



## UvA-DARE (Digital Academic Repository)

### Optimizing Transformer for Low-Resource Neural Machine Translation

Araabi, A.; Monz, C.

**DOI**

[10.18653/v1/2020.coling-main.304](https://doi.org/10.18653/v1/2020.coling-main.304)

**Publication date**

2020

**Document Version**

Final published version

**Published in**

The 28th International Conference on Computational Linguistics

**License**

CC BY

[Link to publication](#)

**Citation for published version (APA):**

Araabi, A., & Monz, C. (2020). Optimizing Transformer for Low-Resource Neural Machine Translation. In D. Scott, N. Bel, & C. Zong (Eds.), *The 28th International Conference on Computational Linguistics: COLING 2020 : Proceedings of the Conference : December 8-13, 2020, Barcelona, Spain (Online)* (pp. 3429-3435). International Committee on Computational Linguistics. <https://doi.org/10.18653/v1/2020.coling-main.304>

**General rights**

It is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), other than for strictly personal, individual use, unless the work is under an open content license (like Creative Commons).

**Disclaimer/Complaints regulations**

If you believe that digital publication of certain material infringes any of your rights or (privacy) interests, please let the Library know, stating your reasons. In case of a legitimate complaint, the Library will make the material inaccessible and/or remove it from the website. Please Ask the Library: <https://uba.uva.nl/en/contact>, or a letter to: Library of the University of Amsterdam, Secretariat, Singel 425, 1012 WP Amsterdam, The Netherlands. You will be contacted as soon as possible.

# Optimizing Transformer for Low-Resource Neural Machine Translation

**Ali Araabi**  
Informatics Institute  
University of Amsterdam  
a.araabi@uva.nl

**Christof Monz**  
Informatics Institute  
University of Amsterdam  
c.monz@uva.nl

## Abstract

Language pairs with limited amounts of parallel data, also known as low-resource languages, remain a challenge for neural machine translation. While the Transformer model has achieved significant improvements for many language pairs and has become the de facto mainstream architecture, its capability under low-resource conditions has not been fully investigated yet. Our experiments on different subsets of the IWSLT14 training data show that the effectiveness of Transformer under low-resource conditions is highly dependent on the hyper-parameter settings. Our experiments show that using an optimized Transformer for low-resource conditions improves the translation quality up to 7.3 BLEU points compared to using the Transformer default settings.

## 1 Introduction

Despite the success of Neural Machine Translation (NMT) (Sutskever et al., 2014; Cho et al., 2014; Bahdanau et al., 2015), for the vast majority of language pairs for which only limited amounts of training data exist (a.k.a. low-resource languages), the performance of NMT systems is relatively poor (Koehn and Knowles, 2017; Gu et al., 2018a). Most approaches focus on exploiting additional data to address this problem (Gülçehre et al., 2015; Sennrich et al., 2016; He et al., 2016; Fadaee et al., 2017). However, Sennrich and Zhang (2019) show that a well-optimized NMT system can perform relatively well under low-resource data conditions. Unfortunately, their results are confined to a recurrent NMT architecture (Sennrich et al., 2017), and it is not clear to what extent these findings also hold for the nowadays much more commonly used Transformer architecture (Vaswani et al., 2017).

Like all NMT models, Transformer requires setting various hyper-parameters but researchers often stick to the default parameters, even when their data conditions differ substantially from the original data conditions used to determine those default values (Gu et al., 2018b; Aharoni et al., 2019).

In this paper, we explore to what extent hyper-parameter optimization, which has been applied successfully to recurrent NMT models for low-resource translation, is also beneficial for the Transformer model. We show that with the appropriate settings, including the number of BPE merge operations, attention heads, and layers up to the degree of dropout and label smoothing, translation performance can be increased substantially, even for data sets with as little as 5k sentence pairs. Our experiments on different corpus sizes, ranging from 5k to 165k sentence pairs, show the importance of choosing the optimal settings with respect to data size.

## 2 Hyper-Parameter Exploration

In this section, we first discuss the importance of choosing an appropriate degree of subword segmentation before we describe the other optimal hyper-parameter settings.

**Vocabulary representation.** In order to improve the translation of rare words, word segmentation approaches such as Byte-Pair-Encoding (BPE) (Sennrich et al., 2016) have become standard practice in NMT. This is especially true for language pairs with small amounts of data where rare words are a common phenomenon. Sennrich and Zhang (2019) show that reducing the number of BPE merge operations

---

This work is licensed under a Creative Commons Attribution 4.0 International Licence. Licence details: <http://creativecommons.org/licenses/by/4.0/>.

Step	Hyper-parameter	Values
1	feed-forward dimension	128, 256, 512, 1024, <u>2048</u> , 4096
2	embedding dimension	256, <u>512</u> , 1024
3	attention heads	1, 2, 4, <u>8</u> , 16
4	dropout	<u>0.1</u> , 0.2, 0.3, 0.4, 0.5
5	number of layers	1, 2, 3, 4, 5, <u>6</u> , 7
6	label smoothing	<u>0.1</u> , 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8
7	enc/dec layer dropout	<u>0</u> , 0.1, 0.2, 0.3, 0.4
8	src/tgt word dropout	<u>0</u> , 0.1, 0.2, 0.3
9	attention dropout	<u>0</u> , 0.1, 0.2, 0.3
10	activation dropout	<u>0</u> , 0.1, 0.2, 0.3, 0.4, 0.5
11	embedding layer normalization	yes, <u>no</u>
12	batch size	512, 1024, 2048, <u>4096</u> , 8192, 12288
13	learning rate scheduler	<u>Transformer standard</u> , inverse square root
14	warm-up steps	2000, <u>4000</u> , 5000, 6000, 8000, 10000
15	learning rate	0.01, 0.001, 0.0001, 0.00001

Table 1: Order in which different hyper-parameters are explored and the corresponding values considered for each hyper-parameter. Underlined values indicate the default value.

can result in substantial improvements of up to 5 BLEU points for a recurrent NMT model. It is natural to assume that reducing the BPE vocabulary is similarly effective for Transformer.

**Architecture tuning.** A current observation in neural networks, and in particular in Transformer architectures, is that increasing the number of model parameters improves performance (Raffel et al., 2019; Wang et al., 2019). However, those findings are mostly obtained for scenarios with ample training data and it is not clear if they are directly applicable to low-resource conditions. While Biljon et al. (2020) show that using fewer Transformer layers improves the quality of low-resource NMT, we expand our exploration towards the effects of using a narrow and shallow Transformer by reducing i) the number of layers in both the encoder and decoder, ii) the number of attention heads, iii) feed-forward layer dimension ( $d_{ff}$ ), and iv) embedding dimensions ( $d_{model}$ ).

**Regularization.** Following Sennrich and Zhang (2019), we analyze the impact of regularization by applying dropouts to various Transformer components (Konda et al., 2015). In addition to regular dropout which is applied to the output of each sub-layer (feed-forward and self-attention) and after adding the positional embedding in both encoder and decoder (Vaswani et al., 2017), we employ attention dropout after the softmax for self-attention and also activation dropout inside the feed-forward sub-layers. Moreover, we drop entire layers using layer dropout (Fan et al., 2020). We further drop words in the embedding matrix using discrete word dropout (Gal and Ghahramani, 2016). We also experiment with larger label-smoothing factors (Müller et al., 2019).

### 3 Experiments

#### 3.1 Experimental setup

Exploring all possible values for several hyper-parameters at once is prohibitively expensive from a computational perspective. Possible ways to circumvent this are random search (Bergstra and Bengio, 2012) or grid search for one hyper-parameter at a time. For simplicity, we opt for the latter. Table 1 shows the order in which the hyper-parameters are tuned. Once the optimal value of a hyper-parameter has been determined, it remains fixed for later steps; see Table 2. Obviously, there are no guarantees that this will result in a global optimum.

To be comparable with Sennrich and Zhang (2019), we take the TED data from the IWSLT 2014 German-English (De-En) shared translation task (Cettolo et al., 2014). We apply punctuation normal-

ID	System	BLEU					
		5k	10k	20k	40k	80k	165k
1	Transformer-big	3.3	3.4	4.3	4.7	5.1	5.5
2	Transformer-base	8.3	11.9	16.8	23.2	28.0	32.1
3	2 + feed-forward dimension (2048 $\rightarrow$ 512)	8.8	12.0	16.7	22.3	27.7	31.7
4	3 + attention heads (8 $\rightarrow$ 2)	9.2	12.7	19.0	23.6	28.7	32.3
5	4 + dropout (0.1 $\rightarrow$ 0.3)	10.6	17.0	21.9	26.7	<b>31.0</b>	<b>33.4</b>
6	5 + layers (6 $\rightarrow$ 5)	10.9	16.9	21.9	26.0	30.2	33.0
7	6 + label smoothing (0.1 $\rightarrow$ 0.6)	11.3	16.5	22.0	26.9	30.4	33.3
8	7 + decoder layerDrop (0 $\rightarrow$ 0.3)	12.9	17.3	22.5	26.9	30.3	33.1
9	8 + target word dropout (0 $\rightarrow$ 0.1)	13.7	18.1	23.1	27.0	30.7	33.0
10	9 + activation dropout (0 $\rightarrow$ 0.3)	<b>14.3</b>	<b>18.3</b>	<b>23.6</b>	<b>27.4</b>	30.4	32.6

Table 2: Results of Transformer optimized on the 5k dataset for different subsets and full corpus of IWSLT14 German  $\rightarrow$  English. Averages over three runs from three different samples are reported.

ization, tokenization, data cleaning, and truecasing using the Moses scripts (Koehn et al., 2007). We also limit the sentence length to a maximum of 175 tokens during training. Our pre-processing pipeline results in 165,667 sentence pairs for training and 1,938 sentence pairs for development. In order to create smaller training sets, we randomly sample 5k, 10k, 20k, 40k, and 80k sentence pairs from the training data. Similar to Sennrich and Zhang (2019), we use the concatenation of the IWSLT 2014 dev sets (tst2010–2012, dev2010, dev2012) as our test set, which consists of 6,750 sentence pairs.

For actual low-resource languages, we evaluate our optimized systems on the original test sets of Belarusian (Be), Galician (Gl), and Slovak (Sk) TED talks (Qi et al., 2018) and also Slovenian (Sl) from IWSLT2014 (Cettolo et al., 2012) with training sets ranging from 4.5k to 55k sentence pairs.

We use Transformer-base and Transformer-big as our baselines, with the hyper-parameters and optimizer settings described in (Vaswani et al., 2017). We use the Fairseq library (Ott et al., 2019) for our experiments and sacreBLEU (Post, 2018) as evaluation metric.

### 3.2 Results and discussions

**BPE effect.** To evaluate the effect of different degrees of BPE segmentation on performance, we consider merge operations ranging from 1k to 30k, training BPE on the full training corpus instead of subsets and also removing infrequent subword units when applying the BPE model. In contrast to earlier results for an RNN model, we observe that discarding infrequent subword units under extreme low-resource conditions is detrimental to the performance of Transformer. Sennrich and Zhang (2019) report that reducing BPE merge operations from 30k to 5k improves performance (+4.9 BLEU). We find that the same reduction in merge operations affects the Transformer model far less (+0.6 BLEU). We observe no significant differences between training BPE on the full training corpus and training on subsets. Thus, we always train BPE on subsets with an optimized number of merge operations (see Table 3).

**Architecture effect.** Table 2 shows the results of our system optimizations alongside the performance of our baselines. We notice that Transformer-big performs poorly on all datasets, which is most likely due to the much larger number of parameters requiring substantially larger training data. The system column in Table 2 shows our optimization steps on the 5k dataset, which are also applied to the larger datasets.

We gain substantial improvements over Transformer-base for various subset sizes. For the smallest dataset, as expected, reducing Transformer depth and width, including number of attention heads, feed-forward dimension, and number of layers along with increasing the rate of different regularization techniques is highly effective (+6 BLEU). The largest improvements are obtained by increasing the dropout rate (+1.4 BLEU), adding layer dropout to the decoder (+1.6 BLEU), and adding word dropout to the target side (+0.8 BLEU). Most of these findings also hold for the 10k and 20k datasets, but differ

	default	5k	10k	20k	40k	80k	165k
BPE operations	37k	5k	10k	10k	12k	15k	20k
feed-forward dimension	2048	512	1024	1024	2048	2048	2048
attention heads	8	2	2	2	2	2	4
dropout	0.1	0.3	0.3	0.3	0.3	0.3	0.3
layers	6	5	5	5	5	5	5
label smoothing	0.1	0.6	0.5	0.5	0.5	0.4	0.3
enc/dec layerDrop	0.0/0.0	0.0/0.3	0.0/0.2	0.0/0.2	0.0/0.1	0/0.1	0/0.1
src/tgt word dropout	0.0/0.0	0.0/0.1	0.0/0.1	0.1/0.1	0.1/0.1	0.2/0.2	0.2/0.2
activation dropout	0.0	0.3	0.3	0.3	0.3	0	0
batch size	4096	4096	4096	4096	4096	8192	12288

Table 3: Default parameters for Transformer-base and optimal settings for different dataset sizes based on the De→En development data.

sentences	words (En)	BLEU	
		T-base	T-opt
De→En			
5k	100k	8.3	14.3
10k	200k	11.9	18.7
20k	410k	16.8	24.1
40k	830k	23.2	28.6
80k	1.6M	28.0	31.9
165k	3.4M	32.1	35.2
Be→En (4.5k)			
90k	90k	5.0	8.1
Gl→En (10k)			
196k	196k	13.1	22.3
Sl→En (13k)			
269k	269k	9.1	15.5
Sk→En (55k)			
1.2M	1.2M	24.8	29.9

Table 4: Results for Transformer-base/optimized. T-opt results for Be, Gl, Sl, and Sk use the optimized settings on De→En development data for 5k, 10k, 10k, and 40k training examples, respectively.

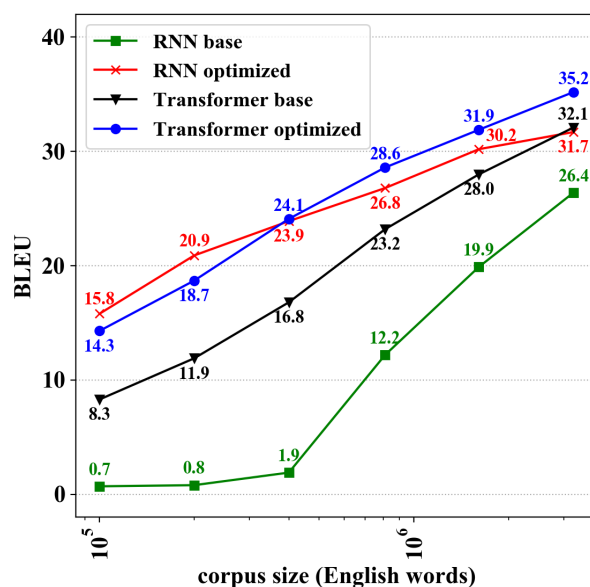


Figure 1: Comparison between RNN and Transformer with base and optimized settings.

for larger subsets. By applying these settings to the 10k, 20k, 40k, 80k, and 165k datasets, BLEU scores increase by +6.4, +6.8, +4.2, +2.4, and +0.5 points, respectively. However, the effect of each adjustment is different for each dataset. For example, reducing the feed-forward layer dimension to 512 is only effective for the two smallest subsets.

We also conducted experiments with different values for the learning rate and warm-up steps using the inverse-square root learning rate scheduler, as implemented within Fairseq (Ott et al., 2019), which is slightly different from the proposed learning rate scheduler in the original Transformer paper. However, we did not observe any improvements over the default Transformer learning rate scheduler.

**Optimized parameter settings.** Table 3 shows the optimal settings for each dataset size, achieved by tuning the parameters on the development data. We observe that a shallower Transformer combined with a smaller feed-forward layer dimension and BPE vocabulary size is more effective under lower-resource conditions. However, as mentioned above, Transformer is not as sensitive to the BPE vocabulary size as RNNs and reducing the embedding dimension size is not effective.

Vaswani et al. (2017) and Chen et al. (2018) show that reducing the number of attention heads decreases the BLEU score under high-resource conditions. Raganato et al. (2020) show that using one

sentences	words (En)	BLEU	
		T-base	T-opt
En→De			
5k	100k	6.4	11.3
10k	200k	9.3	15.6
20k	410k	13.5	20.8
40k	830k	20.2	24.5
80k	1.6M	24.2	27.2
165k	3.4M	27.4	29.8

Table 5: Results for En→De based on the optimal settings for De→En for the corresponding corpus size (see Table 3).

sentences	words (En)	BLEU	
		T-base	T-opt
En→Be (4.5k)	90k	3.6	6.6
En→Gl (10k)	196k	10.6	18.7
En→Sl (13k)	269k	6.8	12.2
En→Sk (55k)	1.2M	19.5	23.5

Table 6: Results for low-resource translation from English using the optimal settings from the De→En system with the closest number of parallel sentence pairs.

attention head does not cause much degradation on moderate-size training data. Our results show that it is even beneficial to use only two attention heads (+0.5 BLEU) under low-resource conditions.

While Sennrich and Zhang (2019) use a high dropout rate of 0.5 for their optimized RNN model, our findings suggest a lower rate of 0.3 for Transformer. In line with their results, we find word dropout effective for most low-resource conditions. Our results show that a higher degree of label smoothing and higher decoder layer dropout rates are beneficial for smaller data sizes and less effective for larger sizes.

Sennrich and Zhang (2019) report substantial gains by using small batch sizes. However, our results show that Transformer still requires larger batches, even under very low-resource conditions. It is worth mentioning that applying attention dropout did not result in improvements in our experiments.

**Optimized Transformer.** The results of our optimized systems for the corresponding subsets are shown in the upper half of Table 4 with improvements of up to 3 BLEU points over the results obtained in Table 2, indicating that under low-resource conditions, the optimal choice of Transformer parameters is highly sensitive with respect to the data size.

The BLEU improvements in the bottom half of Table 4 show that determining the optimal settings on one language pair (De→En) is also effective for actual low-resource language pairs, especially if the size of the training data is taken into account. Furthermore, the results in Table 5 show that the optimal settings for De→En also hold for the opposite translation direction of the same language pair. They even carry over to translating from English for actual low-resource language pairs, see Table 6, which can be considered the more challenging scenario (Aharoni et al., 2019). Note that the results of Tables 5 and 6 and the bottom half of Table 4 are obtained by using the closest systems optimized on De→En subsets with respect to their number of training sentences.

To compare Transformer with an RNN architecture, we replicate the baseline and optimized RNN for low-resource NMT, as described in (Sennrich and Zhang, 2019), on our datasets. Figure 1 shows the BLEU scores for different data sizes. Surprisingly, even without any hyper-parameter optimization, Transformer performs much better than the RNN model under very limited data conditions. However, the optimized Transformer only outperforms the optimized RNN with more than 20k training examples.

## 4 Conclusion

In this paper, we study the effects of hyper-parameter settings for the Transformer architecture under various low-resource data conditions. While our findings are largely in line with previous work (Sennrich and Zhang, 2019) for RNN-based models, we show that very effective optimizations for RNN-based models such as reducing the number of BPE merge operations or using small batch sizes are less effective or even hurt performance. Our experiments show that a proper combination of Transformer configurations combined with regularization techniques results in substantial improvements over a Transformer system with default settings for all low-resource data sizes. However, under extremely low-resource conditions an optimized RNN model still outperforms Transformer.

## References

- Roei Aharoni, Melvin Johnson, and Orhan Firat. 2019. Massively multilingual neural machine translation. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019*, pages 3874–3884.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In *3rd International Conference on Learning Representations, ICLR 2015*.
- James Bergstra and Yoshua Bengio. 2012. Random search for hyper-parameter optimization. *The Journal of Machine Learning Research*, 13(1):281–305.
- Elan Van Biljon, Arnu Pretorius, and Julia Kreutzer. 2020. On optimal transformer depth for low-resource language translation. *CoRR*, abs/2004.04418.
- Mauro Cettolo, Christian Girardi, and Marcello Federico. 2012. Wit3: Web inventory of transcribed and translated talks. In *Conference of European Association for Machine Translation*, pages 261–268.
- Mauro Cettolo, Jan Niehues, Sebastian Stüker, Luisa Bentivogli, and Marcello Federico. 2014. Report on the 11th iwslt evaluation campaign, iwslt 2014. In *Proceedings of the International Workshop on Spoken Language Translation, Hanoi, Vietnam*.
- Mia Xu Chen, Orhan Firat, Ankur Bapna, Melvin Johnson, Wolfgang Macherey, George F. Foster, Llion Jones, Mike Schuster, Noam Shazeer, Niki Parmar, Ashish Vaswani, Jakob Uszkoreit, Lukasz Kaiser, Zhifeng Chen, Yonghui Wu, and Macduff Hughes. 2018. The best of both worlds: Combining recent advances in neural machine translation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018*, pages 76–86.
- Kyunghyun Cho, Bart van Merriënboer, Çağlar Gülçehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using RNN encoder-decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014*, pages 1724–1734.
- Marzieh Fadaee, Arianna Bisazza, and Christof Monz. 2017. Data augmentation for low-resource neural machine translation. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017*, pages 567–573.
- Angela Fan, Edouard Grave, and Armand Joulin. 2020. Reducing transformer depth on demand with structured dropout. In *8th International Conference on Learning Representations, ICLR 2020*.
- Yarin Gal and Zoubin Ghahramani. 2016. A theoretically grounded application of dropout in recurrent neural networks. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems*, pages 1019–1027.
- Jiatao Gu, Hany Hassan, Jacob Devlin, and Victor O. K. Li. 2018a. Universal neural machine translation for extremely low resource languages. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2018*, pages 344–354.
- Jiatao Gu, Yong Wang, Yun Chen, Victor O. K. Li, and Kyunghyun Cho. 2018b. Meta-learning for low-resource neural machine translation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3622–3631.
- Çağlar Gülçehre, Orhan Firat, Kelvin Xu, Kyunghyun Cho, Loïc Barrault, Hui-Chi Lin, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2015. On using monolingual corpora in neural machine translation. *CoRR*, abs/1503.03535.
- Di He, Yingce Xia, Tao Qin, Liwei Wang, Nenghai Yu, Tie-Yan Liu, and Wei-Ying Ma. 2016. Dual learning for machine translation. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems*, pages 820–828.
- Philipp Koehn and Rebecca Knowles. 2017. Six challenges for neural machine translation. In *Proceedings of the First Workshop on Neural Machine Translation, NMT@ACL 2017*, pages 28–39.
- Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondrej Bojar, Alexandra Constantin, and Evan Herbst. 2007. Moses: Open source toolkit for statistical machine translation. In *ACL 2007, Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics*.

- Kishore Reddy Konda, Xavier Bouthillier, Roland Memisevic, and Pascal Vincent. 2015. Dropout as data augmentation. *CoRR*, abs/1506.08700.
- Rafael Müller, Simon Kornblith, and Geoffrey E. Hinton. 2019. When does label smoothing help? In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems*, pages 4696–4705.
- Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. Fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019*, pages 48–53.
- Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers, WMT 2018, Belgium, Brussels, October 31 - November 1, 2018*, pages 186–191.
- Ye Qi, Devendra Singh Sachan, Matthieu Felix, Sarguna Padmanabhan, and Graham Neubig. 2018. When and why are pre-trained word embeddings useful for neural machine translation? In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT*, pages 529–535.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. *CoRR*, abs/1910.10683.
- Alessandro Raganato, Yves Scherrer, and Jörg Tiedemann. 2020. Fixed encoder self-attention patterns in transformer-based machine translation. *CoRR*, abs/2002.10260.
- Rico Sennrich and Biao Zhang. 2019. Revisiting low-resource neural machine translation: A case study. In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019*, pages 211–221.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016*, pages 86–96.
- Rico Sennrich, Orhan Firat, Kyunghyun Cho, Alexandra Birch, Barry Haddow, Julian Hitschler, Marcin Junczys-Dowmunt, Samuel Läubli, Antonio Valerio Miceli Barone, Jozef Mokry, and Maria Nädejde. 2017. Nematus: a toolkit for neural machine translation. In *Proceedings of the Software Demonstrations of the 15th Conference of the European Chapter of the Association for Computational Linguistics*, pages 65–68.
- Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 2014. Sequence to sequence learning with neural networks. In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems*, pages 3104–3112.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems*, pages 5998–6008.
- Qiang Wang, Bei Li, Tong Xiao, Jingbo Zhu, Changliang Li, Derek F. Wong, and Lidia S. Chao. 2019. Learning deep transformer models for machine translation. In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019*, pages 1810–1822.