterwards, image skew is corrected simply by shifting vertical coordinates. Since the height of each of the three writing zones may be different, the vertical position of a point in the image is actually translated according to the scale of the writing zone it is in. For example, the first "f" in Figure 3

	A	B	C	D	E	F	G	H	I	J	K	L	M
Width	1.6	1.8	1.3	1.7	1.4	1.8	1.9	1.7	1.7	1.9	1.7	1.7	2.2
Height	2.1	2.2	1.9	2.1	2.0	2.1	2.2	2.3	2.1	2.3	2.3	2.3	2.0
Top	1.1	1.1	1.2	1.1	1.2	1.2	1.1	1.2	1.1	1.1	1.1	1.1	1.3
	N	O	P	Q	R	S	T	U	V	W	X	Y	Z
Width	1.6	1.5	1.6	2.9	1.8	1.7	2.2	1.7	1.6	2.4	1.0	1.4	1.5
Height	1.9	1.7	2.2	2.1	2.3	2.2	2.0	1.9	2.0	2.0	2.0	2.5	2.6
Top	1.3	1.4	1.1	1.2	1.1	1.1	1.3	1.3	1.1	1.2	1.0	1.2	1.3
	a	b	c	d	e	f	g	h	i	j	k	l	m
Width	1.2	1.2	1.0	1.4	1.0	1.1	1.3	1.4	0.8	1.1	1.5	0.9	1.9
Height	1.2	1.9	1.2	1.9	1.2	2.2	1.9	1.9	1.0	2.0	1.8	1.7	1.2
Top	1.9	1.2	1.9	1.3	1.9	1.4	1.9	1.3	2.0	1.6	1.3	1.4	1.9
	n	o	p	q	r	s	t	u	v	w	x	y	z
Width	1.3	1.0	1.3	1.4	1.0	1.1	1.5	1.3	1.2	1.6	1.1	1.3	1.4
Height	1.1	1.1	2.0	1.9	1.2	1.2	1.7	1.1	1.2	1.3	1.7	2.0	1.5
Top	2.0	2.0	1.8	1.8	1.9	1.9	1.4	2.0	2.0	1.9	1.6	1.8	1.9

Table 1. Expected relative character width, height and top.

has relative top of 0.1, relative bottom of 3.0 and relative height of 2.9. Table 1 also gives expected relative character height and top obtained in the training procedure outlined as follows.

- For each training word image,
 - the recognizer is given an one-entry lexicon containing *only* the truth and produces the best alignment of characters for the truth;
 - reference lines are detected for the word image;
 - the relative height/top of a character in the word is obtained according to the character's vertical position in the image.
- Compute the expected relative height/top of each character by averaging the relative heights/tops of all its occurrences.

As shown in Table 1, the expected relative height of a character is always around 2 for uppercase and around 1 or 2 for lowercase depending on the presence of ascender and descender. The relative bottom of a character can be calculated as relative top plus relative height, thus omitted in the table.

3 Contextual geometric distances

We calculate contextual geometric distances based on contextual character geometry in matching an image segment against a candidate character.

Define functions $\bar{w}(\cdot)$, $\bar{h}(\cdot)$, $\bar{t}(\cdot)$ and $\bar{b}(\cdot)$ to return the expected relative width, height, top and bottom of a character and functions $w(\cdot)$, $h(\cdot)$, $t(\cdot)$ and $b(\cdot)$ to return the relative

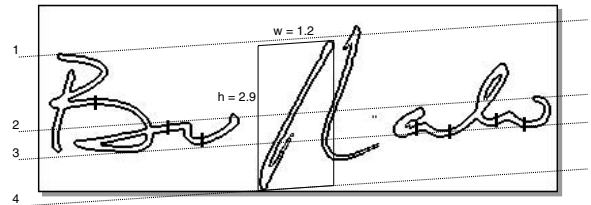


Figure 3. Contextual character geometry.

width, height, top and bottom of an image segment. Also define function $W(\cdot)$ to return the absolute width (in pixels) of an image segment.

Given a word image I and a word candidate $\alpha = (c_1, c_2, \dots, c_n)$ where c_i are characters, the unit width (in pixels) of the word U_α is calculated as the image width (with spaces removed) divided by the sum of expected relative widths of characters in the word, i.e.

$$U_\alpha = W(I) / \sum_{i=1}^n \bar{w}(c_i).$$

As a result, the unit width tends to be large for a short word and small for a long word. Then, the relative width of a segment s is measured as $w(s) = W(s)/U_\alpha$. It should be noticed that this measure depends on the current word candidate and must be re-calculated for each entry in the lexicon. However, for the relative height, top and bottom of a segment, they need to be calculated only once regardless what the current word candidate is.

When a super-segment s is matched against a character c , besides the square distance $d_0(s, c)$ returned by the underlying character recognizer, the following new square distances are calculated.

- Square distance between the segment’s relative width and the character’s expected relative width, normalized by the latter:

$$d_1(s, c) = \left(\frac{w(s) - \bar{w}(c)}{\bar{w}(c)} \right)^2.$$

- Square distance between the segment’s relative height and the character’s expected relative height, normalized by the latter:

$$d_2(s, c) = \left(\frac{h(s) - \bar{h}(c)}{\bar{h}(c)} \right)^2.$$

- Square distance between the segment’s relative top and the character’s expected relative top:

$$d_3(s, c) = (t(s) - \bar{t}(c))^2.$$

- Square distance between the segment’s relative bottom and the character’s expected relative bottom:

$$d_4(s, c) = (b(s) - \bar{b}(c))^2.$$

The new distance between super-segment s and character c is defined as

$$d(s, c) = \sum_{i=0}^4 d_i(s, c)$$

and it will replace $d_0(s, c)$ in the search for the best character alignment. Weighted sum has also been tried to combine distances but no significant difference is observed.

In the calculation of $d_0(s, c)$, each feature contributes at most 1 in distance because it takes a real value in range $[0, 1]$ by design. However, all these new distances take some value in range $[0, +\infty)$ and become overwhelming when some large deviation from the expected geometry is observed. Therefore, the use of contextual character geometry reduces the cases that some very long or very short super-segments are confidently recognized as characters, and, more importantly, characters are now required to have comparable sizes and to align more evenly. Since relative character height, top and bottom have been taken into account, they actually utilize the recognition power of ascenders and descenders as in holistic word recognition approaches [3].

Consider again the example in Figure 2. After penalty for uneven alignment is given to alignment (b) and penalty for ignoring ascenders and descenders to alignment (c), the new ranking becomes “Buffalo”, “Fairfield” and “Bville”, giving the correct top choice.

Unlike character top and bottom, character left and right are not used because the error in predicting them will be accumulated character by character and they become less reliable as the word length increases.

Lexicon size	Without geometry			With geometry		
	10	100	1000	10	100	1000
Top 1	97.03	91.39	79.58	97.80	92.79	84.08
Top 2	98.87	95.16	88.29	98.67	95.96	91.29
Top 5	99.43	97.36	93.29	99.40	97.93	95.30
Top 10		98.53	95.50		98.67	97.50
Top 20		99.07	97.10		99.10	97.80
Top 50		99.50	98.00		99.50	98.50

Table 2. Comparison of recognition rates on US postal words, without and with contextual character geometry.

4 Experimental results

The expected relative width, height and top of characters (Table 1) are obtained from a training set of 5735 US postal word images. The recognizer matches each image against its truth and produces the best alignment from which the contextual geometry of each character in the word is derived and averaged.

Then the recognizer incorporated with new contextual geometric features is tested on a different set of 3,000 US postal words, including city names, firm names, street names and state names. For each word image, lexicons of size 10, 100, and 1000 are supplied with the truth always included.

For comparison, we gather results of the recognizer without new features and those of the recognizer with new features in Table 2. It should be noticed that the results without new features are already improved by using Gaussian mixtures in character recognition [7] when compared against the very original results (96.80%, 88.23% and 73.80% top choice accuracy rate for lexicon size 10, 100 and 1000, respectively) published in [2]. We can see that the addition of new features further improves recognition accuracy in almost all cases, with 0.77%, 1.40% and 4.50% gain in the top choice for lexicon size 10, 100 and 1000, respectively.

To confirm the effectiveness of using contextual character geometry, we also test the recognizer on two sets of UK postal words: 1137 county names and 1720 post-town names. A fixed lexicon of size 1789 including all county names and all post-town names is always used. The results are given in Table 3 including cases when images are enhanced and not enhanced. It is very clear that contextual geometric features contribute around 3%-6% gain in the recognition accuracy.

Since the calculation of new features is very simple, it only increases recognition time by 3% in our implementation.

		Without geometry	With geometry	improvement
Counties	Top 1	66.75	70.36	+3.61
	Top 2	74.14	78.19	+4.05
Counties enhanced	Top 1	66.40	68.16	+1.76
	Top 2	74.23	77.04	+2.81
Post-towns	Top 1	65.87	71.05	+5.18
	Top 2	72.91	78.02	+5.11
Post-towns enhanced	Top 1	68.84	74.77	+5.93
	Top 2	76.22	80.64	+4.42

Table 3. Comparison of recognition rates on UK postal words, without and with contextual character geometry.

5 Conclusions

We improve word recognition by using contextual character geometry, which is only available when a character presents in the context of a word, to refine the match between an image segment and a character candidate. With contextual character geometry incorporated in to emphasize on utilizing ascenders and descenders and aligning character more evenly, a typical high-performance word recognizer achieves about 3%-6% improvement of recognition accuracy on large lexicons (1000-2000 entries) with negligible loss of speed.

References

- [1] G. Kim and V. Govindaraju. Bankcheck recognition using cross validation between legal and courtesy amounts. In S. Impedovo, P. Wang, and H. Bunke, editors, *Automatic Bankcheck Processing*, volume 28 of *Machine Perception and Artificial Intelligence*, pages 195–212. World Scientific, 1997.
- [2] G. Kim and V. Govindaraju. A lexicon driven approach to handwritten word recognition for real-time applications. *IEEE Transactions on PAMI*, 19(4):366–379, April 1997.
- [3] S. Madhvanath and V. Govindaraju. Holistic verification of handwritten phrases. *IEEE Transactions on PAMI*, 21(12):1344–1356, 1999.
- [4] R. Plamondon and S. Srihari. On-line and off-line handwriting recognition: A comprehensive survey. *IEEE Transactions on PAMI*, 22(1):63–84, January 2000.
- [5] P. Slavik and V. Govindaraju. An overview of run-length encoding of handwritten word images. Technical report, CEDAR, SUNY at Buffalo, 1999.
- [6] S. Srihari and E. Keubert. Integration of hand-written address interpretation technology into the united states postal service remote computer reader system. In *International Conference on Document Analysis and Recognition*, pages 892–896, Ulm, Germany, August 1997.
- [7] S. Tulyakov and V. Govindaraju. Probabilistic model for segmentation based word recognition with lexicon. In *International Conference on Document Analysis and Recognition*, pages 154–158, Seattle, Washington, September 2001.