

# Improving Sequence to Sequence Neural Machine Translation by Utilizing Syntactic Dependency Information

An Nguyen Le, Ander Martinez, Akifumi Yoshimoto\* and Yuji Matsumoto

Graduate School of Information Science, Nara Institute of Science and Technology,

8916-5, Takayama, Ikoma, Nara 630-0192, Japan

{nguyen.an.mr9, ander.martinez.zy4, akifumi-y, matsu}@is.naist.jp

## Abstract

Sequence to Sequence Neural Machine Translation has achieved significant performance in recent years. Yet, there are some existing issues that Neural Machine Translation still does not solve completely. Two of them are translation of long sentences and “over-translation”. To address these two problems, we propose an approach that utilize more grammatical information such as syntactic dependencies, so that the output can be generated based on more abundant information. In addition, the output of the model is presented not as a simple sequence of tokens but as a linearized tree construction. Experiments on the Europarl-v7 dataset of French-to-English translation demonstrate that our proposed method improves BLEU scores by 1.57 and 2.40 on datasets consisting of sentences with up to 50 and 80 tokens, respectively. Furthermore, the proposed method also solved the two existing problems, ineffective translation of long sentences and over-translation in Neural Machine Translation.

## 1 Introduction

Our task is to construct a model which learns input in sequence form and decodes output as a linearized dependency tree. In this work, we propose an approach in which dependency labels are incorporated into the model to represent more grammatical information in the output sequence. As we know, the Sequence to Sequence (Seq2Seq) Learning model (Sutskever et al., 2014; Aharoni et al., 2016) is extremely effective on a va-

riety of tasks that require a mapping between a sequence to sequence. Therefore, it is used to solve many tasks in natural language processing. The Seq2Seq model consists of an encoder-decoder neural network which encodes a variable-length input sequence into a vector and decodes it into a variable-length output. Since the model uses the information of the source representation and the previously generated words to produce the next-word token, this distributed representation allows the Seq2Seq model to generate appropriate mapping between the input and the output (Li et al., 2016). For specific tasks, Neural Machine Translation (NMT) model, which is based on the Seq2Seq learning, has achieved excellent translation performance in recent years (Sutskever et al., 2014; Bahdanau et al., 2015; Luong et al., 2015; Firat et al., 2016). In particular, the NMT model which is built upon an encoder-decoder framework with attention mechanism (Bahdanau et al., 2015) can also pay attention and its decoder knows which part of the input is relevant for the word that is currently being translated. Therefore, it has shown competitive results and outperformed conventional statistical methods (Bentivogli et al., 2016). Despite of these advantages, NMT model still has a couple particular issues to be solved such as dealing with fixed vocabulary, not applicable to small-data, additional phrases, wrong lexical choice errors, long sentence translation, over and under translation, etc. In this paper, we touch upon the following two major problems:

- Translation of long sentences
- Over-translation

Since the decoder of the Seq2Seq model produces the target language word by word simply based on the previous target words and the source-side representation vector until it reaches the spe-

This author’s present affiliation is CyberAgent, Inc., Tokyo, Japan, yoshimoto.akifumi\_xa@cyberagent.co.jp

cial end token, it is incapable in capturing long-distance dependencies in history, so ineffective for long sentences translation (Zhang et al., 2016; Toral and Sánchez-Cartagena, 2017). Even with an attention mechanism, the Seq2Seq model just pays attention to the current alignment information between the inputs and the output at the current position but ignores past alignments information. Therefore, it cannot keep track of the attention history when it updates information at each current time step, leading to the over-production (Tu et al., 2016a,c; Mi et al., 2016; Tu et al., 2016b).

In order to address the above two issues, it is worth considering that using syntactic dependency information and representing the output as a tree structure would be effective. This approach allows the next tokens to be output based on not only the previous tokens but also the syntactic dependencies so far, thereby conditioning them on more abundant information so it has the ability to make smarter predictions. Basically, in this paper, we train the model with an encoder-decoder neural network and using dependencies in which the input of the source language is in sequence form and the output of the target language will be generated in a linearized dependency-based tree structure. That is, instead of predicting only words at each time step, the model trains the network to predict both words and their grammatical dependencies as a dependency tree at each time step. Therefore, it is hoped that the accuracy of output will be improved.

The major contributions of this work are as follows:

1. To utilize the information of both “head” words and syntactic dependencies between them to produce better output.
2. To settle the problems in the NMT task. In this paper, we desire to solve two tasks. First is the ineffective translation for long sentences. Second is the over-translation in NMT task.

Empirically, to assess the performance of the proposed method, we used Conditional Gated Recurrent Unit with Attention mechanism model of Bahdanau (2015) on the French-English portions of the Europarl-v7 dataset. As a result, the BLEU score is improved by 1.57 and 2.40 points for sentences of length up to 50 and 80 tokens, respec-

tively. Also, we compare and analyze the results of attention-based Seq2Seq model and the proposed approach.

## 2 Related Work

In fact, the effectiveness of using dependency information of words has been reported in some previous NLP tasks, for example, in dependency-based word embeddings, relation classification and sentence classification tasks (Liu et al., 2015; Socher et al., 2014; Levy and Goldberg, 2014; Komnios, 2016; Ono and Hatano, 2014). It has been shown that the combination of words and their dependency information can boost performance. Besides, in the work of Vinyals et al. (Vinyals et al., 2014), they also represent output as a linearized tree structure, but their work showed that generic sequence-to-sequence approaches can achieve excellent results on syntactic constituency parsing. At a glance, our proposed method is a little similar to the works of Dyer et al., Aharoni et al., Eriguchi et al., Wu et al. (Dyer et al., 2016; Aharoni and Goldberg, 2017; Eriguchi et al., 2017; Wu et al., 2017) in use of parse tree and generation. However, Dyer et al. and Aharoni et al.’s works concern predicting constituent trees. Eriguchi et al.’s model employs syntactic dependency parsing but their model is hybridized the decoder of NMT and the Recurrent Neural Network Grammars, and the target sentences are parsed in transition-based parsing. Wu et al.’s model also employs dependency parsing but their model separately predicts the target translation sequence and parsing action sequence which maps to translation. On the other hand, our proposed model’s decoder directly predicts the linearized dependency tree itself in a single neural network in *Depth-first pre-order* order so that the next-word token is generated based on syntactic relations and tree construction itself. In other words, our model is able to learn and produce a tree of words and their dependency relations by itself.

## 3 Sequence-to-Dependency Model

In our proposed approach, the neural network model is trained to map the target-side output in a linearized dependency tree construction from the source-side input in a sequence. Thus, we call this model Sequence-to-Dependency (Seq2Dep) model. The problem is defined as follows: Given a source sequence  $\mathbf{X} = (x_1, x_2, \dots, x_N)$  of length

$N$ , we want the model to encode the input sequence  $X$  and decode it to a tree structure with both words and dependency information conditioned on the encoded vector. Therefore, the output will be represented in the form  $(LY) = (\mathbf{ly}_1, \mathbf{ly}_2, \dots, \mathbf{ly}_M)$ . The conditional probability  $p(\mathbf{ly}|x)$  is decomposed as:

$$p(\mathbf{ly}|x) = \prod_{i=1}^{\infty} p(\mathbf{ly}_i | \mathbf{ly}_{<i}, x), \quad (1)$$

in which  $(\mathbf{ly}_1, \mathbf{ly}_2, \dots, \mathbf{ly}_M)$  are words or dependency labels.

Therefore, the hidden state  $\mathbf{s}_j$  at time step  $j$  is computed as follows:

$$\mathbf{s}_j = \text{cGRU}_{\text{att}}(\mathbf{s}_{j-1}, \mathbf{ly}_{j-1}, \mathbf{C}_j), \quad (2)$$

and the next token  $\mathbf{ly}_j$ , which may be a word or dependency label, will be generated as follows:

$$\mathbf{ly}_j = f(\mathbf{s}_j, \mathbf{ly}_{j-1}, \mathbf{C}_j), \quad (3)$$

In this paper, dependencies are defined as the dependency labels which are achieved from the Stanford Dependency Parser (Chen and Manning, 2014). The decoder will decode the next output based on relations between governors and dependents in a linearized tree structure. In regards to the order of generating the dependency labels and the words, the decoder will produce these symbols in a manner called *Depth-first pre-order* traversal. In this section, we will describe the model step-by-step as follows:

### 3.1 Processing Data

Since there is no parallel corpus in which the source-side is represented in sequence and target-side is represented in linearized dependency tree, we have to prepare data for training by doing dependency parsing for the target-side language.

#### 3.1.1 Dependency Parsing

In this paper, we do experiments on a French-English language pair so we use the Stanford Dependency Parser to obtain dependency parsing results for English. The Stanford Dependency Parser produces results in the form of a tree structure in which each word of the sentence is the dependent of exactly one token, either another word in the sentence or the distinguished “ROOT-0” token. The parsing result is represented in the format “*abbreviated relation name(governor, depen-*

*dent)*” in which a governor is a head word and dependency is a syntactic relation between a governor and a dependent. The governor and the dependent are words in the sentence. This dependency parsing result will be transformed in another step for traversing the tree, which will be described in the next section to create a dependency tree. The dependency tree represents the target language as an ordered tree structure which is necessary for training. The reason we chose the Stanford Dependency Parser for the parsing portion of this method is because it can represent the order of words in sentence. This information of the order is useful to traverse tree in the following step.

#### 3.1.2 Transformation and Tree Traversal

In this section, we describe the Tree Transform and Tree Traversal process in which output in a linearized dependency tree form is created from the Stanford Dependency Parsing tree. For example, given a sentence “*She ate an apple today .*”, after obtaining dependency parsing tree from the above dependency parsing phase, we move the rooted “*ate*” and “*apple*” headwords to the same layers of their *dependents* which are directly connected to the headwords. We also concurrently make consideration to their positions in order while shifting headwords. The headwords are shifted in such a manner that the word order of sentence can be preserved, so we can evaluate the translated output afterwards. Next, the tree structure obtained in the first step will be transformed into another tree structure for the next tree traversal step. Then we traverse this tree in a *Depth-first pre-order* traversal, which is the search tree in which tree is traversed from its left subtree to right subtree recursively until current node is empty, to create output with a linearized tree structure to train the model. That is, for each rooted subtree, governors and dependency labels of the sentence are predicted first, and their information will be used to predict the next dependent words. In other words, the model can capture the dependency information between label-word and word-word pairs to predict the next tokens. This means that the model is capable of modeling grammatical dependencies in the output symbols. Also, in Seq2Dep model, we define the *Nonterminal* “{*DEPENDENCY LABEL*”, and *Node-closing* “}” tokens. *Nonterminal* indicates subtree (Dong and Lapata, 2016), which means open subtree to visit its children nodes. *Node-closing* indicates end-of-

---

**Algorithm 1** Tree Transform

---

```
1: procedure TRANSFORM TREE
2:   Transform(T,Labels):
3:     for label in Labels do
4:       if label.children.size! = 0 then
5:         Recur Transform(T,Labels)
6:       else
7:         Compare the order of current
8:           label's parent & children
9:         if (label's children order is larger
10:            than label's parent order) then
11:           INSERT label's parent first
12:         else
13:           INSERT label's children
```

---

subtree, that means finishing subtree traversal and returning to the upper layer to continue the next subtree traversal. And these defined tokens do not appear in original source and target datasets. Algorithms 1 and 2 show the definition of transformation and tree traversal in more detail respectively. The purpose of using *Depth-first pre-order* traversal is as follows:

1. To keep the words of the target language sequence in order when they are generated. With this generating order, the word order of the sentence is preserved, thus, we do not have to do any post-processing subsequently.
2. To utilize both information of the words and the dependency labels generated in the previous rooted subtree to predict the tokens of the next rooted subtree.

Figures 1, 2 and 3 show the Stanford dependency parsing tree, tree structure after the positions of “head” words are shifted and *Depth-first pre-order* Tree Traversal.

### 3.2 Sequence-to-Dependency Model

The proposed (Seq2Dep) model consists of an encoder which is a bidirectional GRU layer as in Bahdanau’s model (2015)<sup>1</sup>. The input embeddings of the source sentences are shared by the forward and backward GRU, and the hidden states of the corresponding forward and backward GRU are added to obtain the hidden representation for that time step. The decoder of the model will decode the output as words and dependency labels in a linearized dependency tree structure in

---

<sup>1</sup><https://github.com/nyu-dl/dl4mt-tutorial>

---

**Algorithm 2** Tree *Depth-first pre-order* traversal

---

```
Input: Sentence
Output: Linearized Dependency Tree
1: Stanford Dependency Parsing
2: Make Tree from Dependency Parsing Result
3: Tree transform
4: procedure TRAVERSE TREE
5: Traverse(T,N):
6:   N as discovered
7:   for all Node not in N do
8:     if Node.children.size! = 0 then
9:       Recursively call Traverse(T,N)
10:      in pre-order traverse
11:     else
12:       if Node is Nonterminal then
13:         OUTPUT Node-opening
14:         VISIT children
15:         OUTPUT Node-closing
16:       else
17:         OUTPUT Node
```

---

a *Depth-first pre-order* traversal. Figure 4 shows the decoder which generates both dependency labels and words in the Seq2Dep model. In Figure 4, the previous token and context vector feeding are omitted for simplicity.

## 4 Experiments

### 4.1 Dataset

In our experiment, the proposed model was trained on the French-English parallel corpus of the *Europarl-v7* dataset. We used *newstest2011* and *newstest2012* of WMT16 as development and test data respectively. To confirm translation for long sentences, the whole test set was used without removing any sentences with a maximum length of 50 or 80. We performed experiment on the following two datasets:

- *Europarl-v7* dataset consisting of sentences with a maximum length of 50.
- *Europarl-v7* dataset consisting of sentences with a maximum length of 80.

For preprocessing data, we filtered out sentences which were longer than the above maximum lengths and cleaned the special symbols or characters which were not strings. We also omitted sentences which had multiple sentences in one line. The reason is that the parsing results obtained from the Stanford Dependency Parser in parsing



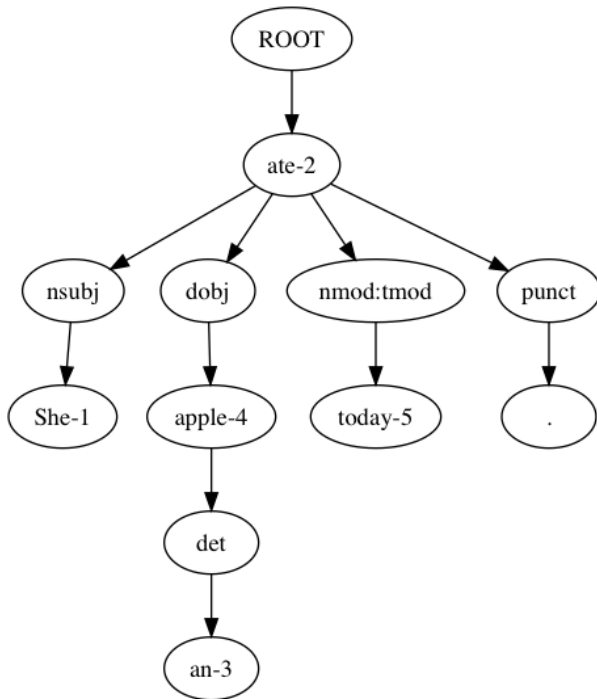


Figure 1: Stanford Dependency Parsing Tree

step would contain multi “{ROOT}” tokens for sentences which have multiple sentences in one line, while it is necessary to generate the next child nodes starting from just one top {ROOT} of a tree. Next, we tokenize and lowercase this dataset and perform dependency parsing. After that, we traverse the tree in a *Depth-first pre-order* to create the parallel corpus for the training model in which the source language, French is in sequence form, and the target language, English is in a linearized dependency tree structure form. The longer sentences are (particularly sentences with a maximum length of 80 tokens), the more CPU’s memory and time cost for this processing data step.

In addition, we built a dictionary of the target language (English) that consists of both words and dependency labels. In this dictionary, we define 74 dependency labels based on the current representation of grammatical relations of the Stanford Dependency Parser.

## 4.2 Settings

In order to evaluate the performance of the proposed method, we set the same hyperparameters as the attention-based cGRU model in DL4MT-Tutorial and compare the obtained results of both Seq2Seq and Seq2Dep models.

The recurrent transformation weights for gates

and hidden state proposal matrices were initialized as random orthogonal matrices. Weights were optimized using the Adadelta algorithm and were updated with a mini-batch size of 32 sentences. The vocabulary sizes of both source and target languages were set at 30k words, the beam size was set to 5, dropout was not applied and the gradients were clipped at 1.0. Moreover, because the generated tokens are not only words but also dependency labels in Seq2Dep model, the maxlen parameter was set up so that dependency labels are not counted, therefore long sentences will not be removed in training.

## 4.3 Model Training

In the experiments, we trained the following 2 models on 1.65M sentences with a maximum length of 50 and 1.89M sentences with a maximum length of 80 from the Europarl-v7 French-English bitext.

### Baseline Model

This model is a Seq2Seq model with attention mechanism as in Firat (2016) that consists of an encoder that encodes the source language input in sequence form and a decoder that decodes target language output in sequence form.

**Seq2Dep Model** The proposed method. In this model, the model architecture is the same as the attention-based Seq2Seq model but the input is in sequence form and the output is in linearized dependency tree structure.

## 5 Results

In the Seq2Dep model, because the output consists of both words and dependency labels, we evaluated the result with post-processing, which is the process that removes the dependency labels from the translated result. From this section onwards, we will refer to the Seq2Seq and Seq2Dep models with sentences of maximum length 50 and 80 tokens as Seq2Seq-50, Seq2Dep-50, Seq2Seq-80 and Seq2Dep-80. As a result, the BLEU score of Seq2Dep-50 with post-processing was 20.88, which is higher than the BLEU score of 19.31 obtained by the attention-based Seq2Seq-50 model with a gain of up to 1.57 points. Similarly, the BLEU score improved by 2.40 points for datasets with maximum sentence lengths of 80.

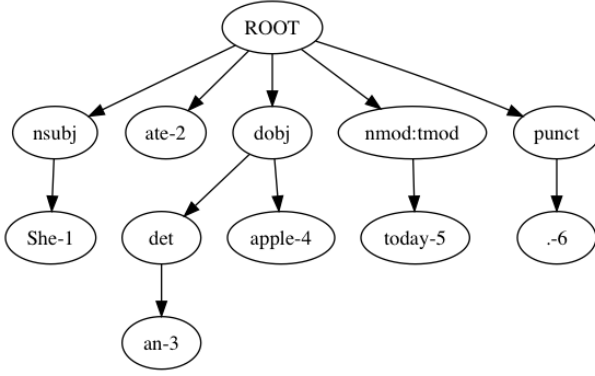


Figure 2: Dependency tree after shifting the positions of “head” words

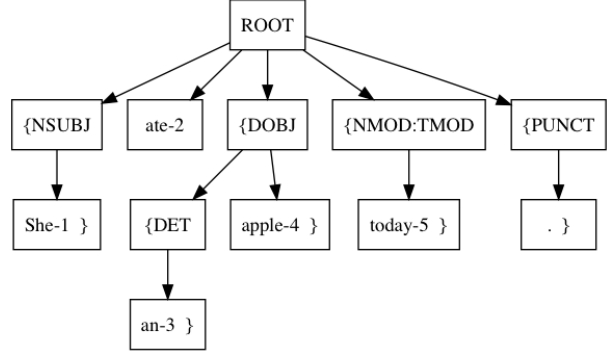


Figure 3: *Depth-first pre-order* Tree Traversal

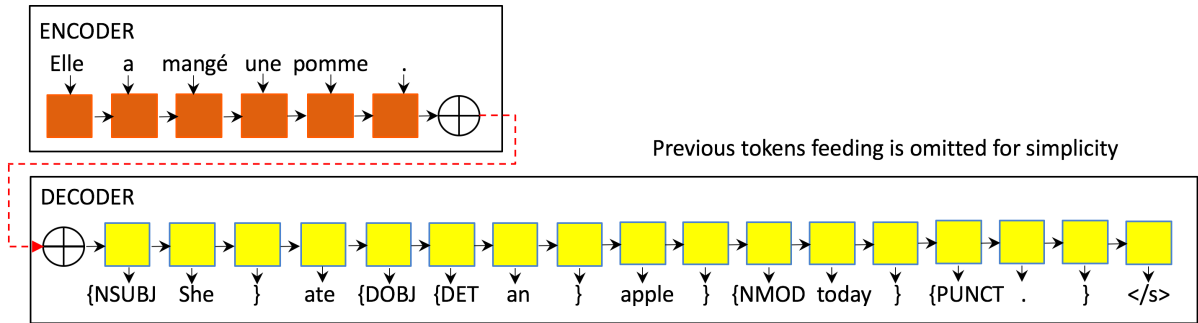


Figure 4: Encoder and decoder of *Seq2Dep* model

Table 1 shows BLEU and METEOR scores and TER error of the attention-based *Seq2Seq* and *Seq2Dep* models. Figure 5 shows the relation between BLEU score and the length of sentence.

Moreover, when we made a trial to evaluate the translation results without post-processing, the BLEU scores without post-processing were 42.76 and 43.41 for both datasets. From these scores, it is thought that the model can predict not only word-based tokens but also dependency labels well.

## 6 Additional Experiments

In order to verify the ability of the proposed approach to solve the repetition problem of NMT, over-translation, we measured the repetition of words in the translation results of attention-based *Seq2Seq* and *Seq2Dep* learnings in this section. The repetition rate is measured by the following formula:

$$rep\_rat = \sum_{i=1}^{T(y)} \frac{1 + r(\tilde{y}_i)}{1 + r(Y)}, \quad (4)$$

in which  $\tilde{y}_i$  and  $Y_i$  are the  $i^{th}$  hypothesis sentence and  $i^{th}$  reference sentence respectively, and  $r$  is the number of the repeated words and is computed by:

$$r(X) = len(X) - len(set(X)) \quad (5)$$

in which  $len(X)$  is the length of the sentence  $X$  and  $len(set(X))$  is the number of words that are not repeated in sentence  $X$ . For example, given the sentence  $X = \text{“The big fish ate the smaller fish”}$ , in this case,  $set(X) = \{\text{The, big, fish, ate, smaller}\}$ ,  $len(X) = 7$ ,  $len(set(X)) = 5$ . Figure 6 shows the comparison of repetition rate in both models in which the horizontal axis is the length of sentences, vertical axis is the repetition rate respectively. In Figure 6, the repetition rate in both *Seq2Seq* and *Seq2Dep* learnings decreases as the length of the sentences increases. From Figure 6, we can see that the more tokens the model learns, the more the repetition rate decreases. Also, the repetition rate is reduced in the *Seq2Dep* model compared to the attention-based *Seq2Seq* model.

Table 1: Translation quality as measured by different metrics.

Model	Post-processing		
	BLEU	METEOR	TER
Seq2Seq-50	19.31	26.3	66.1
Seq2Dep-50	<b>20.88</b>	<b>27.0</b>	<b>62.5</b>
Seq2Seq-80	16.97	25.5	78.5
Seq2Dep-80	<b>19.37</b>	<b>25.6</b>	<b>65.6</b>

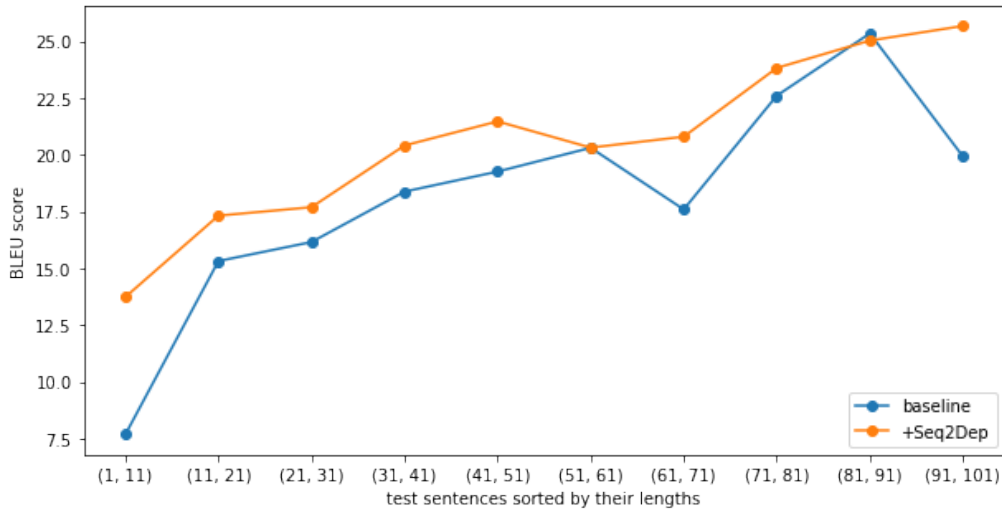


Figure 5: Comparison of BLEU score with respect to the length of sentences

## 7 Analysis and Discussion

In figure 5, except the span in which the sentence length is between 41 and 51 words, the BLEU score of the Seq2Dep model goes up gradually and almost overcomes that of the attention-based Seq2Seq model. The BLEU score falls from 19.31 to 16.97 with a 2.34 points difference for the attention-based Seq2Seq model while the point difference is 1.51 in the Seq2Dep model. From the experiments, we confirm that by using the syntactic dependency information, the Seq2Dep model can learn well and reduce the drop in BLEU score compared to the baseline model even if the sentence is very long. Besides, we can see the BLEU score is low for short sentences which have a length of 10 words or less. This is because of the brevity penalty on short sentences in BLEU (Papineni et al., 2002).

With regards to the BLEU score without post-processing, we see that the score of the Seq2Dep-80 model is higher than that of the Seq2Dep-50 model. The reason could be: The longer the sentences are, the more syntactic de-

pendencies the models require for generating better outputs.

Also, in terms of the over-translation problem, Figure 6 shows that the repetition rates of the two models decrease gradually with respect to the length of the sentences and the Seq2Dep model has a lower repetition rate. When we checked the translation results, we saw that *Node-closing* token “}” was almost generated after each subtree. Moreover, we saw that there were some very long sentences which the over-generation of “UNK”’s occurred in the translation result of Seq2Seq model while that did not occur in translation results of Seq2Dep model. Our assumption is that after generating subtree, the Seq2Dep model can learn that it should generate the *Node-closing* token “}” next, instead of a chain of words. In other words, as mentioned in Kuncoro et al.’s work (Kuncoro et al., 2016) in which modeling of composition can achieve better performance, the Seq2Dep model which learns about the syntactic dependencies and tree structure performance is probably able to learn the blocks of the form “*Non-terminal word }*” like a phrase-structure in sen-

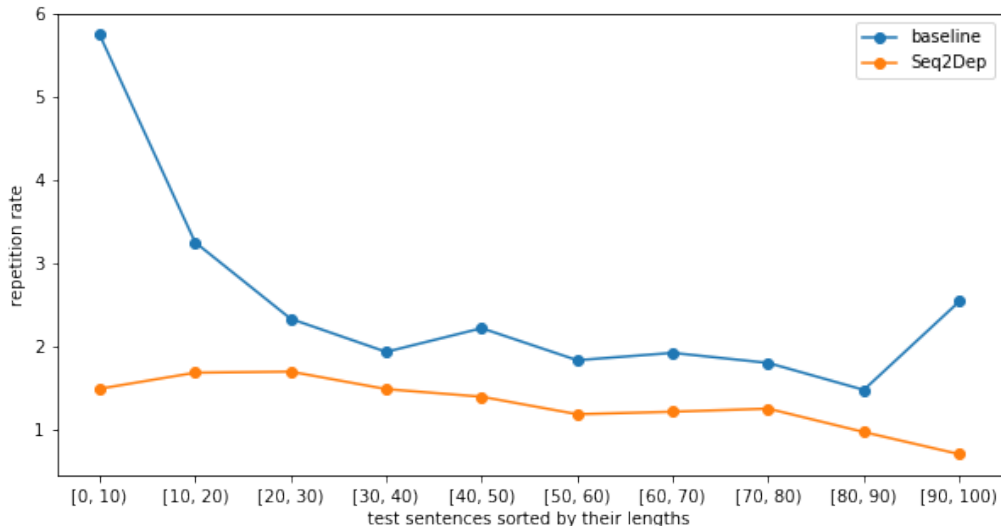


Figure 6: Comparison of the repetition rate of the baseline and Seq2Dep models

tences, so it is unlikely to generate the same word repeatedly. Therefore, it is possible to prevent the long repeated words in long sentences. Usually, because the block of the form “*Nonterminal word }*” is seen as a phrase in sentence or a subtree in tree structure, and it is rare for a phrase to occur repeatedly in sentence or for a subtree to repeat in a tree structure, so it is assumed that repetition of the blocks of form “*Nonterminal word }*” are also rare.

## 8 Conclusion

In this work, we proposed a method in which the Seq2Dep NMT model is trained by utilizing syntactic dependencies to provide the model more abundant information. In other words, Seq2Dep model learns the potential internal relative connections among tokens and their long term syntactic dependencies to predict the next-word tokens. Furthermore, the Seq2Dep model can also generate output as a linearized dependency tree structure in a *Depth-first pre-order* tree traversal over words and dependencies. The purpose of this work is to alleviate issues of translating long sentences and repetitive translation. We conduct experiments on the French-English parallel corpus of the Europarl-v7 dataset to compare the performance of the proposed method with the attention-based Seq2Seq model. The results demonstrated that the proposed model achieved a 1.57 and 2.40 points BLEU score improvement for sentences of length at most 50 and 80 tokens re-

spectively. Moreover, experiments verify that the proposed model also reduces the over-translation, particularly long sentences with over-generation of “*UNK*”s.

## 9 Future work

- Confirm how accurate the Seq2Dep model generates the dependency labels and the whole tree structure as well.
- In this paper, to compare performance of the proposed method with the baseline model, we set the same hyperparameters as the attention-based cGRU model in dl4mt-tutorial and trained the Seq2Dep model on only Europarl-v7 dataset. Since experiments were done on small vocabulary size and dataset, we plan to train the model on larger vocabulary and datasets with subword units segmentation.
- For future work, we plan to train models on datasets which consist of only long sentences with more than 50 or 80 tokens to compare the performance of long-sentences translation of the approach and baseline model.

## Acknowledgments

We thank Assistant Professor Shindo Hiroyuki, Ouchi Hiroki, Michael Wentao Li of the NAIST Computational Linguistics Laboratory, and the reviewers for their valuable and constructive comments. Part of this work was supported by JSPS KAKENHI Grant Number JP17H06101.



## References

- Roei Aharoni and Yoav Goldberg. 2017. Towards string-to-tree neural machine translation. *CoRR*, abs/1704.04743.
- Roei Aharoni, Yoav Goldberg, and Yonatan Belinkov. 2016. Improving sequence to sequence learning for morphological inflection generation: The biu-mit systems for the sigmorphon 2016 shared task for morphological reinflection. In *Proceedings of the 14th Annual SIGMORPHON Workshop on Computational Research in Phonetics, Phonology, and Morphology*, pages 41–48.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In *Proceedings of ICLR 2015*.
- Luisa Bentivogli, Arianna Bisazza, Mauro Cettolo, and Marcello Federico. 2016. Neural versus phrase-based machine translation quality: a case study. *CoRR*, abs/1608.04631.
- Danqi Chen and Christopher D. Manning. 2014. A fast and accurate dependency parser using neural networks. In *Proceedings of EMNLP 2014*.
- Li Dong and Mirella Lapata. 2016. Language to logical form with neural attention. *Proceedings of ACL 2016*, pages 33–43.
- Chris Dyer, Adhiguna Kuncoro, Miguel Ballesteros, and Noah A. Smith. 2016. Recurrent neural network grammars. *CoRR*, abs/1602.07776.
- Akiko Eriguchi, Yoshimasa Tsuruoka, and Kyunghyun Cho. 2017. Learning to parse and translate improves neural machine translation. *CoRR*, abs/1702.03525.
- Orhan Firat, Kyunghyun Cho, and Yoshua Bengio. 2016. Multi-way, multilingual neural machine translation with a shared attention mechanism. *arXiv.org/1601.01073*.
- Alexandros Komninos. 2016. Dependency based embeddings for sentence classification tasks. In *Proceedings of NAACL-HLT 2016*, pages 1490–1500.
- Adhiguna Kuncoro, Miguel Ballesteros, Lingpeng Kong, Chris Dyer, Graham Neubig, and Noah A. Smith. 2016. What do recurrent neural network grammars learn about syntax? *CoRR*, abs/1611.05774.
- Omer Levy and Yoav Goldberg. 2014. Dependency-based word embeddings. In *Proceedings of ACL 2014*, pages 302–308.
- Xiaoqing Li, Jiajun Zhang, and Chengqing Zong. 2016. Towards zero unknown word in neural machine translation. In *Proceedings of IJCAI 2016*.
- Yang Liu, Furu Wei, Sujian Li, Heng Ji, Ming Zhou, and Houfeng Wang. 2015. A dependency-based neural network for relation classification. In *Proceedings of ACL 2015*, pages 285–290.
- Thang Luong, Hieu Pham, and Christopher D. Manning. 2015. Effective approaches to attention-based neural machine translation. In *Proceedings of EMNLP 2015*, pages 1412–1421.
- Haitao Mi, Baskaran Sankaran, Zhiguo Wang, and Abe Ittycheriah. 2016. Coverage embedding models for neural machine translation. In *Proceedings of EMNLP 2016*, pages 955–960.
- Kazuki Ono and Kenji Hatano. 2014. Dependency parsing and its application using hierarchical structure in Japanese language. *International Journal on Advances in Internet Technology*, vol 7 no 3, 4.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. *Proceedings of ACL 2002*, pages 311–318.
- Richard Socher, Andrej Karpathy, Quoc V. Le, Christopher D. Manning, and Andrew Y. Ng. 2014. Grounded compositional semantics for finding and describing images with sentences. In *Transactions of the Association for Computational Linguistics*, pages 2: 207–218.
- Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 2014. Sequence to sequence learning with neural networks. *Proceedings of NIPS 2014*.
- Antonio Toral and Víctor M. Sánchez-Cartagena. 2017. A multifaceted evaluation of neural versus phrase-based machine translation for 9 language directions. *CoRR*, abs/1701.02901.
- Zhaopeng Tu, Yang Liu, Lifeng Shang, Xiaohua Liu, and Hang Li. 2016a. Neural machine translation with reconstruction. In *Proceedings of Association for the Advancement of Artificial Intelligence 2016*.
- Zhaopeng Tu, Zhengdong Lu, Yang Liu, Xiaohua Liu, and Hang Li. 2016b. Coverage-based neural machine translation. *CoRR*, abs/1601.04811.
- Zhaopeng Tu, Zhengdong Lu, Yang Liu, Xiaohua Liu, and Hang Li. 2016c. Modeling coverage for neural machine translation. In *Proceedings of ACL 2016*, pages 76–85.
- Oriol Vinyals, Lukasz Kaiser, Terry Koo, Slav Petrov, Ilya Sutskever, and Geoffrey E. Hinton. 2014. Grammar as a foreign language. *CoRR*, abs/1412.7449.
- Shuangzhi Wu, Dongdong Zhang, Nan Yang, Mu Li, and Ming Zhou. 2017. Sequence-to-dependency neural machine translation. *Proceedings of ACL 2017*, pages 698–707.
- Biao Zhang, Deyi Xiong, and Jinsong Su. 2016. Recurrent neural machine translation. *CoRR*, abs/1607.08725.