

# Retrieval-Based Language Model Adaptation for Handwritten Chinese Text Recognition

Shuying Hu<sup>1</sup>, Qiufeng Wang<sup>2\*</sup>, Kaizhu Huang<sup>3</sup>, Min Wen<sup>1</sup> and Frans Coenen<sup>4</sup>

<sup>1</sup>Department of Applied Mathematics, Xi'an Jiaotong-Liverpool University, 111 Ren'ai Road, Suzhou, 215123, Jiangsu, China.

<sup>2\*</sup>School of Advanced Technology, Xi'an Jiaotong-Liverpool University, 111 Ren'ai Road, Suzhou, 215123, Jiangsu, China.

<sup>3</sup>Data Science Research Center, Duke Kunshan University, No.8 Duke Avenue, Kunshan, 215316, Jiangsu, China.

<sup>4</sup>Department of Computer Science, University of Liverpool, Liverpool, L69 3BX, United Kingdom.

\*Corresponding author(s). E-mail(s): [Qiufeng.Wang@xjtlu.edu.cn](mailto:Qiufeng.Wang@xjtlu.edu.cn);

Contributing authors: [Shuying.Hu20@student.xjtlu.edu.cn](mailto:Shuying.Hu20@student.xjtlu.edu.cn);

[kaizhu.huang@dukekunshan.edu.cn](mailto:kaizhu.huang@dukekunshan.edu.cn); [Min.Wen@xjtlu.edu.cn](mailto:Min.Wen@xjtlu.edu.cn); [Coenen@liverpool.ac.uk](mailto:Coenen@liverpool.ac.uk);

## Abstract

In handwritten text recognition, compared to human, computers are far short of linguistic context knowledge, especially domain-matched knowledge. In this paper, we present a novel retrieval-based method to obtain an adaptive language model for offline recognition of unconstrained handwritten Chinese texts. The content of handwritten texts to be recognized is varied and usually unknown a priori. Therefore we adopt a two-pass recognition strategy. In the first pass, we utilize a common language model to obtain initial recognition results, which are used to retrieve the related contents from Internet. In the content retrieval, we evaluate different types of semantic representation from BERT output and the traditional TF-IDF representation. Then, we dynamically generate an adaptive language model from these related contents, which will consequently be combined with the common language model and applied in the second-pass recognition. We evaluate the proposed method on two benchmark unconstrained handwriting datasets, namely CASIA-HWDB and ICDAR-2013. Experimental results show that the proposed retrieval-based language model adaptation yields improvements in recognition performance, despite the reduced Internet contents hereby employed.

**Keywords:** Recognition, Handwritten Chinese Text Recognition, Internet Content, Information Retrieval, Language Model Adaptation

# 052 1 Introduction

053  
 054 Documents comprising handwritten or printed  
 055 characters are one of the most popular tools for  
 056 our communication and archiving [1]. To dig-  
 057 itize these documents, optical character recog-  
 058 nition (OCR) has been widely researched and  
 059 applied [1, 2]. Solid progress has been made in  
 060 many areas, e.g. from isolated character recogni-  
 061 tion to character string recognition, from printed  
 062 character recognition to unconstrained handwrit-  
 063 ing recognition, and from documents with clear  
 064 background to scene text recognition with com-  
 065 plex background. While related tasks are getting  
 066 more and more complicated, recent advancement  
 067 in OCR has lead to great success in real applica-  
 068 tions. Chinese handwriting recognition has been  
 069 an important branch of OCR since 1970s [3,  
 070 4]. Powered by deep learning, handwritten iso-  
 071 lated Chinese character recognition has achieved  
 072 tremendous advance [4, 5, 6, 7, 8]. Remarkably, the  
 073 reported accuracy rate can even be higher than  
 074 that of human recognition: 97.64% was reported  
 075 in [7] whilst human only get the accuracy of  
 076 96.13%. Nevertheless, automated unconstrained  
 077 handwritten Chinese text recognition still remains  
 078 unsatisfactory and actually far behind human  
 079 recognition, since human can effectively lever-  
 080 age sufficient linguistic context knowledge [4, 5,  
 081 6, 9, 10]. Concretely, there are huge challenges  
 082  
 083  
 084  
 085  
 086  
 087  
 088  
 089  
 090  
 091  
 092  
 093  
 094  
 095  
 096  
 097  
 098  
 099  
 100  
 101  
 102

in automated unconstrained handwritten Chi-  
 nese text recognition including the low-quality  
 of text images, flexible handwriting styles, and  
 unusual topics with possible specific terminologies,  
 and shortage of linguistic context. Figure 1 illus-  
 trates one example from the benchmark dataset of  
 CASIA-HWDB with the cursive handwriting style  
 and specific terminologies, where the airport name  
 ‘白’ (white) is often incorrectly recognized as ‘自’  
 (self).

截止到现在下午17时左右机务几天下来故障原因不在旅客滞  
 留广州,为做好航班延误情况下的航班保障工作,白云机务  
 这部地区正启动航班不正常预案,客舱部临时成立了航班  
 不正常指挥小组,组建了由值机员、乘务员担任的航班不正常  
 处理小组,专人负责航班延误信息的发布、旅客信息的传  
 达,旅客后续工作安排跟进与航空公司相关部门沟通协调工  
 作,设立了由新员工组成的引导岗位,负责带领旅客由值机大  
 厅到航班不正常柜台到西二餐外、酒店候车。此外,白云  
 机务客舱部提醒提醒旅客,航班延误情况下,请旅客注  
 意观察指示牌,注意航班信息;办理登机手续时可选择同航  
 司工作人员航班延误情况;由航班更改调整早班处理延误;航班  
 延误听从机务工作人员的统一安排。

**Fig. 1:** Example of a handwritten cursive Chinese text page where the airport name is often incorrectly recognized.

Generally, handwritten Chinese text recogni-  
 tion (HCTR) is a sequence pattern recognition  
 problem, which can be translated to searching for  
 the optimal path in a complicated candidate lat-  
 tice by over-segmentation under certain path eval-  
 uation criterion [9, 10]. Inspired by how humans  
 read texts, handwritten text recognition usually

utilizes language models to represent linguistic context knowledge, which characterizes the statistical dependency between characters and assigns the prior probability of a sequence of characters. Unfortunately, in lack of sufficient linguistic context knowledge, current automated methods are still far away from humans' ability, thus limiting the ceiling point especially for HCTR.

Language models (LMs) play a very important role in HCTR. There are a fast-growing body of methods which have explored how to apply LMs in HCTR recently [9, 13, 14]. In some literature, LM is merely exploited as post-processing for correcting recognition errors [23, 24]. Despite its simplicity and plug-and-play property, such practice may actually limit LM's full play in the learning process. Instead, as indicated in many investigations, LM needs to be seamlessly integrated in the path search so that the learning process can be well guided to the optimum [9, 22].

On the other hand, there are sufficient and diverse contents on Internet, where linguistic context might be mined. Indeed, Internet contents have been exploited in many fields, e.g. never-ending machine learning [15, 16], image recognition [17], speech recognition [18], postcards recognition [19], and scene text recognition [20, 21]. However, to the best of our knowledge, all of these works either directly utilize the contents or simply correct the recognized result as post-processing. There have been rare studies investigating

how to integrate Internet contents with LM for HCTR.

However, if the domain of LM does not match the handwritten texts to be recognized, its effect could be limited. In fact, one can usually observe such phenomenon, since characters or words usually enjoy various statistics in the text corpus of different domains [26]. To deal with this mismatch problem, language model adaptation has been widely investigated. This is particularly the case in unsupervised adaptation where no prior information exists about the domain of handwritten texts [14]. In [14], a large text corpus was collected with different domains downloaded from the Internet in advance. A set of LMs were then trained for those pre-defined domains. However, these LMs are unchanged for all handwritten texts, which is not flexible.

In this paper, we propose a method to dynamically retrieve the related content from Internet for HCTR, then train an adaptive LM from this content, which is integrated in the whole recognition process. For the retrieval of the related content, we adopt the Transformer-based language model BERT [58], which maps each word to a semantic space with the context information. Since we do not have any prior information for the handwritten texts to be recognized, we utilize a two-pass recognition strategy. In the first-pass recognition, we apply a common language model to get an initial text result. We then retrieve the related

103  
104  
105  
106  
107  
108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153

154 contents based on this transcription and build  
155 a domain-matched  $N$ -gram model, which will be  
156 integrated with the common LM in the second-  
157 pass recognition. The retrieval of the related  
158 contents is based on a criterion that measures  
159 the similarity between the recognized text and  
160 contents from the Internet corpus.

165 The rest of this paper is organized as follows:  
166 Section 2 reviews some related work on handwrit-  
167 ting recognition and language models; Section 3  
168 gives an overview of our HCTR system; Section 4  
169 describes in details the proposed retrieval-based  
170 language model; Section 5 presents the experimen-  
171 tal results, and Section 6 concludes this paper with  
172 final remarks.

## 180 2 Related works

181 In the context of retrieval-based handwritten Chi-  
182 nese text recognition from internet contents, much  
183 work has contributed to the related issues, includ-  
184 ing internet-based pattern recognition, language  
185 model, and handwritten Chinese text recognition.

186 In the following, we will give an overview of these  
187 previous work.

### 194 2.1 Internet-based pattern 195 recognition

196 An LM is generally developed using a text cor-  
197 pus of millions or even billions of sentences, so  
198 crawling sources of online text data is a simple  
199

way for building LMs. Given the popularization of  
the Internet in recent decades, some researchers  
have considered to utilize the idea of web search  
on Internet contents to improve the performance  
of pattern recognition tasks. These methods can  
be roughly categorized into three groups accord-  
ing to how they utilize the Internet: (i) collecting  
and labelling the data set, (ii) training the classi-  
fier, and (iii) integrating in the recognition system  
directly. In the first category, people have used the  
Internet to work on the data set cooperatively. For  
example, Russel et al. designed a web-based tool  
for Image annotation [27], and the famous large-  
scale ImageNet image dataset was collected and  
labelled using Internet crowdsourcing [28, 29]. The  
second category of methods utilizes the Internet to  
get a large set of related data, which is then used  
to train the classifier. Fergus et al. [17] automat-  
ically learned the categories of objects from the  
images retrieved on Google, and Hanzk et al. [30]  
utilized Internet contents in the transfer learn-  
ing between different domains for action-model  
learning. The last category is to use the Internet  
content as the linguistic context in the recognition.  
For example, Nishizaki et al. [18] utilized Internet  
content to correct the errors in the post-processing  
of speech recognition, Chen et al. [31] applied  
Internet content to language model adaptation,  
and White law et al. [32] explored Internet content  
for spell checking and autocorrection. All of these

work has promoted the development of pattern recognition, and improved the performance.

In the area of OCR, much work has also been reported on exploring Internet content to improve recognition performance, and most of these work used the idea of the aforementioned third category. Clemens et al. [19] utilized Internet content to verify and correct text recognition for post-cards. Bassil et al. [33] applied the autocorrection function of Google in the post-processing of OCR. Donoser et al. [20, 21] re-scored candidate characters based on the Internet retrieval results in scene text recognition, where the authors assumed that the web search results of correct recognition text were many more than those of the wrong results. However, all of these work only leveraged the Internet contents in the post-processing of OCR instead of being integrated in the whole recognition process. Recently, Oprean et al. [34, 35] utilized the contents of Wikipedia to construct the dictionary dynamically to overcome the Out-of-Vocabulary (OOV) issue in handwritten English word recognition, and successfully extended it to the handwritten English text recognition using a deep learning framework [36].

## 2.2 Language model

LMs have been widely used in speech recognition, handwriting recognition, and machine translation [25, 37]. The  $N$ -gram model is the most popular LM in handwriting recognition, which

characterizes the statistical dependency between the neighbour  $N$  characters or words [9, 22, 23, 24, 38, 39]. However, the  $N$ -gram model usually has two issues: zero-probability for the unseen  $N$ -gram sequences (though various smoothing methods have been proposed) and the local context limit (considering the moderate model size and decoding time where  $n$  is usually 2 or 3). In recent decades, the neural network based language model was developed, quickly to overcome the zero-probability issue due to the distributed representation of all words [40, 41, 42, 43, 44]. Some of these models have been successfully applied in the HCTR [10]. For the issue of local context limit, many topic model based language models have been proposed [45, 46, 47], and Xie et al. [48] developed an implicit LM to integrate the global linguistic context in the online HCTR.

Language model adaptation (LMA) is an important technique used to adapt a common language model to match the domain of each recognition task, which can be categorized into supervised LMA and unsupervised LMA [26]. Supervised LMA assumes that the domain of the recognition task is known in advance. Therefore, a large set of related texts can be obtained to train a domain-matched LM [49]. However, the domain information is usually unknown a priori, which requires unsupervised LMA. The basic idea is to use a common LM to get an initial recognition result, which is then used to either search related

205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215  
216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255

256 text to train an adaptive LM [31] or learn the  
 257 weights to combine various LMs to get a balanced  
 258 LM [50]. Most work on LMA has been conducted  
 259 in speech recognition and natural language pro-  
 260 cessing (NLP), while there are few examples where  
 261 LMA has been used for HCTR. Recently, our  
 262 previous work [14] has investigated three unsuper-  
 263 vised LMA methods based on a pre-defined LM  
 264 set, validating the effectiveness of LMA in HCTR.

## 272 **2.3 Handwritten Chinese text** 273 **recognition**

274 HCTR has attracted lots of attention since  
 275 two databases (HIT-MW [51] and CASIA-  
 276 HWDB [12]) were released and two competitions  
 277 were organized at ICDAR 2011 [5] and ICDAR  
 278 2013 [52]. HCTR has achieved great progress  
 279 in recent years [9, 10, 38, 39, 48, 53, 54, 63,  
 280 64]. The approaches adopted can be divided  
 281 into two categories: over-segmentation based and  
 282 segmentation-free.

283 In the over-segmentation based approaches,  
 284 the text line is over-segmented into a sequence of  
 285 primitive segments, each corresponding to a char-  
 286 acter or a part of a character. Then a candidate  
 287 segmentation-recognition lattice is constructed by  
 288 combining neighboring segments to be recognized  
 289 via the character classifier, where each candi-  
 290 date path represents one segmentation-recognition  
 291 result. This approach can take advantage of the

character shape and overlapping characteristic to  
 better separate the characters at their boundaries.  
 Most reported HCTR work has been based on this  
 framework [9, 10, 38, 39]. However, a high recall  
 over-segmentation algorithm is usually necessary,  
 which might be however difficult for touching  
 characters [55]. To improve the recall of segmenta-  
 tion, Wu et al. [10] proposed an over-segmentation  
 algorithm based on the convolutional neural net-  
 work (CNN). To avoid difficulties of finding exact  
 boundaries in character-level annotation, Wang  
 et al. [11] proposed a weakly supervised learn-  
 ing method to optimize the character classifier  
 by string-level training such that strong annota-  
 tion is not required. The recent works [62, 63]  
 formulated the character segmentation as a char-  
 acter detection based on a fully convolutional  
 network, which demonstrated great potentials for  
 the segmentation-based HCTR.

The segmentation-free approach is also called  
 the implicit segmentation approach<sup>1</sup>, which first  
 uses the technique of a sliding window to split a  
 text image into a sequence, then uses the Hid-  
 den Markov Model (HMM) or Recurrent Neural  
 Networks (RNN) based models with Connectionist  
 temporal classification (CTC) or attention strate-  
 gies to get the recognition result. Most reported  
 work of English text recognition is based on  
 this approach [56]. Su et al. [53] was the first

---

<sup>1</sup>Accordingly, the over-segmentation is also called as explicit segmentation.

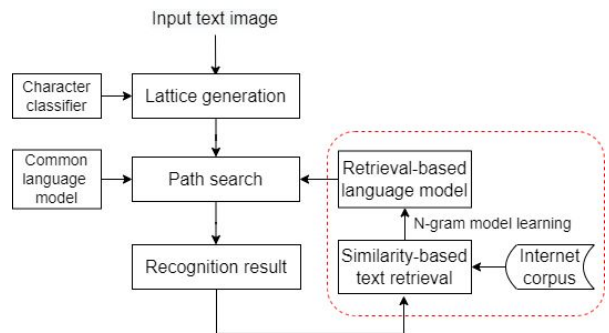
to apply this approach to HCTR, and Wang et al. [54] estimated the observation probability by using a deep CNN model instead of the Gaussian mixture model (GMM). Recently, deep neural network based approaches have been widely used in the segmentation-free recognition, namely, RNN, especially Long-Short Term Memory (LSTM) and Bidirectional LSTM. The work in [57] proposed a multi-dimensional LSTM and CTC framework for end-to-end HCTR. Additionally, the work reported in [48] utilized the fully connected RNN model for the online text recognition. Wang et al. [64] further proposed a writer-aware CNN based on parsimonious HMM to address the issues of large vocabulary and diversity of writing styles in offline HCTR. Under the powerful seq-to-seq framework, the segmentation-free approach has shown the great progress in the HCTR. However, this approach is usually difficult to obtain the character boundary information.

### 3 System overview

In this paper, we design the over-segmentation based framework of [11] as the baseline recognition system, then propose to integrate a retrieval-based language model to improve the recognition performance. As the text to recognize is usually not known in advance, we adopt the two-pass recognition strategy. In the first-pass, we utilize the baseline system to obtain an initial recognition

result, which will be used to search the related contents from an Internet corpus and build a retrieval-based language model integrated in the path search as shown in the dashed-line box in Figure 2.

In the general over-segmentation based system, we usually regard the text recognition problem as searching the optimal path in a candidate lattice. In the following sections, we will describe the lattice generation and path search, respectively. Further details are provided in the following two sections.



**Fig. 2:** System diagram for handwritten Chinese text recognition with the retrieval-based language model.

#### 3.1 Lattice generation

Given an input text line image, we first over-segment it to a sequence of primitive segments (corresponding to a character or a part of a character), then combine several consecutive segments to form a candidate character pattern. Finally, we



engage a pre-trained CNN-based character classifier to output the top  $K$  character classes for each character pattern. By combining all the candidate character patterns and the corresponding character classes, we obtain a candidate segmentation-recognition lattice, where each path represents a candidate recognition result. Since we focus on the language model in this paper, we will omit an exhaustive description of the lattice generation and refer readers to the work [10].

### 3.2 Path search

Once the candidate lattice is generated, our target is to find the optimal path under a path evaluation function. The lattice is very complicated as it contains the uncertainties of both segmentation and recognition. Therefore, it is not possible to utilize the exhaustive search method. To improve the efficiency, we adopt the refined beam search algorithm [9], where the beam width is the same as the setting in the baseline work [9].

In the path evaluation function, we consider both the character recognition score and the linguistic context score as shown in Eqn. 1, where the weight  $\lambda_{LM}$  is used to balance two scores,

$$\mathcal{L} = \mathcal{L}_{\text{char}} + \lambda_{LM} \cdot \mathcal{L}_{LM}. \quad (1)$$

In our system, the character recognition score  $\mathcal{L}_{\text{char}}$  is calculated by a deep CNN on the corresponding character pattern as shown in Eqn. 2,

where the coefficient  $w_i$  represents the normalized character width to overcome the bias issue of the short-length path [10]. In our framework, we refer to the same CNN structure in the baseline work [10, 11].

$$\mathcal{L}_{\text{char}} = -w_i \cdot \log[\mathbf{CNN}(\cdot)]. \quad (2)$$

In the evaluation function, the LM plays an important role to provide the linguistic context information, and we utilize the popular  $N$ -gram model as shown in Eqn. 3.

$$\mathcal{L}_{LM} = -\log[\mathbf{Ngram}(\cdot)]. \quad (3)$$

Once we obtain the scores of both character recognition and language model, we aim to minimize the total function value to get the optimal path in the lattice.

## 4 Retrieval-based language model

In general, one common LM is used to evaluate the linguistic context expressed by Eqn. 1, however, if this LM does not match the document to be recognized, it may give a wrong score. To overcome this issue, we consider integrating an adaptive language model in Eqn. 1, which is constructed from the texts matched to the document to be recognized.



As the document to be recognized is variable, we first utilize the information retrieval techniques to obtain the matched texts from a large corpus (e.g., Internet resources), then learn an LM from these retrieved texts. This is called a retrieval-based language model. Finally, we integrate the loss function of the retrieval-based language model,  $\mathcal{L}_{\text{rLM}}$ , into the path evaluate function to obtain a more accurate score with a balanced weight  $\lambda_{\text{rLM}}$  as shown in Eqn. 4:

$$\mathcal{L} = \mathcal{L}_{\text{char}} + \lambda_{\text{LM}} \cdot \mathcal{L}_{\text{LM}} + \lambda_{\text{rLM}} \cdot \mathcal{L}_{\text{rLM}}. \quad (4)$$

In information retrieval, the similarity metric plays an important role, and we adopt the widely used Cosine similarity as shown in Eqn. 5:

$$\cos \theta = \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|}. \quad (5)$$

In the above, the symbols  $\mathbf{u}$  and  $\mathbf{v}$  denote the vector representation of query document (i.e., the document to be recognized) and retrieved content in the prepared corpus, respectively. In the following, we will describe two methods to obtain the vector representation of a document in our system.

#### 4.1 TF-IDF based content retrieval

The first approach to generate vectors for documents is based on the Term Frequency-Inverse Document Frequency (TF-IDF). TF-IDF is a technique for vectorizing documents based on the Bag

of words (BoW) model. It measures how important a term is within a document relative to a collection of documents. Term Frequency (TF) measures how frequently a term occurs in a document. It is the ratio of the number of times a term  $w$  appear in the document to the document length (total number of terms in the document):

$$\text{TF}(w) = \frac{N_w}{N}, \quad (6)$$

where  $N_w$  represents the number of times for the item  $w$  in a document, and  $N$  represents the total number of terms in the document. Inverse Data Frequency (IDF) measures how important a term is. It is defined as the log of the ratio of total number of documents in a collection to number of documents that contain a particular word:

$$\text{IDF}(w) = \log\left(\frac{D}{D_w}\right), \quad (7)$$

where the variable  $D$  represents the total number of documents in the corpus, and the variable  $D_w$  represents the number of documents that contain  $w$ . As a result, it weighs down the frequent terms and scales up the rare ones. Finally the TF-IDF value is the product of TF and IDF:

$$\text{TF-IDF}(w) = \text{TF}(w) \cdot \text{IDF}(w), \quad (8)$$

To obtain the vector representation of a document, we first segment all documents in the corpus into

409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431  
432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459

word levels using the Jieba toolkit<sup>2</sup> with a dictionary, then calculate the TF-IDF value of each word in the dictionary for each document (if the word does not show in this document, we set the TF-IDF value as zero), and finally we concatenate all TF-IDF values to be the vector of each document (i.e., the dimension of each vector is the size of the dictionary).

## 4.2 BERT based content retrieval

In our system, we also utilize the BERT model [58] to extract the vector of each document. BERT is a method for pre-training language representations, which is able to capture the contextual information of each word. For NLP tasks like semantic textual similarity (STS), BERT has achieved new state-of-the-art performance.

In our experiment, we adopt the pre-trained Bert-Base-Chinese model<sup>3</sup>, which contains 12 layers with 68 hidden states for each layer. We first segment each document into tokens using the BERT tokenizer as the input of the multi-layer Transformer, then obtain the document vector by four different methods. Specifically, **BE** stands for Bert-based embedding, and these four methods are (1) **BE1**: the representation of [CLS] of the last layer; (2) **BE2**: the average of the sequence of token vectors from the last layer; (3) **BE3**: the average of the sequence of token vectors from the

second last layer; (4) **BE4**: the average of the sequence of token vectors from the last two layers.

## 4.3 Language model construction

After retrieving the top-N related news articles, the SRI Language Modeling Toolkit (SRILM) [59] is utilized to build the retrieval-based language model. SRILM can provide frequency counts for  $N$ -grams after processing the related news corpus. It is expected that some phrases, terminologies, and names of places or celebrities can be better recognized with the dynamically generated retrieval-based language model. For example, if the document is related to tennis, then in the small Internet corpus, tennis-related terminologies and names of famous players may appear multiple times, increasing their possibility to be chosen during the search algorithm.

## 4.4 Additional improvements on HCTR

In addition to the retrieval-based model, we also implemented two extra improvements for HCTR.

**Concatenating Adjacent Lines.** Since the text line image is recognized one by one, the character at the beginning and the end of a line is recognized without context information from its adjacent lines in the same document. As a result, the language model has incomplete stencils at the beginning and the end of a text line, which will

---

<sup>2</sup><https://pypi.org/project/jieba/>

<sup>3</sup><https://github.com/google-research/bert>

lead to less accurate recognition. A simple idea is to concatenate two adjacent text lines for recognition. Thus, both previous and current text lines can benefit from it.

**Adjusting Weights of Punctuation Marks.** It is observed that some punctuation marks are often incorrectly recognized, and sometimes their neighboring Chinese characters are affected and not correctly recognized as well. This is partially because that some punctuation marks, such as commas and enumeration commas, are relatively simple shapes and resemble the strokes of Chinese characters. The framework may confuse them with some of the over-segmented parts of a Chinese character, or vice versa. For these scenarios, we proposed to adjust the weight for the loss function of the language model such that the importance of the context over the shape of the character pattern can be modified accordingly. When the top candidate character pattern belongs to a certain punctuation mark, the weight of the language model will be scaled by a factor of  $\alpha_{\text{mark}}$ ,

$$\lambda_{\text{LM}} = \alpha_{\text{mark}} \cdot \lambda_{\text{LM}}. \quad (9)$$

## 5 Experiments

### 5.1 Dataset and experimental setting

In this paper, we evaluate the proposed method on two benchmark Chinese handwriting recognition datasets: CASIA-HWDB [12] and ICDAR-2013 competition dataset [6]. The CASIA-HWDB database contains both isolated characters and unconstrained handwritten texts, where the training set contains 3,118,447 isolated character samples of 7,356 classes and 4076 pages of handwritten pages (including 41,781 textline samples). We tested our system on the test set containing 1,015 pages (including 10,449 textline samples). The dataset ICDAR-2013 contains 300 pages (including 3,432 textline samples) for testing only.

The values of the hyper-parameters in this paper are set by following the baseline system, which make a trade-off between the accuracy and efficiency [10]. Specifically, we set the maximum number of concatenated segments as 4 (e.g., a character pattern can contain at most 4 segments) in the candidate lattice generation, and the candidate number of character classification as 20 (i.e., the top 20 character classes with the high classification scores). In the path search, we set the beam width as 10. The common language model  $\mathcal{L}_{\text{LM}}$  is a character-level 2-gram language model, and the weight  $\lambda_{\text{LM}}$  is set as 0.1 by the trial-and-error method.

511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561

562 For the retrieval-based language model, we  
 563 take the top 200 documents for each recognized  
 564 document from the retrieval in the Internet cor-  
 565 pus. For each document to be recognized, the arti-  
 566 cles that have similar content are selected from a  
 567 large Internet corpus from Sogou<sup>4</sup> which consists  
 568 of more than 70,000 news articles. The selection  
 569 is based on the sentence similarity and the top-N  
 570 most related news articles that are stored to form a  
 571 small Internet corpus. In Eqn. 4, we set the weight  
 572  $\lambda_{\text{LM}}$  as 0.05 by the trial-and-error method.

580 Furthermore, we concatenate two adjacent  
 581 text lines during the path search to add the lin-  
 582 guistic context for the beginning characters. In  
 583 addition, we adjust the language model weight  
 584  $\lambda_{\text{LM}}$  on the punctuation marks as they have weak  
 585 linguistic context. Both tricks are effective for the  
 586 improvement of the recognition performance.

592 We evaluate the recognition performance in  
 593 terms of two character-level metrics, i.e., Correct  
 594 Rate (CR) and Accurate Rate (AR):

$$595 \text{CR} = \frac{N_t - N_{\text{de}} - N_{\text{se}}}{N_t}, \quad (10)$$

$$602 \text{AR} = \frac{N_t - N_{\text{de}} - N_{\text{se}} - N_{\text{ie}}}{N_t}. \quad (11)$$

604 In the above equations,  $N_t$  is the total number  
 605 of characters in the transcript of test documents.  
 606 The numbers of substitution errors  $N_{\text{se}}$ , deletion  
 607 errors  $N_{\text{de}}$  and insertion errors  $N_{\text{ie}}$  are calculated

by aligning the recognition result string with the  
 transcript using dynamic programming.

## 5.2 Experimental results

For the evaluation, we mainly evaluate the effec-  
 tiveness of the proposed retrieval-based language  
 model on the benchmark datasets of both CASIA-  
 HWDB and ICDAR-2013, and the results are  
 shown in Tables 1 and 2, respectively. In our base-  
 line recognition system, we only use a common  
 language model (i.e., character bi-gram) as shown  
 in the first row of each table. In the retrieval-based  
 language model, we also utilize the character bi-  
 gram and take top 200 relevant texts retrieved by  
 the TF-IDF method. The last line in each table  
 shows the proposed method with the proposed  
 two further improvement tricks (i.e., concatenat-  
 ing adjacent lines and adjusting language model  
 weights on punctuation marks).

As shown in Tables 1 and 2, we can see that  
 the performance of only using the retrieved-based  
 model decreases significantly in comparison to  
 the baseline recognition. The reason is that the  
 retrieved-based model is only built on a small  
 size of related contents, resulting in serious spar-  
 sity in the  $N$ -grams. However, we can see that  
 the recognition performance is increased by the  
 combination of the common language model and  
 retrieval-based model. Since the retrieval-based  
 model can provide complimentary linguistic infor-  
 mation from the relevant contents, it increases

---

611 <sup>4</sup>[http://www.sogou.com/labs/resource/list\\_news.php](http://www.sogou.com/labs/resource/list_news.php)  
 612

**Table 1:** Comparison of different models on CASIA-HWDB. (The bold values indicate the highest performance)

Experiment	CR(%)	AR(%)	Ch(%)	Sb(%)	Dg(%)	Lt(%)
Common LM	94.85	94.09	96.79	85.01	89.20	60.37
Retrieval LM	92.70	90.78	94.96	81.65	83.06	62.92
Common + Retrieval LM	94.86	94.16	96.78	84.72	90.93	67.07
Common + Retrieval LM+Improvements	<b>94.91</b>	<b>94.25</b>	<b>96.85</b>	<b>84.78</b>	<b>91.15</b>	<b>67.43</b>

**Table 2:** Comparison of different models on ICDAR-2013. (The bold values indicate the highest performance)

Experiment	CR(%)	AR(%)	Ch(%)	Sb(%)	Dg(%)	Lt(%)
Common LM	94.28	93.28	96.19	81.90	85.30	42.79
Retrieval LM	90.65	87.87	92.85	76.63	80.16	45.15
Common+Retrieval LM	94.68	93.79	96.50	82.55	87.64	47.04
Common+Retrieval LM+Improvements	<b>94.87</b>	<b>94.00</b>	<b>96.57</b>	<b>83.92</b>	<b>88.15</b>	<b>48.46</b>

the potential that characters or phrases that often appear in certain fields are recognized correctly. Based on this, we also implemented two further improvement tricks and found they increase the recognition accuracy to the CR 94.91% and the AR 94.25% on CASIA-HWDB, the CR 94.87% and the AR 94.00% on ICDAR-2013. In summary, we can see that both CR and AR values are improved by using the retrieval-based model, connecting two adjacent textlines and adjusting the language model weights of punctuation marks.

截止到目前为止，白云机场几天下来依然有近千名旅客滞留。截止到目前为止，白云机场几天下来依然有近千名旅客滞留。截止到目前为止，白云机场几天下来依然有近千名旅客滞留。

**Fig. 3:** An example of recognition with/without the retrieval-based language model. The first row is the text line image; the second row is the transcript (ground-truth); the third row is the result without the retrieval-based language model; the last row is result with the retrieval-based language model.

死了”让制片方恼羞成怒，立即对所有演员下达了封口令。  
死了”让制片方恼羞成怒，立即对所有演员，下达了封口令。  
死了”让制片方恼羞成怒，立即对所有演员下达了封口令。

**Fig. 4:** An example of recognition with/without adjusted weights for punctuation marks. The first row is the text line image; the second row is the transcript (ground-truth); the third row is the result without adjusted weights for punctuation marks; the last row is result with adjusted weights for punctuation marks.

In order to demonstrate the effectiveness of the proposed method intuitively, we show some recognition examples in Figures 3, 4 and 5. In Figure 3, we show that the airport name is recognized correctly after the retrieval-based language model is integrated. The possible reason is that this airport name frequently appeared in the retrieved corpus, which increases the estimation of language model. As shown in Figure 4, the recognition was improved by adjusting the language model weight

664 立了航班不正常处理小组, 组建由 值机员、服务受理  
 665 任的航班不正常处理小组, 专人负责航班延误信息的公布、  
 666 任的航班不正常处理小组, 专人负责航班延误信息的公布、  
 667 任的航班不正常处理小组, 专人负责航班延误信息的公布、  
 668 任的航班不正常处理小组, 专人负责航班延误信息的公布、  
 669 **Fig. 5:** An example of recognition with/without  
 670 concatenating adjacent lines. The first two rows  
 671 are two adjacent text line images; the third row is  
 672 the transcript (ground-truth) for the second text  
 673 line image; the fourth row is the result for the  
 674 second text line image without concatenating two  
 675 adjacent lines; the last row is result for the second  
 676 text line image with concatenating two adjacent  
 677 lines.

678  
 679 for punctuation marks, where the original recog-  
 680 nition is misclassified by a comma. In Figure 5, we  
 681 show one example of recognition by concatenating  
 682 adjacent lines. We can see that the character at  
 683 the beginning of the current line is misrecognized  
 684 to ‘住’ due to the high similarity to that char-  
 685 acter image, which is corrected by concatenating  
 686 the previous line recognition as the last character  
 687 provides more contexts for this correction.

694 Taking the paragraph in Figure 3 as an exam-  
 695 ple, we examined how many the retrieved news are  
 696 related to the target text. As shown in Table 3,  
 697 158 out of the top 200 retrieved news hit the tar-  
 698 get of air transportation, while the remaining 42  
 699 news also had strong overlapping with air trans-  
 700 portation. In 19 sports news and 4 entertainment  
 701 news, sports teams and celebrities encountering  
 702 flight delay are often reported. It can be expected  
 703 that news of other types of transportation may  
 704 also be retrieved such as railway/road/boat trans-  
 705 portation. The 6 social news are mainly using air  
 706  
 707  
 708  
 709  
 710  
 711  
 712  
 713  
 714

**Table 3:** Statistics on the top 200 news retrieved for a paragraph related to air transportation.

Theme	Number of news
Air transportation	158
Sports	19
Other transportation	13
Entertainment	4
Society	6

transportation for natural disaster rescue. The air-  
 port name is widely used in these retrieved texts,  
 that’s why the recognition error in the Figure 3 is  
 corrected by adding this retrieval LM.

### 5.3 Comparison of different Bert-based embedding

Four different approaches to represent the docu-  
 ment vector were compared: **BE1**, **BE2**, **BE3** and  
**BE4**. As shown in Tables 4 and 5, we can see that  
 the performance of **BE1** was the worst because  
 [CLS] token appears at the start of the text for  
 classification tasks, and [CLS] token embedding  
 does not convey much semantic information as a  
 sentence representation. On the other hand, the  
 difference of performance among **BE2**, **BE3** and  
**BE4** was negligible, and the reason is that the size  
 of the news corpus is not large enough. As a con-  
 sequence, we adopt **BE2** to obtain the document  
 embedding in the following experiments.

**Table 4:** Comparison of different Bert-based embedding on CASIA-HWDB.

Experiment	CR(%)	AR(%)
BE1	94.82	94.15
BE2	94.89	94.24
BE3	94.89	94.23
BE4	94.89	94.23

**Table 5:** Comparison of different Bert-based embedding on ICDAR-2013.

Experiment	CR(%)	AR(%)
BE1	94.74	93.88
BE2	94.85	93.99
BE3	94.86	94.01
BE4	94.86	94.01

## 5.4 Comparison of different retrieval-based models

In this section, we evaluate different patterns in the retrieval-based language model, including two retrieval methods and different orders of  $N$ -gram models. The results are shown in Table 6. By comparison of different  $N$ -gram models, we can see that the 2-gram model performs the worst for both the TF-IDF and the BERT retrieval methods on both datasets, due to capturing very short contexts. After increasing to the 3-gram model, the accuracy is increased considerably, especially on the ICDAR-2013 dataset, where the AR is increased from 94.00% to 94.18%, and from 93.99% to 94.18% for TF-IDF and BERT methods respectively. However, the 4-gram model does not boost the accuracy further because the number of retrieved text is very limited and leads to very sparse 4-gram items. In other words, most of

the 4-gram context scores will be learned by the back-off to the 3-gram estimation [59]. Comparing TF-IDF with BERT, we find that their performance is very similar on the ICDAR-2013 dataset, and note only a little improvement by BERT for the 3-gram model on CASIA-HWDB. The possible reason is that our query (i.e., transcript from the first-pass recognition) is not reliable due to some recognition errors.

## 5.5 Comparison with existing methods

Table 7 shows the comparison of existing methods and ours on the ICDAR-2013 dataset. In our approach, we apply common LM, retrieval LM and the two additional improvements (see Sec.4.4). For the language model, we utilize a 5-gram Common LM and 4-gram Retrieval LM. As the number of retrieved text is very limited leading to very sparse 4-gram (or higher order) items, it is not necessary to utilize higher order LM (see Sec.5.4). For the character model, we apply a VGG-style CNN model which is the same as the baseline system [10, 11]. As shown in Sec. 3, we adopt the recognition system of [11] as the baseline, which is a weakly supervised learning method with only transcript level annotation under the over-segmentation framework. From the table, we can observe that our method achieves a much higher accuracy than the baseline work [11],

715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
760  
761  
762  
763  
764  
765



**Table 6:** Results of different retrieval-based models on CASIA-HWDB and ICDAR-2013.

Retrival methods	N-gram	CASIA-HWDB		ICDAR-2013	
		CR(%)	AR(%)	CR(%)	AR(%)
TF-IDF	2-gram	94.91	94.28	94.87	94.00
	3-gram	94.93	94.27	95.02	94.18
	4-gram	94.92	94.26	95.03	94.19
BERT	2-gram	94.89	94.24	94.85	93.99
	3-gram	94.95	94.29	95.02	94.18
	4-gram	94.94	94.28	95.02	94.18

**Table 7:** Comparison with existing methods on the ICDAR-2013 dataset.

Methods	CR(%)	AR(%)
Wang et al. [60]	95.53	94.02
Wu et al. [10]	96.32	96.20
Xie et al. [61]	96.70	96.22
Peng et al. [62]	95.51	94.88
Wang et al. [11]	95.73	95.11
Peng et al. [63]	97.32	96.79
Ours (Common+Retrieval LM+Improvements)	96.13	95.48

which demonstrates the effectiveness of the proposed retrieval-based language model adaptation. By comparison with the other methods, our performance is even competitive to the accuracy in the work [10]. Note that [10] adopted the same over-segmentation based recognition framework, but optimized the recognition model under strong supervision with character-level annotation and elaborately integrated geometric context models. Although the recent work [63] optimized the model under the segmentation-free framework with transcript level annotation only, it utilized a large set of synthetic data to boost the accuracy.

## 6 Conclusion

In this paper, we proposed a retrieval-based language model for handwritten Chinese text recognition, which obtains the adaptive linguistic context during the recognition. Since the document to be recognized is unknown a priori, we engaged a two-pass recognition strategy. In the first-pass recognition, we take a common language model to recognize the document to output an initial transcript, which was used to retrieve related contents from a large text corpus. For the retrieval method, we evaluated both TF-IDF and four BERT-based embedding methods in our experiments. Finally, we built an adaptive language model from the retrieved contents, and combined with a common language model in the second-pass recognition to obtain the final transcript. We evaluated the proposed method on two benchmark datasets, and the extensive experimental results demonstrated the effectiveness of the proposed retrieval-based language model. In the future, we will consider an online Internet-based retrieval method to obtain

the related common sense knowledge to build the adaptive language model.

**Acknowledgments.** The work was funded by National Natural Science Foundation of China under no.61876154 and no.61876155; Jiangsu Science and Technology Programme (Natural Science Foundation of Jiangsu Province) under no.BE2020006-4.

## References

- [1] G. Nagy. Disruptive developments in document recognition. *Pattern Recognition Letters*, vol. 79, pp. 106–112, 2016.
- [2] H. Fujisawa. Forty Years of Research in Character and Document Recognition—An Industrial Perspective. *Pattern Recognition*, vol. 41, pp. 2435-2446, 2008.
- [3] R.-W. Dai, C.-L. Liu, B.-H. Xiao. Chinese Character Recognition: History, Status and Prospects. *Frontiers of Computer Science in China*, vol. 1(2), pp. 126-136, 2007.
- [4] C.-L. Liu, Y. Lu, editors. *Advances in Chinese Document and Text Processing*. book in Series on Language Processing, Pattern Recognition, and Intelligent Systems, vol. 2, World Scientific, 2017.
- [5] Cheng-Lin Liu, Fei Yin, Qiu-Feng Wang, Da-Han Wang. ICDAR 2011 Chinese Handwriting Recognition Competition. 11th International Conference on Document Analysis and Recognition (ICDAR), pp. 1464-1469, 2011.
- [6] Fei Yin, Qiu-Feng Wang, Xu-Yao Zhang, Cheng-Lin Liu. ICDAR 2013 Chinese Handwriting Recognition Competition. 12th International Conference on Document Analysis and Recognition (ICDAR), pp. 1464-1470, 2013.
- [7] Cheng Cheng, Xu-Yao Zhang, Xiaohu Shao, Xiang-Dong Zhou. Handwritten Chinese Character Recognition by Joint Classification and Similarity Ranking. 15th International Conference on Frontiers in Handwriting Recognition (ICFHR), pp. 507-511, 2016.
- [8] Xu-Yao Zhang, Yoshua Bengio, Cheng-Lin Liu. Online and offline handwritten Chinese character recognition: A comprehensive study and new benchmark. *Pattern Recognition*, vol. 61, pp. 348-360, 2017.
- [9] Qiu-Feng Wang, Fei Yin, Cheng-Lin Liu. Handwritten Chinese Text Recognition by Integrating Multiple Contexts. *IEEE Trans. Pattern Analysis and Machine Intelligence (PAMI)*, vol. 34(8), pp. 1469-1481, 2012.
- [10] Yi-Chao Wu, Fei Yin, Cheng-Lin Liu. Improving handwritten Chinese text recognition using neural network language models and convolutional neural network shape models. *Pattern Recognition*, vol. 65, pp. 251-264, 2017.
- [11] Zhen-Xing Wang, Qiu-Feng Wang, Fei Yin,

817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863  
864  
865  
866  
867

- 868 Cheng-Lin Liu. Weakly Supervised Learn-  
869 ing for Over-Segmentation Based Handwrit-  
870 ten Chinese Text Recognition. 17th Inter-  
871 national Conference on Frontiers in Hand-  
872 writing Recognition (ICFHR), pp. 157-162,  
873 2020.
- 874  
875  
876  
877
- [12] C.-L. Liu, F. Yin, D.-H. Wang, Q.-F. Wang.  
878 CASIA Online and Offline Chinese Handwrit-  
879 ing Databases. 11th International Confer-  
880 ence on Document Analysis and Recognition  
881 (ICDAR), pp. 37-41, 2011.
- 882  
883  
884  
885
- [13] Qiu-Feng Wang, Erik Cambria, Cheng-Lin  
886 Liu, Amir Hussain. Common Sense Knowl-  
887 edge for Handwritten Chinese Text Recogni-  
888 tion. *Cognitive Computation*, vol. 5 (2), pp.  
889 234-242, 2013.
- 890  
891  
892  
893
- [14] Qiu-Feng Wang, Fei Yin, Cheng-Lin Liu.  
894 Unsupervised Language Model Adaptation  
895 for Handwritten Chinese Text Recognition.  
896 *Pattern Recognition*, vol. 47, pp. 1202-1216,  
897 2014.
- 898  
899  
900  
901  
902
- [15] A. Carlson, J. Betteridge, B. Kisiel, B. Set-  
903 tles, E.R. Hruschka Jr. and T.M. Mitchell.  
904 Toward an Architecture for Never-Ending  
905 Language Learning. In *Proceedings of the*  
906 *24th Conference on Artificial Intelligence*  
907 (AAAI), 2010.
- 908  
909  
910  
911  
912
- [16] T. Mitchell, W. Cohen, E. Hruschka, et al.  
913 Never-Ending Learning. In *Proceedings of*  
914 *the 29th Conference on Artificial Intelligence*  
915 (AAAI), pp. 2302-2310, 2015.
- 916  
917  
918
- [17] Fergus R, Fei-Fei L, Perona P, et al. Learn-  
ing object categories from Google's image  
search. 10th IEEE International Conference  
on Computer Vision (ICCV), pp. 1816-1823,  
2005.
- [18] Nishizaki H., Sekiguchi Y. Word Error Cor-  
rection of Continuous Speech Recognition  
Using WEB Documents for Spoken Docu-  
ment Indexing. In: *Computer Processing of*  
*Oriental Languages. Beyond the Orient: The*  
*Research Challenges Ahead (ICCPOL)*, vol.  
4285, pp. 213-221, 2006.
- [19] Clemens Oertel, Shauna O'Shea, Adam Bod-  
nar, D. Blostein. Using the web to validate  
document recognition results: experiments  
with business cards. In *Proceedings of SPIE,*  
*Document Recognition and Retrieval XII*,  
vol. 5676, pp. 17-27, 2005.
- [20] Donoser M, Bischof H, Wagner S. Using  
web search engines to improve text recog-  
nition. *International Conference on Pattern*  
*Recognition (ICPR)*, pp. 1-4, 2008.
- [21] Donoser M., Wagner S., Bischof H.. Context  
information from search engines for docu-  
ment recognition. *Pattern Recognition Let-*  
*ters*, vol. 31, pp. 750-754, 2010.
- [22] Qiu-Feng Wang, Fei Yin, Cheng-Lin Liu.  
Integrating language model in handwrit-  
ing Chinese text recognition. 10th Interna-  
tional Conference on Document Analysis and

- Recognition (ICDAR), pp. 1036-1040, 2009.
- [23] Y.X. Li, C.L. Tan, X.Q. Ding. A hybrid postprocessing system for offline handwritten Chinese Script recognition. *Pattern Analysis and Applications*, vol. 8, pp. 272-286, 2005.
- [24] R.F. Xu, D.S. Yeung, D.M. Shi. A hybrid postprocessing system for offline handwritten Chinese character recognition based on a statistical language model. *International Journal of Pattern Recognition and Artificial Intelligence*, vol. 19(3), pp. 415-428, 2005.
- [25] Rosenfeld R. Two decades of statistical language modeling: where do we go from here? *IEEE*, vol. 88(8), pp. 1270–8, 2000.
- [26] J.R.Bellegarda. Statistical language model adaptation: review and perspectives. *Speech Communication*, vol. 42(1), pp. 93-108, 2004.
- [27] B.C. Russell, A. Torralba, K.P. Murphy, et al. LabelMe: A Database and Web-Based Tool for Image Annotation. *International Journal of Computer Vision*, vol. 77(1-3), pp. 157–173, 2008.
- [28] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. *IEEE Computer Vision and Pattern Recognition (CVPR)*, pp. 248-255, 2009.
- [29] L. Fei-Fei. ImageNet: crowdsourcing, benchmarking & other cool things. *CMU VASC Seminar*, March, 2010.
- [30] Hankz Hankui Zhuo, Qiang yang, Rong Pan, Lei Li. Cross-Domain Action-Model Acquisition for Planning Via Web Search. *21th International Conference on Automated Planning and Scheduling*, pp. 298-305, 2011.
- [31] Chen L, Lamel L, Gauvain J L, et al. Dynamic language modeling for broadcast news. *8th International Conference on Spoken Language Processing*, 2004.
- [32] Whitelaw C, Hutchinson B, Chung G Y, et al. Using the web for language independent spellchecking and autocorrection. *Conference on Empirical Methods in Natural Language Processing*, Association for Computational Linguistics, pp. 890-899, 2009.
- [33] Bassil Y, Alwani M. OCR Post-Processing Error Correction Algorithm Using Google’s Online Spelling Suggestion. *Emerging Trends in Computing and Information Sciences*, vol. 3(1), pp. 90-99, 2012.
- [34] Oprean C, Likforman-Sulem L, Popescu A, et al. Using the Web to Create Dynamic Dictionaries in Handwritten Out-of-Vocabulary Word Recognition. *12th International Conference on Document Analysis and Recognition (ICDAR)*, pp. 989-993, 2013.
- [35] Oprean C, Popescu A, Popescu A, et al. Handwritten word recognition using Web resources and recurrent neural networks.

919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969

- 970 International Journal on Document Analy-  
 971 sis and Recognition (IJ DAR), vol. 18(4), pp.  
 972 287-301, 2015.
- 973  
 974 [36] Oprean C, Likformansulem L, Mokbel C, et  
 975 al. BLSTM-based handwritten text recog-  
 976 nition using Web resources. 13th Interna-  
 977 tional Conference on Document Analysis and  
 978 Recognition (ICDAR), pp. 466-470, 2015.
- 979  
 980 [37] Marti U. V., Bunke H. Using a statistical lan-  
 981 guage model to improve the performance of  
 982 an HMM-based cursive handwriting recogni-  
 983 tion systems. International Journal of Pat-  
 984 tern Recognition and Artificial Intelligence,  
 985 vol. 15(01), pp. 65-90, 2001.
- 986  
 987 [38] N.-X. Li, L.-W. Jin. A Bayesian-Based Prob-  
 988 abilistic Model for Unconstrained Handwrit-  
 989 ten Offline Chinese Text Line Recognition.  
 990 IEEE International Conference on Systems,  
 991 Man and Cybernetics, pp. 3664-3668, 2010.
- 992  
 993 [39] X.-D. Zhou, D.-H. Wang, F. Tian, C.-L. Liu,  
 994 M. Nakagawa. Handwritten chinese/japanese  
 995 text recognition using semi-markov condi-  
 996 tional random fields. IEEE Trans. Pattern  
 997 Analysis and Machine Intelligence (PAMI),  
 998 vol. 35(10), pp. 2413-2426, 2013.
- 999  
 1000 [40] Yoshua Bengio, Réjean Ducharme, Pascal  
 1001 Vincent, Christian Jauvin. A Neural Proba-  
 1002 bilistic Language Model. Journal of Machine  
 1003 Learning Research, vol. 3, pp. 1137-1155,  
 1004 2003.
- 1005  
 1006 [41] T. Mikolov, M. Karafiat, L. Burget, J. H.  
 1007 Cernocky, S. Khudanpur. Recurrent neural  
 1008 network based language model. Interspeech  
 1009 2010, 11th Annual Conference of the Inter-  
 1010 national Speech Communication Association,  
 1011 pp. 1045-1048, 2010.
- 1012  
 1013 [42] Yann N. Dauphin, Angela Fan, Michael  
 1014 Auli, David Grangier. Language Modeling  
 1015 with Gated Convolutional Networks. arXiv  
 1016 preprint arXiv:1612.08083v3, 2017.
- 1017  
 1018 [43] Irie K, Tüske Z, Alkhouli T, et al. LSTM,  
 1019 GRU, Highway and a Bit of Attention: An  
 1020 Empirical Overview for Language Model-  
 ing in Speech Recognition. Interspeech 2016,  
 17th Annual Conference of the Interna-  
 tional Speech Communication Association,  
 pp. 3519-3523, 2016.
- [44] Luong T, Kayser M, Manning C D. Deep  
 Neural Language Models for Machine Trans-  
 lation. 19th Conference on Computational  
 Natural Language Learning, pp. 305-309,  
 2015.
- [45] J.R.Bellegarda. Exploiting latent semantic  
 information in statistical language modeling.  
 IEEE, vol. 88(8), pp. 1279-1296, 2000.
- [46] Hofmann T. Unsupervised Learning by Prob-  
 abilistic Latent Semantic Analysis. Machine  
 Learning, vol. 42(1-2), pp. 177-196, 2001.
- [47] Blei D M, Ng A Y, Jordan M I. Latent dirich-  
 let allocation. Journal of Machine Learning  
 Research, vol. 3, pp. 993-1022, 2003.



- 1072 (ICSLP), pp. 901-904, 2002.
- 1073 [60] S. Wang, L. Chen, L. Xu, W. Fan, J. Sun,  
1074 and S. Naoi. Deep knowledge training and  
1075 heterogeneous cnn for handwritten chinese  
1076 text recognition. 15th International Confer-  
1077 ence on Frontiers of Handwriting Recognition  
1078 (ICFHR), pp. 84-89, 2016.
- 1083 [61] Z.-C. Xie, Y.-X. Huang, Y.-Z. Zhu, L.-W.  
1084 Jin, Y.-L. Liu, and L.-L. Xie. Aggregation  
1085 cross-entropy for sequence recognition. IEEE  
1086 Conference on Computer Vision and Pattern  
1087 Recognition (CVPR), pp. 6538-6547, 2019.
- 1091 [62] D.-Z. Peng, L.-W. Jin, Y.-Q. Wu, Z.-P.  
1092 Wang, and M.-X. Cai. A Fast and Accu-  
1093 rate Fully Convolutional Network for End-to-  
1094 End Handwritten Chinese Text Segmentation  
1095 and Recognition. 15th International Confer-  
1096 ence on Document Analysis and Recognition  
1097 (ICDAR), pp. 25-30, 2019.
- 1103 [63] Peng, Dezhi and Jin, L. and Ma, Wei-  
1104 hong and Xie, Canyu and Zhang, Hesuo and  
1105 Zhu, Shenggao and Li, Jing. Recognition of  
1106 Handwritten Chinese Text by Segmentation:  
1107 A Segment-annotation-free Approach. IEEE  
1108 Trans. on Multimedia, 2022.
- 1113 [64] Z.-R. Wang, J. Du, and J.-M. Wang. Writer-  
1114 aware CNN for parsimonious HMM-based  
1115 offline handwritten Chinese text recognition.  
1116 Pattern Recognition, vol. 100, pp. 107-102,  
1117 2020.
- 1118  
1119  
1120  
1121  
1122