

# Korean Neural Machine Translation Using Hierarchical Word Structure

Jeonghyeok Park and Hai Zhao\*

Department of Computer Science and Engineering, Shanghai Jiao Tong University

Key Laboratory of Shanghai Education Commission for Intelligent Interaction

and Cognitive Engineering, Shanghai Jiao Tong University, Shanghai, China

MoE Key Lab of Artificial Intelligence, AI Institute, Shanghai Jiao Tong University, Shanghai, China

117033990011@sjtu.edu.cn, zhaohai@cs.sjtu.edu.cn

**Abstract**—Korean neural machine translation may significantly suffer from low-resource issues. We thus propose an enhancement method that fully exploits the hierarchical Korean word embedding structure from source representation. To our best knowledge, this is the first attempt for such Korean NMT tasks. Every Korean word can be decomposed into character- and jamo-level (sub-character) units. Therefore, We merge the character- and jamo-level representations with word embeddings to capture important combining word meaning. And then the merged representations are fed into NMT model. Our simple and novel method achieves BLEU improvements (up to 0.6) compared to word-based NMT baselines on Korean-to-Chinese and Korean-to-English translation tasks.

**Index Terms**—Machine Translation, Hierarchical Word Structure, Korean Language.

## I. INTRODUCTION

For different language pairs, traditional word-based neural machine translation (NMT) models suffer from the out-of-vocabulary (OOV) issue as they can only model a limited number of words. Character-based representations were proposed as a solution to overcome OOV problems, but it may be too fine-grained to miss some important information. And Byte Pair Encoding (BPE) [16] demonstrated extremely competitive performance by providing effective subword segmentation for NMT systems. Though the technique has solved the OOV problem efficiently, it still misses the semantic and syntactic information of the word itself. In this paper, we introduce a simple and novel method of supplementing additional information to the encoder of the NMT model by utilizing a unique compositional structure of the Korean language.

In this work, we focus on the Korean language. Unlike other languages, the Korean language has a unique compositional structure because it has both the features of the alphabetic and syllabic writing systems. Moreover, the decomposition of Korean syllables is deterministic [19]. Korean word is constructed by a regular hierarchy, so no special markers or measures are needed for the decomposition. *Jamos*, the Korean alphabet, are a set of the smallest unit that forms the language.

\*Corresponding author. This paper was partially supported by National Key Research and Development Program of China (No. 2017YFB0304100), Key Projects of National Natural Science Foundation of China (U1836222 and 61733011), Huawei-SJTU long term AI project, Cutting-edge Machine reading comprehension and language model.

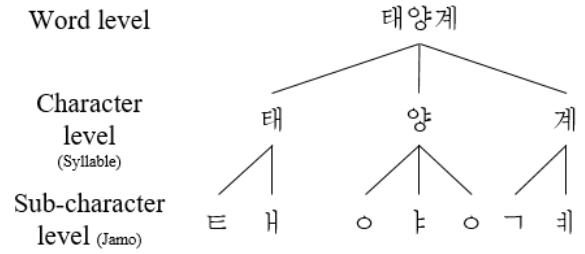


Fig. 1. The hierarchical structure of Korean language.

Every Korean word can be formed in the following hierarchy: A Korean word consists of a sequence of syllables, and a syllable (character) is arranged in a square, two-dimensional space with 2 or 3 jamo letters (sub-character). In other words, Korean words can be decomposed into two stages sequentially. As shown in Figure 1, the word 태양계 (solar system, *taeyanggye*) is a composition of three characters {태 (*tae*), 양 (*yang*) and 계 (*gye*)}, and the character 태 is a composition of {ㄸ (*t*), 애 (*ae*)}. The unique compositional structure of the Korean language are often ignored in most Korean NLP. Recently there have been studies that have proved that both characters and jamo letters contain important semantic and syntactic information in Korean language [19], [20]. Motivated by these works, we propose a hierarchical representation that merges word-/character-/jamo-level representation for Korean NMT. In our experiment, we demonstrate that the proposed method improves about 0.5 BLEU score of Korean NMT on average.

## II. RELATED WORKS

Since the dawn of NLP, there have been a variety of studies that encode more information that extracts from subword-level units into the representation of neural models. Many morpheme-based and character-based models have been proposed [5]–[10]. Among them, [9] proposed a neural language model that uses a convolutional neural network (CNN) to produce high-quality character-level representation. Partially motivated by this work, we further extend the range from character level to the sub-character level. For Korean, some previous studies exploit information that extracts from jamo





TABLE II  
EXPERIMENTAL RESULTS ON KOREAN-TO-CHINESE TRANSLATION AND  
KOREAN-TO-ENGLISH TRANSLATION.

Features	SacreBLEU Score	
	Dong-A (KO→ZH)	AIHub (KO→EN)
jamo	40.2	28.6
char	39.7	29.3
word	41.8	30.6
Merge methods	Concat/Aver/Gate	Concat/Aver/Gate
jamo, char	40.9/41.8/41.6	29.6/30.0/29.8
jamo, word	<b>42.3</b> /42.0/41.9	<b>31.2</b> /30.7/30.6
char, word	<b>42.3</b> /42.2/ <b>42.3</b>	31.0/31.0/30.9
jamo, char, word	42.2/ <b>42.3</b> /41.8	<b>31.2</b> /31.0/30.9

We use the 6-layer base Transformer architecture as a baseline model. For all experiments, the size of the embedding vector of source-/target-side is set to 512, and the initial learning rate is set to 0.2. We optimize the model parameters using Adam optimizer [4] with  $\beta_1 = 0.9$  and  $\beta_2 = 0.998$ , and Noam learning rate decay [21] with 8000 warm-up steps. All our experiments were performed using the open source Torch-based toolkit OpenNMT [18].

### B. Experimental Results

Table II shows the sacreBLEU evaluation of our systems. For both KO→ZH and KO→EN, the results show similar patterns: (1) The performance of transformer models trained with only the jamo- or character-level representation is lower than with only the word-level representation. (2) Merging word-level representation and other representations led to the transformer model’s performance improvement. (3) Among the merging methods, the method of concatenating the representations generally shows the highest improvement in the result. The results show that the merged representation for Korean achieves small but consistent sacreBLEU improvements over the baseline (only use word-level representation) on all parallel corpora.

The proposed method requires up to 10% more parameters and longer training time than the base Transformer model to generate low-level representations, but it can encode richer information than the word-based NMT model.

## VI. ANALYSIS AND DISCUSSIONS

### A. Ablation Study

Our composition model consists of three different networks: CNN, Highway Network, and LSTM. We perform an ablation study on Korean-to-Chinese translation task (Dong-A dataset) to understand the importance of these networks in producing low-level representations. The baseline model uses word-/character-/jamo-level representations and merge them through the Averaging method. As Table III shows, if we remove Highway Network, the translation performance decreases by 0.4 BLEU point. When LSTM is excluded from the composition model, it results in a significant drop of 0.7 BLEU point. When using only CNN models, translation performance dropped by

TABLE III  
ABLATION STUDY OF COMPOSITION MODEL ON KOREAN-TO-CHINESE  
TRANSLATION TASK (DONG-A DATASET).

Model	SacreBLEU Score
Baseline	42.3
-LSTM	41.6
-Highway Network	41.9
-LSTM&Highway Network	41.5

0.8 BLEU point. Similar to our model, [17] proposed syllable-based word embedding for Korean using CharCNN, and demonstrated good performance on the word similarity and relatedness task. However, in our preliminary experiment, we found that low-level representations produced using CharCNN only dropped performance on translation tasks. Therefore, we mitigated this issue by adding LSTM to the model.

### B. Why use the Composition Model?

We use the composition model to integrate low-level unit representation into the NMT model. It requires the character-/jamo-level representation for each word and each word consists of units of different lengths. Hence, we have adopted a model that can extract information effectively without being constrained by the length of the word. Similar to our model, [28] (in Chinese) proposed a new approach that integrates both word embeddings for Chinese words and Chinese character stroke sequence information into NMT system. However, They average a sequence of Chinese character stroke vectors which are induced from Chinese word without using any specific model and merge with the word embedding. In the preliminary experiment, we found that the method did not work well for the Korean translation task.

### C. Is jamo presentation useful in translation task?

Another observation is that jamo-level representation has a relatively lower contribution to performance either individually or in a merged manner than character-level representation. There may be several causes: (1) Jamo information is useful in speech recognition/synthesis task but may not in translation task. (2) All the datasets used in the experiment consist of formal sentences. Jamo representation may show better performance in the Ill-formed sentences commonly seen on the internet (i.e.,  $\text{ㅎㅇ}$  is shorthand for  $\text{하ㅇ(hi)}$ ) in formal sentences. As a future work, we have a plan to expand our model so that it can cover real-word data that are actually used on the Internet.

## VII. CONCLUSION

In this paper, we have presented the merged representations that exploit the hierarchical structure of the Korean language for Korean NMT model. In our experiment, we have demonstrated that the proposed model requires additional parameters, but is competitive and efficient in capturing additional semantic and syntactic information. Furthermore, the merged representations can identify unseen verbs more effectively than

simple word-level representations. Although this proposed method seems only suitable for the Korean language, it could be easily applied to similar scenarios.

## REFERENCES

- [1] Sepp Hochreiter, Jürgen Schmidhuber.: Long short-term memory. In: *Neural Computation*, pp. 1735–1780.
- [2] Kishore Papineni, Salim Roukos, Todd Ward, Wei-Jing Zhu.: BLEU: a method for automatic evaluation of machine translation. In: *Proceedings of the 40th annual meeting on association for computational linguistics. Association for Computational Linguistics (2002)*, pp. 311–318.
- [3] Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, et al.: Moses: Open source toolkit for statistical machine translation. In: *Proceedings of the 45th annual meeting of the association for computational linguistics companion volume proceedings of the demo and poster sessions*, pp. 177–180
- [4] Diederik P. Kingma, Jimmy Ba.: Adam: A method for stochastic optimization. In: *CoRR*, abs/1412.6980.
- [5] Thang Luong, Richard Socher, Christopher D Manning.: Better word representations with recursive neural networks for morphology. In: *Proceedings of CoNLL*, pp. 104–113.
- [6] Siyu Qiu, Qing Cui, Jiang Bian, Bin Gao, Tie-Yan Liu.: Co-learning of word representations and morpheme representations. In: *Proceedings of CoNLL*.
- [7] Xinxiong Chen, Lei Xu, Zhiyuan Liu, Maosong Sun, Huanbo Luan.: Joint learning of character and word embeddings. In: *Proceedings of IJCAI*. (2015)
- [8] Rupesh Kumar Srivastava, Klaus Greff, and Jurgen Schmidhuber.: Highway networks. *arXiv preprint arXiv:1505.00387* (2015).
- [9] Yoon Kim, Yacine Jernite, David Sontag, Alexander M. Rush. Character-aware neural language models. In: *Proceedings of AAAI*, pp. 2741–2749.
- [10] Zhen Yang, Wei Chen, Feng Wang, Bo Xu.: A character-aware encoder for neural machine translation. In: *Proceedings of the 26th International Conference on Computational Linguistics (COLING 2016)*, pp. 3063–3070.
- [11] Rupesh Kumar Srivastava, Klaus Greff, and Jurgen Schmidhuber.: Highway networks. *arXiv preprint arXiv:1505.00387* (2015).
- [12] Junyoung Chung, Kyunghyun Cho, Yoshua Bengio.: A character-level decoder without explicit segmentation for neural machine translation. In: *Proceedings of Association for Computational Linguistics (2016)*
- [13] Jason Lee, Kyunghyun Cho, Thomas Hofmann. :Fully character-level neural machine translation without explicit segmentation. *arXiv preprint arXiv:1610.03017* (2016)
- [14] Austin Matthews, Eva Schlinger, Alon Lavie, Chris Dyer.: Synthesizing compound words for machine translation. In: *Proceedings of Association for Computational Linguistics (2016)*
- [15] Rico Sennrich, Barry Haddow.: Linguistic input features improve neural machine translation. In: *Proceedings of WMT (2016)*
- [16] Rico Sennrich, Barry Haddow, Alexandra Birch.: Neural machine translation of rare words with subword units. In: *Proceedings of ACL*.
- [17] Sanghyuk Choi, Taeuk Kim, Jinseok Seol, Sangwoo Lee. A syllable-based technique for word embeddings of Korean words. In: *Proceedings of the First Workshop on Subword and Character Level Models in NLP*, pp. 36–40.
- [18] Guillaume Klein, Yoon Kim, Yuntian Deng, Jean Senellart, Alexander M. Rush.: OpenNMT: Open-Source Toolkit for Neural Machine Translation. In: *Proceedings of Association for Computational Linguistics (2017)*
- [19] Karl Stratos.: A sub-character architecture for Korean language processing In: *the 2017 Conference on Empirical Methods in Natural Language Processing*, pp. 721–726. Copenhagen, Denmark.(2017)
- [20] Sungjoon Park, Jeongmin Byun, Sion Baek, Yongseok Cho, Alice Oh.: Subword-level Word Vector Representations for Korean. In: *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2429–2438
- [21] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, Illia Polosukhin.: Attention is all you need. In: *Proceedings of the 31st Conference on Neural Information Processing Systems (NIPS 2017)*, pp. 5998–6008.
- [22] Joshua Coates, Danushka Bollegala.: Frustratingly easy meta-embedding—computing meta embeddings by averaging source word embeddings. In: *Proceedings of NAACL-HLT*. (2018)
- [23] Longtu Zhang, Mamoru Komachi.: Neural machine translation of logographic language using subcharacter level information. In: *Proceedings of the 3th Conference on Machine Translation: Research Papers, WMT 2018.*, pp. 17–25.
- [24] Longtu Zhang, Mamoru Komachi.: Chinese–Japanese Unsupervised Neural Machine Translation Using Sub-character Level Information. In: *Proceedings of the 33rd Pacific Asia Conference on Language, Information and Computation*, pp. 309–315.
- [25] Shaohui Kuang, Lifeng Han.: Apply Chinese radicals into neural machine translation: Deeper than character level. In: *30Th European Summer School In Logic, Language And Information*. (2018)
- [26] Matt Post.: A call for clarity in reporting BLEU scores. In: *Proceedings of the Third Conference on Machine Translation: Research Papers*, pp. 186–191.
- [27] Jeonghyeok Park and Hai Zhao.: Korean-to-Chinese Machine Translation using Chinese Character as Pivot Clue. In: *Proceedings of 33rd Pacific Asia Conference on Language, Information and Computation (PACLIC 33)*, pp. 558–566. Hakodate, Japan.
- [28] TAN Xin and KUANG Shaohui and ZHANG Longyin and XIONG Deyi.: Integration of Chinese character stroke sequence into neural machine translation. In: *Journal of Xiamen University(Natural Science)* (2019), pp. 164–169.