

Neural System Combination for Machine Translation

Long Zhou[†], Wenpeng Hu[†], Jiajun Zhang^{†*}, Chengqing Zong^{†‡}

[†]University of Chinese Academy of Sciences, Beijing, China

National Laboratory of Pattern Recognition, CASIA, Beijing, China

[‡]CAS Center for Excellence in Brain Science and Intelligence Technology, Shanghai, China

{long.zhou, wenpeng.hu, jjzhang, cqzong}@nlpr.ia.ac.cn

Abstract

Neural machine translation (NMT) becomes a new approach to machine translation and generates much more fluent results compared to statistical machine translation (SMT). However, SMT is usually better than NMT in translation adequacy. It is therefore a promising direction to combine the advantages of both NMT and SMT. In this paper, we propose a neural system combination framework leveraging multi-source NMT, which takes as input the outputs of NMT and SMT systems and produces the final translation. Extensive experiments on the Chinese-to-English translation task show that our model archives significant improvement by 5.3 BLEU points over the best single system output and 3.4 BLEU points over the state-of-the-art traditional system combination methods.

1 Introduction

Neural machine translation has significantly improved the quality of machine translation in recent several years (Kalchbrenner and Blunsom, 2013; Sutskever et al., 2014; Bahdanau et al., 2015; Junczys-Dowmunt et al., 2016a). Although most sentences are more fluent than translations by statistical machine translation (SMT) (Koehn et al., 2003; Chiang, 2005), NMT has a problem to address translation adequacy especially for the rare and unknown words. Additionally, it suffers from over-translation and under-translation to some extent (Tu et al., 2016). Compared to NMT, SMT, such as phrase-based machine translation (PBMT, (Koehn et al., 2003)) and hierarchical phrase-based machine translation (HPMT,

(Chiang, 2005)), does not need to limit the vocabulary and can guarantee translation coverage of source sentences. It is obvious that NMT and SMT have different strength and weakness. In order to take full advantages of both NMT and SMT, system combination can be a good choice.

Traditionally, system combination has been explored respectively in sentence-level, phrase-level, and word-level (Kumar and Byrne, 2004; Feng et al., 2009; Chen et al., 2009). Among them, word-level combination approaches that adopt confusion network for decoding have been quite successful (Rosti et al., 2007; Ayan et al., 2008; Freitag et al., 2014). However, these approaches are mainly designed for SMT without considering the features of NMT results. NMT opts to produce diverse words and free word order, which are quite different from SMT. And this will make it hard to construct a consistent confusion network. Furthermore, traditional system combination approaches cannot guarantee the fluency of the final translation results.

In this paper, we propose a neural system combination framework, which is adapted from the multi-source NMT model (Zoph and Knight, 2016). Different encoders are employed to model the semantics of the source language input and each best translation produced by different NMT and SMT systems. The encoders produce multiple context vector representations, from which the decoder generates the final output word by word. Since the same training data is used for NMT, SMT and neural system combination, we further design a smart strategy to simulate the real training data for neural system combination.

Specifically, we make the following contributions in this paper:

- We propose a neural system combination method, which is adapted from multi-source

*Corresponding author.

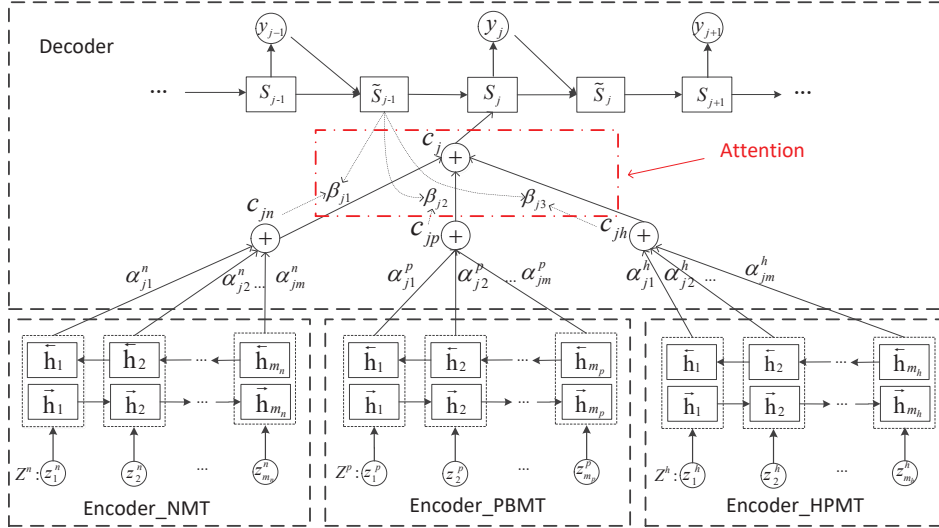


Figure 1: The architecture of Neural System Combination Model.

NMT model and can accommodate both source inputs and different system translations. It combines the fluency of NMT and adequacy (especially the ability to address rare words) of SMT.

- We design a good strategy to construct appropriate training data for neural system combination.
- The extensive experiments on Chinese-English translation show that our model archives significant improvement by 3.4 BLEU points over the state-of-the-art system combination methods and 5.3 BLEU points over the best individual system output.

2 Neural Machine Translation

The encoder-decoder NMT with an attention mechanism (Bahdanau et al., 2015) has been proposed to softly align each decoder state with the encoder states, and computes the conditional probability of the translation given the source sentence.

The encoder is a bidirectional neural network with gated recurrent units (GRU) (Cho et al., 2014) which reads an input sequence $X = (x_1, x_2, \dots, x_m)$ and encodes it into a sequence of hidden states $H = h_1, h_2, \dots, h_m$.

The decoder is a recurrent neural network that predicts a target sequence $Y = (y_1, y_2, \dots, y_n)$. Each word y_j is predicted based on a recurrent hidden state s_j , the previously predicted word y_{j-1} , and a context vector c_j . c_j is obtained from the

weighted sum of the annotations h_i . We use the latest implementation of attention-based NMT¹.

3 Neural System Combination for Machine Translation

Macherey and Och (2007) gave empirical evidence that these systems to be combined need to be almost uncorrelated in order to be beneficial for system combination. Since NMT and SMT are two kinds of translation models with large differences, we attempt to build a neural system combination model, which can take advantage of the different systems.

Model: Figure 1 illustrates the neural system combination framework, which can take as input the source sentence and the results of MT systems. Here, we use MT results as inputs to detail the model.

Formally, given the result sequences $Z(Z^n, Z^p, \text{ and } Z^h)$ of three MT systems for the same source sentence and previously generated target sequence $Y_{<j} = (y_1, y_2, \dots, y_{j-1})$, the probability of the next target word y_j is

$$p(y_j|Y_{<j}, Z) = \text{softmax}(f(c_j, y_{j-1}, s_j)) \quad (1)$$

Here $f(\cdot)$ is a non-linear function, y_{j-1} represents the word embedding of the previous prediction word, and s_j is the state of decoder at time step j , calculated by

$$s_j = \text{GRU}(\tilde{s}_{j-1}, c_j) \quad (2)$$

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

System	MT03	MT04	MT05	MT06	Ave
PBMT	37.47	41.20	36.41	36.03	37.78
HPMT	38.05	41.47	36.86	36.04	38.10
NMT	37.91	38.95	36.02	36.65	37.38
Jane (Freitag et al., 2014)	39.83	42.75	38.63	39.10	40.08
Multi	40.64	44.81	38.80	38.26	40.63
Multi+Source	42.16	45.51	40.28	39.03	41.75
Multi+Ensemble	41.67	45.95	40.37	39.02	41.75
Multi+Source+Ensemble	43.55	47.09	42.02	41.10	43.44

Table 1: Translation results (BLEU score) for different machine translation and system combination methods. Jane is a open source machine translation system combination toolkit that uses confusion network decoding. **Best** and **important** results per category are highlighted.

$$\tilde{s}_{j-1} = GRU(s_{j-1}, y_{j-1}) \quad (3)$$

where s_{j-1} is previous hidden state, \tilde{s}_{j-1} is an intermediate state. And c_j is the context vector of system combination obtained by attention mechanism, which is computed as weighted sum of the context vectors of three MT systems, just as illustrated in the middle part of Figure 1.

$$c_j = \sum_{k=1}^K \beta_{jk} c_{jk} \quad (4)$$

where K is the number of MT systems, and β_{jk} is a normalized item calculated as follows:

$$\beta_{jk} = \frac{\exp(\tilde{s}_{j-1} \cdot c_{jk})}{\sum_{k'} \exp(\tilde{s}_{j-1} \cdot c_{jk'})} \quad (5)$$

Here, we calculate k th MT system context c_{jk} as a weighted sum of the source annotations:

$$c_{jk} = \sum_{i=1}^m \alpha_{ji}^k h_i \quad (6)$$

where $h_i = [\vec{h}_i; \overleftarrow{h}_i]$ is the annotation of z_i from a bi-directional GRU, and its weight α_{ji}^k is computed by

$$\alpha_{ji}^k = \frac{\exp(e_{ji})}{\sum_{l=1}^m \exp(e_{jl})} \quad (7)$$

where $e_{ji} = v_a^T \tanh(W_a \tilde{s}_{j-1} + U_a h_i)$ scores how well \tilde{s}_{j-1} and h_i match.

Training Data Simulation: The neural system combination framework should be trained on the outputs of multiple translation systems and the gold target translations. In order to keep consistency in training and testing, we design a strategy

to simulate the real scenario. We randomly divide the training corpus into two parts, then reciprocally train the MT system on one half and translate the source sentences of the other half into target translations. The MT translations and the gold target reference can be available.

4 Experiments

We perform our experiments on the Chinese-English translation task. The MT systems participating in system combination are PBMT, HPMT and NMT. The evaluation metric is case-insensitive BLEU (Papineni et al., 2002).

4.1 Data preparation

Our training data consists of 2.08M sentence pairs extracted from LDC corpus. We use NIST 2003 Chinese-English dataset as the validation set, NIST 2004-2006 datasets as test sets. We list all the translation methods as follows:

- **PBMT:** It is the start-of-the-art phrase-based SMT system. We use its default setting and train a 4-gram language model on the target portion of the bilingual training data.
- **HPMT:** It is a hierarchical phrase-based SMT system, which uses its default configuration as PBMT in Moses.
- **NMT:** It is an attention-based NMT system, with the same setting given in section 2.

4.2 Training Details

The hyper-parameters used in our neural combination system are described as follows. We limit both Chinese and English vocabulary to 30k in our experiments. The number of hidden units is 1000

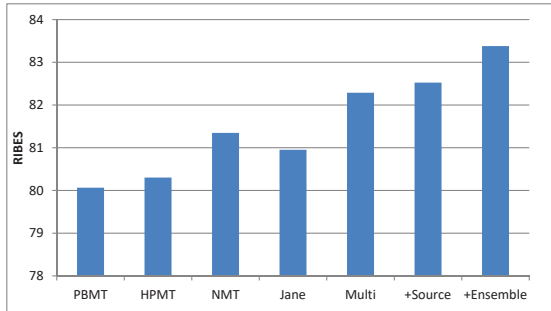


Figure 2: Translation results (RIBES score) for different machine translation and system combination methods.

and the word embedding dimension is 500 for all source and target word. The network parameters are updated with Adadelta algorithm. We adopt beam search with beam size $b=10$ at test time.

As to confusion-network-based system Jane, we use its default configuration and train a 4-gram language model on target data and 10M Xinhua portion of Gigaword corpus.

4.3 Main Results

We compare our neural combination system with the best individual engines, and the state-of-the-art traditional combination system Jane (Freitag et al., 2014). Table 1 shows the BLEU of different models on development data and test data. The BLEU score of the multi-source neural combination model is 2.53 higher than the best single model HPMT. The source language input gives a further improvement of +1.12 BLEU points.

As shown in Table 1, Jane outperforms the best single MT system by 1.92 BLEU points. However, our neural combination system with source language gets an improvement of 1.67 BLEU points over Jane. Furthermore, when augmenting our neural combination system with ensemble decoding², it leads to another significant boost of +1.69 BLEU points.

4.4 Word Order of Translation

We evaluate word order by the automatic evaluation metrics RIBES (Isozaki et al., 2010), whose score is a metric based on rank correlation coefficients with word precision. RIBES is known to have stronger correlation with human judgements than BLEU for English as discussed in Isozaki et al. (2010).

²We use four neural combination models in ensemble model.

System	MT03	MT04	MT05	MT06	Ave
NMT	1086	1145	1020	708	989.8
Ours	869	1023	909	609	852.5

Table 2: The number of unknown words in the results of NMT and our model.

System	MT03	MT04	MT05	MT06	Ave
E-NMT	39.14	40.78	37.31	37.89	38.78
Jane	40.61	43.28	39.05	39.18	40.53
Ours	43.61	47.65	42.02	41.17	43.61

Table 3: Translation results (BLEU score) when we replace original NMT with strong E-NMT, which uses ensemble strategy with four NMT models. All results of system combination are based on strong outputs of E-NMT.

Figure 2 illustrates experimental results of RIBES scores, which demonstrates that our neural combination model outperforms the best result of single MT system and Jane. Additionally, although BLEU point of Jane is higher than single NMT system, the word order of Jane is worse in terms of RIBES.

4.5 Rare and Unknown Words Translation

It is difficult for NMT systems to handle rare words, because low-frequency words in training data cannot capture latent translation mappings in neural network model. However, we do not need to limit the vocabulary in SMT, which are often able to translate rare words in training data. As shown in Table 2, the number of unknown words of our proposed model is 137 fewer than original NMT model.

Table 4 shows an example of system combination. The Chinese word zǔzhīwǎng is an out-of-vocabulary(OOV) for NMT and the baseline NMT cannot correctly translate this word. Although PBMT and HPMT translate this word well, they does not conform to the grammar. By combining the merits of NMT and SMT, our model gets the correct translation.

4.6 Effect of Ensemble Decoding

The performance of candidate systems is very important to the result of system combination, and we use ensemble strategy with four NMT models to improve the performance of original NMT system. As shown in Table 3, the E-NMT with

Source	海珊也与恐怖组织网建立了联系。
Pinyin	<i>hǎishān yě yǔ kǒngbù zǔzhīwǎng jiànlì le liánxì .</i>
Reference	hussein has also established ties with terrorist networks .
PBMT	hussein also has established relations and terrorist group .
HPMT	hussein also and terrorist group established relations .
NMT	hussein also established relations with UNK .
Jane	hussein also has established relations with .
Multi	hussein also has established relations with the terrorist group .

Table 4: Translation examples of single system, Jane and our proposed model.

ensemble strategy outperforms the original NMT system by +1.40 BLEU points, and it has become the best system in all MT systems, which is +0.68 BLEU points higher than HPMT.

After replacing original NMT with strong E-NMT, Jane outperforms original result by +0.45 BLEU points, and our model gets an improvement of +3.08 BLEU points over Jane. Experiments further demonstrate that our proposed model is effective and robust for system combination.

5 Related Work

The recently proposed neural machine translation has drawn more and more attention. Most of the existing approaches and models mainly focus on designing better attention models (Luong et al., 2015a; Mi et al., 2016a,b; Tu et al., 2016; Meng et al., 2016), better strategies for handling rare and unknown words (Luong et al., 2015b; Li et al., 2016; Zhang and Zong, 2016a; Sennrich et al., 2016b), exploiting large-scale monolingual data (Cheng et al., 2016; Sennrich et al., 2016a; Zhang and Zong, 2016b), and integrating SMT techniques (Shen et al., 2016; Junczys-Dowmunt et al., 2016b; He et al., 2016).

Our focus in this work is aiming to take advantage of NMT and SMT by system combination, which attempts to find consensus translations among different machine translation systems. In past several years, word-level, phrase-level and sentence-level system combination methods were well studied (Bangalore et al., 2001; Rosti et al., 2008; Li and Zong, 2008; Li et al., 2009; Heafield and Lavie, 2010; Freitag et al., 2014; Ma and Mckeown, 2015; Zhu et al., 2016), and reported state-of-the-art performances in benchmarks for SMT. Here, we propose a neural system combination model which combines the advantages of NMT and SMT efficiently.

Recently, Niehues et al. (2016) use phrase-

based SMT to pre-translate the inputs into target translations. Then a NMT system generates the final hypothesis using the pre-translation. Moreover, multi-source MT has been proved to be very effective to combine multiple source languages (Och and Ney, 2001; Zoph and Knight, 2016; Firat et al., 2016a,b; Garmash and Monz, 2016). Unlike previous works, we adapt multi-source NMT for system combination and design a good strategy to simulate the real training data for our neural system combination.

6 Conclusion and Future Work

In this paper, we propose a novel neural system combination framework for machine translation. The central idea is to take advantage of NMT and SMT by adapting the multi-source NMT model. The neural system combination method cannot only address the fluency of NMT and the adequacy of SMT, but also can accommodate the source sentences as input. Furthermore, our approach can further use ensemble decoding to boost the performance compared to traditional system combination methods.

Experiments on Chinese-English datasets show that our approaches obtain significant improvements over the best individual system and the state-of-the-art traditional system combination methods. In the future work, we plan to encode n-best translation results to further improve the system combination quality. Additionally, it is interesting to extend this approach to other tasks like sentence compression and text abstraction.

Acknowledgments

The research work has been funded by the Natural Science Foundation of China under Grant No. 61333018 and No. 61673380, and it is also supported by the Strategic Priority Research Program of the CAS under Grant No. XDB02070007.

References

- Necip Fazil Ayan, Jing Zheng, and Wen Wang. 2008. *Improving alignments for better confusion networks for combining machine translation systems*. In Proceedings of COLING 2008.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. *Nuural machine translation by jointly learning to align and translate*. In Proceedings of ICLR 2015.
- Srinivas Bangalore, German Bordel, and Giuseppe Richardi. 2001. *Computing consensus translation from multiple machine translation systems*. In Proceedings of IEEE ASRU.
- Boxing Chen, Min Zhang, Haizhou Li, and Aiti Aw. 2009. *A comparative study of hypothesis alignment and its improvement for machine translation system combination*. In Proceedings of ACL 2009.
- Yong Cheng, Wei Xu, Zhongjun He, Wei He, Hua Wu, Maosong Sun, and Yang Liu. 2016. *Semi-supervised learning for neural machine translation*. In Proceedings of ACL 2016.
- David Chiang. 2005. *A hierarchical phrase-based model for statistical machine translation*. In Proceedings of ACL 2005.
- Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. *Learning phrase representations using RNN encoder - decoder for statistical machine translation*. In Proceedings of EMNLP 2014.
- Yang Feng, Yang Liu, Haitao Mi, Qun Liu, and Yajuan Lu. 2009. *Lattice-based system combination for statistical machine translation*. In Proceedings of ACL 2009.
- Orhan Firat, Kyunghyun Cho, and Yoshua Bengio. 2016a. *Multi-way, multilingual neural machine translation with a shared attention mechanism*. In Proceedings of NAACL-HLT 2016.
- Orhan Firat, Baskaran Sankaran, Yaser Al-Onaizan, Fatos T. Yarman Vural, and Kyunghyun Cho. 2016b. *Zero-resource translation with multi-lingual neural machine translation*. In Proceedings of EMNLP 2016.
- Markus Freitag, Matthias Huck, and Hermann Ney. 2014. *Jane: open source machine translation system combination*. In Proceedings of EACL 2014.
- Ekaterina Garmash and Christof Monz. 2016. *Ensemble learning for multi-source neural machine translation*. In Proceedings of COLING 2016.
- Wei He, Zhongjun He, Hua Wu, and Haifeng Wang. 2016. *Improved neural machine translation with SMT features*. In Proceedings of AAAI 2016.
- Kenneth Heafield and Alon Lavie. 2010. *Combining machine translation output with open source*. The Prague Bulletin of Machematical Linguistics.
- Hideki Isozaki, Tsutomu Hirao, Kevin Duh, Katsuhito Sudoh, and Hajime Tsukada. 2010. *Automatic evaluation of translation quality for distant language pairs*. In Proceedings of EMNLP 2010.
- Marcin Junczys-Dowmunt, Tomasz Dwojak, and Hieu Hoang. 2016a. *Is neural machine translation ready for deployment? A case study on 30 translation directions*. In Proceedings of IWSLT 2016.
- Marcin Junczys-Dowmunt, Tomasz Dwojak, and Rico Sennrich. 2016b. *The AMU-UEDIN submission to the WMT16 news translation task: attention-based NMT models as feature functions in phrase-based SMT*. In Proceedings of WMT 2016.
- Nal Kalchbrenner and Phil Blunsom. 2013. *Recurrent continuous translation models*. In Proceedings of EMNLP 2013.
- Philipp Koehn, Franz J. Och, and Daniel Marcu. 2003. *Statistical phrase-based translation*. In Proceedings of ACL NAACL 2013.
- Shankar Kumar and William Byrne. 2004. *Minimum bayes-risk decoding for statistical machine translation*. In Proceedings of HLT-NAACL 2004.
- Maoxi Li, Jiajun Zhang, Yu Zhou, and Chengqing Zong. 2009. *The CASIA statistical machine translation system for IWSLT 2009*. In Proceedings of IWSLT2009.
- Maoxi Li and Chengqing Zong. 2008. *Word reordering alignment for combination of statistical machine translation systems*. In Proceedings of the International Symposium on Chinese Spoken Language Processing.
- Xiaoqing Li, Jiajun Zhang, and Chengqing Zong. 2016. *Towards zero unknown word in neural machine translation*. In Proceedings of IJCAI 2016.
- Minh-Thang Luong, Hieu Pham, and Christopher D. Manning. 2015a. *Effective approaches to attention-based neural machine translation*. In Proceedings of EMNLP 2015.
- Minh-Thang Luong, Ilya Sutskever, Quoc V Le, Oriol Vinyals, and Wojciech Zaremba. 2015b. *Addressing the rare word problem in neural machine translation*. In Proceedings of ACL 2015.
- Wei-Yun Ma and Kathleen Mckeown. 2015. *System combination for machine translation through paraphrasing*. In Proceedings of EMNLP 2015.
- Wolfgang Macherey and Franz Josef Och. 2007. *An empirical study on computing consensus translations from multiple machine translation systems*. In Proceedings of EMNLP 2007.

- Fandong Meng, Zhengdong Lu, Hang Li, and Qun Liu. 2016. *Interactive attention for neural machine translation*. In Proceedings of COLING 2016.
- Haitao Mi, Baskaran Sankaran, Zhiguo Wang, Niyu Ge, and Abe Ittycheriah. 2016a. *A coverage embedding model for neural machine translation*. In Proceedings of EMNLP 2016.
- Haitao Mi, Zhiguo Wang, Niyu Ge, and Abe Ittycheriah. 2016b. *Supervised attentions for neural machine translation*. In Proceedings of EMNLP 2016.
- Jan Niehues, Eunah Cho, Thanh-Le Ha, and Alex Waibel. 2016. *Pre-translation for neural machine translation*. In Proceedings of COLING 2016.
- Franz Josef Och and Hermann Ney. 2001. *Statistical multi-source translation*. In Proceedings of MT Summit.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. *Bleu: a method for automatic evaluation of machine translation*. In Proceedings of ACL 2002.
- Antti-Veikko I. Rosti, Spyros Matsoukas, and Richard Schwartz. 2007. *Improved word-level system combination for machine translation*. In Proceedings of ACL 2007.
- Antti-Veikko I. Rosti, Bing Zhang, Spyros Matsoukas, and Richard Schwartz. 2008. *Incremental hypothesis alignment for building confusion networks with application to machine translation systems combination*. In Proceedings of the Third ACL Workshop on Statistical Machine Translation.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016a. *Improving neural machine translation models with monolingual data*. In Proceedings of ACL 2016.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016b. *Neural machine translation of rare words with subword units*. In Proceedings of ACL 2016.
- Shiqi Shen, Yong Cheng, Zhongjun He, Wei He, Hua Wu, Maosong Sun, and Yang Liu. 2016. *Minimum risk training for neural machine translation*. In Proceedings of ACL 2016.
- Ilya Sutskever, Oriol Vinyals, and Quoc VV Le. 2014. *Sequence to sequence learning with neural networks*. In Proceedings of NIPS 2014.
- Zhaopeng Tu, Zhengdong Lu, Yang Liu, Xiaohua Liu, and Hang Li. 2016. *Modeling coverage for neural machine translation*. In Proceedings of ACL 2016.
- Jiajun Zhang and Chengqing Zong. 2016a. *Bridging neural machine translation and bilingual dictionaries*. arXiv preprint arXiv:1610.07272.
- Jiajun Zhang and Chengqing Zong. 2016b. *Exploiting source-side monolingual data in neural machine translation*. In Proceedings of EMNLP 2016.
- Junguo Zhu, Muyun Yang, Sheng Li, and Tiejun Zhao. 2016. *Sentence-level paraphrasing for machine translation system combination*. In Proceedings of ICYCSEE 2016.
- Barret Zoph and Kevin Knight. 2016. *Multi-source neural translation*. In Proceedings of NAACL-HLT 2016.