

Phrase Based Direct Model for Improving Handwriting Recognition Accuracies

Faisal Farooq, Damien Jose and Venu Govindaraju
Center for Excellence in Document Analysis and Recognition (CEDAR)
Center for Unified Biometrics and Sensors (CUBS)
SUNY at Buffalo, Amherst, NY 14228
{ffarooq2,djose,govind}@cedar.buffalo.edu

Abstract

We propose a method for increasing word recognition accuracies by correcting the output of a handwriting recognition system. We treat the handwriting recognizer as a black-box, such that there is no access to its internals. This enables us to keep our algorithm general and independent of any particular system. We use a novel method for correcting the output based on a direct “phrase-based” system in contrast to traditional source-channel models. We report the accuracies of an in-house handwritten word recognizer before and after the correction. We achieve highly encouraging results for a large dataset.

1 Introduction

Any handwriting system, such as handwritten information retrieval, is composed of a preprocessing component, a handwriting recognition(HR) system and a search engine. For appreciable user experience, such components must be robust against noise as the input to the the next component in the pipeline is the output of the previous. Component systems are often developed independently and possibly by different groups. In most cases, the internals of one component are not accessible to the developers of the next component in the pipeline. In such cases, these components (e.g. HR) must be treated as black boxes where only their output is observed (Figure 1a). Unfortunately, the output of these systems is error-prone. Table 1 shows performance of pre-processing (line separation and word segmentation) and word recognition on 20 (control) document images of different writing styles [6]. As can be seen from the illustration, word recognition is the most difficult piece of the solution.

Thus, the performance of any usable system that involves a handwriting recognition system as one of its component will highly depend on its accuracy. Unfortunately, the state-of-art in handwriting recognition for a moderate lexicon of size $\sim 10k$ is $\sim 40\%$ [15]. Any improvement in

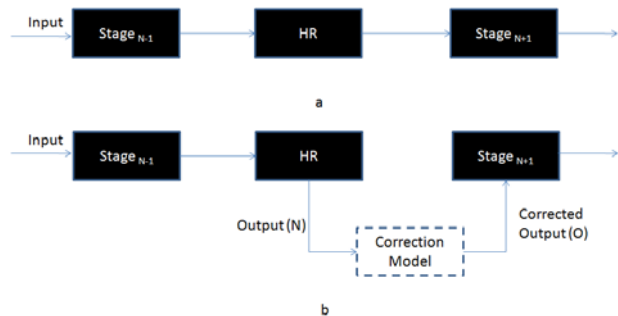


Figure 1. a) Typical NLP architecture b) NLP architecture with proposed correction model.

the performance of the recognizer will improve the overall user experience of the systems that use them as one of the components. In this paper we describe a novel “phrase-based” system for the correction of the output of a word recognizer that was developed in CEDAR [17]. We treat the recognizer as a black box that requires three inputs viz. i) a word image, ii) a lexicon and iii) an integer $n \leq 10$. As outputs we can acquire top- n hypothesis from the lexicon for that word image. We process the output from the HR and correct it using a model trained for that particular HR system. Thus, we add an additional module in the pipeline. Figure 1 describes in a block diagram the outline of the proposed system.

In section 2 we describe the related work. In section 3 we introduce our proposed technique. We describe our training data, test data and experimental settings in section 4. We describe our model in section 5 and the evaluation of our system performance in section 6. We finally conclude in section 7.

2 Related Work

There has been some work on post-processing of the OCR output. A general survey of the research in this area

Table 1. Performance of line separation, word segmentation and word recognition modules.

	writing style # of images	discrete 6	cursive 5	mixed 9	total 20
Line Separation	# of lines	118	101	210	429
	# of lines separated	114 96.6%	97 96.0%	198 94.3%	409 95.3%
Word Segmentation	# of words	641	692	1,427	2,760
	# of words segmented	597 93.1%	631 91.2%	1,379 96.6%	2,607 94.5%
Word Recognition	top 1	432 72.4%	268 42.5%	750 54.4%	1,450 55.6%
	top 10	541 90.6%	459 72.7%	1,081 78.4%	2,081 79.8%

can be found in [9]. Jones et al. [4] describe a multi-pass OCR post-processing system which carries out individual word corrections, combined edit distance corrections and bigram probability based correction in different passes. Perez-Cortes et al. [11] use a stochastic finite state machine to test hypothesis of words. If the machine accepts the word, then no correction is made, otherwise the smallest set of transitions that could not be traversed show the most similar string in the model. Pal et al. [10] describe a method for OCR error correction of Devanagiri script using morphological parsing. Some of these techniques are highly dependent on the language and use features that are specific to the language in question, and so do not scale to a different language. Secondly, they have been applied on machine print OCR that conventionally use “character-models” as opposed to the handwriting recognition systems that follow the word-based “multiple choice” paradigm. Thus, training the character confusion matrices is not straight forward. The method we propose here is trainable. Although we only show results on the in-house recognizer, the technique is adaptable to other recognizers and even other scripts where training data is available. Okan et al. [8] introduce a generative probabilistic character OCR model that describes an end-to-end process in a noisy channel framework. The technique they describe is implemented as a cascade of finite state transducers. We propose a different approach to this problem. We cast this problem as a simplified “translation” task where the source language is the error-full HR output and the target language is the corrected output (truth).

3 Proposed Technique

The erroneous output from the OCR can be thought of as an output from a noisy channel problem[1]. It can be considered as a black-box through which the signal (truth) when passed gets corrupted and emerges out as

the degraded output. We can learn the corruption models (correction model in the direct model framework) and then use this knowledge to correct the output. We leverage the research in Statistical Machine Translation (SMT) systems by casting this correction task as a simplified translation task. We propose a direct phrase-based translation approach that has proven effective in machine translation (e.g., [14], [7]). The phrase correction model has various advantages over a simple word model with post-processing. One, that this model adapts to the behaviour of the word recognizer instead of only post processing the output based on linguistic and dictionary rules. Second, due to modelling longer “ngrams”, it captures the context in which the errors were made and thus yields different correction options based on different context. This strength can be utilized to model the effects of co-articulation in handwriting. In addition to this, as we shall show in Section 5, we could potentially model segmentation errors as well.

As is done in a typical SMT system, given *sentence pairs* in the source (*foreign*) and the target (*english*) languages, we first align the source and target words. Word alignments can be obtained using one of several techniques such as HMM alignments [16], Maximum Entropy alignments [3], or Maximum-Posterior alignments [18]. Once a sentence pair is word-aligned, we extract phrase pairs (also referred to as block in the statistical machine translation literature) in a manner similar to [14]. These phrases are then used in combination with an *n*-gram language model to translate the source language into the target language.

We model the correction as a simplified translation task. The source language in our case is the OCR output (error-full) and the target language is the truth (correct). Since there is no word re-ordering as is the case of a typical MT task, the alignments can be obtained eas-

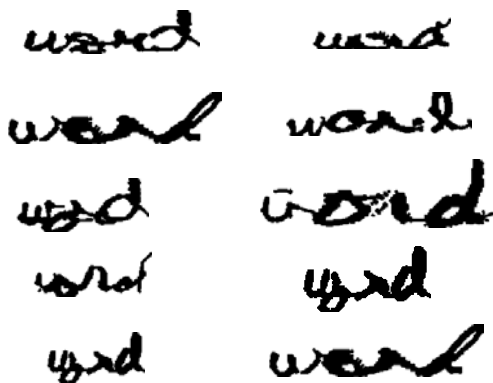


Figure 2. Samples generated for the same word 'word'.

ily by minimum edit distance alignments or monotone HMM alignments. We used a simple dynamic programming minimum edit distance alignment.

4 Training and Test Data

Due to lack of large truthed handwritten datasets, we automatically generated images from the english text used in the Conference on Computational Natural Language Learning (CONLL 2003) shared task (<http://www.cnts.ua.ac.be/conll2003/ner/>). Since our model required the learning the behaviour of the recognizer and the most common corrections for the mistakes that it makes, it is critical that the same input to the recognizer is provided in various styles and forms. This was not possible from any existing sentence database. Thus, we had to rely on generating data that would simulate handwriting. In case such databases are available, they can easily be utilized or even added to the current model with absolute ease. The handwriting is generated from text by concatenating character templates generated by the Blum's Medial Axis Transform [2] from actual handwritten samples. This is followed by character auto-scaling, automatic baseline determination, ligature modelling, ligature joining, skeleton thickening and smoothing. The visual analysis and accuracies of the recognizers on this data support our simulation technique for generating handwriting samples. Details of this method can be found in [13]. In order to model the multiple writer scenario we sampled the characters at random from a set. Thus, during generation of a word e.g. 'word' at different positions in the data we sampled uniformly for i from different handwritten character samples $\{w_i, o_i, r_i, d_i\}$. Figure 2 shows the different samples generated for the same word.

We split the data into a training set (75%) and a test set

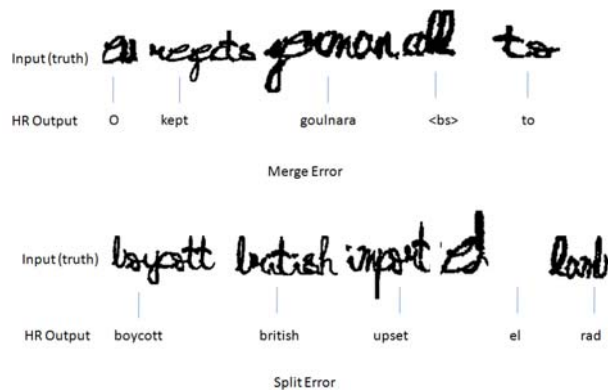


Figure 3. Examples of merge/split errors.

(25%). The training data consisted of 127690 words from 1055 documents. Similarly, the test data consisted 40789 words from 340 documents. Table 2 gives the details of the two sets.

The in-house handwriting recognizer used for recognition of these images is a lexicon driven HMM based word model recognizer [17]. The recognizer takes input as a lexicon and the word image. A lexicon of size $\sim 21K$ was generated from the training data. Thus, during the recognition of the training data we used a closed lexicon such that the truth of the image being recognized was present in the lexicon. We used the same lexicon for recognizing the test data and observed that it covered 95% of the words and the rest were treated as OOVs (Out Of Vocabulary).

5 Phrase-based Direct Model

In order to extract the phrases from the training data we align the truth and the OCR output character streams obtained above. Since there is no word re-ordering, we use a dynamic programming based minimum edit distance alignment. The edit distance was calculated using the regular Levenshtein edit distance formula with equal cost for deleting, inserting or substituting a unit. Using this, one-to-one alignments are obtained for the paired truth and output datasets. An unaligned character in the truth models a character deletion. Similarly an unaligned character in the OCR output models an insertion. Moreover, an unaligned *blank* in the truth models a merge error where the OCR merges two words and an unaligned *blank* in the OCR output models a word split. Figure 3a shows an example where two words were merged by the OCR into one and the *blank* is unaligned. Similarly, Figure 3b shows an example where a word was split and a word inserted (*el*). This time the *el* on the OCR side is unaligned.

Using this alignment we extract all phrase pairs of lengths 1 through n ($n = 6$ in our case) words. Figure 4 shows example phrases extracted from the training data

Table 2. Details of the training and test data.

	Documents	Lines	Words
Training	1055	7878	127690
Test	340	2540	40789
Total	1395	10418	168479

with the possible corrections alongwith their correction probabilities for a given word in the left column.

Recognized words	Hidden words	Probability
pitcher	pitcher	1.00
pascolo	financially	0.40
	pascolo	0.20
	speculates	0.20
	poisonous	0.20
notation	protection	0.50
	invitation	0.40
	motivation	0.05
	notation	0.05
experts	experts	0.88
	expired	0.13
updated	injunction	0.40
	uprooted	0.20
	infrastructural	0.40

Figure 4. Phrases extracted from training data.

5.1 Decoding

To correct the OCR output for a given test sentence, we “translate” this sentence by decoding using two weighted components - the phrases obtained above and the language model. We trained a trigram language model on the Reuters data (<http://trec.nist.gov/data/reuters/reuters.html>) using the CMU language modelling toolkit [12].

The decoder is a trigram viterbi decoder. For the correction problem, the decoding is carried out in a monotone fashion such that words at the beginning of the sentence must be decoded first. Decoding hypotheses are extended in cardinality-synchronous fashion such that at any given time step n , only hypotheses that cover exactly n words (i.e., of cardinality n) are extended. For example, at time step 0 only hypothesis of cardinality 0 (in the special initial hypothesis) are extended. The decoding is continued until all words in the sentence are covered. The final decoding is the hypothesis of cardinality l with the least cost, where l is the length of the source text. More formally, the final decoding e for the source f is the one that satisfies the following equation:

$$\hat{e} = \arg \max_e [-w_{ph} \times \log_{10} P(e|f) - w_{lm} \times \log_{10} P(e)]$$

where $P(e)$ is the trigram character language model

Recognized words	Hidden words	Log probability
om	the	0
official	officials	7.89912
sort	said	15.2251
om	the	24.5661
accord	amanda	54.0033
had	had	86.9505
Gerg	seized	143.344
tra	two	239.679
Roberta	kalashnikov	389.403
dealt	assault	639.717
aples	wales	1043.42
ad	and	1679.26
disappears	disappointed	2728.21
	attackers	50.6336
	batter	99.1431
	leaving	152.173
	nagatsuka	491.44
	decades	649.419
	maybe	1043.54
	cash	1691.52
	disappeared	2721.85
	antara	57.4298
	stated	97.6456
	tsang	152.516
	consistent	650.351
	rifles	1036.11
	seed	1691.66
	disappears	2731.83
	freeze	156.673
	nationwide	650.707
	engineers	1046.72

Decoded Sentence:
the official said the attackers had seized two kalashnikov assault rifles and disappeared

Figure 5. Example of Sentence Decoding.

probability and $P(e|f)$ is the phrase-based direct model. We used equal weights w_{ph} and w_{lm} ($w_{ph} + w_{lm} = 1$) for both mixture components. However, in general these weights need refinement (can also be trained) based on the prior evaluation of the word recognizer. For example, a recognizer with better accuracies would require a slightly larger w_{ph} and vice-versa.

The drawback of this approach, however, would be that in order to get features for all possible words, we need much more data. In order to overcome this problem, we adopted the following procedure:

1. Extract phrasepairs of length=1 pairs $P = \bigcup p_w$, where p_w is a phrase pair containing source word w and all target word(s) from the training data.
2. Get the vocabulary V_t of the test set.
3. $\forall w_i \in V_t$, if $\exists p_{w_i} \in P$ such that $count(p_{w_i}) < T$ then $P = P \cup \{w_i, recog(w_i)\}$, i.e. add w_i and $recog(w_i)$ to P .

Given the test lexicon V_t , the function $recog(w_i)$ yields top-10 words $w'_i \in V_t$ from the word recognizer. This technique was motivated by the back-off language modelling [5]. Thus, wherever, the unigram is not found in the training model ($T < 1$), we simply utilize the top-10 unigram outputs from the word recognizer for that word image.

Table 3. Accuracy of the HR with and without correction.

	Raw	Top-10+3-gram LM	Corrected
Accuracy	8.1%	13.5%	71.3%

6 Results

Figure 5 shows an example of sentence decoding. The matrix shows the correction model options for the “observed” output from the recognizer with the corresponding probabilities. The arrows show the decoding path in conjunction with the language model. The final output sentence is shown at the bottom of the figure.

Table 3 shows the accuracies of the recognizer on the test set before and after the correction. Before correction, both raw accuracies as well as accuracy with post processing using top-10 choices in conjunction with a trigram language model are reported. We observe that there is a considerable increase in the accuracy after the correction. Since, there has been no related work in handwriting domain, the increase in accuracies can only be compared with methods from machine print OCRs. The relative improvement due to our method is far superior than many methods as in [8, 10, 11].

7 Conclusion

We have presented a method for improving the accuracies of a handwriting recognizer. This method is statistical and highly adaptable. It can be adapted to machine print OCRs and other Handwriting Recognizers with absolute ease. In fact, due to its statistical nature, given the training data it can be extended to other scripts as well. The increase in the accuracy can be reflected in the final outcome of any system that utilizes these components.

References

- [1] L. Bahl, F. Jelinek and R. Mercer, “A Maximum Likelihood Approach to Continuous Speech Recognition”, *IEEE Transactions on PAMI*, 5(2):179–190, 1983.
- [2] H. Blum, “A Transformation for Extracting New Descriptors of Shape”, *Models for the perception of Speech and Visual Form*, MIT Press, 1967, pp 362–380, Cambridge, MA.
- [3] A. Ittycheriah and S. Roukos, “A Maximum Entropy Word Aligner for Arabic-English Machine Translation”, *Proceedings of the Human Language Technology Conference (HLT-NAACL)*, 2005, Vancouver, Canada.
- [4] M. Jones, G. Story and B. Ballard, “Integrating multiple knowledge sources in a Bayesian OCR post-processor”, *International Conference on Document Analysis and Recognition*, 1991, pp 925–933, St. Malo, France.
- [5] S. M. Katz, “Estimation of probabilities from sparse data for the language model component of a speech recognizer”, *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 35(3):400–401, 1987.
- [6] G. Kim, V. Govindaraju and S. Srihari, “Architecture for handwriting recognition systems”, *International Journal of Document Analysis and Recognition*, 2(1):37–44, 1999.
- [7] P. Koehn, F. Och and D. Marcu, “Statistical Phrase-Based Translation”, *Proceedings of the Human Language Technology Conference (HLT-NAACL)*, 2003, Edmonton, Canada.
- [8] O. Kolak, W. Byrne and P. Resnik, “A Generative Probabilistic OCR Model for NLP Applications”, *Proceedings of the Human Language Technology Conference (HLT-NAACL)*, 2003, pp 55–62, Edmonton, Canada.
- [9] K. Kukich, “Techniques for automatically correcting words in text”, *ACM Computing Surveys*, 24(4):377–439, 1992.
- [10] U. Pal, P. Kundu and B. Chaudhuri, “OCR error correction of an inflectional Indian language using morphological parsing”, *Journal of Information Science and Engineering*, 16(6):903–922, 2000.
- [11] J. Perez-Cortes, J. Amerngual, J. Arlandis and R. Llobet, “Stochastic error-correcting parsing for OCR post-processing”, *International Conference on Pattern Recognition*, 2000, pp 4405–4408, Barcelona, Spain.
- [12] R. Rosenfeld, “The CMU Statistical Language Modeling Toolkit and its use in the 1994 ARPA CSR Evaluation”, *Proceedings of the ARPA Spoken Language Systems Technology Workshop*, 1995, San Mateo, California.
- [13] A. O. Thomas, A. Rusu, S. Mukund and V. Govindaraju, “Non-Writer Specific Handwriting Generation for the CAPTCHA Application”, *Proceedings of the IEEE Western New York Image Processing Workshop*, 2007, Rochester, NY.
- [14] C. Tillmann and H. Ney, “Word Re-ordering and a DP Beam Search Algorithm for Statistical Machine Translation”, *Computational Linguistics*, 29(1):97–133, 2003.
- [15] A. Vinciarelli, S. Bengio and H. Bunke, “Offline recognition of unconstrained handwritten texts using HMMs and statistical language models”, *IEEE transactions on Pattern analysis and Machine intelligence*, 26(6):709–720, 2004.
- [16] S. Vogel, H. Ney and C. Tillmann, “HMM Based Word Alignment in Statistical Translation”, *Proceedings of the 16th International Conference on Computational Linguistics (COLING)*, 1996, pp 836–841, Copenhagen, Denmark.
- [17] H. Xue and V. Govindaraju, “Hidden Markov Models Combining Discrete Symbols and Continuous Attributes in Handwriting Recognition”, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(3):458–462, 2006.
- [18] B. Zhao, N. Ge and K. Papineni, “Inner-Outer Bracket Models for Word Alignments using Hidden Blocks”, *Proceedings of the Human Language Technology Conference (HLT-NAACL)*, 2005, Vancouver, Canada.