# AN END-TO-END CHINESE TEXT NORMALIZATION MODEL BASED ON RULE-GUIDED FLAT-LATTICE TRANSFORMER

*Wenlin Dai*[1,†], *Changhe Song*[1,†], *Xiang Li*[1], *Zhiyong Wu*[1,2,*], *Huashan Pan*[3], *Xiulin Li*[3], *Helen Meng*[1,2]

[1] Tsinghua-CUHK Joint Research Center for Media Sciences, Technologies and Systems,
Shenzhen International Graduate School, Tsinghua University, Shenzhen, China
[2] Department of Systems Engineering and Engineering Management,
The Chinese University of Hong Kong, Hong Kong SAR, China
[3] Databaker (Beijing) Technology Co., Ltd, Beijing, China

{dwl20, sch19, xiang-li20}@mails.tsinghua.edu.cn, {zywu, hmmeng}@se.cuhk.edu.hk, {panhuashan, lixiulin}@data-baker.com

## ABSTRACT

Text normalization, defined as a procedure transforming non-standard words to spoken-form words, is crucial to the intelligibility of synthesized speech in text-to-speech system. Rule-based methods without considering context can not eliminate ambiguation, whereas sequence-to-sequence neural network based methods suffer from the unexpected and uninterpretable errors problem. Recently proposed hybrid system treats rule-based model and neural model as two cascaded sub-modules, where limited interaction capability makes neural network model cannot fully utilize expert knowledge contained in the rules. Inspired by **F**lat-**LA**ttice **T**ransformer (FLAT), we propose an end-to-end Chinese text normalization model, which accepts Chinese characters as direct input and integrates expert knowledge contained in rules into the neural network, both contribute to the superior performance of proposed model for the text normalization task. We also release a first publicly accessible large-scale dataset for Chinese text normalization. Our proposed model has achieved excellent results on this dataset.

***Index Terms***— Chinese text normalization, rule-based, none-standard word, flat-lattice Transformer, relative position encoding

## 1. INTRODUCTION

Text normalization (TN) is crucial to the intelligibility of synthesized speech in text-to-speech (TTS) system. It is defined as a procedure that transforms non-standard words (NSWs), e.g. written-form numbers, symbols or characters, to spoken-form words (SFWs), such as transforming "3.4" to "three point four" and "2021/10" to "October Twenty Twenty-one". To deal with the ambiguity problem in transforming NSWs to SFWs, context information and NSW's inherent special construct should be considered. For example, context can decide whether to read "2021" as year or number, whereas special construct of "172.0.0.1" can be determined as IP address. Furthermore, to form a context, word information is crucial. Take "2021 光年" as example, the "2021" will be read as number only if the word "光年 (light-years)" is correctly identified; otherwise, "2021" might be read as year if the keyword "年 (year)" is inaccurately matched.

Based on the taxonomy approach for NSW [1], the TN tasks can be resolved by rule-based approaches which utilize handcrafted regular expressions and/or keywords [2, 3, 4] to determine the category of NSWs and then convert to corresponding SFWs with predefined conversion functions. However, the selection of keywords as well as the construction of regular expression rules are time-consuming and labor-intensive. Several machine learning methods have been proposed for the disambiguation task of NSWs, including finite state automata (FSA) [2], maximum entropy (ME) [5], conditional random fields (CRF) [6], etc.

With the development of deep learning technologies, using neural network to model contextual information has achieved impressive progress for TN task. Sequence-to-sequence (seq2seq) models typically encode the written-form text representation into a state vector, and decode it into a sequence of spoken-form text output directly [7]. Long short-term memory (LSTM) and attention-based recurrent neural network (RNN) sequence-to-sequence models are well applied in English and Russian text normalization [8, 9]. Bidirectional LSTM or gated recurrent unit (GRU) are further utilized in both encoder and decoder [10, 11, 9]. However, directly applying sequence-to-sequence models to TN task may cause unexpected and uninterpretable errors caused by the model or data bias.

Recently, a hybrid TN system for Mandarin has been proposed, which combines a rule-based model based on pattern match and a multi-head self-attention based non-seq2seq neural network model, to address the corresponding shortcomings mentioned above [12]. According to the priority, NSWs are sent to the rule-based and neural models respectively. If the result of the neural model is of mismatched format, the NSW will be processed by the rule-based model again. However, the hybrid system simply treats rule-based model and neural model as cascaded sub-modules serially, which may cause error accumulation. It is easy to tell that rule-based model and neural network model can supplement each other. But limited interaction capability of these two cascaded sub-modules makes neural network model cannot fully utilize the expert knowledge included in the rules.

Inspired by the superior performance and the flexibility of the latest **F**lat-**LA**ttice **T**ransformer (FLAT) [13], we propose a FLAT based end-to-end Chinese TN model, named ***FlatTN***, which can directly incorporate the expert knowledge in predefined rules into the network, providing a novel way of leveraging the complementary advantages of the two models.

The advantages of using ***FlatTN*** for the text normalization task falls into two aspects. First, there is no need of the prerequisite

---

† Equal contributions.
* Corresponding author.

word segmentation module. FLAT can obtain all potential words in the sentence that match the specific lexicon, organize all characters and matched words to a lattice structure and flatten the lattice structure into spans [13], then send them into Transformer encoder. The method of combining lexicon is fully independent of word segmentation, and more effective in using word information thanks to the freedom of choosing lexicon words in a context [14]. Second, the NSW matching rules can be easily incorporated into the model and the definition of rules is greatly simplified. In the proposed model, rules are only adopted for the purpose of pattern match to derive all possible candidate NSWs in the input sentence. There is no need for the rules to account for complex context matching for disambiguation task as in the conventional method. We also release a large-scale dataset for the Chinese text normalization task, which will be open-sourced for public access soon. Experimental results on this dataset demonstrate that our proposed model has achieved an excellent performance.

The contributions of our work are:

1). For the first time, we propose the use of flat-lattice transformer (FLAT) for the task of Chinese TN problem, enhancing the controllability and scalability of TN task.

2). We come up with a novel rule-guided FLAT model, that can directly incorporate the expert knowledge on the predefined rules into the neural network based model. The proposed model is also an end-to-end model which predicts the NSW categories directly from the raw Chinese characters with NSWs.

3). We release, as open-sourced resources, a Chinese text normalization dataset[1] with standard NSW taxonomies to eliminate the ambiguity in pronunciation of NSWs. It is a first publicly accessible dataset for the Chinese TN task.

## 2. METHODOLOGY

We propose a fully end-to-end Chinese text normalization model based on FLAT, which accepts characters as direct input and can conveniently incorporate the expert knowledge from NSW matching rules. As shown in Fig.1, the model is made up of 4 parts: (i) Lexicon and rules matching that processes the input text and outputs a flat-lattice. (ii) An embedding presentation layer that generates embeddings for each token in the lattice. (iii) A Transformer encoder that produces lattice representations based on the generated embeddings and relative positional encodings of all tokens. (iv) A linear and CRF layer that predicts NSW category labels, given the lattice representations.

For an input sentence, all the potential words that match the lexicon are obtained. Furthermore, all the possible candidate NSWs are derived by matching the input sentence against the NSW matching rules. The sequence of characters, the potential words and the potential NSWs in the sentence are then organized as a sequence, named flat-lattice, where each character, word or NSW is called a token. For each token, its head and tail positions in the original sentence are also recorded. The character embeddings are then derived from a pretrained BERT model [15], from which word embeddings are obtained using a pooling layer [16] for each potential word and NSW. The character and word embeddings, together with the relative positional encoding derived from the head and tail positions for each token, are then fed into the transformer-based neural network. The self-attention mechanism of Transformer enables embeddings
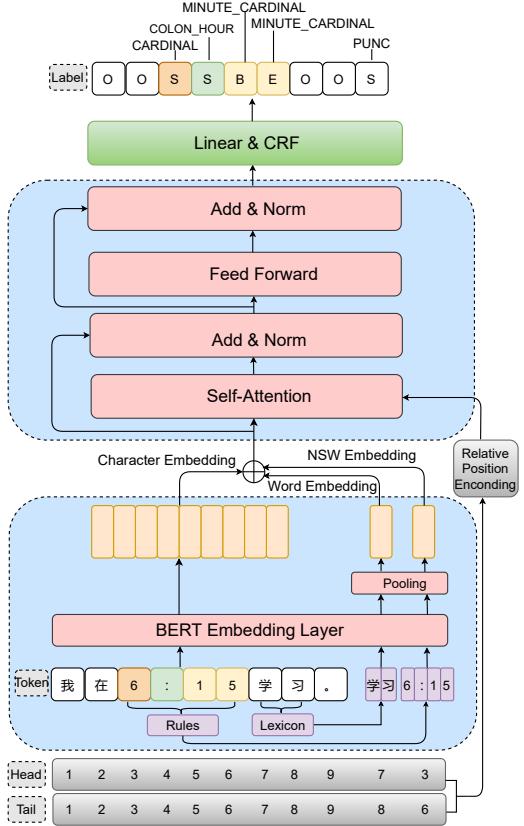
---
[1] https://github.com/thuhcsi/FlatTN



**Fig. 1**. Proposed end-to-end model structure.

to directly interact with each other [17], from which the contextual feature representation of the lattice is obtained. Finally, a linear layer and a conditional random field (CRF) layer are adopted to generate the categories of NSWs in the sentence.

### 2.1. Lexicon and rules matching

The model takes original characters in the sentence as its input. The characters are then handled by two pattern matching processes. On one hand, a lexicon is utilized for word matching. On the other hand, regular expressions, keywords, or any other form of rules could be incorporated to search through the text for potential NSWs. Each matched words, NSWs, as well as the original characters are taken as individual tokens, and then combined into a flat-lattice.

### 2.2. Embedding representation layer

Each token of the lattice is processed by a pre-trained BERT model [15] separately to obtain the character embeddings. Since BERT produces character level representations, for words and NSWs, a pooling layer is adopted to obtain the final word and NSW embeddings.

### 2.3. Transformer encoder with relative position encoding

We employ a Transformer model as our encoder, considering its strength in modeling the dependence between arbitrary nodes. This advantage is used to make bridges for information exchanging

among all tokens, which is expected to improve NSWs disambiguation by recognizing the context of individual tokens. Moreover, the calculation process of Transformer is static, which indicates it is agnostic to the structure of input lattice.

The Transformer encoder takes the token embeddings and relative position embeddings as its input. The token embeddings are offered by the embedding presentation layer. While the relative position embeddings are obtained through a sophisticatedly designed process proposed in FLAT [13].

As shown in Fig.1, the absolute position of each token in the lattice is represented by its start and end location (character index) in the input sentence, named as head and tail. The heads and tails of all tokens are then utilized to calculate four relative distances between every two nodes $x_i$ and $x_j$:

$$d_{ij}^{(hh)} = \text{head}[i] - \text{head}[j] \tag{1}$$

$$d_{ij}^{(ht)} = \text{head}[i] - \text{tail}[j] \tag{2}$$

$$d_{ij}^{(th)} = \text{tail}[i] - \text{head}[j] \tag{3}$$

$$d_{ij}^{(tt)} = \text{tail}[i] - \text{tail}[j] \tag{4}$$

where $d_{ij}^{(hh)}$ denotes the distance between head of $x_i$ and head of $x_j$, and other $d_{ij}^{(ht)}$, $d_{ij}^{(th)}$, $d_{ij}^{(tt)}$ have similar meanings. Fig.2 shows an example of four relative position matrixes calculated with head and tail information of tokens in lattice "学 习 学习".
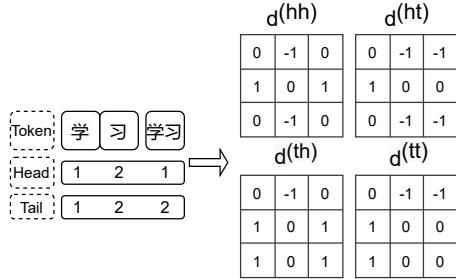


**Fig. 2**. Relative position encoding of "学 习 学习".

To get the final relative position encoding, a non-linear transformation is applied to the four relative distances:

$$\mathbf{R}_{ij} = \text{ReLU} \left| \left( W_r \left| (\mathbf{p}_{d_{ij}^{(hh)}} \oplus \mathbf{p}_{d_{ij}^{(th)}} \oplus \mathbf{p}_{d_{ij}^{(ht)}} \oplus \mathbf{p}_{d_{ij}^{(tt)}} | \right) \right| \right) \tag{5}$$

where $W_r$ is a learnable parameter, $\oplus$ denotes the concatenation operator, and $\mathbf{p}_d$ is calculated by following equations [17]:

$$\mathbf{p}_d^{(2k)} = \sin \left( d/10000^{2k/d_{\text{model}}} \right) \tag{6}$$

$$\mathbf{p}_d^{(2k+1)} = \cos \left( d/10000^{2k/d_{\text{model}}} \right) \tag{7}$$

where $d$ is one of $d_{ij}^{(hh)}$, $d_{ij}^{(ht)}$, $d_{ij}^{(th)}$, $d_{ij}^{(tt)}$ and $k$ denotes the index of dimension of position encoding.

Then the relative position encoding with fused information are sent to the Transformer encoder, together with token embeddings. Original Transformer calculates self-attention with absolute position coding [18]. In FLAT, Transformer encoder calculates self-attention with token embeddings and relative position encoding using following equation [19]:

$$\mathbf{A}_{i,j}^* = \mathbf{W}_q^\top \mathbf{E}_{x_i}^\top \mathbf{E}_{x_j} \mathbf{W}_{k,E} + \mathbf{W}_q^\top \mathbf{E}_{x_i}^\top \mathbf{R}_{ij} \mathbf{W}_{k,R} \\ + \mathbf{u}^\top \mathbf{E}_{x_j} \mathbf{W}_{k,E} + \mathbf{v}^\top \mathbf{R}_{ij} \mathbf{W}_{k,R} \tag{8}$$

where $\mathbf{W}_q, \mathbf{W}_{k,R}, \mathbf{W}_{k,E} \in \mathbb{R}^{d_{\text{model}} \times d_{\text{head}}}$ and $\mathbf{u}, \mathbf{v} \in \mathbb{R}^{d_{\text{head}}}$ are learnable parameters, and $\mathbf{E}_{x_i}, \mathbf{E}_{x_j} \in \mathbb{R}^{L \times d_{\text{model}}}$ are token embeddings ($d_{\text{model}} = H \times d_{\text{head}}$, $H$ is the number of attention heads, $d_{\text{head}}$ is the dimension of each head).

### 2.4. Linear and CRF layer

With the contextual feature representation from Transformer encoder, a linear and CRF layer are included as decoder to predict entity labels as the final output of our model. CRF can obtain an optimal prediction sequence through the relationship between adjacent labels and help to reduce error occurrence.

## 3. CHINESE TN DATASET

In this work, we release a large-scale Chinse TN dataset, which is the first large-scale open source Chinse text normalization dataset to the best of our knowledge. A well-designed Chinese NSW classification standard is proposed along with the data, which is made up by a total of 29 categories. As shown in Table.1, each category corresponds to a handcrafted conversion function for SFW generation. For example, a NSW likes "2021" belongs to the "DIGIT" category, and should be converted into SFW as "二零二一" ("two-zero-two-one") by the Read_DIGIT function. Specifically, ordinary characters that do not require transformation are labeled as "O" category, which are kept intact during conversion. And punctuation marks are labeled as "PUNC" category, which are simply dropped during conversion. The original text data in the dataset are extracted from Chinese Wikipedia, which come out as 30,000 sentences with average length of 50 characters. Detailed NSW category labels distribution of our dataset is depict as Fig.3.

The proposed Chinese TN dataset and NSW classification standard build up a benchmark for Chinese text normalization task. Future researches are relieved from the burden of NSW label and SFW converson function design, but focus on making steady improvements based on the same foundation.
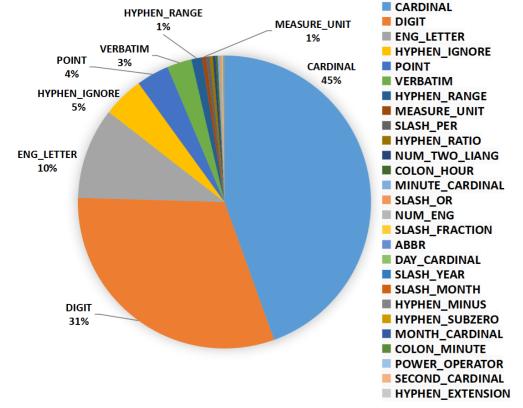


**Fig. 3**. Proportion of NSW categories in the proposed dataset. ("O" and "PUNC" excluded)

## 4. EXPERIMENT

### 4.1. Experiment setup

We split all sentences into characters and label them with "BMESO" format, and then randomly separate the dataset into training set, ver-

**Table 1**. Category set of the Chinese TN dataset.

| Category | How to read | Example |
|---|---|---|
| O | Self-reading | 你好. (你) |
| CARDINAL | Read cardinal | 11 pears. (11) |
| DIGIT | One by one | Call 911. (911) |
| PUNC | No reading | See you. (.) |
| ENG_LETTER | One by one | NBA. (NBA) |
| HYPHEN_IGNORE | No reading | See-you. (-) |
| POINT | Read as "diǎn" (点) | PI is 3.14. (.) |
| VERBATIM | One by one | I like C++. (++) |
| HYPHEN_RANGE | Read as "dào" (到) | July 12-20. (-) |
| MEASURE_UNIT | Read the unit | It is 24cm. (cm) |
| SLASH_PER | Read as "měi" (每) | 100$/Year. (/) |
| HYPHEN_RATIO | Read as "bǐ" (比) | Score is 3:2. (:) |
| NUM_TWO_LIANG | Read as "liǎng" (两) | 2个人. (2) |
| COLON_HOUR | Read as "diǎn" (点) | At 9:10am. (:) |
| MINUTE_CARDINAL | Add "fēn" (分) | At 9:10am. (10) |
| SLASH_OR | Read as "huò" (或) | Apple/pear. (/) |
| NUM_ENG | Read as English | Seq2seq. (2) |
| SLASH_FRACTION | Read as fraction | 3/4. (/) |
| ABBR | Abbreviation | Mr Smith. (Mr) |
| DAY_CARDINAL | Add "rì" (日) | 2021/09/06. (06) |
| SLASH_YEAR | Read as "nián" (年) | 2021/09. (/) |
| SLASH_MONTH | Add "yuè" (月) | 09/06. (/) |
| HYPHEN_MINUS | Read as "fù" (负) | -5. (-) |
| HYPHEN_SUBZERO | Read "líng xià" (零下) | -20°C. (-) |
| MONTH_CARDINAL | Add "yuè" (月) | 2021/09. (09) |
| COLON_MINUTE | Read as "fēn" (分) | 59:20. (:) |
| SECOND_CARDINAL | Add "miǎo" (秒) | 23:59:20. (20) |
| HYPHEN_EXTENSION | Read as "zhuǎn" (转) | 12345-678. (-) |
| POWER_OPERATOR | Read "cì fāng" (次方) | 2^3. (^) |

**Table 2**. Accuracy and F1 of different models.

| Method | Accuracy | F1 |
|---|---|---|
| Rule-based | 0.8775 | 0.8729 |
| BERT-MLP | 0.9869 | 0.9580 |
| BERT-LSTM | 0.9885 | 0.9638 |
| **FlatTN** | **0.9907** | **0.9708** |

text are more likely to obtain good performance, while the categories rarely appear in Chinese text may get worse performance than others.

**Table 3**. Model performance on the test set.

| Category | Precision | Recall | F1 |
|---|---|---|---|
| PUNC | 0.9952 | 0.9978 | 0.9965 |
| MINUTE_CARDINAL | 0.9706 | 1.0000 | 0.9851 |
| POINT | 0.9610 | 0.9769 | 0.9689 |
| CARDINAL | 0.9676 | 0.9607 | 0.9641 |
| DIGIT | 0.9487 | 0.9567 | 0.9527 |
| SLASH_PER | 0.8889 | 1.0000 | 0.9412 |
| HYPHEN_RATIO | 0.9677 | 0.9091 | 0.9375 |
| VERBATIM | 0.9385 | 0.8750 | 0.9057 |
| HYPHEN_RANGE | 0.7876 | 0.9468 | 0.8599 |
| HYPHEN_IGNORE | 0.8386 | 0.8469 | 0.8428 |

*4.2.3. Ablation study*

Ablation study of the proposed FlatTN model is established by removing lexicon, or rules, or both of them from the input process of the proposed model. As shown in Table.4, both lexicon and rules can help to improve results on accuracy and F1-score , and the performance promotion of the each part is complementary to the other. This is consistent with our expectation that rule-guided FlatTN should present better text normalization result, since it is embodied with both context information modeling and expert knowledge.

**Table 4**. Results of ablation experiment.

| Method | Accuracy | F1 |
|---|---|---|
| FlatTN | **0.9907** | **0.9708** |
| - Lexicon | 0.9886 | 0.9658 |
| - Rules | 0.9880 | 0.9632 |
| - Lexicon | 0.9879 | 0.9596 |

fication set, and test set at a ratio of 8:1:1. To compare our proposed model with rule based model and mainstream models based on neural network, three baseline models are constructed for comparison:

1). **Rule-based**: The adopted rule-based system is a commercial system by Databaker. Regular expressions are firstly used to match the NSW candidates, and then predefined rules are used to disambiguate the category for each NSW.

2). **BERT-MLP**: A neural network that obtains character embeddings using BERT, and predicts NSW labels using multi-layer perception [20] followed by CRF.

3). **BERT-LSTM**: A neural network similar to BERT-MLP, but with perception layer replaced by LTSM layer [20].

### 4.2. Model performance

*4.2.1. Overall performance*

Experiment results confirm that the proposed FlatTN model reaches the best results on both accuracy and F1-score. As revealed in Table.2, BERT-LSTM outperforms all other baseline models, due to its sequential structure. And there exists a noticeable gap between the Rule-based model and other models, which indicates the significance of pre-trained language models and neural modules in the promotion of text normalization performance.

*4.2.2. Model performance regarding different categories*

To explore the classification performance on different categories, our model is evaluated further in terms of precision, recall and F1-score on all categories. The results of 10 typical categories are shown in Table 3. It is clear that the categories frequently appear in Chinese

## 5. CONCLUSION

In this paper, we propose an end-to-end Chinese text normalization model based on rule-guided flat-lattice Transformer. Our model combines the scalability and flexibility of rules and the ability to model context and utilize data efficiently of Transformer encoder. We also release a first publicly accessible large-scale dataset for Chinese text normalization task. Our proposed model achieves the accuracy of 99.1% on NSW classification on the test set, achieving better performance compared with the methods using rules or neural networks alone, and the ablation experiment proves the importance of rules and lexicon in our model.

## 6. REFERENCES

[1] Richard Sproat, Alan W Black, Stanley Chen, Shankar Kumar, Mari Ostendorf, and Christoper Richards, "Normalization of non-standard words," *Computer speech & language*, vol. 15, no. 3, pp. 287–333, 2001.

[2] Z. Tao, D. Yuan, D. Huang, L. Wu, and H. Wang, "A three-stage text normalization strategy for mandarin text-to-speech systems," in *Chinese Spoken Language Processing, 2008. ISC-SLP '08. 6th International Symposium on*, 2009.

[3] X. Zhou, Z. Wu, C. Yuan, and Y. Zhong, "Document structure analysis and text normalization for chinese putonghua and cantonese text-to-speech synthesis," in *Intelligent Information Technology Application, 2008. IITA '08. Second International Symposium on*, 2009.

[4] G. T. Liou, Y. R. Wang, and C. Y. Chiang, "Text normalization for mandarin tts by using keyword information," in *2016 Conference of The Oriental Chapter of International Committee for Coordination and Standardization of Speech Databases and Assessment Techniques (O-COCOSDA)*, 2016.

[5] Yuxiang Jia, Dezhi Huang, Wu Liu, Shiwen Yu, and Haila Wang, "Text normalization in mandarin text-to-speech system," in *2008 IEEE International Conference on Acoustics, Speech and Signal Processing*, 2008, pp. 4693–4696.

[6] Guan-Ting Liou, Yih-Ru Wang, and Chen-Yu Chiang, "Text normalization for mandarin tts by using keyword information," in *The Oriental Chapter of International Committee for Coordination and Standardization of Speech Databases and Assessment Techniques (O-COCOSDA)*, 2016, pp. 73–78.

[7] Massimo Lusetti, Tatyana Ruzsics, Anne Göhring, Tanja Samardzic, and Elisabeth Stark, "Encoder-decoder methods for text normalization," in *Proceedings of the Fifth Workshop on NLP for Similar Languages, Varieties and Dialects (VarDial 2018)*, 2018, pp. 18–28.

[8] R. Sproat and N. Jaitly, "Rnn approaches to text normalization: A challenge," *arXiv preprint arXiv:1611.00068*, 2016.

[9] Courtney Mansfield, Ming Sun, Yuzong Liu, Ankur Gandhe, and Björn Hoffmeister, "Neural text normalization with subword units," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Industry Papers)*, 2019, pp. 190–196.

[10] Hao Zhang, Richard Sproat, Axel H Ng, Felix Stahlberg, Xiaochang Peng, Kyle Gorman, and Brian Roark, "Neural models of text normalization for speech applications," *Computational Linguistics*, vol. 45, no. 2, pp. 293–337, 2019.

[11] Richard Sproat and Navdeep Jaitly, "Rnn approaches to text normalization: A challenge," *arXiv preprint arXiv:1611.00068*, 2016.

[12] Junhui Zhang, Junjie Pan, Xiang Yin, Chen Li, Shichao Liu, Yang Zhang, Yuxuan Wang, and Zejun Ma, "A hybrid text normalization system using multi-head self-attention for mandarin," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 6694–6698.

[13] Xiaonan Li, Hang Yan, Xipeng Qiu, and Xuanjing Huang, "Flat: Chinese ner using flat-lattice transformer," *arXiv preprint arXiv:2004.11795*, 2020.

[14] Yue Zhang and Jie Yang, "Chinese ner using lattice lstm," *arXiv preprint arXiv:1805.02023*, 2018.

[15] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.

[16] Nikita Kitaev and Dan Klein, "Constituency parsing with a self-attentive encoder," *arXiv preprint arXiv:1805.01052*, 2018.

[17] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin, "Attention is all you need," in *Advances in neural information processing systems*, 2017, pp. 5998–6008.

[18] R. Al-Rfou, D. Choe, N. Constant, M. Guo, and L. Jones, "Character-level language modeling with deeper self-attention," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, pp. 3159–3166, 2019.

[19] Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V Le, and Ruslan Salakhutdinov, "Transformer-xl: Attentive language models beyond a fixed-length context," *arXiv preprint arXiv:1901.02860*, 2019.

[20] H. Y. Zhao, "Research and implementation of named entity recognition of electronic medical records based on deep learning," *Computer Engineering & Software*, 2019.