# Neural Chinese Word Segmentation as Sequence to Sequence Translation

Xuewen Shi, Heyan Huang, Ping Jian, Yuhang Guo,
Xiaochi Wei, and Yikun Tang

Beijing Engineering Research Center of High Volume Language Information
Processing and Cloud Computing Applications,
School of Computer Science and Technology,
Beijing Institute of Technology, Beijing 100081, China
{xwshi,hhy63,pjian,guoyuhang,wxchi,tangyk}@bit.edu.cn

**Abstract.** Recently, Chinese word segmentation (CWS) methods using neural networks have made impressive progress. Most of them regard the CWS as a sequence labeling problem which construct models based on local features rather than considering global information of input sequence. In this paper, we cast the CWS as a sequence translation problem and propose a novel sequence-to-sequence CWS model with an attention-based encoder-decoder framework. The model captures the global information from the input and directly outputs the segmented sequence. It can also tackle other NLP tasks with CWS jointly in an end-to-end mode. Experiments on Weibo, PKU and MSRA benchmark datasets show that our approach has achieved competitive performances compared with state-of-the-art methods. Meanwhile, we successfully applied our proposed model to jointly learning CWS and Chinese spelling correction, which demonstrates its applicability of multi-task fusion.

**Keywords:** Chinese word segmentation, sequence-to-sequence, Chinese spelling correction, natural language processing

## 1 Introduction

Chinese word segmentation (CWS) is an important step for most Chinese natural language processing (NLP) tasks, since Chinese is usually written without explicit word delimiters. The most popular approaches treat CWS as a sequence labelling problem [21,14] which can be handled with supervised learning algorithms, e.g. Conditional Random Fields [11,14,24,18]. However the performance of these methods heavily depends on the design of handcrafted features.

Recently, neural networks for CWS have gained much attention as they are capable of learning features automatically. Zheng et al. [25] adapted word embedding and the neural sequence labelling architecture [7] for CWS. Chen et al. [4] proposed gated recursive neural networks to model the combinations of context characters. Chen et al. [5] introduced Long Short-Term Memory (LSTM) into neural CWS models to capture the potential long-distance dependencies. The

aforementioned methods predict labels of each character in the order of the sequence by considering context features within a fixed-sized window and limited tagging history [2]. In order to eliminate the restrictions of previous approaches, we cast the CWS as a sequence-to-sequence translation task.

The sequence-to-sequence framework has successful applications in machine translation [19,1,20,9], which mainly benefits from (i) distributed representations of global input context information, (ii) the memory of outputs dependencies among continuous timesteps and (iii) the flexibilities of model fusion and transfer.

In this paper, we conduct sequence-to-sequence CWS under an attention-based recurrent neural network (RNN) encoder-decoder framework. The encoder captures the whole bidirectional input information without context window limitations. The attention based decoder directly outputs the segmented sequence by simultaneously considering the global input context information and the dependencies of previous outputs. Formally, given an input characters sequence $\mathbf{x}$ with $T$ words i.e. $\mathbf{x} = (x_1, x_2, ..., x_{T_x})$, our model directly generates an output sequence $\mathbf{y} = (y_1, y_2, ..., y_T)$ with segmentation tags inside. For example, given a Chinese sentence " 我爱夏天" (I love summer) , the input $\mathbf{x} = ($ 我, 爱, 夏, 天 $)$ and the output $\mathbf{y} = ($ 我, </s>, 爱, </s>, 夏, 天 $)$ where the symbol '</s>' denotes the segmentation tag. In the post-processing step, we replace '</s>' with word delimiters and join the characters sequence into a sentence as $\mathbf{s} = $ "我爱夏天".

In addition, considering that the sequence-to-sequence CWS is an end-to-end process of natural language generation, it has the capacity of jointly learning with other NLP tasks. In this paper, we have successfully applied our proposed method to jointly learning CWS and Chinese spelling correction (CSC) in an end-to-end mode, which demonstrates the applicability of the sequence-to-sequence CWS framework.

We evaluate our model on three benchmark datasets, Weibo, PKU and MSRA. The experimental results show that the model achieves competitive performances compared with state-of-the-art methods.

The main contributions of this paper can be summarized as follows:

- We first treat CWS as a sequence-to-sequence translation task and introduce the attention-based encoder-decoder framework into CWS. The encoder-decoder captures the whole bidirectional input information without context window limitations and directly outputs the segmented sequence by simultaneously considering the dependencies of previous outputs and the input information.
- We let our sequence-to-sequence CWS model simultaneously tackle other NLP tasks, e.g., CSC, in an end-to-end mode, and we well validate its applicability in our experiments.
- We propose a post-editing method based on longest common subsequence (LCS) [12] to deal with the probable translation errors of our CWS system. This method solves the problem of missing information in the translation process and improves the experiment results.
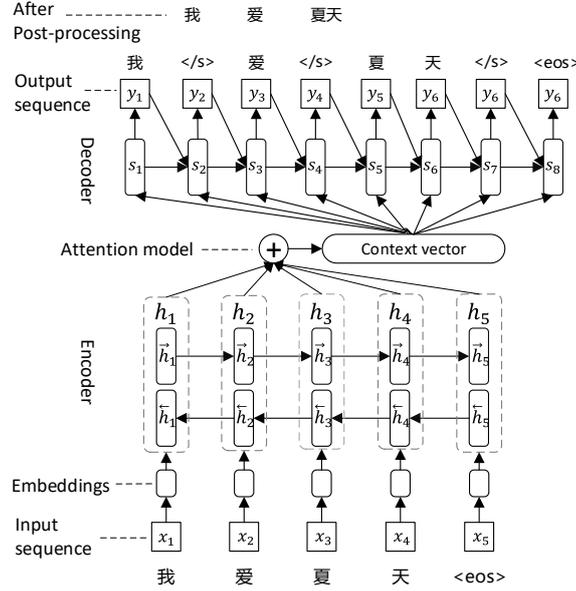
**Fig. 1.** Illustration of the presented model for CWS. The tag '<eos>' refers to the end of the sequence.

## 2 Method

### 2.1 Attention based RNN Encoder-decoder Framework for CWS

Our approach uses the attention based RNN encoder-decoder architecture called RNNsearch [1]. From a probabilistic perspective, our method is equivalent to finding a character sequence $\mathbf{y}$ with segmentation tags inside via maximizing the conditional probability of $\mathbf{y}$ given a input character sequence $\mathbf{x}$, i.e., $argmax_y p(\mathbf{y}|\mathbf{x})$.

The model contains (i) an bidirectional RNN encoder to maps the input $\mathbf{x} = (x_1, x_2, ..., x_{T_x})$ into a sequence of annotations $(h_1, h_2, ..., h_{T_x})$, and (ii) an attention based RNN decoder to generate the output sequence $\mathbf{y} = (y_1, y_2, ..., y_T)$. Fig. 1 gives an illustration of the model architecture.

### 2.2 bidirectional RNN Encoder

The bidirectional RNN encoder consists of forward and backward RNNs. The forward RNN $\overrightarrow{f}$ reads the input sequence in the order of (from $x_1$ to $x_{T_x}$) and calculates the sequence $(\overrightarrow{h}_1, \overrightarrow{h}_2, ..., \overrightarrow{h}_{T_x})$, while the backward RNN $\overleftarrow{f}$ reads the input sequence in the reverse order of (from $x_{T_x}$ to $x_1$) and calculates the

sequence ($\overleftarrow{h}_1, \overleftarrow{h}_2, ..., \overleftarrow{h}_{T_x}$). Finally, the annotation $h_j$ for each $x_j$ is obtained by $h_j = \left[ \overrightarrow{h}_j^T ; \overleftarrow{h}_j^T \right]^T$.

### 2.3 Attention-based RNN Decoder

The attention-based RNN decoder estimates the conditional probability $p(\mathbf{y}|\mathbf{x})$ as

$$p(\mathbf{y}|\mathbf{x}) = \prod_t^T p(y_t|y_1, ..., y_{t-1}, \mathbf{x}). \tag{1}$$

In Eq.(1), each conditional probability is defined as:

$$p(y_t|y_1, ..., y_{t-1}, \mathbf{x}) = g(y_{t-1}, s_t, c_t), \tag{2}$$

where $s_t$ is the RNN hidden state for time $t$ and computed by

$$s_t = f(s_{t-1}, y_{t-1}, c_t). \tag{3}$$

The $c_t$ in Eq.(2) and Eq.(3) is the context vector computed as a weighted sum of the annotations $(h_1, h_2, ..., h_{T_x})$:

$$c_t = \sum_{j=1}^{T_x} \alpha_{t,j} h_j.$$

The weight $\alpha_{t,j}$ is computed by:

$$\alpha_{t,j} = \frac{exp(e_{t,j})}{\sum_{k=1}^{T_x} exp(e_{t,k})},$$

where $e_{t,j} = a(s_{t-1}, h_j)$, therein, $a(\cdot)$ is an attention model constructed with a feedforward neural network.

### 2.4 Post-editing Method for Sequence-to-sequence CWS

We found some negative outputs of our model caused by translation errors such as missing words and extra words. The cause of the errors is mostly due to out-of-vocabulary or rare Chinese characters of input sequence.

Table 1 shows an example with translation errors of our sequence-to-sequence CWS system. The original input comes from the Weibo dataset (seen in Section 3.1). The output missed three Japanese characters ""(extreme), "の"(of) and " "(parent), and introduced three extra characters "UNK" instead which means 'unknown word' in the vocabulary.

Inspired by Lin and Och [12], we proposed an LCS based post-editing algorithm[1] (seen in Algorithm 1) to alleviate the negative impact to CWS. In the

---

[1] Executable source code is available at
https://github.com/SourcecodeSharing/CWSpostediting

**Table 1.** An example of translation errors in our CWS system and post-editing results.

| Original input | 岛国一超精分的小品《道の子》，看完之后我想说，为什么我没有这么 "通情达理" 的老爸呢? |
|---|---|
| System output | 岛国一超精分的小品《<u>UNK</u>道<u>UNK</u> <u>UNK</u>子》，看完之后我想说，为什么我没有这么 "通情达理" 的老爸呢? |
| After post-editing | 岛国一超精分的小品《道の 子》，看完之后我想说，为什么我没有这么 "通情达理" 的老爸呢? |
| Gold standard | 岛国一超精分的小品《道の 子》，看完之后我想说，为什么我没有这么 "通情达理" 的老爸呢? |

---

**Algorithm 1** Post-editing algorithm for our CWS system

---

**Input:**

  The original input character sequence: $s_{ori}$;

  The segmented word sequence with translation errors from sequence-to-sequence CWS system: $s_{seg}$;

**Output:**

  Segmentation labels set: $L \leftarrow \{B, M, E, S\}$;

  Labeling characters in $s_{seg}$ with labels in $L$ gets $lab_{seg}$;

  $length_{ori} \leftarrow getLength(s_{ori})$, $length_{seg} \leftarrow getLength(s_{seg})$

  **if** $length_{ori} \neq length_{seg}$ **then**

    Labeling characters in $s_{ori}$ with position labels;

    Extracting the longest common subsequences between $s_{ori}$ and $s_{seg}$ using longest common subsequence (LCS) algorithm: $s_{sub} = LCS(s_{seg}, s_{ori})$;

    Taking $s_{ori}$ as a reference, filling the missing characters in $s_{sub}$ and labeling them with label $X$;

    Replacing label $X$ with labels in $L$ according to manually prepared rules;

  **else**

    Labeling $s_{ori}$ according to $lab_{seg}$;

  **end if**

  Merging the characters in $s_{ori}$ into word sequence $s_{pe}$ according to their segmentation labels;

  **return**  $s_{pe}$;

---

algorithm, we define an extended word segmentation labels set $\{B, M, E, S, X\}$. $\{B, M, E\}$ represent begin, middle, end of a multi-character segmentation respectively, and $S$ represents a single character segmentation. The additional label $X$ in $L$ can be seen as any other labels according to its context. For example, given a CWS label sequence $(S, S, B, E, B, X, E)$, the $X$ should be transformed into label $M$ and in the other case of $(S, X, B, E, B, M, E)$, the $X$ should be treated as label $S$. The above transformation strategy can be based on handcraft rules or machine learning methods. In this paper, we use the transformation rules written manually. Table 1 also gives an example of post-editing results.

**Table 2.** Statistics of different datasets. The size of training/testing datasets are given in number of sentences (Sents), words (Words) and characters (Chars).

| Datasets | Training | | | Testing | | |
|---|---|---|---|---|---|---|
| | Sents | Words | Chars | Sents | Words | Chars |
| Weibo | 20,135 | 421,166 | 688,743 | 8,592 | 187,877 | 315,865 |
| PKU | 43,475 | 1,109,947 | 1,826,448 | 4,261 | 104,372 | 172,733 |
| MSRA | 86,924 | 2,368,391 | 4,050,469 | 3,985 | 106,873 | 184,355 |

## 3 Experiments

### 3.1 Datasets

We use three benchmark datasets, Weibo, PKU and MSRA, to evaluate our CWS model. Statistics of all datasets are shown in Table 2.

**Weibo**[2]: this dataset is provided by NLPCC-ICCPOL 2016 shared task of Chinese word segmentation for Micro-blog Texts [16]. The data are collected from Sina Weibo[3]. Different with the popular used newswire dataset, the texts of the dataset are relatively informal and consists various topics. Experimental results on this dataset are evaluated by eval.py scoring program[1], which calculates standard precision (P), recall (R) and F1-score (F) and weighted precision (P), recall (R) and F1-score (F) [15] [16] simultaneously.

**PKU and MSRA**[4] these two datasets are provided by the second International Chinese Word Segmentation Bakeoff [8]. We found that the PKU dataset contains many long paragraphs consisting of multiple sentences, which has negative impacts on the training of the sequence translation models. To solve this problem, we divide the long paragraphs in the PKU dataset into sentences. Experiment results on those two datasets are evaluated by the standard Bakeoff scoring program[3], which calculates P, R and F scores.

### 3.2 Model Setup and Pre-training

We use the RNNsearch[5] model [1] to achieve our sequence-to-sequence CWS system. The model is set with embedding size 620, 1000 hidden units and an alphabet with the size of 7190. We also apply the Moses' phrase-based (Moses PB) statistical machine translation system [10] with 3-gram or 5-gram language model as sequence-to-sequence translation baseline systems.

Since our sequence-to-sequence CWS model contains large amount numbers (up to ten million) of free parameters, it is much more likely to be overfitting when training on small datasets [17]. In fact, we make an attempt to train

---

[2] All data and the program are available at
  `https://github.com/FudanNLP/NLPCC-WordSeg-Weibo`
[3] `http://www.weibo.com`
[4] All data and the program are available at
  `http://sighan.cs.uchicago.edu/bakeoff2005/:`
[5] Implementations are available at `https://github.com/lisa-groundhog/GroundHog`

**Table 3.** Experimental results on benchmark datasets w/o pre-training.

| Datasets | P | R | F |
|---|---|---|---|
| Weibo | 89.8 | 89.5 | 89.6 |
| PKU | 87.0 | 88.6 | 87.8 |
| MSRA | 95.1 | 93.2 | 94.1 |

**Table 4.** Experimental results on the CWS dataset of Weibo. The contents in parentheses represent the results of comparison with other systems.

| Groups | Models | Standard Scores | | | Weighted Scores | | |
|---|---|---|---|---|---|---|---|
| | | P | R | F | P | R | F |
| A | LTP [3] | 83.98 | 90.46 | 87.09 | 69.69 | 80.43 | 74.68 |
| B [16] | S1 | 94.13 | 94.69 | 94.41 | 79.29 | 81.62 | 80.44 |
| | S2 | 94.21 | 95.31 | 94.76 | 78.18 | 81.81 | 79.96 |
| | S3 | 94.36 | 95.15 | 94.75 | 78.34 | 81.34 | 79.81 |
| | S4 | 93.98 | 94.78 | 94.38 | 78.43 | 81.20 | 79.79 |
| | S5 | 93.93 | 94.80 | 94.37 | 76.24 | 79.32 | 77.75 |
| | S6 | 93.90 | 94.42 | 94.16 | 75.95 | 78.20 | 77.06 |
| | S7 | 93.82 | 94.60 | 94.21 | 75.08 | 77.91 | 76.47 |
| | S8 | 93.74 | 94.31 | 94.03 | 74.90 | 77.14 | 76.00 |
| | S9 | 92.89 | 93.65 | 93.27 | 71.25 | 73.92 | 72.56 |
| M | Moses PB w/ 3-gram LM | 92.42 | 92.26 | 92.34 | 76.74 | 77.23 | 76.98 |
| | Moses PB w/ 5-gram LM | 92.37 | 92.26 | 92.31 | 76.58 | 77.25 | 76.91 |
| | RNNsearch w/o fine-tuning | 86.10 | 88.82 | 87.44 | 68.88 | 75.20 | 71.90 |
| | RNNsearch | 92.09 | 92.79 | 92.44 | 75.00 | 78.27 | 76.60 |
| | RNNsearch w/ post-editing | 93.48 (>S9) | 94.60 (>S6) | 94.04 (>S8) | 76.30 (>S5) | 79.99 (>S5) | 78.11 (>S5) |

the model on the benchmark datasets directly and get poor scores as shown in Table 3. To deal with this problem, a large scale pseudo data is utilized to pre-train our model. The Weibo, PKU and MSRA datasets are then used for fine-tuning. To construct the pseudo data, we label the UN1.0 [26] with LTP[6] [3] Chinese segmentor. The pseudo data contains 12,762,778 sentences in the training set and 4,000 sentences in the validation set and the testing set. The testing set of the pseudo data is used to evaluate the pre-training performance of the model, and the result P, R and F scores are **98.2**, **97.1** and **97.7** respectively w.r.t the LTP label as the ground truth.

### 3.3   CWS Experiment Results

**Weibo:** for Weibo dataset, we compare our models with two groups of previous works on CWS as shown in Table 4. The LTP [3] in group A is a general CWS tool

---

[6] Available online at `https://github.com/HIT-SCIR/ltp`

**Table 5.** Experimental results on the CWS benchmark datasets of PKU and MSRA.

| Groups | Models | PKU | | | MSRA | | |
|---|---|---|---|---|---|---|---|
| | | P | R | F | P | R | F |
| A | LTP [3] | 95.9 | 94.7 | 95.3 | 86.8 | 89.9 | 88.3 |
| B | Zheng et al., 2013 [25] | 93.5 | 92.2 | 92.8 | 94.2 | 93.7 | 93.9 |
| | Pei et al., 2014 [13] | 94.4 | 93.6 | 94.0 | 95.2 | 94.6 | 94.9 |
| | Chen et al., 2015 [4] | 96.3 | 95.9 | 96.1 | 96.2 | 96.3 | 96.2 |
| | Chen et al., 2015 [5] | 96.3 | 95.6 | 96.0 | 96.7 | 96.5 | 96.6 |
| | Cai and Zhao, 2016 [2] | 95.8 | 95.2 | 95.5 | 96.3 | 96.8 | 96.5 |
| C | Zhang et al., 2013 [23] | - | - | 96.1 | - | - | 97.4 |
| M | Moses PB w/ 3-gram LM | 92.9 | 93.0 | 93.0 | 96.0 | 96.2 | 96.1 |
| | Moses PB w/ 5-gram LM | 92.7 | 92.8 | 92.7 | 95.9 | 96.3 | 96.1 |
| | Moses PB w/ 3-gram LM w/ CSC | 92.9 | 93.0 | 92.9 | 95.3 | 96.5 | 95.9 |
| | Moses PB w/ 5-gram LM w/ CSC | 92.6 | 93.2 | 92.9 | 95.9 | 96.3 | 96.1 |
| | RNNsearch w/o fine-tuning | 93.1 | 92.7 | 92.9 | 84.1 | 87.9 | 86.0 |
| | RNNsearch | 94.7 | 95.3 | 95.0 | 96.2 | 96.0 | 96.1 |
| | RNNsearch w/ post-editing | 94.9 | 95.4 | 95.1 | 96.3 | 96.1 | 96.2 |
| | RNNsearch w/ CSC | 95.2 | 94.6 | 94.9 | 96.1 | 96.1 | 96.1 |
| | RNNsearch w/ CSC and post-editing | 95.3 | 94.7 | 95.0 | 96.2 | 96.1 | 96.2 |

which we use to label pseudo data. S1 to S8 in Group B are submitted systems results of NLPCC-ICCPOL 2016 shared task of Chinese word segmentation for Micro-blog Texts [16]. Our works are shown in Group M. Since the testing set of Weibo dataset has many out-of-vocabulary (OOV) words, our post-editing method shows its effective for enhancing our CWS results for its abilities to recall missing words and replace extra words.

**PKU and MSRA:** for the two popular benchmark datasets, PKU and MSRA, we compare our model with three groups of previous models on CWS task as shown in Table 5. The LTP [3] in group A is same as Table 4. Group B presents a series of published results of previous neural CWS models with pre-trained character embeddings. The work proposed by Zhang et al. [23] in group C is one of the state-of-the-art methods. Our post-editing method dose not significantly enhance the CWS results for PKU and MSRA datasets comparing with Weibo dataset. The reason is that the text style in the two datasets is formal and the OOV words are less common than Weibo dataset. In addition, the sequence translation baselines of Moses PB also gained decent results without pre-training or any external data.

According to all experimental results, our approaches still have gaps with the state-of-the-art methods. Considering the good performance (F1-score 97.7) on the pseudo testing data, the sequence-to-sequence CWS model has shown its capacity on this task and the data scale may be one of main limitations for enhancing our model.

**Table 6.** An example of modified data. The character with double underline is wrong and the characters with single underlines are correct.

| | |
|---|---|
| Original input | 在这个基础上，公安机关还从<u>原</u>料采购等方面加以严格控制，统一发放"准购证"。 |
| Modified input | 在这个基础上，公安机关还从<span style="color:red">源</span>料采购等方面加以严格控制，统一发放"准购证"。 |
| Gold standard | 在这个基础上，公安机关还从<u>原</u>料采购等方面加以严格控制，统一发放"准购证"。 |

**Table 7.** Experimental results on modified PKU data. The numbers in parentheses represent the changes compared with the normal CWS results shown in Table 5.

| Models | P | R | F |
|---|---|---|---|
| Modified testing data | 99.0 | 99.0 | 99.0 (-1.0) |
| LTP [3] | 94.0 | 93.2 | 93.6 (-1.7) |
| Moses PB w/ 3-gram LM | 90.8 | 91.5 | 91.2 (-1.8) |
| Moses PB w/ 3-gram LM w/ CSC | 92.7 | 92.9 | 92.8 (-0.1) |
| Moses PB w/ 5-gram LM | 90.6 | 91.3 | 91.0 (-1.7) |
| Moses PB w/ 5-gram LM w/ CSC | 92.3 | 93.0 | 92.6 (-0.3) |
| RNNsearch | 93.2 | 93.2 | 93.2 (-1.8) |
| RNNsearch w/ CSC | **95.0** | **94.5** | **94.8** (-0.1) |

### 3.4 Learning CWS and Chinese Spelling Correction Jointly

As a sequence translation framework, the model can achieve any expected kinds of sequence-to-sequence transformation with the reasonable training. It hence leaves a lot of space to tackle other NLP tasks jointly.

In this paper, we apply the model to jointly learning CWS and Chinese spelling correction (CSC). To evaluate the performance of spelling correction, we use automatic method to build two datasets, modified PKU and MSRA, based on assumptions that (i) most spelling errors are common with fixed pattern and (ii) the appearance of spelling errors are randomly. The details are as follows: firstly, we construct a correct-to-wrong word pair dictionary counting from the Chinese spelling check training dataset of SIGHAN 2014 [22] as a fixed pattern of spelling errors; secondly, we randomly select 50% sentences from PKU and MSRA training set respectively and replace one of the correct words with the wrong one according to the dictionary for each selected sentence. The testing set is generated in the same way.

We treat the modified sentences and the original segmented sentences as the input sequence and the golden standard respectively in the training procedure. Table 6 gives an example of the modified data. In the testing procedure, we send the sentence with wrong words into the model, and expect to get the segmented sentence with all correct words. The results are shown in Table 7 and Table 8. Since the general CWS tool LTP does not have the ability to correct spelling

**Table 8.** Experimental results on modified MSRA data. The numbers in parentheses represent the changes compared with the normal CWS results shown in Table 5.

| Models | P | R | F |
|---|---|---|---|
| Modified testing data | 98.5 | 98.5 | 98.5 (-1.5) |
| LTP [3] | 84.8 | 88.4 | 86.6 (-1.7) |
| Moses PB w/ 3-gram LM | 93.7 | 94.6 | 94.2 (-1.9) |
| Moses PB w/ 3-gram LM w/ CSC | 95.0 | 96.3 | 95.6 (-0.3) |
| Moses PB w/ 5-gram LM | 93.7 | 94.7 | 94.2 (-1.9) |
| Moses PB w/ 5-gram LM w/ CSC | 94.6 | 95.9 | 95.3 (-0.7) |
| RNNsearch | 93.8 | 94.7 | 94.2 (-1.9) |
| RNNsearch w/ CSC | **96.0** | **96.0** | **96.0** (-0.1) |

mistakes, the performance decreases. Whereas, the impact of the wrong words is limited in our models trained to do CWS and CSC jointly.

## 4   Related Work

CWS using neural networks have gained much attention in recent years as they are capable of learning features automatically. Collobert et al. [7] developed a general neural architecture for sequence labelling tasks. Zheng et al. [25] adapted word embedding and the neural sequence labelling architecture [7] for CWS. Pei et al. [13] improved upon Zheng et al. [25] by modeling complicated interactions between tags and context characters. Chen et al. [4] proposed gated recursive neural networks to model the combinations of context characters. Chen et al. [5] introduced LSTM into neural CWS models to capture the potential long-distance dependencies. However, the methods above all regard CWS as sequence labelling with local input features. Cai and Zhao [2] re-formalize CWS as a direct segmentation learning task without the above constrains, but the maximum length of words is limited.

**Sequence-to-sequence Machine Translation Models.** Neural sequence-to-sequence machine translation models have rapid developments since 2014. Cho et al. [6] proposed an RNN encoder-decoder framework with gated recurrent unit to learn phrase representations. Sutskever et al. [19] applied LSTM for RNN encoder-decoder framework to establish a sequence-to-sequence translation framework. Bahdanau et al. [1] improved upon Sutskever et al. [19] by introducing an attention mechanism. Wu et al [20] presented Google's Neural Machine Translation system which is serving as an online machine translation system. Gehring et al. [9] introduce an architecture based entirely on convolutional neural networks to sequence-to-sequence learning tasks which improved translation accuracy at an order of magnitude faster speed. Other efficient sequence-to-sequence models will be introduced into this task and compared with existing works in our future work.

## 5   Conclusion

In this paper, we re-formalize the CWS as a sequence-to-sequence translation problem and apply an attention based encoder-decoder model. We also make an attempt to let the model jointly learn CWS and CSC. Furthermore, we propose an LCS based post-editing algorithm to deal with potential translating errors. Experimental results show that our approach achieves competitive performances compared with state-of-the-art methods both on normal CWS and CWS with CSC.

In the future, we plan to apply other efficient sequence-to-sequence models for CWS and study an end-to-end framework for multiple natural language pre-processing tasks.

## References

1. Bahdanau, D., Cho, K., Bengio, Y.: Neural machine translation by jointly learning to align and translate. arXiv preprint arXiv:1409.0473 (2014)
2. Cai, D., Zhao, H.: Neural word segmentation learning for chinese. arXiv preprint arXiv:1606.04300 (2016)
3. Che, W., Li, Z., Liu, T.: Ltp: A chinese language technology platform. In: Proceedings of the 23rd International Conference on Computational Linguistics: Demonstrations. pp. 13–16. Association for Computational Linguistics (2010)
4. Chen, X., Qiu, X., Zhu, C., Huang, X.: Gated recursive neural network for chinese word segmentation. In: ACL (1). pp. 1744–1753 (2015)
5. Chen, X., Qiu, X., Zhu, C., Liu, P., Huang, X.: Long short-term memory neural networks for chinese word segmentation. In: EMNLP. pp. 1197–1206 (2015)
6. Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., Bengio, Y.: Learning phrase representations using rnn encoder-decoder for statistical machine translation. arXiv preprint arXiv:1406.1078 (2014)
7. Collobert, R., Weston, J., Bottou, L., Karlen, M., Kavukcuoglu, K., Kuksa, P.: Natural language processing (almost) from scratch. Journal of Machine Learning Research **12**(Aug), 2493–2537 (2011)
8. Emerson, T.: The second international chinese word segmentation bakeoff. In: Proceedings of the fourth SIGHAN workshop on Chinese language Processing. vol. 133 (2005)
9. Gehring, J., Auli, M., Grangier, D., Yarats, D., Dauphin, Y.N.: Convolutional Sequence to Sequence Learning. ArXiv e-prints (May 2017)
10. Koehn, P., Hoang, H., Birch, A., Callison-Burch, C., Federico, M., Bertoldi, N., Cowan, B., Shen, W., Moran, C., Zens, R., et al.: Moses: Open source toolkit for statistical machine translation. In: Proceedings of the 45th annual meeting of the ACL on interactive poster and demonstration sessions. pp. 177–180. Association for Computational Linguistics (2007)

11. Lafferty, J., McCallum, A., Pereira, F., et al.: Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In: Proceedings of the eighteenth international conference on machine learning, ICML. vol. 1, pp. 282–289 (2001)

12. Lin, C.Y., Och, F.J.: Automatic evaluation of machine translation quality using longest common subsequence and skip-bigram statistics. In: Proceedings of the 42nd Annual Meeting on Association for Computational Linguistics. p. 605. Association for Computational Linguistics (2004)

13. Pei, W., Ge, T., Chang, B.: Max-margin tensor neural network for chinese word segmentation. In: ACL (1). pp. 293–303 (2014)

14. Peng, F., Feng, F., McCallum, A.: Chinese segmentation and new word detection using conditional random fields. In: Proceedings of the 20th international conference on Computational Linguistics. p. 562. Association for Computational Linguistics (2004)

15. Qiu, P.Q.X., Huang, X.: A new psychometric-inspired evaluation metric for chinese word segmentation (2016)

16. Qiu, X., Qian, P., Shi, Z.: Overview of the NLPCC-ICCPOL 2016 shared task: Chinese word segmentation for micro-blog texts. In: Proceedings of The Fifth Conference on Natural Language Processing and Chinese Computing & The Twenty Fourth International Conference on Computer Processing of Oriental Languages (2016)

17. Srivastava, N., Hinton, G.E., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.: Dropout: a simple way to prevent neural networks from overfitting. Journal of Machine Learning Research **15**(1), 1929–1958 (2014)

18. Sun, X., Li, W., Wang, H., Lu, Q.: Feature-frequency–adaptive on-line training for fast and accurate natural language processing. Computational Linguistics **40**(3), 563–586 (2014)

19. Sutskever, I., Vinyals, O., Le, Q.V.: Sequence to sequence learning with neural networks. In: Advances in neural information processing systems. pp. 3104–3112 (2014)

20. Wu, Y., Schuster, M., Chen, Z., Le, Q.V., Norouzi, M., Macherey, W., Krikun, M., Cao, Y., Gao, Q., Macherey, K., et al.: Google's neural machine translation system: Bridging the gap between human and machine translation. arXiv preprint arXiv:1609.08144 (2016)

21. Xue, N., Shen, L.: Chinese word segmentation as lmr tagging. In: Proceedings of the Second SIGHAN Workshop on Chinese Language Processing. pp. 11212–179 (2003)

22. Yu, L.C., Lee, L.H., Tseng, Y.H., Chen, H.H., et al.: Overview of sighan 2014 bake-off for chinese spelling check. In: Proceedings of the 3rd CIPSSIGHAN Joint Conference on Chinese Language Processing (CLP'14). pp. 126–132 (2014)

23. Zhang, L., Wang, H., Mansur, X.S.M.: Exploring representations from unlabeled data with co-training for chinese word segmentation (2013)

24. Zhao, H., Huang, C.N., Li, M., Lu, B.L.: A unified character-based tagging framework for chinese word segmentation. ACM Transactions on Asian Language Information Processing (TALIP) **9**(2), 5 (2010)

25. Zheng, X., Chen, H., Xu, T.: Deep learning for chinese word segmentation and pos tagging. In: EMNLP. pp. 647–657 (2013)

26. Ziemski, M., Junczys-Dowmunt, M., Pouliquen, B.: The united nations parallel corpus v1. 0. In: Proceedings of the Tenth International Conference on Language Resources and Evaluation LREC. pp. 23–28 (2016)