

Chinese Spelling Corrector Is Just Language Learner

Anonymous ACL submission

Abstract

This paper emphasizes the Chinese spelling correction of self-supervised learning, which means there are no annotated errors within the training data. This setting is a pivotal issue that has received broad attention in the community. Our intuition is that humans are naturally good correctors with exposure to monolingual sentences, which contrasts with current unsupervised methods that strongly rely on the usage of confusion sets to produce parallel sentences. In this paper, we demonstrate that learning a spelling correction model is identical to learning a language model from monolingual data alone, with decoding it in a greater search space. We propose *Denoising Decoding Correction (D²C)*, which selectively imposes noise upon the source sentence to solve out the underlying correct characters. Our method largely inspires the ability of language models to perform correction, including both BERT-based models and large language models (LLMs). We show that the self-supervised learning manner generally outstrips the confusion set in specific domains because it bypasses the need to introduce error characters to the training data which can impair the patterns in the target domains. We evaluate our methods on multi-domain datasets Syn-LEMON and ECSpell.

1 Introduction

Chinese spelling correction (CSC) stands as a fundamental task in natural language processing, supporting many downstream applications, e.g. web search (Martins and Silva, 2004; Gao et al., 2010), named entity recognition (Yang et al., 2023b), optical character recognition (Afi et al., 2016; Gupta et al., 2021). Recent studies (Wu et al., 2023a; Liu et al., 2024) show that simply using the supervised signals within parallel sentences to fine-tune pre-trained language models (PLMs) achieves notable results across a series of benchmarks.

However, the great cost of annotation is blamed for the low accessibility of parallel sentences.

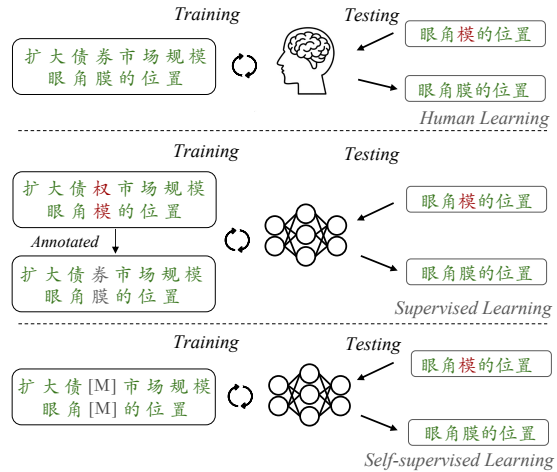


Figure 1: Comparison of human learning, supervised learning, and proposed self-supervised learning process for spelling correction. [M] refers to the mask token.

Therefore, these models remain mediocre in handling massive domains in real applications. This paper thus emphasizes the value of self-supervised learning, where only monolingual data is used to adapt models to specific target domains, which still achieved marginal progress in recent years.

Previous unsupervised methods (Zhao and Wang, 2020; Liu et al., 2021; Li, 2022) focus on synthesizing pseudo parallel sentences, while the supervised signals do not derive from the real distribution but from the confusion set (an empirically constructed word set of common misspelled cases). By replacing certain characters in the original sentences with those in the confusion set, parallel sentences are obtained for fine-tuning the models. However, the gap between the confusion set and the real error patterns in the target domain can induce a high false positive rate (Wu et al., 2023a). This paper raises a bold idea: *Can machine spelling correction learn from monolingual data alone?*

Intriguingly, humans naturally learn to rectify mistakes in a sentence with minimal exposure to parallel data. We give an illustration in Figure 1,

066 which shows that humans only learn to use the
067 correct sentences (monolingual data) in daily life.
068 When encountering a sentence with an error character
069 “模” (*modal*), they can correct it to “膜” (*cornea*)
070 with ease based on their knowledge. In contrast,
071 the machine spelling correction models cannot do
072 this if it isn’t exposed to annotated edit pairs like
073 “模” → “膜” in the training process.

074 In this paper, we demonstrate that a machine
075 spelling corrector can also be learned from solely
076 monolingual data as illustrated at the bottom of Fig-
077 ure 1. The key is to have the model learn semantics
078 rather than character-to-character editing. In light
079 of this, we find that rephrasing models (Liu et al.,
080 2024), where the source sentence will first be en-
081 coded into the semantic space, and then rephrased
082 to the correct sentence, demonstrate this ability. We
083 call this manner self-supervised spelling correction.
084 However, the resultant models still exhibit a low
085 recall tendency.

086 To this end, we propose a novel decoding al-
087 gorithm *Denoising Decoding Correction* (D^2C),
088 which selectively imposes noise upon the source
089 sentence to solve out the underlying correct char-
090 acters. We apply D^2C to two architectures:
091 bidirectional models (represented by ReLM (Liu
092 et al., 2024), the state-of-the-art model in Chinese
093 spelling correction) and auto-regressive models
094 (represented by a series of LLMs (OpenAI, 2023;
095 Touvron et al., 2023; Yang et al., 2023a)). D^2C
096 achieves a significant performance boost over raw
097 language models, trained with monolingual data.

098 To evaluate our method’s performance across
099 different domains, we propose a LEMON (Wu
100 et al., 2023a) training set synthesized by GPT3.5,
101 which only contains monolingual sentences. This
102 dataset permits the fine-tuning and evaluation of
103 self-supervised models in different domains.

104 We summarize the contributions of this paper.

- 105 • We demonstrate that spelling correction can be
106 intrinsically transferred by language modeling on
107 monolingual data.
- 108 • With the proposed novel decoding algorithm,
109 we build an effective self-supervised learning man-
110 ner, allowing the spelling correction models to
111 adapt to target domains at a minimal expense.
- 112 • We build synthetic monolingual training data
113 from LEMON to benchmark the unsupervised do-
114 main adaption in the community.

2 Related Work 115

116 Correcting spelling errors poses a challenging yet
117 crucial task in natural language processing. Early
118 endeavors primarily relied on unsupervised tech-
119 niques, assessing sentence perplexity as a key met-
120 ric (Yeh et al., 2013; Yu and Li, 2014; Xie et al.,
121 2015). Recent methods model spelling correction
122 as a sequence tagging problem that maps each char-
123 acter in a given sentence to its accurate counterpart
124 (Wang et al., 2018, 2019). On top of pre-trained
125 language models (PLMs), some BERT-based mod-
126 els with the sequence tagging training objective are
127 proposed. Zhang et al. (2020) identify the potential
128 error characters by a detection network and then
129 leverage the soft masking strategy to enhance the
130 eventual correction decision. Zhu et al. (2022a) use
131 a multi-task network to minimize the misleading
132 impact of the misspelled characters (Cheng et al.,
133 2020). There is also a line of work that incorpo-
134 rates phonological and morphological knowledge
135 through data augmentation and enhances the BERT-
136 based encoder to assist mapping the error to the
137 correct one (Guo et al., 2021; Li et al., 2021; Liu
138 et al., 2021; Cheng et al., 2020; Huang et al., 2021;
139 Zhang et al., 2021). Recent studies (Wu et al.,
140 2023a; Liu et al., 2024) focus on the rephrasing
141 training objective, which achieves notable results.

142 While in the unsupervised spelling correction do-
143 main, previous works focus on generating pseudo
144 annotated data or detecting error characters with
145 confusion dataset (Zhao and Wang, 2020; Liu et al.,
146 2021; Li, 2022). We are the first to raise a notable
147 self-supervised method with pure monolingual Chi-
148 nese spelling correction data in the community. Our
149 method inherits the ability of PLMs and presents a
150 transferability from language modeling to spelling
151 correction.

3 From Language Modeling to Spelling Correction 152

153 This section serves as the preliminary of our work.
154 The basic effort is to learn spelling correction
155 from monolingual data. We call it self-supervised
156 spelling correction. We first discuss the transfer-
157 ability between language modeling and spelling
158 correction. Second, we point out that rephrasing is
159 the primary training objective for self-supervised
160 spelling correction.
161

162 We discuss the transferability from two angles:
163 (1) The coherence of training objectives between
164 rephrasing spelling correction and language mod-

eling. (2) The knowledge in spelling correction is included in the pre-training process.

3.1 Language Modeling

First, we introduce the training objectives of language modeling.

Given an input sentence $Y = \{y_1, y_2, \dots, y_n\}$ of n characters, (auto-regressive) language modeling seeks to solve the character y_i based on its left context, namely $P(y_i|y_1, y_2, \dots, y_{i-1})$. A spelling correction model can be learned by two dominant objectives, sequence tagging and rephrasing.

3.2 Spelling Correction

Second, we introduce the training objectives of spelling correction.

Spelling correction aims to rectify the underlying misspelled characters in the source sentence. Denote the source sentence as $X = \{x_1, x_2, \dots, x_n\}$ and the target sentence as $Y = \{y_1, y_2, \dots, y_n\}$ and suppose x_i is one of the typos in X , the model learns to correct x_i to y_i based on the entire source sentence, namely $P(y_i|x_1, x_2, \dots, x_n)$.

Tagging The above modeling process can also be viewed as sequence tagging from X to Y . While this has been widely adopted in previous work, a recent study (Liu et al., 2024) shows that tagging-based spelling correction models will lean towards point-to-point editing, thus ignoring the specific context. The final training objective degenerates into $P(y_i|x_i)$.

Rephrasing In comparison, rephrasing (Liu et al., 2024) is shown to be a more effective training objective for spelling correction. It specifically seeks to rewrite the entire sentence after it, namely $P(y_i|x_1, x_2, \dots, x_n, y_1, y_2, \dots, y_{i-1})$. To ensure that the rephrasing process is based on semantics instead of copying, a ratio of noise (e.g. masking with an unused token) is introduced to the source sentence, written as $P(y_i|\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n, y_1, y_2, \dots, y_{i-1})$.

3.3 Self-supervised Spelling Correction

The unsupervised learning setting is naturally akin to language modeling, where the model is trained on monolingual data. Comparing the above two training objectives with language modeling, we find that rephrasing and language modeling are formally the same. In rephrasing, the input sentence is the concatenation of the source and target. This

implies that the spelling correction model can better utilize the knowledge in a pre-trained language model and be transferred from it more easily.

3.4 Knowledge in Vanilla PLMs

The second tiny experiment is to probe the pre-trained knowledge in pre-trained language models. We hypothesize that, after large-scale pre-training, the language model already contains the literal knowledge needed for spelling correction. What we do is to mask the error characters in the source sentence and have the vanilla model (non-fine-tuned one) predict that. From Table 1, we see that the vanilla model can already recall the correct characters in its top- k candidates without any fine-tuning on spelling correction. For example, in about 90% of the cases, the model’s top 10 predictions have covered the correct answer.

	LAW	MED	ODW
Top-20	93.8	88.8	93.8
Top-10	90.8	86.0	90.6
Top-5	86.9	82.0	88.7
Top-1	69.5	66.3	76.8

Table 1: Accuracy of the top- k predictions of MLM from the vanilla BERT model.

It indicates that the pre-trained language models have already possessed the needed knowledge for spelling correction in the form of mask infilling.

3.5 Tagging Model vs. Rephrasing Model

In this section, we evaluate the tagging model and rephrasing model on self-supervised spelling correction and present empirical results in Table 2, which reveals that the rephrasing model hugely surpasses the performances of the tagging models in self-supervised spelling correction.

Monolingual Data It shows that the tagging model trained on monolingual data is powerless. We conjecture that the model only learns point-to-point copying since the source is always the same as its target, thus losing the ability to make modifications to the source sentence. In contrast, the rephrasing model can learn well even with monolingual data. It paves the way for us that pre-trained language models can learn spelling correction from solely monolingual data.

Shuffling of Characters We conduct a second tiny experiment to compare that the rephrasing (Liu

	Method	LAW	MED	ODW
Mono.	Tagging	0.5	0.6	0.5
	Tagging-MFT	10.1	5.3	10.5
	Rephrasing	71.3	68.6	71.9
Shuf.	Tagging	29.5	15.3	16.7
	Tagging-MFT	34.0	17.3	18.9
	Rephrasing	27.6	12.3	13.3

Table 2: Comparison (F1) of tagging and rephrasing on monolingual (self-supervised) / shuffled characters. The details of the models and dataset are in Sec. 6. Mono. means monolingual and Shuf. means shuffled.

et al., 2024) and tagging training objective for self-supervised spelling correction. Specifically, we shuffle the characters in the source and target sentences in correspondence to spoil their semantics. We use these highly noisy samples to fine-tune the rephrasing and tagging models. From Table 2 (Shuf.), we find that the tagging model outperforms the rephrasing model on samples that do not convey semantic information. It inversely verifies that the tagging model learns more of point-to-point editing at the expense of semantics. As aforementioned, it is the semantics that are the key to learn spelling correction from monolingual data. In this paper, therefore, we pick rephrasing as the primary training objective for self-supervised spelling correction.

4 Synthetic LEMON Train Set

LEMON (Wu et al., 2023b) is a multi-domain benchmark that allows us to evaluate the multi-domain generalization of CSC models. However, for each domain, it only has a test set without a train set. We release a synthetic LEMON train set, which GPT3.5 generates. This synthetic train set allows us to evaluate self-supervised models’ performance across multiple domains.

The synthetic data is generated in two steps: (1) Extract the words in each domain. (2) Randomly select words and request GPT3.5 to generate monolingual sentences mimicking the style of specific domains. See our prompts in the appendix 10.

The train set’s information is in the table 3.

GAM	ENC	COT	MEC	CAR	NOV	NEW
2389	2489	1707	2222	2381	3669	4273

Table 3: Number of sentences in each training set.

5 Method

In this section, we propose an enhanced decoding method to unleash the potential of pre-trained language models further. We also propose a method that uses a confusion dataset to upgrade the recall score as another option.

5.1 Two Rephrasing Architectures

Our method focuses on rephrasing-based spelling correction, which can be achieved in two architectures, non-auto-regressive rephrasing, and auto-regressive rephrasing.

Auto-regressive models Auto-regressive model is the primary choice to generate the rephrasing following the input sentence, represented by GPT-like models (Brown et al., 2020) and large language models (LLMs).

To improve the quality of rephrasing, it is an easy yet effective way to mask a ratio of characters in the source sentence with an unused token. In this paper, we denote the masked source sentence as $\tilde{Y} = \{\tilde{y}_1, \tilde{y}_2, \dots, \tilde{y}_n\}$.

ReLM Rephrasing Language Model (ReLM) (Liu et al., 2024) is the current state-of-the-art spelling correction model based on BERT (Devlin et al., 2019). It rephrases the source sentence by infilling the mask slots. Specifically, the model is fed with the concatenation of the source sentence and a sequence of mask tokens. Due to the bidirectional nature of BERT, the rephrasing process can be written as $P(y_i | \tilde{y}_1, \tilde{y}_2, \dots, \tilde{y}_n, m_1, m_2, \dots, m_n)$, where m_i refers to the mask token. As opposed to auto-regressive models, ReLM predicts all characters at once.

5.2 Denoising Decoding Correction

In self-supervised spelling correction, where the source sentence equals the target sentence, the resultant model trained with rephrasing still suffers from a low recall when testing on real sentences because there is no mask token. Therefore, we can introduce mask tokens to test sentences. A more severe situation happens when there are multiple errors in one sentence. The cascade effect of errors makes it even harder to correct the error characters. To this end, we propose a novel decoding algorithm, where we actively introduce noise to the source sentence and encourage the model to recall more candidates. Since the mask operation in the inference stage is consistent with that in the

training stage of rephrasing, the model’s correction capability can be boosted. We call this method *Denosing Decoding Correction (D²C)*.

Concretely, we first mask the characters in the source sentence from the left side during each iteration if the character’s confidence is lower than β (0.995). The character in such a position is regarded as a potential error. To determine which character to update, we send this sentence to the model and figure out whether the original character appears in its **top- k** candidates. If it does, we remain the original character, else we record the new character and its confidence if this confidence is bigger than a **threshold** ϵ . After each iteration, we choose the character with the biggest confidence recorded before and update the original sentence with it. We do iterations continually until there is nothing to update after an iteration. Note that once a character is updated, the confidence of the other characters will change correspondingly. As the error number decreases, the challenges associated with multi-typo spelling correction will also diminish. So this iterative decoding is robust to multiple errors.

We notice that picking a character with the biggest confidence in each iteration costs a large decoding overhead. Given that there is always a small number of errors in a sentence, we rank the characters in the sentence by their confidence from the lowest to highest, mask the top α of them respectively, and send the sentence to the model. Figure out whether the original character appears in its top- k candidates. If it does, we remain the original character (same as original D²C strategy), else we update it with a new character that has the highest confidence if this confidence is bigger than a threshold ϵ .

Pseudo code The overall procedure of D²C is described in Algorithm 1.

5.3 Fine-tune with Confusion Set

Considering the low recall rate, we can also change some tokens with a confusion set instead of mask tokens during the fine-tuning process. The confusion set is built based on the pronunciations and fonts in Chinese. We can use a confusion set to produce a parallel dataset as a train set. For a monolingual sentence, we randomly choose one character and replace it with a character in the confusion set.

Concretely, our method to use the confusion set is as follows: we initially train our self-supervised

Algorithm 1: D²C

Input: Source sentence Y ; threshold ϵ , top- k .
Output: predict result Z

- 1 Sort the characters in Y on their confidences ascendingly and record the indices I ;
- 2 **for** $i \in I$ **do**
- 3 Mask y_i ;
- 4 Get top- k predictions $\{y_i^1, y_i^2, \dots, y_i^k\}$;
- 5 Get confidences $\{p_i^1, p_i^2, \dots, p_i^k\}$;
- 6 **if** $y_i \notin \{y_i^1, \dots, y_i^k\}$ **and** $p_i^1 > \epsilon$ **then**
- 7 Replace y_i with y_i^1 ;
- 8 Decode the new Y and update it;
- 9 **else**
- 10 Keep y_i unchanged;
- 11 **end**
- 12 **end**
- 13 $Z = Y$;

model using entirely monolingual data and then use a rate of the monolingual data to continually fine-tune the model with the confusion set.

6 Experiments

In this section, we report the empirical results of a series of spelling correction benchmarks.

We concentrate on two benchmarks:

- *ECSpell* (Lv et al., 2023): a small-scale multi-domain Chinese spelling correction dataset of Law (LAW), medical treatment (MED), and official document writing (ODW), which is particular in that there are a large number of errors in the test set that do not appear in the training set;

- *Syn-LEMON* (Lv et al., 2023): it is generated from LEMON (Wu et al., 2023b) which spans 7 different domains with a total of 19,130 synthetic train samples.

We consider the following methods:

- *BERT* (Devlin et al., 2019): the fine-tuned tagging model based on BERT;

- *MDCSpell* (Zhu et al., 2022b): the strongest tagging model with a multi-task network of error detection and correction;

- *Masked-FT (MFT)* (Wu et al., 2023a): a simple yet effective fine-tuning technique on tagging models to uniformly mask the non-error characters in the source sentence;

- *ReLM* (Liu et al., 2024): the newly released state-of-the-art models on spelling correc-

	Method	EC-LAW (%)				EC-MED (%)				EC-ODW (%)			
		F1	P	R	FPR	F1	P	R	FPR	F1	P	R	FPR
Supervised	BERT	38.6	42.1	35.7	12.2	24.2	27.1	21.9	10.5	24.9	29.9	21.3	13.9
	BERT-MFT	74.6	73.2	76.1	14.3	61.7	62.4	60.9	10.5	60.8	59.7	62.0	18.9
	MDCSpell-MFT	81.5	77.2	86.3	15.9	65.1	62.3	68.1	16.8	64.1	61.3	67.2	21.4
	Baichuan2	86.0	85.1	87.1	4.5	73.2	72.6	79.3	5.5	82.6	86.1	79.3	4.0
	ReLM	95.8	93.6	98.0	5.7	89.9	86.6	93.5	7.4	92.2	93.3	91.1	2.5
Self-supervised	BERT	0.5	0.7	0.4	9.0	0.6	0.9	0.4	8.0	0.5	0.8	0.4	12.4
	BERT-MFT	10.1	14.1	7.8	9.4	5.3	7.7	4.0	9.1	10.5	15.1	8.0	12.8
	MDCSpell-MFT	36.2	45.3	30.2	9.4	20.9	28.7	16.4	8.8	25.9	33.7	21.7	13.7
	Baichuan2	23.5	25.5	21.6	26.5	17.4	25.2	13.3	13.5	24.4	27.2	22.2	20.9
	Baichuan2-UD	26.9	30.8	23.9	20.4	18.3	27.4	13.7	11.7	28.0	32.7	24.4	14.5
	Baichuan2-D ² C	27.6	30.6	25.1	22.4	20.2	26.2	16.4	12.4	30.5	33.8	27.8	17.5
	ReLM	71.3	78.1	75.7	0.4	68.6	70.8	66.5	7.02	71.9	79.7	65.5	0.8
	ReLM-UD	89.5	89.2	89.9	4.7	79.3	74.1	85.4	18.5	84.6	88.5	81.0	2.3
	ReLM-Conf.(10%)	83.8	79.1	89.0	15.6	70.8	67.5	74.4	14.7	75.5	71.5	79.8	18.5
	ReLM-Conf.(100%)	84.1	77.7	91.8	19.7	69.7	57.6	88.4	41.1	73.4	68.5	79.1	19.3
	ReLM-D ² C	90.2	87.7	92.9	8.6	75.7	66.8	87.4	25.5	85.9	85.7	86.1	7.3

Table 4: Results on ECSpell, where F1, P, R, FPR refers to the F1 score, precision, recall, and false positive rate. Conf. (10%) means continually fine-tuning the self-supervised model with 10% confusion data. Conf. (100%) means continually fine-tuning the self-supervised model with 100% confusion data.

tion, which rephrases the sentence in a non-auto-regressive manner;

- *Baichuan2-7b* (Yang et al., 2023a): one of the strongest Chinese LLMs following the auto-regressive architecture;

- *User Dictionary (UD)* (Lv et al., 2023): an enhanced decoding method that leverages an expertise dictionary (law, medical treatment, and official document writing) to bias the beam search.

6.1 Training Settings

For BERT-based models, we set the batch size to 128 and the learning rate to $5e-5$, swept from grid search. For Baichuan2, we set the batch size to 32 and the learn rate to $3e-4$, and use LoRA (Hu et al., 2022) to reduce the training budget. For supervised spelling correction, the masking ratio is chosen from {0.2, 0.3}, while for self-supervised spelling correction, it is set to 0.5.

When fine-tuning with the confusion set, we set the batch size to 64 and the learning rate to $5e-5$.

6.2 Results on ECSpell

Table 4 summarizes the performances of different training methods on ECSpell and we also report the supervised performances for reference. For self-supervised spelling correction, we first find that ReLM outperforms MDCSpell-MFT by 35.1, 47.7, and 46.0 absolute points of F1 respectively on LAW, MED, and ODW, suggesting the great promise of rephrasing models.

When empowered with D²C, it further signifi-

cantly produces the increase of 18.9, 7.1, and 14.0 absolute points. The biggest increase is in the recall rate, which is consistent with the design of D²C. Furthermore, we find that D²C is competitive against using a user dictionary (UD), or even more powerful. It suggests that some of the domain knowledge in the user dictionary has already been stored in the pre-trained language models, and D²C plays a key role in unlocking the great power of pre-training.

When empowered with the confusion set, the increase is weaker. The confusion set method increases the FPR score, which even meets 41.1 on MED. A higher confusion rate is related to a higher FPR score.

6.3 Results on Syn-LEMON

Table 5 summarizes the results of self-supervised methods on Syn-LEMON. It indicates that except GAM(game) and NOV(novel), using the confusion set outperformed D²C’s F1 score. These variances reveal that different domains always have distinct data properties, and these properties play a role in varying performance outcomes between employing the confusion set and D²C.

7 Discussion

7.1 D²C v.s using Confusion Set

We compare D²C and the data augmentation method using the confusion set, a widely used technique in previous work. In Table 4 we find that

Method	GAM	ENC	COT	MEC	CAR	NOV	NEW
<i>Previous SotA (Wu et al., 2023a)</i>	33.8	48.6	67.2	54.3	53.1	38.6	58.7
ReLM	63.7	51.5	69.3	57.6	55.3	43.9	58.6
ReLM-D ² C	65.5	53.7	69.6	58.4	58.6	50.0	63
ReLM-Conf.(100%)	46.5	58.6	75.5	65.8	63.3	49.7	70.0
ReLM-Conf.(10%)	52.2	51.7	71.1	55.1	55.0	40.4	59.2

Table 5: Results on LEMON. Conf. (100%) means 100% data is trained as a confusion set and Conf. (10%) means we use 90% data as self-supervised training data and 10% data as continued fine-tuning confusion data.

D²C outperforms using the confusion set on two of the chosen datasets. Table 5 indicates that D²C surpasses using confusion set on GAM and NOV.

First, the results indicate that both D²C and using the confusion set can increase the recall rate. The common phenomenon is caused by different reasons. The confusion set introduces the character-to-character correction during the training process that is similar to test examples. While D²C introduces mask tokens to the test examples, which is inherited from the fine-tuning process. However, using the confusion set has a disadvantage compared with D²C. The non-matching segments in the confusion set can cause gaps in the real error patterns in the testing time. Therefore, using the confusion set always has lower P scores and higher FPR scores. D²C is a more suitable choice when it comes to domains that contain professional knowledge.

Second, compared to D²C, using the confusion set is relatively straightforward and efficient. Employing the confusion set presents an alternative approach in various application scenarios, offering efficiency but potentially posing a risk to performance.

	Models	F1(%)		
		LAW	MED	ODW
Supervised	MDCSpell (I)	71.8	51.3	54.9
	MDCSpell (E)	7.5	4.0	0.8
	MDCSpell-MFT (I)	94.3	78.4	81.7
	MDCSpell-MFT (E)	76.0	60.7	57.8
Self-supervised	MDCSpell-MFT (I)	52.6	32.9	32.1
	MDCSpell-MFT (E)	48.0	26.0	33.7
	ReLM (I)	93.2	73.5	82.2
	ReLM (E)	92.5	74.7	73.1
	ReLM-D ² C (I)	98.2	79.2	88.3
	ReLM-D ² C (E)	97.0	81.5	82.7

Table 6: Performances on seen (I) and unseen (E) errors, measured by F1 scores.

7.2 Seen and Unseen Errors

To take a closer look at the correction ability, we divide the test set into two subsets, exclusive (E) and inclusive (I) sets, which refer to the test errors that occur or do not occur in the training set.

From table 6, it is discernible that supervised models fit the internal error set well but the performances drop sharply on the external error set. While models trained with monolingual data have a high degree of similarity between the performance on the external error set and the internal error set. Besides, D²C boosts the performance on the external and internal sets simultaneously.

Surprisingly, MDCSpell-MFT performs even better on self-supervised learning than supervised on the exclusive set. It suggests that the tagging objective degenerates the learned representation in the pre-trained language model, incurring a drop in generalizability.

7.3 Effect of Mask Rate

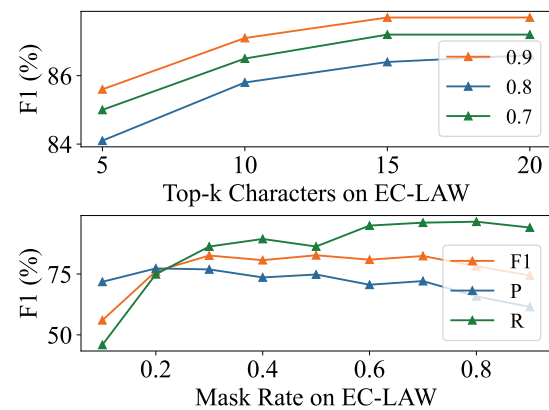


Figure 2: Self-supervised performances with different mask rates on Law of EcsPELL.

We also investigate the impact of mask rate. From Figure 2 it is apparent that the F1 scores on EcsPELL’s Law keep improving when the mask rate grows from 0% to about 30%, and then drop

slightly. To dig further, an increase in the mask rate uplifts recall (R) scores more apparently than precision (P) scores while P scores either lean to unchanged or even decline. Because the monolingual fine-tuning process introduces noise solely through mask tokens, the models are more inclined to preserve the source sentences without modification, which means lower R scores. During the evaluation stage, error characters serve as noise for the model, therefore a higher mask rate boosts models’ performances on R scores.

7.4 Effect of Hyperparameters

We access the effect of hyperparameters in D²C. As a representative, we depict the curves on ReLM in Figure 2.

Threshold Figure 2 shows that different datasets are suitable with different thresholds (ϵ). For example, D²C with higher ϵ (0.9) gains better performances on LAW, MED, and ODW domains. It reveals that ϵ should be set based on different datasets.

Top- k There is a common phenomenon in Figure 2 that a higher top- k character uplifts the F1 score under different thresholds ϵ .

7.5 Efficiency

We compare the decoding efficiency of D²C and normal decoding in Table 7. We can observe that compared with decoding each sentence directly, D²C requires about twice the time on ReLM and three times the time on Baichuan.

	Dataset	Normal (s)	D ² C (s)
ReLM	EC-MED	0.024	0.048
	EC-LAW	0.022	0.038
	EC-ODW	0.022	0.044
Baichuan	EC-MED	1.0	3.2
	EC-LAW	0.6	1.6
	EC-ODW	0.7	2.2

Table 7: Comparison between D²C and normal decoding on ReLM and Baichuan, by second per sample.

8 Case Study

We further showcase some cases to illustrate how D²C improves the decoding process.

Multi-typo In this case, (How does calcification (钙化) of the meniscus (半月板) occur), error characters are (钙→改) and (半→伴), which

SRC	伴月板改化的病因有哪些
TRG	半月板钙化的病因有哪些
ReLM	伴月板改化的病因有哪些
ReLM-D ² C	半月板钙化的病因有哪些

Table 8: Multi-typo case can be better corrected by D²C. Blue characters are right and red are wrong.

are very similar in pronunciation but meaningless as words in the sentence. We noticed in the experiment that ReLM without D²C failed to correct this sentence with two error characters while successful with a single error character if one of the two errors has been corrected before. Therefore, with D²C we introduce noise into the source sentence to correct “伴” and “改” step by step.

Not recall Considering sentences in spelling correction sometimes have short lengths, models receive limited semantics information and tend to under-correct error characters just like the case in Table 9. This case (How to calculate children’s weight (体重)) has the error pattern of (体→休), which are similar in terms of their visual appearance. In the presence of semantics limitations, D²C directs models to reword specified positions to incorporate more suitable characters and effectively mitigate the issue of under-correction.

SRC	小孩休重怎么计算
TRG	小孩体重怎么计算
ReLM	小孩休重怎么计算
ReLM-D ² C	小孩体重怎么计算

Table 9: D²C improves the recall rate.

9 Conclusion

This paper studies self-supervised spelling correction based on the rephrasing-based models. We demonstrate that machine spelling correction does not necessitate parallel data, and can be learned from monolingual data alone. We propose a novel decoding algorithm named D²C to effectively enhance the recall ability of the self-supervised model. We also compare the D²C method with using the confusion set method. Results on Chinese spelling correction showcase the significant improvement brought by our method. We hope this paper can bring new insight and vigor to future research on self-supervised spelling correction.

586
587
588
589
590

591

592
593
594
595
596
597

598
599
600
601
602
603
604
605
606
607
608
609
610
611
612

613
614
615
616
617
618
619
620

621
622
623
624
625
626
627
628
629

630
631
632
633
634
635

636
637
638
639
640

Limitations

Our work focuses on Chinese. Other languages such as Korean have not been studied in this work. D²C costs a decline in the speed of single sentence processing.

References

Haithem Afli, Zhengwei Qiu, Andy Way, and Pádraic Sheridan. 2016. Using smt for ocr error correction of historical texts. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC 2016)*, Paris, France. European Language Resources Association (ELRA).

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.

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*, pages 871–881. Association for Computational Linguistics.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Jianfeng Gao, Xiaolong Li, Daniel Micol, Chris Quirk, and Xu Sun. 2010. [A large scale ranker-based system for search query correction](#). In *Proceedings of the 23rd International Conference on Computational Linguistics (Coling 2010)*, pages 358–366, Beijing, China. Coling 2010 Organizing Committee.

Zhao Guo, Yuan Ni, Keqiang Wang, Wei Zhu, and Guotong Xie. 2021. [Global attention decoder for chinese spelling error correction](#). In *Findings of the Association for Computational Linguistics: ACL/IJCNLP 2021, Online Event, August 1-6, 2021*, volume

ACL/IJCNLP 2021 of *Findings of ACL*, pages 1419–1428. Association for Computational Linguistics. 641
642

Harsh Gupta, Luciano Del Corro, Samuel Broscheit, Johannes Hoffart, and Eliot Brenner. 2021. [Unsupervised multi-view post-OCR error correction with language models](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 8647–8652, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics. 643
644
645
646
647
648
649
650

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. [Lora: Low-rank adaptation of large language models](#). In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net. 651
652
653
654
655
656

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, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 5958–5967. Association for Computational Linguistics. 657
658
659
660
661
662
663
664
665
666

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*, pages 441–446. Association for Computational Linguistics. 667
668
669
670
671
672
673
674
675

Piji Li. 2022. [uChecker: Masked pretrained language models as unsupervised Chinese spelling checkers](#). In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 2812–2822, Gyeongju, Republic of Korea. International Committee on Computational Linguistics. 676
677
678
679
680
681

Linfeng Liu, Hongqiu Wu, and Hai Zhao. 2024. [Chinese spelling correction as rephrasing language model](#). In *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024*. AAAI Press. 682
683
684
685

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, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 2991–3000. Association for Computational Linguistics. 686
687
688
689
690
691
692
693
694
695

Qi Lv, Ziqiang Cao, Lei Geng, Chunhui Ai, Xu Yan, and Guohong Fu. 2023. [General and domain-adaptive](#) 696
697

698	chinese spelling check with error-consistent pretraining . <i>ACM Trans. Asian Low-Resour. Lang. Inf. Process.</i> , 22(5).	
699		
700		
701	Bruno Martins and Mário J. Silva. 2004. Spelling correction for search engine queries . In <i>Advances in Natural Language Processing, 4th International Conference, ESTAL 2004, Alicante, Spain, October 20-22, 2004, Proceedings</i> , volume 3230 of <i>Lecture Notes in Computer Science</i> , pages 372–383. Springer.	
702		
703		
704		
705		
706		
707	OpenAI. 2023. GPT-4 technical report . <i>CoRR</i> , abs/2303.08774.	
708		
709	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models . <i>CoRR</i> , abs/2307.09288.	
710		
711		
712		
713		
714		
715		
716		
717		
718		
719		
720		
721		
722		
723		
724		
725		
726		
727		
728		
729		
730		
731		
732	Dingmin Wang, Yan Song, Jing Li, Jialong Han, and Haisong Zhang. 2018. A hybrid approach to automatic corpus generation for chinese spelling check . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018</i> , pages 2517–2527. Association for Computational Linguistics.	
733		
734		
735		
736		
737		
738		
739		
740	Dingmin Wang, Yi Tay, and Li Zhong. 2019. Confusionset-guided pointer networks for chinese spelling check . In <i>Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers</i> , pages 5780–5785. Association for Computational Linguistics.	
741		
742		
743		
744		
745		
746		
747	Hongqiu Wu, Shaohua Zhang, Yuchen Zhang, and Hai Zhao. 2023a. Rethinking masked language modeling for chinese spelling correction . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 10743–10756. Association for Computational Linguistics.	
748		
749		
750		
751		
752		
753		
754		
	Hongqiu Wu, Shaohua Zhang, Yuchen Zhang, and Hai Zhao. 2023b. Rethinking masked language modeling for Chinese spelling correction . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 10743–10756, Toronto, Canada. Association for Computational Linguistics.	755 756 757 758 759 760 761
	Weijian Xie, Peijie Huang, Xinrui Zhang, Kaiduo Hong, Qiang Huang, Bingzhou Chen, and Lei Huang. 2015. Chinese spelling check system based on n-gram model . In <i>Proceedings of the Eighth SIGHAN Workshop on Chinese Language Processing, SIGHAN@IJCNLP 2015, Beijing, China, July 30-31, 2015</i> , pages 128–136. Association for Computational Linguistics.	762 763 764 765 766 767 768 769
	Aiyuan Yang, Bin Xiao, Bingning Wang, Borong Zhang, Ce Bian, Chao Yin, Chenxu Lv, Da Pan, Dian Wang, Dong Yan, Fan Yang, Fei Deng, Feng Wang, Feng Liu, Guangwei Ai, Guosheng Dong, Haizhou Zhao, Hang Xu, Haoze Sun, Hongda Zhang, Hui Liu, Jiaming Ji, Jian Xie, JunTao Dai, Kun Fang, Lei Su, Liang Song, Lifeng Liu, Liyun Ru, Luyao Ma, Mang Wang, Mickel Liu, MingAn Lin, Nuolan Nie, Peidong Guo, Ruiyang Sun, Tao Zhang, Tianpeng Li, Tianyu Li, Wei Cheng, Weipeng Chen, Xiangrong Zeng, Xiaochuan Wang, Xiaoxi Chen, Xin Men, Xin Yu, Xuehai Pan, Yanjun Shen, Yiding Wang, Yiyu Li, Youxin Jiang, Yuchen Gao, Yupeng Zhang, Zenan Zhou, and Zhiying Wu. 2023a. Baichuan 2: Open large-scale language models .	770 771 772 773 774 775 776 777 778 779 780 781 782 783 784
	Yifei Yang, Hongqiu Wu, and Hai Zhao. 2023b. Attack named entity recognition by entity boundary interference . <i>CoRR</i> , abs/2305.05253.	785 786 787
	Jui-Feng Yeh, Sheng-Feng Li, Mei-Rong Wu, Wen-Yi Chen, and Mao-Chuan Su. 2013. Chinese word spelling correction based on n-gram ranked inverted index list . In <i>Proceedings of the Seventh SIGHAN Workshop on Chinese Language Processing, SIGHAN@IJCNLP 2013, Nagoya, Japan, October 14-18, 2013</i> , pages 43–48. Asian Federation of Natural Language Processing.	788 789 790 791 792 793 794 795
	Junjie Yu and Zhenghua Li. 2014. Chinese spelling error detection and correction based on language model, pronunciation, and shape . In <i>Proceedings of The Third CIPS-SIGHAN Joint Conference on Chinese Language Processing, Wuhan, China, October 20-21, 2014</i> , pages 220–223. Association for Computational Linguistics.	796 797 798 799 800 801 802
	Ruiqing Zhang, Chao Pang, Chuanqiang Zhang, Shuo-huan Wang, Zhongjun He, Yu Sun, Hua Wu, and Haifeng Wang. 2021. Correcting chinese spelling errors with phonetic pre-training . In <i>Findings of the Association for Computational Linguistics: ACL/IJCNLP 2021, Online Event, August 1-6, 2021</i> , volume ACL/IJCNLP 2021 of <i>Findings of ACL</i> , pages 2250–2261. Association for Computational Linguistics.	803 804 805 806 807 808 809 810 811

812	Shaohua Zhang, Haoran Huang, Jicong Liu, and Hang	given word sets	864
813	Li. 2020. Spelling error correction with soft-masked	4. Your answer should be abundant and	865
814	BERT . In <i>Proceedings of the 58th Annual Meeting of</i>	include details, but not too long	866
815	<i>the Association for Computational Linguistics, ACL</i>	5. Try to generate realistic and fluent	867
816	<i>2020, Online, July 5-10, 2020</i> , pages 882–890. Asso-	sentences like a human writer	868
817	ciation for Computational Linguistics.	6. Your answer should be in Chinese in	869
818	Zewei Zhao and Houfeng Wang. 2020. Maskgec: Im-	JSON format	870
819	proving neural grammatical error correction via dy-	7. Your generated sentence should follow	871
820	namic masking . In <i>The Thirty-Fourth AAAI Con-</i>	the style of my given example sentences	872
821	<i>ference on Artificial Intelligence, AAAI 2020, The</i>	This is my given word sets:	873
822	<i>Thirty-Second Innovative Applications of Artificial</i>	"words":	874
823	<i>Intelligence Conference, IAAI 2020, The Tenth AAAI</i>	["word1", "word2",...],	875
824	<i>Symposium on Educational Advances in Artificial In-</i>	["word1", "word2",...]	876
825	<i>telligence, EAAI 2020, New York, NY, USA, February</i>	...	877
826	<i>7-12, 2020</i> , pages 1226–1233. AAAI Press.]}]	878
827	Chenxi Zhu, Ziqiang Ying, Boyu Zhang, and Feng Mao.	This is my given example sentences:	879
828	2022a. Mdcspell: A multi-task detector-corrector	{sentences}	880
829	framework for chinese spelling correction . In <i>Find-</i>	Your answer:	881
830	<i>ings of the Association for Computational Linguistics:</i>	[sentence1,sentence2,...]	882
831	<i>ACL 2022, Dublin, Ireland, May 22-27, 2022</i> , pages		883
832	1244–1253. Association for Computational Linguis-		
833	tics.		
834	Chenxi Zhu, Ziqiang Ying, Boyu Zhang, and Feng Mao.		
835	2022b. MDCSpell: A multi-task detector-corrector		
836	framework for Chinese spelling correction . In <i>Find-</i>		
837	<i>ings of the Association for Computational Linguis-</i>		
838	<i>tics: ACL 2022</i> , pages 1244–1253, Dublin, Ireland.		
839	Association for Computational Linguistics.		

840 10 Prompts

841 10.1 Extract Words

842 *Sentences* are extracted from the original LEMON
843 dataset.

844 1. Please extract the words in the given
845 sentences

846 2. Your answer should be in Chinese and
847 JSON format

848 {sentences}

849 Your answer format:

850 "words":

851 ["word1", "word2",...],

852 ["word1", "word2",...]

853 ...

854]}]

855 10.2 Generate Data

856 We propose the GAM domain’s prompt as an ex-
857 ample.

858 1. You are a professional game writer.
859 Try to use your professional knowledge
860 and think step by step.

861 2. Please make your answers diverse in
862 formats, words, and expressions.

863 3. Generate 5 smooth sentences Using the