# General and Domain Adaptive Chinese Spelling Check with Error Consistent Pretraining

QI LV*, School of Computer Science and Technology, Soochow University, China

ZIQIANG CAO*†, School of Computer Science and Technology, Institution of Artificial Intelligence, Soochow University, China

LEI GENG, School of Computer Science and Technology, Soochow University, China

CHUNHUI AI, School of Computer Science and Technology, Soochow University, China

XU YAN, School of Computer Science and Technology, Soochow University, China

GUOHONG FU, School of Computer Science and Technology, Institution of Artificial Intelligence, Soochow University, China

The lack of label data is one of the significant bottlenecks for Chinese Spelling Check (CSC). Existing researches use the automatic generation method by exploiting unlabeled data to expand the supervised corpus. However, there is a big gap between the real input scenario and automatically generated corpus. Thus, we develop a competitive general speller **ECSpell** which adopts the **E**rror **C**onsistent masking strategy to create data for pretraining. This error consistency masking strategy is used to specify the error types of automatically generated sentences consistent with the real scene. The experimental result indicates that our model outperforms previous state-of-the-art models on the general benchmark.

Moreover, spellers often work within a particular domain in real life. Due to many uncommon domain terms, experiments on our built domain specific datasets show that general models perform terribly. Inspired by the common practice of input methods, we propose to add an alterable user dictionary to handle the zero-shot domain adaption problem. Specifically, we attach a **U**ser **D**ictionary guided inference module (**UD**) to a general token classification based speller. Our experiments demonstrate that ECSpell$^{UD}$, namely ECSpell combined with UD, surpasses all the other baselines broadly, even approaching the performance on the general benchmark[1].

CCS Concepts: • **Computing methodologies → Natural language processing**.

Additional Key Words and Phrases: Chinese spelling check, domain adaptive, user dictionary

---

*Both authors contributed equally to this research.
†Corresponding author.
[1]https://github.com/Aopolin-Lv/ECSpell

---

Authors' addresses: Qi Lv, 20205227047@stu.suda.edu.cn, School of Computer Science and Technology, Soochow University, China; Ziqiang Cao, zqcao@suda.edu.cn, School of Computer Science and Technology, Institution of Artificial Intelligence, Soochow University, China; Lei Geng, School of Computer Science and Technology, Soochow University, China; Chunhui Ai, School of Computer Science and Technology, Soochow University, China; Xu Yan, School of Computer Science and Technology, Soochow University, China; Guohong Fu, ghfu@hotmail.com, School of Computer Science and Technology, Institution of Artificial Intelligence, Soochow University, China.

---

# 1 INTRODUCTION

Chinese spelling check (CSC) [26] aims to identify and correct the misspelling characters in Chinese sentences. Since spelling errors may cause deviations in the semantics of sentences, the CSC task is significant to most downstream NLP applications such as named entity recognition, machine translation, and text summarization. CSC is usually regarded as a token classification problem [11] as the lengths of input and output are the same.

Recently, CSC has made remarkable progress with the help of the large-scale pretrained language models like BERT [4]. Hong et al. [11] firstly modeled CSC as a BERT token classification task. Since almost all the spelling errors are related to phonological or visual similarity [17], many subsequent studies incorporated the similarity knowledge into spellers. For example, Nguyen et al. [22] utilized glyph information while Cheng et al. [2], Zhang et al. [36], and Xu et al. [33] employed phonetic features. The fusion method of these similarities has also been explored, such as Graph Convolution Network (GCN) [2] and multi-modal [12, 33].

One of the main challenges of utilizing supervised learning for CSC is the lack of high-quality annotated corpora. Most previous CSC studies concentrated on the general benchmark SIGHAN [26, 32, 35]. To alleviate this problem, large-scale unlabelled corpus was used to enhance the spelling check ability of models. Some works [6, 18, 29, 36] adopted multimodal-based methods to build corpora automatically for weakly supervised learning. Although some of these methods [29] work well, they still lack the consideration of the consistency between the real datasets and the automatically generated corpora, which will bring a big gap.

To build a practical competitive speller, we develop a novel general speller **ECSpell** which uses the **E**rror **C**onsistent masking strategy to exploit unlabeled data for pretraining. Previous work [18, 36] ignored the significant fact that the error source of a sentence could often be the only one due to the input method. The proposed masking strategy uses this property to simulate the natural scene as much as possible to construct the corpus automatically. Concretely, this masking strategy first specifies an error type of the given sentence, then masks single or continuous characters according to this type to create spelling errors. It can be simply applied to large-scale corpus for unsupervised CSC learning. Meanwhile, we fuse the glyph information and fine-grained phonetic features. See details in following Section 3.2 and Section 3.3.

In addition, spellers often work within a specific domain. Considering most domain terms, our experiments show that general spellers perform terribly in this case. Many such errors can not be corrected, and some correct parts are even wrongly changed. For example, with no medical knowledge, it is difficult to detect that "Xian Bing" means "癣病 (Tinea)" instead of the common word "馅饼 (pie)". Under the circumstance of the document writing, it is common for a speller to wrongly change "金字村 (Jinzi Village)" into "金子村 (Golden Village)" as the lack of the particular external terms knowledge.

It is unacceptable in time and economy to annotate data in each domain and fine-tune models individually, especially for domains with increasing new terms. Considering the actual application, we extend the traditional CSC task to the domain adaptive CSC task. Following the common practice of the industrial word segmentation tools (e.g., Jieba[2]) and input methods (e.g., Sougou[3]), we propose to add an alterable user dictionary to handle the adaption problem of CSC in the zero-shot schema. Since the output probability in token classification is independent token-by-token, the user dictionary can not be directly adopted in such spellers. Hence we propose a novel **U**ser **D**ictionary guided inference module (UD) to post-process the prediction results of a token classification based speller. Specifically, given the predicted token path candidates, UD rewards the

---

paths containing more dictionary terms. Notably, this method applies to both general and domain adaptive CSC tasks.

Due to the lack of existing evaluation datasets, we annotate a CSC dataset of the law, medical treatment and official document writing domains to represent these three scenarios. We both conduct experiments on the general benchmark [26] and our domain specific evaluation datasets. Our experimental results show that our ECSpell achieves new state-of-the-art performance on the general benchmark and UD can improve the performance of all tests with no extra fine-tuning. We combine our ECSpell and UD to form our CSC extensive speller ECSpell$^{UD}$. Related experiments demonstrate that ECSpell$^{UD}$ surpasses all the other baselines.

Our contributions can be summarized as follows:

- We develop an error consistent masking strategy for unsupervised CSC learning.
- We annotate domain specific datasets for follow-up CSC research.
- We use the user dictionary guided inference module (UD) to enhance the domain adaptive for any token classification based speller.

We organize the rest of this paper as follows. First, in Section 2, we introduce the related work. Following Section 3 and Section 4, we describe our work for the general CSC task and user dictionary guided domain adaptive CSC task, respectively. Next, Section 5 presents the experimental result and related analysis. Finally, we conclude our work in Section 6.

## 2 RELATED WORK

### 2.1 Chinese Spelling Check

Previous research on CSC can be divided into three distinct categories: rule based methods, machine learning based methods and deep learning based methods. Jiang et al. [13] utilized the China National Matriculation Examinations (NME) rules to classify spelling errors into Idiom Error and Word Error. Yu and Li [34] used character-level n-gram language models and a word vocab to detect and correct potential misspelled characters. For machine learning based methods, CRF-based word segmentation/part of speech tagger was integrated into a tri-gram language model after the rule-based fronted [7, 31].

With the rapid development of deep learning techniques, the process of CSC has moved a big step forward. Wang et al. [30] leveraged the pointer network by picking the correct character from the confusion set. Hong et al. [11] firstly modeled CSC as a BERT token classification task. As most spelling errors come from similar pronunciations or glyphs, many successive studies merged the similarity knowledge into spellers. Nguyen et al. [22] adopted glyph features while Cheng et al. [2], Zhang et al. [36], and Xu et al. [33] employed phonetic information. Researcher also explored the mixture method of these similarities, such as an adaptive gating module [33], GCN [2] and multi-modal [12]. To refine the learning object, Li et al. [14] applied an adversarial strategy to enhance the robustness of the model while Li and Shi [15] adopted the focal loss penalty strategy to alleviate the class imbalance problem. Further, Li et al. [16] refined the knowledge representation of pretrained language models to narrow the gap between it and the essential of CSC task via the contrastive learning method. Zhang et al. [37] combined the losses of detection and correction with a soft-mask strategy. Different from the above mentioned approaches, Bao et al. [1] applied a non-autoregressive model to improve the phrase correction performance.
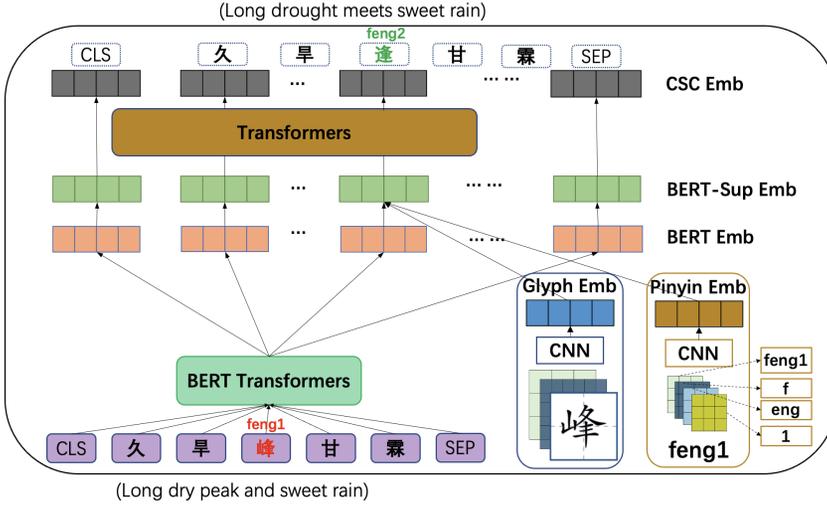
Fig. 1. Overview of our general model.

## 2.2 Domain Terminology Injection

In the development of the neural machine translation (NMT), significant studies proposed methods to integrate external specialized domain terminologies. It can be roughly divided into the following three categories [21]: 1). Placeholders. Placeholders replace terms appearing in the sentence pair in pre- and post-processing [3]. It is obviously that this method lacks flexibility as the model will always replace the placeholder with the same phrase irrespective of grammatical context. 2). Learning to apply constraints. It learns a copy behavior of terminology at training time [5, 24]. For example, "疟疾疟疾是一种由寄生虫引起的威胁生命的疾病" (Malaria Malaria is a life-threatening disease caused by parasites). The model is trained to incorporate terminology translations when provided as additional input in the source sentence. It also lacks generalization power as it simply "copy" the term found in the terminology base on the source sentence, irrespective of the target hypothesis context [5]. 3). Constrained decoding. The model takes the translation terms as the decoding constraints applied in the inference stage [9, 10, 25]. However, a source term may have multiple translation term inflections among which the MT engine should on-the-fly select the best one depending on the source and target context.

As far as we known, there is little study devoted to the domain adaptive spelling check although it is quite important in real life.

## 3 GENERAL CSC TASK

### 3.1 Problem Formulation

Given an input sentence $X = \{x_1, x_2, \cdots, x_n\}$ of length $n$, the model needs to generate its corresponding correct sentence $Y = \{y_1, y_2, \cdots, y_n\}$. For CSC, most tokens in output sentence $Y$ are the same as that in the input sentence $X$ and the rest is the target of error correction. Since the length of input sentence $X$ is equal to the output sentence $Y$, this task is usually formed token classification task [2, 11].

### 3.2 Our ECSpell Overview

Since spelling errors focus on shape and sound similarities, we enrich the BERT token classification

model with related features as shown in Figure 1. In order to make more effective use of these two knowledge mentioned above, our general model ECSpell has made two improvements. On one hand, the knowledge of glyph including the meaning of radical and the frame structure are considered. On the other hand, we explore the fusion of pinyin[4] knowledge.

### 3.2.1 Embedding Layer.

**BERT Embedding**. Following previous works [11, 12, 33], we employ BERT [4] encoder as the semantic encoder. Benefit from the pretraining on large-scale corpus, BERT embedding of a sentence contains its rich contextual information.

For the input sentence $X = \{x_1, x_2, \ldots, x_n\}$, the semantic embedding $E^S = \{e_1^S, e_2^S, \ldots, e_n^S\}$ is formulated as follows:

$$e_i^S = BERTEncoder(x_i) \tag{1}$$

where $E^S \in \mathbb{R}^{l \times 768}$ and $l$ is the length of the input sentence.

**Glyph Embedding**. Compared with recent research, we pay more attention to the radical meaning and the frame structure of Chinese characters. The Glyce encoder [20] which utilizes a Tianzige-CNN structure is adopted as the visual information encoder. To some extent, it fits the origin of Chinese characters better than other methods, such as stroke sequence [8, 18] or object detection [12]. Taking account of the evolution of Chinese characters and their current form, we select two fonts (楷书，kǎishū) in both traditional and simplified Chinese finally.

Given an input sentence $X = \{x_1, x_2, \ldots, x_n\}$, we define its glyph embedding $E^G = \{e_1^G, e_2^G, \ldots, e_n^S\}$ as follows:

$$e_i^G = GlyphEncoder(x_i) \tag{2}$$

where $E^G \in \mathbb{R}^{l \times d_g}$ and $d_g$ is the dimension of glyph embedding.

**Pinyin Embedding**. For phonetic errors, following recent research [12, 28], we use *Pinyin* to represent pronunciation. Previous works usually treated the pinyin as a whole token [18, 36], ignoring its internal components of the initial, final, and tone. Although two pinyin strings are different, they may still be similar in phonetic due to the same initials or finals. For example, the pinyin string of "插 (insert)" and "擦 (wipe)" are "chā" and "cā" respectively. Thus, we separate the integral pinyin of each character into its initial, final and tone.

When given an input sentence $X = \{x_1, x_2, \ldots, x_n\}$, we first convert each pinyin of character $x$ in $X$ to its corresponding separation form. Then we concatenate pinyin of each character's integral form and separation form as its pinyin representation. Finally, we apply a CNN network of width 2 and a max pooling function to extract the phonetic information as:

$$p_i^S = [p_{x_i}^{Initial} \cdot p_{x_i}^{Final} \cdot p_{x_i}^{Tone}] \tag{3}$$

$$p_i = [p_i^I \cdot p_i^S] \tag{4}$$

$$e_i^P = CNN(p_i) \tag{5}$$

where $[\cdot]$ means concatenate operation between embeddings, $p_i^I$ and $p_i^S$ is the integral pinyin string and separation pinyin form of $x_i$ respectively, $p_i$ is the final pinyin representation, $e_i^P \in \mathbb{R}^{l \times d_p}$ and $d_p$ is the dimension of pinyin embedding.

---

[4]https://en.wikipedia.org/wiki/Pinyin

*3.2.2   Output Layer.* The original BERT embedding integrated with glyph embedding and pinyin embedding is fed into the 2-layer Transformer encoder [27]:

$$e_i = \text{MLP}([e_i^S \cdot e_i^G \cdot e_i^P]) \tag{6}$$

$$\tilde{e}_i = \text{LayerNorm}(e_i) \tag{7}$$

$$h_i = \text{Transformer}(\tilde{e}_i) \tag{8}$$

where MLP is a linear layer, $h_i \in \mathbb{R}^{d_t}$ and $d_t$ is the output dimension of the Transformer encoder.

Following Meng et al. [20], we combine the loss of token classification task and glyph classification task as the final training objective. The training objective $\mathcal{L}$ is given as follows:

$$\mathcal{L}_{glyph} = -\log p(z|x) \tag{9}$$

$$= -\log \text{softmax}(W \times h_{image}) \tag{10}$$

$$\mathcal{L}_{csc} = -\sum_{i=1}^{n} \log P(\hat{y}_i = y_i | X) \tag{11}$$

$$\mathcal{L} = (1 - \lambda)\mathcal{L}_{csc} + \lambda \mathcal{L}_{glyph} \tag{12}$$

where $z$ is the label of font image $x$, $h_{image}$ is the hidden state from CNNs in the Glyce. $\hat{y}_i, y_i$ are the prediction and label of $X$ respectively. $\lambda$ controls the trade-off between the CSC classification objective and the auxiliary image classification objective, where $\lambda \in [0, 1]$.

## 3.3   Error Consistent Masking Strategy

We propose an Error Consistent Masking (ECM) strategy to create misspelling. In Zhang et al. [36] and Liu et al. [18], candidates to replace the original character are independently selected from the confusion set. For example, "报导 (bào dǎo)" (report) could be changed into "爆异 (bào yì) (explosive difference)", where "报 (bào)" (newspaper) and "爆 (bào)" (burst) hold the same pronunciation while "导 (dǎo)" (guide) and "异 (yì)" (difference) look similar.

However, people use only one input method when inputting sentences in fact. Especially under the pinyin input method, people are primarily used to inputting continuous characters. Previous masking strategies break the consistency of spelling errors in one sentence.

In order to better simulate the real input scene, we simulate the idea of generating candidates by input method based on N-gram method and build an N-gram level confusion set:

- Step 1. We collect 2-gram, 3-gram, 4-gram spans from large-scale corpus[5].
- Step 2. We match spans in different gram sets using the single character pinyin confusion set which only includes phonetic-similar candidates.
- Step 3. To further robust, we segment large-scale corpus sentences, collecting medium and high-frequency phrases. These phrases are converted to pinyin[6] and reconverted to characters[7] by tools. We collect the candidate phrases provided in this process.

Therefore we get an N-gram confusion set of which size is 57,363 and it can provide candidates if given a fragment. Notably, this N-gram confusion set contains not only phrases with similar pinyin but also high frequency input segments with similar voices. For example, "一年 (yī nián) (one year) - 意念 (yì niàn) (mind)", both of which are phrases and "四类 (sì lèi) (four types) - 室内 (shì nèi) (indoor)", where the former is a high-frequency input segment, and the latter is a phrase.

In addition to the N-gram confusion set, we divide the prior confusion set of a single character into phonetic-similar confusion set and morphological-similar confusion set, according to the error

---

[5]https://github.com/brightmart/nlp_chinese_corpus
[6]https://pypi.org/project/pypinyin/
[7]https://github.com/letiantian/Pinyin2Hanzi

type annotated in the prior one. Overall, we get three confusion sets and then do the following operations:

(1) determining the type of error in a sentence according to 30% of the pronunciation, 30% of the shape, 20% of the random and 20% of the unchanged, referring to Liu et al. [17].
(2) selecting single or continuous characters randomly and replacing them with the same type from confusion set or the N-gram confusion set. Continuous errors are selected only when the error type is sound.

The total number of chosen characters is limited to less than 15% of the sentence length according to the mask ratio in Devlin et al. [4].

## 4 USER DICTIONARY GUIDED DOMAIN ADAPTIVE CSC TASK

Most of the existing benchmark SIGHAN [26, 32, 35] focuses on the general domain. However, it is common for spellers to function within a specific domain. Unfortunately, based on our information, there is no domain specific CSC evaluation benchmark. Hence we annotate three domain specific CSC datasets. Meanwhile, to extend the CSC task to the practical application, we introduce the zero-shot CSC task which aims to be adaptive to specific domains without extra training data. Accordingly, we also propose a simple but effective framework to handle this problem with the help of the user dictionary.

### 4.1 Domain Specific Dataset Construction

We collect the raw sentences from publicly available websites, consisting of three typical domains: Law[8], medical treatment (Med)[9] and official document writing (Odw)[10]. The legal data is composed of both the question stems and options of the multiple-choice questions in the judicial examination. The medical treatment data consists of sentences in QA pairs from online network consultations. The official document writing data comprises news, policies and national conditions officially reported by the state. We use all sentences in the public dataset as candidates for legal and medical treatment domains. As for the official writing domain, we crawl 1,000 official documents and split their content into single sentences using natural separators such as periods and question marks.

We only retain sentences with more than 5 Chinese characters from these raw data and sample them randomly as the candidates. Then we enlist five native volunteers to copy the raw sentences and create possible spelling errors in the three domains. The guideline for manufacturing errors is as follows:

- Annotators should fully consider the context and the meaning of the selected characters when making spelling errors.
- Apart from single-character errors, continuous phrase-level errors should also be considered.
- As SIGHAN15, error characters in each sentence should be no more than 15%, with a maximum total of 7.

Finally, the annotation results are merged and manually proofread. The dataset holds the following three appealing properties.

***Reduce Subjective Bias.*** Five annotators label each sentence and the results are then proofread manually to minimize the annotation deviation caused by subjective factors.

---

[8]http://cail.cipsc.org.cn:2020/
[9]https://github.com/CBLUEbenchmark/CBLUE
[10]http://www.gov.cn/

|                          | Law         | Med         | Odw         | SIGHAN15  |
|--------------------------|-------------|-------------|-------------|-----------|
| # Error sents/Sents      | 1,314/2,460 | 1,699/3,500 | 1,259/2,220 | 542/1,100 |
| Min.Len                  | 12          | 11          | 9           | 5         |
| Max.Len                  | 120         | 127         | 161         | 108       |
| Avg.Len                  | 30.5        | 50.1        | 41.2        | 30.7      |
| # Continuous error sents | 229         | 253         | 265         | 51        |
| # Annotators per sent    | 5           | 5           | 5           | 1         |
| PPL                      | 30.71       | 32.45       | 26.67       | 30.76     |

Table 1. Statistics of different datasets.

***Reduce Bias of Input Methods***. To simulate people's input as much as possible, these annotators are required to use different input methods including Microsoft Pinyin[11], Google Pinyin[12], Tencent Pinyin[13], Baidu Pinyin[14] and Sougou Handwriting[3]. Meanwhile, we ensure that an annotator keeps one input method to label a sentence, just like the real input scenario.

***Continuous Misspelling***. Continuous spelling errors are also common in real life. We create this type of misspelling via phrase-level phonetic or visual similar replacement, Internet buzzword replacement and different input intervals in continuous input.

We supplement correct sentences to make the correction rate close to SIGHAN15. The statistic of this dataset is shown in Table 1. We introduce the SIGHAN15 test dataset to make an intuitive comparison. Compared with the general domain, we choose some domains with representative terms, which can be seen from the differences in the PPL metric in the above table. As can be seen, the sentences in Law, Med and Odw datasets are much more than that in the SIGHAN15 test dataset. The proportion of error sentences in our built datasets is similar to that in the SIGHAN15 test dataset. The average length of sentences in the Med dataset is 50.1, which is longer than others. In addition, the number of continuous error sents is more than that in the SIGHAN15, which can better restore the real input scene.

## 4.2 User Dictionary Guided Framework

We propose a simple and effective framework named **U**ser **D**ictionary (UD) for this domain adaptive CSC task, which is also suitable for classified-based check spellers. As Figure 2 shows, UD gives the final prediction according to the probability matrix and user dictionary jointly. The ECSpell could be replaced with other classified-based check spellers. Compared with obtaining training corpus of different domains, collecting their corresponding dictionaries saves more time and economy. In addition, UD can obtain the ability to adapt to various application scenarios by switching user dictionaries. Our experiment shows that the error correction result will be greatly improved with the help of these dictionaries.

We can first achieve the top-$k$ rank log-probability matrix $O$ generated by any token classification based speller. Then UD will update the prediction result according to the given user dictionary. Specifically, we can select a token $y_j^{\mathcal{P}}$ in each column $j$ of $O$ and connect them to form a candidate correction path $\mathcal{P} = \{y_1^{\mathcal{P}}, y_2^{\mathcal{P}}, \cdots, y_n^{\mathcal{P}}\}$, where $y_j^{\mathcal{P}} \in O_j$. In this way, the score of $\mathcal{P}$ computed by

---

[11]https://en.wikipedia.org/wiki/Microsoft_Pinyin_IME
[12]https://en.wikipedia.org/wiki/Google_Pinyin
[13]http://qq.pinyin.cn/
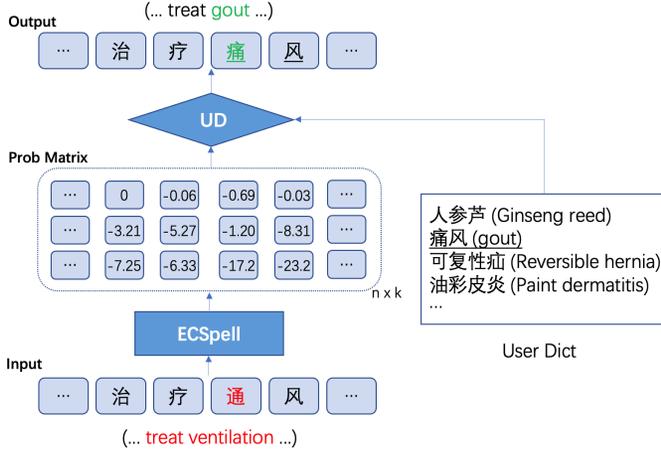[14]https://srf.baidu.com/

Fig. 2. Overview of the UD framework flow chart. Underline characters are included in the user dictionary.

the original speller is:

$$score_O^{\mathcal{P}} = \sum_{j=1}^{n} score(y_j^{\mathcal{P}}) \tag{13}$$

where $score(y_j^{\mathcal{P}})$ is the related log-probability value from $O$. Since the output probability of token classification based models is independent token-by-token, the optimal path from Eq. 13 is simply to pick the first candidate of each position. Our UD module adopts the following rules to encourage paths including more dictionary terms:

**Raw Span Match (RSM).** If the span in the input sentence can be found in the dictionary, we regard this part as correct and keep it unchanged.

**Altered Span Match (ASM).** In the aspect of alternative paths containing dictionary terms, we reward them as follows:

$$score_D^{\mathcal{P}} = l_{match}^{\mathcal{P}} \tag{14}$$

where $l_{match}^{\mathcal{P}}$ stands for the character number of altered span matched the dictionary in $\mathcal{P}$.

Finally, the optimal candidate path is selected as follows:

$$\mathcal{P}_{opt} = \arg\max_{\mathcal{P}}(score_O^{\mathcal{P}} + \eta \times score_D^{\mathcal{P}}) \tag{15}$$

where $\eta$ is a hyper-parameter and we set it to 4 empirically.

For a sentence $X$ of length $n$, the total number of candidate paths is $k^n$. Besides RSM, we set the minimum and maximal threshold for each token candidate to prune. On one hand, when the probability of a token candidate exceeds the maximal threshold, we fix that token as the predicted result. On the other hand, we discard a token if its score is below the minimum threshold. In the following experiment, the minimum and maximal thresholds are set to -11 and -0.001, respectively.

After pruning, the average numbers of candidate paths in different spellers are shown in Table 2. We can find that the number varies greatly. For weak models, there are still too many candidate paths. It will consume much time if the greedy algorithm obtains the global optimum. Hence we use beam search to reduce the search space. To balance the speed and effect, we set the beam size to 20 and the hyperparameter $k$ to 5. As can be seen from Table 2, for our model ECSpell[#], it is
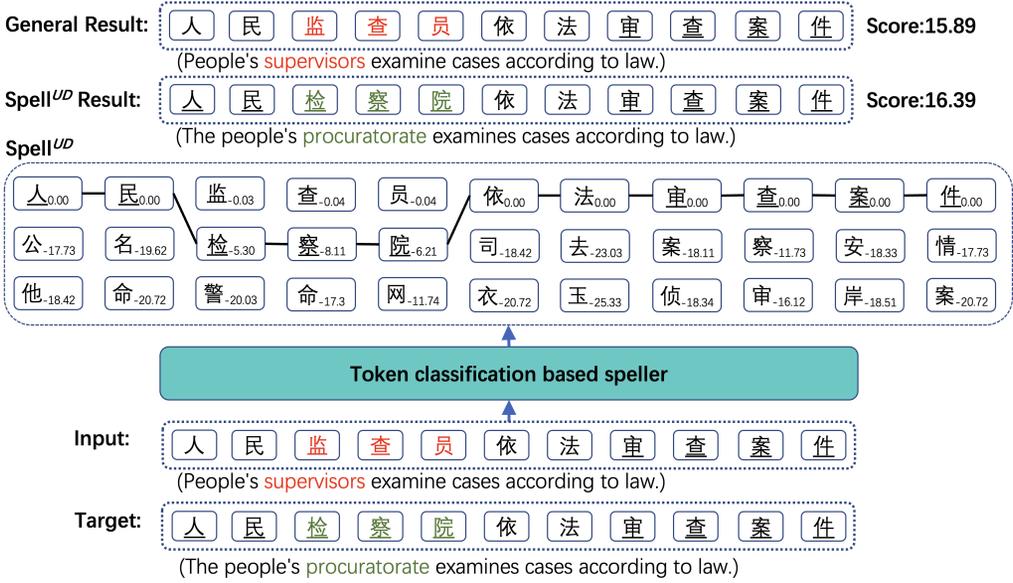
General Result: 人 民 监 查 员 依 法 审 查 案 件　Score:15.89
(People's supervisors examine cases according to law.)

Spell$^{UD}$ Result: 人 民 检 察 院 依 法 审 查 案 件　Score:16.39
(The people's procuratorate examines cases according to law.)

Spell$^{UD}$

人$_{0.00}$ 民$_{0.00}$ 监$_{-0.03}$ 查$_{-0.04}$ 员$_{-0.04}$ 依$_{0.00}$ 法$_{0.00}$ 审$_{0.00}$ 查$_{0.00}$ 案$_{0.00}$ 件$_{0.00}$

公$_{-17.73}$ 名$_{-19.62}$ 检$_{-5.30}$ 察$_{-8.11}$ 院$_{-6.21}$ 司$_{-18.42}$ 去$_{-23.03}$ 案$_{-18.11}$ 察$_{-11.73}$ 安$_{-18.33}$ 情$_{-17.73}$

他$_{-18.42}$ 命$_{-20.72}$ 警$_{-20.03}$ 命$_{-17.3}$ 网$_{-11.74}$ 衣$_{-20.72}$ 玉$_{-25.33}$ 侦$_{-18.34}$ 审$_{-16.12}$ 岸$_{-18.51}$ 案$_{-20.72}$

**Token classification based speller**

Input: 人 民 监 查 员 依 法 审 查 案 件
(People's supervisors examine cases according to law.)

Target: 人 民 检 察 院 依 法 审 查 案 件
(The people's procuratorate examines cases according to law.)

Fig. 3. An example of inference process with UD. "监查员" (supervisors) in this sentence should be corrected to "检察院" (procuratorate). Underline characters, "人民检察院" (people's procuratorate) and "审查案件" (review cases), are included in the user dictionary. The scores of general result and UD inference result are calculated with parameters by Eq.15.

| Model | Law | Med | Odw |
|---|---|---|---|
| SM BERT | 783.45 | 1258.83 | 106.67 |
| SpellGCN | 77.29 | 507.04 | 107.85 |
| BERT$^{\#}$ | 4.50 | 149.85 | 2.89 |
| ECSpell$^{\#}$ | 3.89 | 38.49 | 2.64 |

Table 2. The average candidate paths per sentence. "$^{\#}$" denotes the model is additionally pretrained.

usually able to select the optimal candidate path while this chance in weaker models like SM BERT and SpellGCN is reduced seriously.

Figure 3 shows an example. The input sentence is "人民监查员依法审查案件" (People's supervisors examine cases according to law) where "监查员" (supervisors) should be corrected to "检察院" (procuratorate). These potential errors are too confusing to capture by general models. Nevertheless, UD framework rescores each path according to Eq. 15 and then gives the prediction consistent with the target sentence.

## 5 EXPERIMENTS

### 5.1 Data

*5.1.1 Pretraining.* We collect 38.1 million sentences from news2016[5] and wiki2019zh[5] for pre-training. Compared with Liu et al. [18], our pre-training is quite lightweight.

*5.1.2 General Task.* We evaluate our general model on the widely-used SIGHAN15 benchmark. The corpus for fine-tuning consists of the SIGHAN training data (6,476 samples) [26, 32, 35] and

|  | Law | Med | Odw |
|---|---|---|---|
| Num | 9896 | 18749 | 12509 |
| Min.Len | 2 | 2 | 2 |
| Max.Len | 29 | 13 | 14 |
| Avg.Len | 7.9 | 4.2 | 2.8 |
| # Error Related / All | 403/2,460 | 1,155/3,500 | 1,158/2,220 |

Table 3. Statistics of dictionaries of law, medicine and official document writing. "# Error Related / All", represents the number of dictionary related error sentences / all sentences in each dataset.

automatically generated pseudo data (271,329 samples) [29]. Follow previous works [2, 11, 14, 37], we convert the traditional Chinese character in SIGHAN dataset to simplified Chinese form using OpenCC[15].

*5.1.3 Domain Adaptive Task.* Regarding the domain specific CSC evaluation, the dataset we built is employed. There is no extra data for training. Any dictionaries can be used to guide the general speller. In this paper, we adopt the public Tsinghua University open Chinese dictionaries[16] for law and medical treatment domain. For the official document writing domain, we first crawl 40,000 additional official documents and collect all phrases via the word segment tool. Then we do the same operations on our pretraining data and obtain general phrases. Next, we retain the top 10k/15k phrases with the most occurrences in general/Odw phrases as the candidate general/Odw dictionary. In addition, we adopt the AutoPhrase [23] to mine domain focused phrases to supplement original Odw dictionary. Finally, phrases appearing in the general dictionary are removed from the Odw dictionary. The statistics related to the dictionary is shown in Table 3.

## 5.2 Baselines

We compare our model with these typical baselines:

**BERT[4]** The basic BERT token classification model.
**FASPell[11]** This model utilizes a denoising autoencoder to generate candidates from context.
**SM BERT[37]** This model uses a soft-masked strategy to combined the detection module and correction module.
**SpellGCN[2]** This model incorporates pronunciation and shape similarity graphs into BERT model via GCN.
**DCN[28]** This model considers connection between two adjacent characters.
**PLOME[18]** This model predicts the token and pinyin at the same time.
**REALISE[33]** This model adopts a adaptive gate mechanism to fuse the phonetic and visual information.
**PHMOSpell[12]** This model integrates pinyin and glyph representations with a multi-modal method.
**2ways[14]** This model uses an adversarial strategy to optimize the robustness.
**TtT[15]** This model adopts focal loss penalty strategy to alleviate the class imbalance problem considering that most of the tokens in a sentence are not changed.
**Copy BERT** This model combines a copy behaviour of terminology [5] with BERT for token classification method at training time to handle the domain adaptive CSC task.

---

[15]https://github.com/BYVoid/OpenCC
[16]http://thuocl.thunlp.org/

| Method | Detection Level | | | | Correction Level | | | |
|---|---|---|---|---|---|---|---|---|
| | Acc. | Pre. | Rec. | F1 | Acc. | Pre. | Rec. | F1 |
| FASPell[†] | 67.6 | 67.6 | 60.0 | 63.5 | 66.6 | 66.6 | 69.1 | 62.6 |
| SM BERT[†] | 80.9 | 73.7 | 73.2 | 73.5 | 77.4 | 66.7 | 66.2 | 66.4 |
| SM BERT | 80.5 | 71.2 | 75.1 | 73.1 | 78.5 | 67.3 | 71.0 | 69.1 |
| SpellGCN[†] | 74.8 | 74.8 | 80.7 | 77.7 | 72.1 | 72.1 | 77.7 | 75.9(74.8) |
| SpellGCN | 84.0 | 76.9 | 78.9 | 77.9 | 82.7 | 74.4 | 76.3 | 75.4 |
| DCN-P[#†] | - | 77.1 | 80.9 | 79.0 | - | 74.5 | 78.2 | 76.3 |
| PLOME[#†] | - | 77.4 | 81.5 | 79.4 | - | 75.3 | 79.3 | 77.2 |
| REALISE[#†] | 84.7 | 77.3 | 81.3 | 79.3 | 84.0 | 75.9 | 79.9 | 77.8 |
| 2ways[#†] | - | - | - | 80.0 | - | - | - | 78.2 |
| BERT | 79.4 | 69.8 | 76.6 | 73.0 | 78.4 | 67.9 | 74.5 | 71.1 |
| BERT[#] | 84.7 | 76.0 | 81.0 | 78.4 | 84.0 | 74.7 | 79.5 | 77.0 |
| ECSpell | 83.4 | 76.4 | 79.9 | 78.1 | 82.4 | 74.4 | 77.9 | 76.1 |
| ECSpell[#] | **86.3** | **81.1** | **83.0** | **81.0** | **85.6** | **77.5** | **81.7** | **79.5** |

Table 4. Performance on the SIGHAN15 test. Best results are in **bold**. "#" denotes the model is additionally pretrained. "†" represents the results we quoted. Specially, if the precision and recall are correct, F-score of SpellGCN cited in its essay should be 74.8.

| Method | Detection Level | | | | Correction Level | | | |
|---|---|---|---|---|---|---|---|---|
| | Acc. | Pre. | Rec. | F1 | Acc. | Pre. | Rec. | F1 |
| SpellGCN[†] | 83.7 | 85.9 | 80.6 | 83.1 | 82.2 | 85.4 | 77.6 | 81.3 |
| DCN-P[#†] | **94.6** | 88.0 | 80.2 | 83.9 | 83.2 | 87.6 | 77.3 | 82.1 |
| PHMOSpell[#†] | - | **90.1** | 72.7 | 80.5 | - | **89.6** | 69.2 | 78.1 |
| TtT[†] | 82.7 | 85.4 | 78.1 | 81.6 | 81.5 | 85.0 | 75.6 | 80.0 |
| ECSpell | 82.9 | 85.7 | 78.4 | 81.9 | 82.0 | 85.4 | 76.6 | 80.7 |
| ECSpell [#] | 86.3 | 88.3 | **83.2** | **85.7** | **85.6** | 88.1 | **81.7** | **84.8** |

Table 5. Performance on the SIGHAN15 test evaluated by the official tools. Best results are in **bold**. "#" denotes the model is additionally pretrained. "†" represents the results we quoted.

## 5.3 Settings

Our model is implemented based on huggingface's pytorch implementation of transformers[17]. Concerning pretraining, the training batch size is 2560 and the training step is 150k, and we utilize AdamW [19] optimizer with the learning rate of 5e-5. Instead of training from scratch, we adopt the parameters of Chinese $BERT_{wwm}$ to initialize the Transformer blocks. As to finetuning, the number of fonts in glyce embeddings is 4452 as Meng et al. [20]. The training batch size and the learning rate are set to 128 and 2e-5, respectively. The optimizer is the same as that of pretraining period. In finetuning, there are two ways of initializing the model parameters. One is adopting the weight of Chinese $BERT_{wwm}$ directly and the other is initialized with the weight obtained from the pretraining period mentioned before.

---

[17]https://github.com/huggingface/transformers

## 5.4 Evaluation Metrics

To evaluate the performance, we adopt the widely used metrics following Hong et al. [11] and Zhang et al. [37]. Compared with character-level metrics, sentence-level metrics are more rigorous since they measures the ability to detect and correct the spelling errors for the entire sentence. Metrics includes accuracy, precision, recall and F1.

Besides, we also report the experimental result of general task evaluated by official tool. In order to facilitate the following research, we also rewritten the original java code version of the official tool into Python version and verify it on all the relevant experiments mentioned in this paper[18].

## 5.5 Main Results

*General Task.* The performance on the SIGHAN15 test set is shown in Table 4 and Table 5. The former is evaluated by the method in Hong et al. [11] while the latter is computed with the official tool. As can be seen, our model ECSpell# outperforms all previous methods. Table 4 shows it achieves 2.5% gain on correction F1 score compared with the basic BERT#. From the results, the effect of pretraining for CSC task is very obvious, especially BERT#. Even without pretraining, ECSpell can compete against DCN-P#. In Table 5, the result is consistent. Our ECSpell# achieves the best performance. It is noticed that PHMOSpell# performs well in precision but not in recall. We deduce that the multi-modal features they incorporated are too strict to mine the potential errors. By contrast, ECSpell# has a better trade-off between precision and recall, and it achieves better performance in terms of the ultimate correction-level F1 score. The overall results indicate that our enhanced embedding representation and error consistent masking strategy take effect indeed. Notably, Li et al. [14] applies a complex adversarial training mechanism to largely strengthen the model robustness. We believe that the ECSpell will be further improved if we use the same training mechanism.

*Domain Adaptive Task.* The upper, middle and lower parts of Table 6 correspond to the results on the Law, Med, and Odw datasets, respectively. As shown in this table, general models worsen to some extent in domain specific scenarios. Relying on the copy mechanism, Copy BERT obtains a certain ability to retain domain terminology. However, its lower correction-F1 score denotes that the copy mechanism can only hold terms rather than find potential errors. With the help of our user dictionary guided inference module UD, the performance of all the tested spellers rises although there is no extra data for fine-tuning. In the detection-f1 level, ECSpell# does not comprehensively exceed BERT#. Nevertheless, in the correction-f1 level, the former performs better. It illustrates that ECSpell# has stronger ability for accurate error correction.

It can be seen that UD is good at identifying and correcting term-in-dictionary errors. Take Table 7 as an example. The dictionary term "剂量" (dose) is wrongly changed without UD while another dictionary term "甲苯咪唑" (mebendazole) can be corrected only with UD. The experimental results also verify that with the help of shifting user dictionaries of different domains, the model could have a good adaptability to its current applying domain field when checking spelling errors.

Moreover, we find that UD tends to work more effectively and efficiently for stronger general models. Stronger general models are more likely to contain terms in candidate paths and the number of these paths would be less than weaker ones as Table 2 shows. As a result, the search space of such a model is less which benefits its efficiency.

Overall, our ECSpell model performs best among all the tested spellers regardless of whether it is guided by UD or not. Extraordinarily, ECSpell$^{UD}$ can approach the performance on the general benchmark.

---

[18]The python version code of official evaluation tool is released along with our source code.

| Model | Detection | | | | Correction | | | |
|---|---|---|---|---|---|---|---|---|
| | Acc. | Pre. | Rec. | F1 | Acc. | Pre. | Rec. | F1 |
| Copy BERT# | 80.2/59.8 | 76.1/72.4 | 65.9/59.8 | 70.6/65.5 | 74.7/52.1 | 64.2/63.1 | 55.6/52.1 | 59.6/57.1 |
| SM BERT + UD | 65.7/41.2 +0.7/+8.2 | 53.5 /52.9 +1.7/+**8.1** | 48.3/41.2 +1.0/+8.2 | 50.8/46.3 +1.3/+8.3 | 58.5/32.0 +0.6/+10.7 | 38.4/41.1 +1.4/+**11.7** | 34.7/32.0 +0.9/+10.7 | 36.5/36.0 +1.1/+11.2 |
| SpellGCN + UD | 64.0/35.7 +1.4/+7.5 | 53.4/50.2 +1.9/+5.9 | 44.0/35.7 +2.4/+7.5 | 48.3/41.7 +2.1/+7.1 | 54.9/23.6 +1.7/+10.6 | 32.6/33.1 +**3.0**/+11.4 | 26.9/23.6 +3.0/+10.6 | 29.5/27.5 +**3.0**/+11.2 |
| BERT# + UD | 80.7/56.8 +1.7/+**9.7** | 76.9/69.6 +**2.1**/+6.8 | 65.5/56.8 +**3.1**/+**9.7** | 70.8/62.6 +**2.6**/+8.7 | 77.1/50.9 +**1.8**/+**11.6** | 69.0/62.3 +2.4/+9.9 | 58.8/50.9 +**3.2**/+**11.6** | 63.5/56.0 +2.8/+**11.0** |
| ECSpell# + UD | 80.2/57.1 +**1.9**/+9.5 | 76.5/70.3 +1.7/+5.5 | 65.0/57.1 +2.8/+9.5 | 70.3/63.0 +2.3/+6.4 | 77.4/52.6 +1.6/+10.9 | 70.5/64.8 +1.7/+8.1 | 59.9/52.6 +2.7/+10.9 | 64.8/58.1 +2.4/+9.8 |
| Copy BERT# | 72.5/46.8 | 59.6/60.4 | 53.9/46.8 | 56.6/52.7 | 67.9/38.4 | 49.1/49.6 | 44.4/38.4 | 46.7/43.3 |
| SM BERT + UD | 61.5/35.6 +0.9/+6.4 | 43.6/49.1 +0.8/+3.4 | 41.0/35.6 +4.1/+6.4 | 42.3/41.3 +2.4/+5.6 | 53.8/20.9 +1.2/+7.4 | 26.7/28.6 +0.1/+3.7 | 25.1/20.8 +4.1/+6.7 | 25.9/24.1 +2.1/+5.6 |
| SpellGCN + UD | 56.7/25.8 +1.3/+6.9 | 35.6/38.2 +0.7/+3.1 | 32.7/25.8 +4.8/+7.0 | 34.1/30.8 +2.8/+5.8 | 48.9/11.6 +1.2/+7.3 | 18.2/17.2 +0.3/+4.1 | 16.8/11.6 +4.0/+7.1 | 17.5/13.9 +2.1/+6.0 |
| BERT# + UD | 79.6/55.0 +2.2/+7.1 | 74.5/72.9 +**1.5**/+5.2 | 61.4/55.0 +**5.8**/+7.1 | 67.3/62.7 +**3.7**/+7.6 | 76.0/48.3 +**1.8**/+8.7 | 65.6/64.1 +**0.8**/+5.8 | 54.0/48.3 +3.6/+8.7 | 59.2/55.1 +**3.7**/+7.8 |
| ECSpell# + UD | 79.1/54.2 +**1.9**/+**9.6** | 75.2/71.7 +0.6/+**5.9** | 60.6/54.2 +5.2/+**9.6** | 67.1/61.7 +3.3/+**8.8** | 76.0/49.6 +1.5/+**10.5** | 67.3/65.6 +0.3/+**6.6** | 54.2/49.6 +**4.4**/+**10.5** | 60.0/56.5 +2.8/+**9.1** |
| Copy BERT# | 79.3/67.7 | 78.6/79.9 | 64.9/67.7 | 71.1/73.3 | 74.5/59.4 | 68.3/70.1 | 56.4/59.4 | 61.8/64.3 |
| SM BERT + UD | 62.4/46.3 +0.7/+6.8 | 52.1/59.9 +1.3/+3.2 | 44.9/46.3 +4.1/+6.7 | 48.2/52.2 +2.9/+5.4 | 55.0/33.2 +1.1/+7.1 | 37.1/42.9 +0.4/+**4.3** | 31.9/33.2 +3.2/+5.5 | 34.3/37.5 +2.0/+5.0 |
| SpellGCN + UD | 60.6/39.5 +0.2/+0.5 | 53.1/59.0 +0.5/+0.6 | 38.4/39.5 +0.5/+0.5 | 44.6/47.3 +0.5/+0.6 | 51.7/23.7 +0.5/+0.8 | 31.6/35.4 +1.0/+0.8 | 22.8/23.7 +0.8/+0.8 | 26.5/28.4 +0.9/+1.0 |
| BERT# + UD | 78.4/64.9 +**4.0**/+7.6 | 79.8/80.0 +**2.8**/+3.4 | 62.6/64.9 +3.2/+7.6 | 70.1/71.7 +**6.8**/+5.9 | 75.9/60.7 +**4.5**/+6.7 | 74.0/74.8 +**1.9**/+2.7 | 58.1/60.7 +4.0/+**6.7** | 65.1/67.0 +3.6/+4.9 |
| ECSpell# + UD | 79.1/66.2 +3.4/+7.0 | 81.4/81.6 +1.0/+2.0 | 63.6/66.2 +**6.5**/+7.0 | 71.4/73.1 +4.4/+5.0 | 76.8/62.5 +3.1/+6.6 | 76.3/77.0 +0.6/+2.7 | 59.6/62.5 +**4.7**/+6.6 | 67.0/69.0 +**3.2**/+**5.2** |

Table 6. Performance on the Law (top), Med (mid) and Odw (bottom) test. "#" denotes the results with additional pretraining. Best results are in **bold**. The values to the left/right of the slash distinguish the result of the entire dataset and the result of the sentences containing dictionary-related errors.

| Input | 患者需要按照剂量服用甲苯米坐片。 |
|---|---|
| | Patients need to take toluidine sitting tablets according to the dose. |
| w/o UD | 患者需要按照计量服用甲苯米坐片。 |
| | Patients need to take toluidine sitting tablets according to the meter. |
| w/ UD | 患者需要按照剂量服用甲苯咪唑片。 |
| | Patients need to take mebendazole tablets according to the dose. |

Table 7. A example of the input and output of ECSpell# with/without UD. We highlight the wrong or correct characters in red or green color. The underlined phrases, "剂量" (dose) and "甲苯咪唑" (mebendazole), are in the dictionary.

| Model | Law | | Med | | Odw | |
|---|---|---|---|---|---|---|
| | D-F1. | C-F1. | D-F1. | C-F1. | D-F1. | C-F1. |
| ECSpell$^{UD}$ | 72.9 | 67.2 | 70.4 | 62.8 | 75.8 | 70.2 |
| w/o glyph | 70.8 | 64.3 | 67.7 | 60.1 | 73.7 | 67.7 |
| w/o pinyin | 71.5 | 65.2 | 68.9 | 61.9 | 74.1 | 68.8 |
| w/o inter-c | 72.2 | 65.8 | 69.1 | 62.2 | 75.0 | 69.6 |
| w/o ECM | 71.3 | 65.1 | 68.5 | 61.3 | 73.5 | 68.5 |
| w/o RSM | 72.8 | 67.0 | 69.1 | 62.4 | 75.5 | 70.0 |
| w/o ASM | 70.2 | 63.9 | 67.5 | 60.7 | 71.6 | 67.2 |

Table 8. Ablation results on our built dataset.

## 5.6 Ablation Study

We explore the contribution of each component in ECSpell$^{UD}$ with pretraining by conducting ablation studies with the following settings: 1) removing the glyph information, 2) removing the pinyin information, 3) removing internal components of pinyin (denoted as inter-c), 4) removing the error consistent masking (ECM) strategy, 5) removing the raw span match (RSM) rule in UD, 6) removing the altered span match (ASM) rule in UD.

The result can be seen in Table 8. For general models, removing glyph, pinyin or internal components of pinyin leads to performance degradation. This proves that the glyph information, as well as the entire pinyin information, i.e., pinyin with its initial, final and tone, could benefit CSC. Then we remove the consistent masking strategy. The drop in performance indicates the importance of this strategy. In order to further study the role of two modules in UD, we remove RSM and ASM, respectively. The decline in experimental results illustrates that they are both effective while ASM is more significant. This phenomenon is also very logical. Although RSM allows UD to reduce the probability of incorrect modification of some terms, the rule is too strict that if a fragment containing spelling errors matches, these errors will persist. By contrast, ASM makes the error correction process more fault tolerant. UD would synthesize the prediction results according to the probability predicted by the model and the user dictionary. To sum up, we can observe that removing any components in ECSpell$^{UD}$ with pretraining brings a performance decline.

## 5.7 Upper Bound Analysis

We attempt to build ideal dictionaries consisting of different proportions of original phrase-level errors in each dataset. Specifically, we use the word segmentation tool to split original error sentences in each dataset. Then we collect different proportions of error phrases from the word segmentation results to form the ideal dictionaries. We use the ideal dictionaries to verify the ability of UD to update in real time according to the dictionary. At the same time, it reflects the upper bound of UD to a certain extent. Figure 4 demonstrates that the performance has been further improved with the increase of dictionary coverage ratio. However, gradually slow growth and the final ceiling indicate that the performance of UD relies on the quality of the dictionary as well. For example, if there are several conflicting terms which is hard to distinguish in the target domain like "检察 (jiǎn chá)" (procurator). Furthermore, "监察 (jiān chá)" (monitor), the effect of our UD will also be limited. Meanwhile, it is also challenging to deal with the independent character error such as "他, 她, 它 (tā, tā, tā)" (he, she, it).

Overall, the proposed method has a particular ability to improve performance in general and domain datasets with a more specific user dictionary.
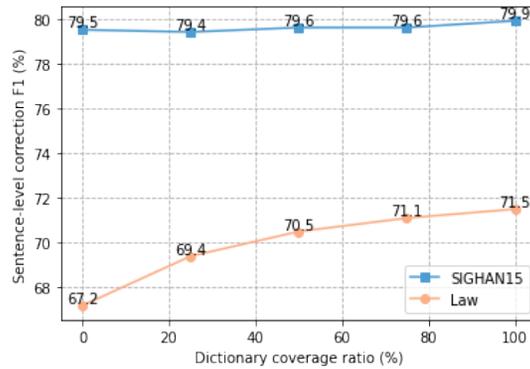
Fig. 4. The upper bound analysis results.

## 6 CONCLUSION

In this paper, we developed an error consistent strategy for unsupervised CSC learning and annotated domain specific datasets. Meanwhile, we proposed a domain adaptive speller ECSpell$^{UD}$ which is composed of a competitive general speller ECSpell and a generic user dictionary guided inference module UD. Experimental results showed that the performance of ECSpell$^{UD}$ on domain-related benchmarks can approach the level of general benchmarks.

We believe our work can be extended in various aspects. On one hand, it is meaningful to collect spelling errors from more different domains especially the emerging areas containing the latest terms. On the other hand, in addition to the domain adaption problem, our model is likely to fit the personalized scenarios.

## ACKNOWLEDGMENTS

## REFERENCES

[1] Zuyi Bao, Chen Li, and Rui Wang. 2020. Chunk-based Chinese Spelling Check with Global Optimization. In *Findings of the Association for Computational Linguistics: EMNLP 2020, Online Event, 16-20 November 2020 (Findings of ACL, Vol. EMNLP 2020)*, Trevor Cohn, Yulan He, and Yang Liu (Eds.). Association for Computational Linguistics, 2031–2040.

[2] Xingyi Cheng, Weidi Xu, Kunlong Chen, Shaohua Jiang, Feng Wang, Taifeng Wang, Wei Chu, and Yuan Qi. 2020. SpellGCN: Incorporating Phonological and Visual Similarities into Language Models for Chinese Spelling Check. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel R. Tetreault (Eds.). Association for Computational Linguistics, 871–881.

[3] Josep Crego, Jungi Kim, Guillaume Klein, Anabel Rebollo, Kathy Yang, Jean Senellart, Egor Akhanov, Patrice Brunelle, Aurelien Coquard, Yongchao Deng, et al. 2016. Systran's pure neural machine translation systems. *arXiv preprint arXiv:1610.05540* (2016).

[4] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *CoRR* abs/1810.04805 (2018).

[5] Georgiana Dinu, Prashant Mathur, Marcello Federico, and Yaser Al-Onaizan. 2019. Training neural machine translation to apply terminology constraints. *arXiv preprint arXiv:1906.01105* (2019).

[6] Jianyong Duan, Lijian Pan, Hao Wang, Mei Zhang, and Mingli Wu. 2019. Automatically Build Corpora for Chinese Spelling Check Based on the Input Method. In *Natural Language Processing and Chinese Computing - 8th CCF International Conference, NLPCC 2019, Dunhuang, China, October 9-14, 2019, Proceedings, Part I (Lecture Notes in Computer Science, Vol. 11838)*, Jie Tang, Min-Yen Kan, Dongyan Zhao, Sujian Li, and Hongying Zan (Eds.). Springer, 471–485.

[7] Lei Gu, Yong Wang, and Xitao Liang. 2014. Introduction to NJUPT Chinese Spelling Check Systems in CLP-2014 Bakeoff. In *Proceedings of The Third CIPS-SIGHAN Joint Conference on Chinese Language Processing*. Association for Computational Linguistics, Wuhan, China, 167–172.

[8] Zijia Han, Chengguo Lv, Qiansheng Wang, and Guohong Fu. 2019. Chinese Spelling Check based on Sequence Labeling. In *International Conference on Asian Language Processing, IALP 2019, Shanghai, China, November 15-17, 2019*, Man Lan, Yuanbin Wu, Minghui Dong, Yanfeng Lu, and Yan Yang (Eds.). IEEE, 373–378. https://doi.org/10.1109/IALP48816.2019.9037652

[9] Eva Hasler, Adrià De Gispert, Gonzalo Iglesias, and Bill Byrne. 2018. Neural machine translation decoding with terminology constraints. *arXiv preprint arXiv:1805.03750* (2018).

[10] Chris Hokamp and Qun Liu. 2017. Lexically constrained decoding for sequence generation using grid beam search. *arXiv preprint arXiv:1704.07138* (2017).

[11] Yuzhong Hong, Xianguo Yu, Neng He, Nan Liu, and Junhui Liu. 2019. FASPell: A Fast, Adaptable, Simple, Powerful Chinese Spell Checker Based On DAE-Decoder Paradigm. In *Proceedings of the 5th Workshop on Noisy User-generated Text (W-NUT 2019)*. Association for Computational Linguistics, Hong Kong, China, 160–169.

[12] Li Huang, Junjie Li, Weiwei Jiang, Zhiyu Zhang, Minchuan Chen, Shaojun Wang, and Jing Xiao. 2021. PHMOSpell: Phonological and Morphological Knowledge Guided Chinese Spelling Check. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. Association for Computational Linguistics, Online, 5958–5967.

[13] Ying Jiang, Tong Wang, Tao Lin, Fangjie Wang, Wenting Cheng, Xiaofei Liu, Chenghui Wang, and Weijian Zhang. 2012. A rule based Chinese spelling and grammar detection system utility. In *2012 International Conference on System Science and Engineering (ICSSE)*. IEEE, 437–440.

[14] Chong Li, Cenyuan Zhang, Xiaoqing Zheng, and Xuanjing Huang. 2021. Exploration and Exploitation: Two Ways to Improve Chinese Spelling Correction Models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 2: Short Papers), Virtual Event, August 1-6, 2021*, Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (Eds.). Association for Computational Linguistics, 441–446.

[15] Piji Li and Shuming Shi. 2021. Tail-to-Tail Non-Autoregressive Sequence Prediction for Chinese Grammatical Error Correction. arXiv:2106.01609 [cs.CL]

[16] Yinghui Li, Qingyu Zhou, Yangning Li, Zhongli Li, Ruiyang Liu, Rongyi Sun, Zizhen Wang, Chao Li, Yunbo Cao, and Hai-Tao Zheng. 2022. The Past Mistake is the Future Wisdom: Error-driven Contrastive Probability Optimization for Chinese Spell Checking. In *Findings of the Association for Computational Linguistics: ACL 2022, Dublin, Ireland, May 22-27, 2022*, Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (Eds.). Association for Computational Linguistics, 3202–3213.

[17] Chao-Lin Liu, Min-Hua Lai, Yi-Hsuan Chuang, and Chia-Ying Lee. 2010. Visually and Phonologically Similar Characters in Incorrect Simplified Chinese Words. In *Coling 2010: Posters*. Coling 2010 Organizing Committee, Beijing, China, 739–747.

[18] Shulin Liu, Tao Yang, Tianchi Yue, Feng Zhang, and Di Wang. 2021. PLOME: Pre-training with Misspelled Knowledge for Chinese Spelling Correction. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. Association for Computational Linguistics, Online, 2991–3000.

[19] Ilya Loshchilov and Frank Hutter. 2019. Decoupled Weight Decay Regularization. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.

[20] Yuxian Meng, Wei Wu, Fei Wang, Xiaoya Li, Ping Nie, Fan Yin, Muyu Li, Qinghong Han, Xiaofei Sun, and Jiwei Li. 2019. Glyce: Glyph-vectors for Chinese Character Representations. *CoRR* abs/1901.10125 (2019).

[21] Elise Michon, Josep Maria Crego, and Jean Senellart. 2020. Integrating Domain Terminology into Neural Machine Translation. In *Proceedings of the 28th International Conference on Computational Linguistics, COLING 2020, Barcelona, Spain (Online), December 8-13, 2020*, Donia Scott, Núria Bel, and Chengqing Zong (Eds.). International Committee on Computational Linguistics, 3925–3937.

[22] Minh Nguyen, Gia H. Ngo, and Nancy F. Chen. 2021. Domain-Shift Conditioning Using Adaptable Filtering Via Hierarchical Embeddings for Robust Chinese Spell Check. *IEEE ACM Trans. Audio Speech Lang. Process.* 29 (2021), 2027–2036.

[23] Jingbo Shang, Jialu Liu, Meng Jiang, Xiang Ren, Clare R Voss, and Jiawei Han. 2018. Automated phrase mining from massive text corpora. *IEEE Transactions on Knowledge and Data Engineering* 30, 10 (2018), 1825–1837.

[24] Kai Song, Yue Zhang, Heng Yu, Weihua Luo, Kun Wang, and Min Zhang. 2019. Code-switching for enhancing NMT with pre-specified translation. *arXiv preprint arXiv:1904.09107* (2019).

[25] Raymond Hendy Susanto, Shamil Chollampatt, and Liling Tan. 2020. Lexically constrained neural machine translation with Levenshtein transformer. *arXiv preprint arXiv:2004.12681* (2020).

[26] Yuen-Hsien Tseng, Lung-Hao Lee, Li-Ping Chang, and Hsin-Hsi Chen. 2015. Introduction to SIGHAN 2015 Bake-off for Chinese Spelling Check. In *Proceedings of the Eighth SIGHAN Workshop on Chinese Language Processing, SIGHAN@IJCNLP 2015, Beijing, China, July 30-31, 2015*, Liang-Chih Yu, Zhifang Sui, Yue Zhang, and Vincent Ng (Eds.). Association for Computational Linguistics, 32–37.

[27] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (Eds.). 5998–6008.

[28] Baoxin Wang, Wanxiang Che, Dayong Wu, Shijin Wang, Guoping Hu, and Ting Liu. 2021. Dynamic Connected Networks for Chinese Spelling Check. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*. Association for Computational Linguistics, Online, 2437–2446.

[29] Dingmin Wang, Yan Song, Jing Li, Jialong Han, and Haisong Zhang. 2018. A Hybrid Approach to Automatic Corpus Generation for Chinese Spelling Check. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun'ichi Tsujii (Eds.). Association for Computational Linguistics, 2517–2527.

[30] Dingmin Wang, Yi Tay, and Li Zhong. 2019. Confusionset-guided Pointer Networks for Chinese Spelling Check. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Florence, Italy, 5780–5785.

[31] Yih-Ru Wang, Yuan-Fu Liao, Yeh-Kuang Wu, and Liang-Chun Chang. 2013. Conditional Random Field-based Parser and Language Model for Tradi-tional Chinese Spelling Checker. In *Proceedings of the Seventh SIGHAN Workshop on Chinese Language Processing, SIGHAN@IJCNLP 2013, Nagoya, Japan, October 14-18, 2013*, Liang-Chih Yu, Yuen-Hsien Tseng, Jingbo Zhu, and Fuji Ren (Eds.). Asian Federation of Natural Language Processing, 69–73.

[32] Shih-Hung Wu, Chao-Lin Liu, and Lung-Hao Lee. 2013. Chinese Spelling Check Evaluation at SIGHAN Bake-off 2013. In *Proceedings of the Seventh SIGHAN Workshop on Chinese Language Processing, SIGHAN@IJCNLP 2013, Nagoya, Japan, October 14-18, 2013*, Liang-Chih Yu, Yuen-Hsien Tseng, Jingbo Zhu, and Fuji Ren (Eds.). Asian Federation of Natural Language Processing, 35–42.

[33] Heng-Da Xu, Zhongli Li, Qingyu Zhou, Chao Li, Zizhen Wang, Yunbo Cao, Heyan Huang, and Xian-Ling Mao. 2021. Read, Listen, and See: Leveraging Multimodal Information Helps Chinese Spell Checking. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*. Online, 716–728.

[34] Junjie Yu and Zhenghua Li. 2014. Chinese Spelling Error Detection and Correction Based on Language Model, Pronunciation, and Shape. In *Proceedings of The Third CIPS-SIGHAN Joint Conference on Chinese Language Processing, Wuhan, China, October 20-21, 2014*, Le Sun, Chengqing Zong, Min Zhang, and Gina-Anne Levow (Eds.). Association for Computational Linguistics, 220–223.

[35] Liang-Chih Yu, Lung-Hao Lee, Yuen-Hsien Tseng, and Hsin-Hsi Chen. 2014. Overview of SIGHAN 2014 Bake-off for Chinese Spelling Check. In *Proceedings of The Third CIPS-SIGHAN Joint Conference on Chinese Language Processing, Wuhan, China, October 20-21, 2014*, Le Sun, Chengqing Zong, Min Zhang, and Gina-Anne Levow (Eds.). Association for Computational Linguistics, 126–132.

[36] Ruiqing Zhang, Chao Pang, Chuanqiang Zhang, Shuohuan Wang, Zhongjun He, Yu Sun, Hua Wu, and Haifeng Wang. 2021. Correcting Chinese Spelling Errors with Phonetic Pre-training. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*. Association for Computational Linguistics, Online, 2250–2261.

[37] Shaohua Zhang, Haoran Huang, Jicong Liu, and Hang Li. 2020. Spelling Error Correction with Soft-Masked BERT. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online, 882–890.