

# Towards Bidirectional Hierarchical Representations for Attention-Based Neural Machine Translation

Baosong Yang<sup>†</sup>, Derek F. Wong<sup>†</sup>, Tong Xiao<sup>‡</sup>, Lidia S. Chao<sup>†</sup> and Jingbo Zhu<sup>‡</sup>

Presenter: Baosong Yang

<sup>†</sup>NLP<sup>2</sup>CT Lab, University of Macau

<sup>‡</sup>NiuTrans Lab, Northeastern University

Aug 16, 2017

# Contents

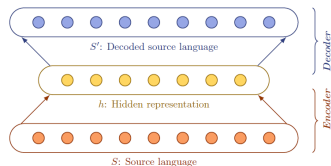
- 1 Introduction & Motivation
- 2 Bidirectional Hierarchical Model
- 3 Evaluations and Analysis

# Contents

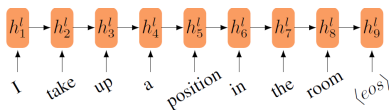
- 1 Introduction & Motivation
- 2 Bidirectional Hierarchical Model
- 3 Evaluations and Analysis

# Neural Machine Translation models

- Encoder-decoder framework
  - Encodes source sentences into a distributed representations.
  - Followed by a decoder generates target translation.



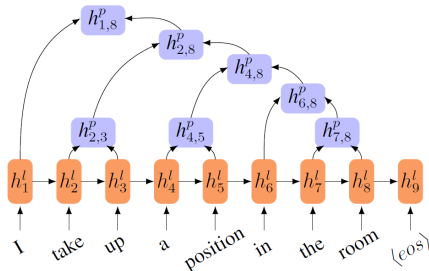
- Traditional sequential encoder
  - Insufficient to fully capture the semantics of a sentence (Tai et al., 2015; Eriguchi et al., 2016).



# Conventional Tree-based Encoder

## ■ Tree-based encoder

- Encodes a source sentence following a syntactic tree (Tai et al., 2015).
- Tree-to-sequence NMT model is outstanding on structurally distant language pair, e.g. en-jp (Eriguchi et al., 2016).

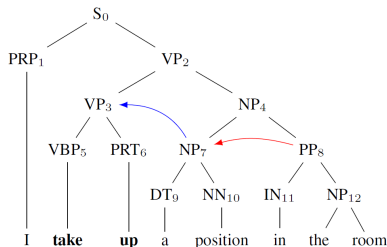


## ■ Problems:

- Recursively generate source representations in a **bottom-up** fashion.
- The learned representations are limited to local information, while failing to capture the global meaning of a sentence.
- Ignoring the neighboring context as well as the remote context.

# Conventional Tree-based Encoder

- An example:
  - **take up** has the meanings of *start doing something new*, *use space/time*, *accept an offer*, etc.
  - **a position** has the meanings of *location*, *job offer*, *rank/status*, etc.



- The differences in meaning arise as a result of ignoring the neighboring context as well as the remote context, i.e:
  - $h_{NP_7} \leftarrow h_{PP_8}$  (sibling)
  - $h_{VP_3} \leftarrow h_{NP_7}$  (child of sibling)

# Goal and Contributions

- Goal:
  - Improving tree-based encoder so that the generated source-side representations cover both local and global semantic information.
- Contributions:
  - Bidirectional tree-based encoder
    - To enhance the source-side hierarchical representations
  - Extending to the sub-word level
    - To alleviate the out-of-vocabulary problem
  - A variant weighted tree-based attention mechanism
    - To effectively leverage hierarchical representations

# Contents

- 1 Introduction & Motivation
- 2 Bidirectional Hierarchical Model**
- 3 Evaluations and Analysis



# Bidirectional Tree-based encoder

- Bidirectional Leaf-Node Encoding

- Jointly take into account both preceding and following annotations.
  - $h_i = [\vec{h}_i, \overleftarrow{h}_i]$

- Bottom-up Encoding

- Recursively propagated the local context to tree nodes.
  - $h_{par}^\uparrow = f_{tree-gru}(h_{l-child}^\uparrow, h_{r-child}^\uparrow)$

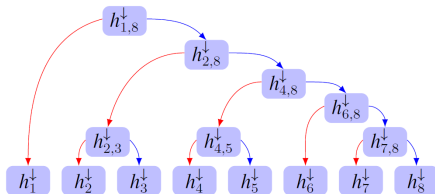
- Top-down Encoding

- Update the representations with global semantic information.

- $h_{l-child}^\downarrow = f_{GRU}^l(h_{l-child}^\uparrow, h_{par}^\downarrow)$

- $h_{r-child}^\downarrow = f_{GRU}^r(h_{r-child}^\uparrow, h_{par}^\downarrow)$

- $f_{GRU}^l$  and  $f_{GRU}^r$  with different parameters are applied to distinguish the left and right structural information.

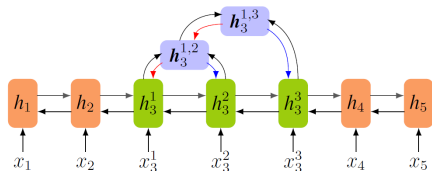


# Handling OOV: Tree-based Rare Word Encoding

## Motivation

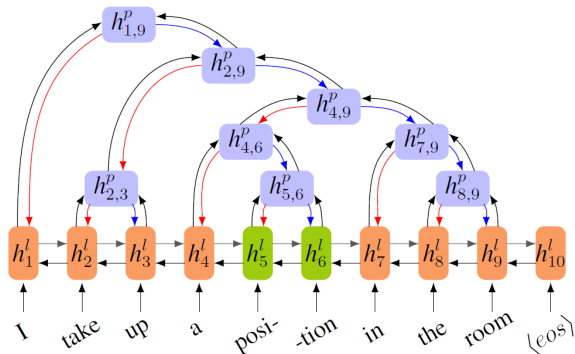
- Sequential sub-word encoding: Representing rare words as a sequence of sub-word units (Sennrich et al. 2016).
- Not applicable to the tree-based NMT model.

## Proposed method



1. Segment the rare word  $x_i$  into a sequence of sub-word units  $(x_i^1, \dots, x_i^n)$  by Byte-pair Encoding (BPE).
2. Built a binary lexical tree by simply composing two nodes in a rightwards fashion,  $((x_i^1, x_i^2), x_i^3) \dots, x_i^n$ .
3. Sub-word units are encoded following the binary lexical tree.

# Bidirectional Hierarchical Encoder



- The vector representations of the sentence, phrases, words as well as sub-word units are therefore based on the global context rather than local information.
- New problem: Attending fairly may cause the problem of over-translation.

# Weighted Variant of Attention Mechanism

- To balance the attentive information between the lexical and phrase vectors in the context vector.

$$d_j = (1 - \beta_j) \sum_{i=1}^n \alpha_j(i) h_i^l + \beta_j \sum_{k=1}^{n-1} \alpha_j(k) h_k^p$$

- $\alpha$  is the attention score which denotes the correspondence between each source annotation and the current target hidden state
- $\beta_j \in [0, 1]$  manually or automatically weights the expected importance of the representations.
- Gating scalar

$$\beta_j = \sigma(W_\beta c_{j-1} + b_\beta),$$

- is dominated by the target composite hidden state alone;
- is a time-dependent scalar;
- enables the attention model to explicitly quantify how far the leaf and no-leaf states contribute to the word prediction at each time step.

# Contents

- 1 Introduction & Motivation
- 2 Bidirectional Hierarchical Model
- 3 Evaluations and Analysis**

# Experimental settings

- Data set:

Training	Dev	Test		
LDC En-Ch	mt08	mt04	mt05	mt06
1.4M	1,357	1,788	1,082	1,664

- The vocabulary size of the training set:

Training set	Before BPE	After BPE
$ V $ in En	120k	40k
$ V $ in Zh	125k	40k

- NMT setting:

$ V $ size of source	40000
$ V $ size of target	40000
Dimension of word embedding	620
Dimension of leaf layer	512
Dimension of other layer	1024
Batch size	16
Beam search	5
Length of sentence	$\leq 40$
Optimizer	AdaDelta

# Evaluations on Hierarchical Encoder

## ■ Bidirectional tree-based encoding

Model	BPE	MT04	MT05	MT06	Dev.
<b>tree-based encoder</b>	no	31.90	24.68	24.40	17.63
+ bidirectional leaf-node encoding	no	32.13	24.94	25.02	18.12
+ top-down encoding	no	32.85	25.37	25.30	18.26
<b>hierarchical encoder</b> ( $\beta = 0.5$ )	no	32.91	25.55	25.52	18.46

- The future context at leaf level can contribute to word prediction.
- The translation quality is improved by considering the global semantic information.

# Evaluations on Hierarchical Encoder

- Tree-based rare word encoding

Model	BPE	MT04	MT05	MT06	Dev.
<b>sequential encoder</b>	no	31.26	23.98	24.02	17.20
+ sequential rare word encoding	yes	32.54	25.09	25.07	18.19
+ tree-based rare word encoding	yes	32.56	25.30	24.96	18.33
<b>tree-based encoder</b>	no	31.90	24.68	24.40	17.63
+ tree-based rare word encoding	yes	33.02	25.62	25.24	18.59
<b>hierarchical encoder</b> ( $\beta = 0.5$ )	no	32.91	25.55	25.52	18.46
<b>hierarchical encoder</b> ( $\beta = 0.5$ )	yes	33.81	26.47	26.31	19.41

- Achieves performance comparable to that of the standard BPE in the sequential model, but is applicable to the tree-based NMT model.



# Evaluations on Weighted Attention Model

## ■ Four cases:

- $\beta = 0.0$ : Ignore the phrase vectors.
- $\beta = 0.5$ : Non-leaf and leaf vectors participate equally.
- $\beta = 1.0$ : Only consider the phrase vectors.
- Gating scalar: Dynamically control the proportion.

Model	BLEU	Perplexity	Avg. Length
$\beta = 1.0$	17.16	98.65	21.13
$\beta = 0.5$	19.41	94.73	23.08
$\beta = 0.0$	19.83	94.68	23.33
Gating scalar	<b>20.10</b>	<b>94.18</b>	23.24

- Phrase representations are unable to fully capture the lexical information of the source sentence. ( $\beta = 1.0$ )
- Phrase representations tends to generate shorter translation. ( $\beta = 1.0$ , the average length of Ref. is 23.19.)
- Global information contributes to distinguishing the differences between word meanings (Compare  $\beta = 0.5$  with  $\beta = 0.0$ ).
- Through the use of the gating scalar, the hierarchical model achieves progressive improvements.

## Totally

Model	BPE	MT04	MT05	MT06	Dev.
<b>sequential encoder</b>	no	31.26	23.98	24.02	17.20
+ sequential rare word encoding	yes	32.54	25.09	25.07	18.19
+ tree-based rare word encoding	yes	32.56	25.30	24.96	18.33
<b>tree-based encoder</b>	no	31.90	24.68	24.40	17.63
+ bidirectional leaf-node encoding	no	32.13	24.94	25.02	18.12
+ top-down encoding	no	32.85	25.37	25.30	18.26
+ tree-based rare word encoding	yes	33.02	25.62	25.24	18.59
<b>hierarchical encoder</b> ( $\beta = 0.5$ )	no	32.91	25.55	25.52	18.46
<b>hierarchical encoder</b> ( $\beta = 0.5$ )	yes	33.81	26.47	26.31	19.41
+ gating scalar	yes	<b>34.33</b>	<b>26.72</b>	<b>26.58</b>	<b>20.10</b>

- Effectively model source-side representations from both the sequential and structural context.
- Outperform conventional models.

# Qualitative Analysis

## ■ A translation example.

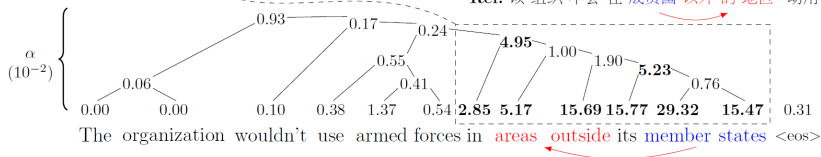
Our: 该组织不会在成员国以外的地区使用武力

$\beta$ : 0.17 0.14 0.22 **0.22** 0.27 0.22 0.19 0.44 0.14 0.56

sq-enc: 该组织不会使用其成员国以外的武装力量

tr-enc: 该组织不会在成员国境外使用武力

Ref: 该组织不会在成员国以外的地区动用军队



## ■ Translation examples of sub-words

Source	Reference	Hierarchical	Sequential
liu/jing/min	刘/敬/民 Liú/jìng/mín	刘/敬/民 Liú/jìng/mín	刘/敬/民 Liú/jìng/mín
adventur/er	探险家 Tàn xiǎn jiā	探险家 Tàn xiǎn jiā	探险者 Tàn xiǎn zhě
hi/k/ed	上调 Shàng tiáo	上升 Shàng shēng	发生 Fā shēng

**Thank you for listening.**