

---

# Appendix for “R-Drop: Regularized Dropout for Neural Networks”

---

Anonymous Author(s)

Affiliation

Address

email

## 1 A Detailed Experimental Settings

2 We provide more detailed settings for the experiments of each task in this part.

### 3 A.1 Neural Machine Translation

4 For all the NMT tasks, we use the public datasets from IWSLT competitions<sup>1</sup> and WMT competitions<sup>2</sup>.  
5 We tokenize all the datasets with byte-pair-encoding (BPE) [10] approach with the dictionary built  
6 jointly upon the source and target sentence pairs except the IWSLT17 En↔Zh translation dataset  
7 that is built separately. After tokenization, the resulted vocabularies for IWSLT datasets are near  $10k$ ,  
8 while for WMT datasets, the vocabulary size is about  $32k$ .

9 To train the Transformer based NMT models, we use `transformer_iwslt_de_en` configuration  
10 for IWSLT translations, which has 6 layers in both encoder and decoder, embedding size 512,  
11 feed-forward size 1,024, attention heads 4, dropout value 0.3, weight decay 0.0001. For the WMT  
12 experiments, the `transformer_vaswani_wmt_en_de_big` setting has 6 layers in encoder and  
13 decoder, embedding size 1,024, feed-forward size 4,096, attention heads 16, dropout value 0.1,  
14 attention dropout 0.1 and relu dropout 0.1. The training is optimized with Adam [5] with  $\beta_1 =$   
15  $0.9$ ,  $\beta_2 = 0.98$ ,  $\epsilon = 10^{-9}$ . The learning rate scheduler is `inverse_sqrt` with default learning  
16 rate 0.0005 and warmup steps 4,000. Label smoothing [11] is adopted with value 0.1. Our code  
17 implementation is based on open-source Fairseq<sup>3</sup>. We train the IWSLT translations on 1 GEFORCE  
18 RTX 3090 card and the WMT translations on 8 GEFORCE RTX 3090 cards.

19 To evaluate the performance, we use `multi-bleu.perl`<sup>4</sup> to evaluate IWSLT14 En↔De and all  
20 WMT tasks for a fair comparison with previous works [15, 8]. For other NMT tasks, we use  
21 `sacre-bleu`<sup>5</sup> [9] for evaluation. When inference, we follow [12] to use beam size 4 and length  
22 penalty 0.6 for WMT14 En→De, beam size 5 and penalty 1.0 for other tasks.

### 23 A.2 Abstractive Summarization

24 For summarization, we take the pre-trained BART [6] model as backbone and fine-tune on the  
25 CNN/DailyMail dataset<sup>6</sup>. BART is a pre-trained sequence-to-sequence model based on the masked  
26 source input and autoregressive target output, which contains 12 layers of Transformer encoder and  
27 12 layers of Transformer decoder, the embedding size is 1,024, and the feed-forward size is 4,096.

---

<sup>1</sup><https://iwslt.org/>

<sup>2</sup><https://www.statmt.org/wmt14/translation-task.html>

<sup>3</sup><https://github.com/pytorch/fairseq/tree/master/examples/translation>

<sup>4</sup><https://github.com/moses-smt/mosesdecoder/blob/master/scripts/generic/multi-bleu.perl>

<sup>5</sup><https://github.com/mjpost/sacrebleu>

<sup>6</sup><https://github.com/abisee/cnn-dailymail>

Hyper-parameter	CoLA	MRPC	RTE	SST-2	MNLI	QNLI	QQP	STS-B
Learning Rate	1e-5	1e-5	1e-5	1e-5	1e-5	1e-5	1e-5	1e-5
Max Update	5336	2296	3120	20935	123873	33112	113272	3598
Max Sentence (Batch)	16	16	8	32	32	32	32	16
Dropout	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1
Coefficient $\alpha$	0.5	1.0	1.0	1.0	0.5	1.0	0.5	1.0

Table 1: Hyper-parameters when fine-tuning our models on GLUE benchmark.

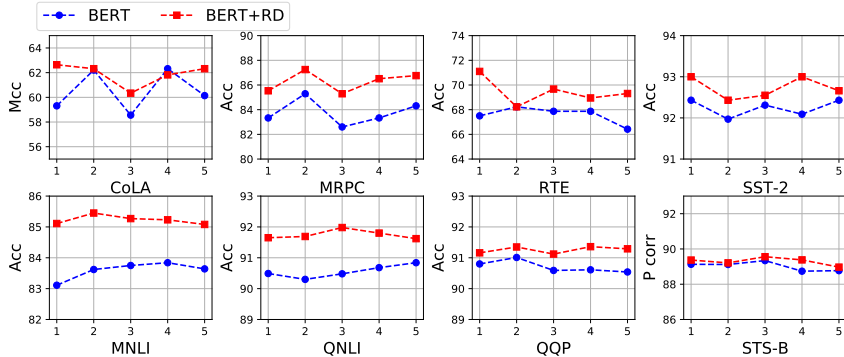


Figure 1: Results on 8 GLUE tasks with different random seed.

28 Dropout value is 0.1. During fine-tuning, we follow the hyper-parameters used in [6]. The pre-trained  
 29 model and the backbone implementations are all from Fairseq<sup>7</sup>. The training is conducted on 8  
 30 GEFORCE RTX 3090 GPU cards.

### 31 A.3 Language Modeling

32 For language modeling, we train on the Transformer decoder [12] and Adaptive Input Transformer [1]  
 33 models. The configuration for Transformer is `transformer_lm_gpt`, which contains 12 layers with  
 34 embedding size 768 and feed-forward size 3,072, attention heads 12. Dropout and attention dropout  
 35 are 0.1. For Adaptive Input Transformer, the configuration is `transformer_lm_wiki103` with 16  
 36 layers, embedding size 1,024, feed-forward size 4,096, attention heads 16, dropout 0.3, attention  
 37 dropout 0.1, gelu dropout 0.1 and adaptive softmax dropout 0.2. We train Transformer model for  
 38 50k steps and Adaptive Input Transformer for 286k steps. The development is based on the code  
 39 base Fairseq<sup>8</sup>. The training is on 8 Tesla V100 GPU cards.

### 40 A.4 Language Understanding

41 For language understanding tasks, we follow the popular pre-training and fine-tuning methodology,  
 42 and the fine-tuned sets are the GLUE [13] benchmark. We follow previous works [2, 7] to work on  
 43 the 8 tasks, including single-sentence classification tasks (CoLA, SST-2), sentence-pair classification  
 44 tasks (MNLI, QNLI, RTE, QQP, MRPC), and sentence-pair regression task (STS-B). The detailed  
 45 data statistics can be found from the original paper [13].

46 The pre-trained BERT-base model is the Transformer [12] encoder network, which contains 12 layers  
 47 with embedding size 768, feed-forward size 3,072 and attention heads 12. Correspondingly, the  
 48 Roberta-large model contains 24 layers with embedding size 1,024, feed-forward size 4,096 and  
 49 attention heads 16. During fine-tuning, we use Adam [5] as our optimizer with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.98$ ,  
 50  $\epsilon = 10^{-6}$ , and  $L_2$  weight decay of 0.01. We select the learning rate in range  $\{5 \times 10^{-6}, 10^{-5}\}$  and  
 51 batch size in  $\{8, 16, 32\}$ . Other hyper-parameter settings are mostly same as previous works [7]. The

<sup>7</sup><https://github.com/pytorch/fairseq/tree/master/examples/bart>

<sup>8</sup>[https://github.com/pytorch/fairseq/tree/master/examples/language\\_model](https://github.com/pytorch/fairseq/tree/master/examples/language_model)

52 pre-trained model and the backbone implementations are all from Huggingface Transformers<sup>9</sup>. We  
 53 report the specific settings of several important hyper-parameters in Table 1, including the dropout  
 54 value. The fine-tuning experiments are conducted on 1 GEFORCE RTX 3090 GPU card.

55 Further, to give a clear comparison of our R-Drop based fine-tuning and vanilla fine-tuning, we  
 56 plot the performance changes from different random seeds over the pre-trained BERT model on  
 57 each GLUE task. The curves are shown in Figure 1. We can see that consistent improvements are  
 58 achieved on different random seeds, which means our R-Drop can robustly help improve the mode  
 59 generalization and model performance.

60 **MSE Regularization** Our R-Drop is presented under the KL-divergence between two distributions.  
 61 To extend our method into the regression task, such as STS-B in GLUE, we introduce the MSE-based  
 62 regularization. For input data  $(x, y)$ , we forward the  $x$  two times similarly as in classification and  
 63 obtain the two predicted values  $y'_1$  and  $y'_2$ . Then we regularize these two predicted values with MSE  
 64 as follow:

$$\mathcal{L}_{mse_r} = \|y'_1 - y'_2\|_2, \quad (1)$$

65 and we add  $\mathcal{L}_{mse_r}$  with conventional MSE loss:  $\mathcal{L}_{mse} = \|y - y'_1\|_2 + \|y - y'_2\|_2$ . The final  
 66 optimization objective is:

$$\mathcal{L} = \mathcal{L}_{mse} + \alpha \mathcal{L}_{mse_r}. \quad (2)$$

## 67 A.5 Image Classification

68 The image classification task is evaluated with the recent popular Vision Transformer (ViT) [4] model,  
 69 which is the same as Transformer but with the image patch data as input. We take the two publicly  
 70 released models<sup>10</sup>, ViT-B/16 and ViT-L/16, which are pre-trained on ImageNet-21k [3] dataset with  
 71 21k classes and 14M images in total. ViT-B/16 is a Transformer model with 12 Transformer encoder  
 72 layers, embedding size 768, feed-forward size 3,072 and attention heads 12, while ViT-L/16 with  
 73 24 layers, 1,024 embedding size, 4,096 feed-forward size and 16 attention heads. We only conduct  
 74 the fine-tuning stage experiments on CIFAR-100 and ImageNet. Note that the ImageNet results  
 75 are computed without additional techniques (Polyak averaging and 512 resolution images) used to  
 76 achieve results in [4]. During fine-tuning, the dropout values are 0.1 for both models. Fine-tuning is  
 77 on 8 GEFORCE RTX 3090 GPU cards.

## 78 B Theoretical Discussion of R-Drop

79 For self-contained, we reformulate the constrained optimization problem as:

$$\min_w \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\xi} [-\log \mathcal{P}_{\xi}^w(y_i|x_i)], \quad (3)$$

$$s.t. \quad \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\xi^{(1)}, \xi^{(2)}} [\mathcal{D}_{KL}(\mathcal{P}_{\xi^{(1)}}^w(y_i|x_i) || \mathcal{P}_{\xi^{(2)}}^w(y_i|x_i))] = 0. \quad (4)$$

80 For a multi-layer perceptron (MLP) with dropout, its output is represented as:

$$\mathcal{P}_{\xi}(y|x) = \text{softmax}(vh_{\xi^L}^L); \quad h_{\xi^l}^l = \xi^l \odot \sigma(w^l h_{\xi^{l-1}}^{l-1}), l = 1, \dots, L; \quad h_{\xi^0}^0 = \xi_0 \odot x, \quad (5)$$

81 where  $\sigma(\cdot)$  is the activation function.

82 We also consider the MLP with skip/residual connections (which is more popular than plain MLP),  
 83 whose output after dropout is represented as:

$$\mathcal{P}_{\xi}^{res}(y|x) = \text{softmax}(vh_{\xi^L}^L); \quad h_{\xi^l}^l = \xi^l \odot \sigma(w^l h_{\xi^{l-1}}^{l-1}) + h_{\xi^{l-1}}^{l-1}, l = 1, \dots, L; \quad h_{\xi^0}^0 = \xi_0 \odot x. \quad (6)$$

<sup>9</sup><https://github.com/huggingface/transformers>

<sup>10</sup><https://github.com/jeonsworld/ViT-pytorch>

84 For a Transformer model, the output can be represented as:

$$z = SA(Q, K, V) = \text{softmax} \left( \frac{(Qw_1)(Kw_2)^T}{\sqrt{d}} \right) Vw_3, \quad (7)$$

$$f(z) = v_2 \max(0, v_1 z), \quad (8)$$

85 where  $v_1, v_2, w_1, w_2, w_3$  are learnable parameters,  $SA(\cdot)$  stands for self-attention operation, and  
 86  $f(\cdot)$  is the position-wise feed-forward network as described in Transformer [12]. We only consider  
 87 dropout inside the self-attention layer, i.e.,

$$z_\xi = SA(Q, K, V) = \left( \xi_2 \odot \text{softmax} \left( \frac{((\xi_1 \odot Q)w_1)((\xi_1 \odot K)w_2)^T}{\sqrt{d}} \right) \right) ((\xi_1 \odot V)w_3), \quad (9)$$

$$f(z_\xi) = v_2 \max(0, (v_1 z_{\xi_1})), \quad (10)$$

$$\mathcal{P}_\xi^{TF}(y|x) = \text{softmax}(f(z_\xi)), \quad (11)$$

88 where the size of random matrix  $\xi_1$  equal to the size of  $Q^{11}$ , and the size of  $\xi_2$  equal to the size of the  
 89 value after softmax.

90 Now we provide the proof for the next proposition. Here, we consider a sampling strategy called  
 91 dropout<sub>2</sub>, which randomly samples  $p$  portion of hidden neurons at each layer. Thus, this strategy can  
 92 keep the width for the sampled sub-structure and avoid extremely thin sub-structures.

93 **Proposition B.1.** For a fully-connected neural network  $\mathcal{P}^w(y|x)$ , e.g., multilayer perceptron and  
 94 Transformer, the constraint in Equation (4) is equivalent to constraining all the parameters of the  
 95 network at the same layer to be equal.

96 *Proof:* We consider two sub-structures  $\mathcal{P}_1$  and  $\mathcal{P}_2$  that differ from each other at only one hidden node.  
 97 Without loss of generality, we suppose the hidden node that is only contained in  $\mathcal{P}_1$  as  $h_i^l$  and the  
 98 hidden node that is only contained in  $\mathcal{P}_2$  as  $h_j^l$ .

99 Since we constrain the  $KL(\mathcal{P}_1||\mathcal{P}_2) = 0, \forall x$ , which is equivalent to  $\mathcal{P}_1(y|x) = \mathcal{P}_2(y|x), \forall x$ . Then  
 100 we have  $h_i^l = h_j^l, \forall x$ . If not, we can construct an  $x \in \mathbb{R}^d$  to make their final output unequal<sup>12</sup>.

101 Because the two sub-structures share the other hidden nodes, the inputs of  $h_i^l$  and  $h_j^l$  are the same,  
 102 i.e.,  $h_i^l = \sigma(\tilde{w}_i^l \tilde{h}^{l-1}), h_j^l = \sigma(\tilde{w}_j^l \tilde{h}^{l-1})$ , where  $w_i$  denotes the ingoing weights of  $h_i^l$  and we use the  
 103 subscript  $\tilde{\cdot}$  to denote the hidden output and weight matrix after the dropout. Then, we have  $\tilde{w}_i^l = \tilde{w}_j^l$ .  
 104 Traversing all the paired sub-structures (produced by dropout<sub>2</sub>) that differ at only one hidden node,  
 105 we can get that all the parameters at the same layer are equal.

106 Similarly, for the Transformer model, we can also get the  
 107 results in the above Proposition, but it is established by re-  
 108 moving the redundancy in weights. In detail, the weights  $w_1$   
 109 and  $w_2$  have redundant freedom due to the scaling invariant  
 110 property of  $Qw_1(Kw_2)^T$  [14], i.e., If we regard  $w_1$  and  $w_2$   
 111 as two linear layers of MLP, the ingoing weights of one hid-  
 112 den node are multiplying a scalar  $c$  and its outgoing weights  
 113 are dividing by  $c$ , the resulted output will not change. To  
 114 remove the redundancy, we fix the first column in  $w_1$  (one  
 115 ingoing weight for each hidden node) to have the same value.  
 116 Then Proposition B.1 is established for  $w_1, w_2, w_3$ .

117 **Remark:** Although Proposition B.1 is established for a  
 118 different sampling strategy with original dropout, it will  
 119 approach to the result for original dropout **with high proba-  
 120 bility** because the mean portion of the sampled neurons will  
 121 be closed to  $p$  according to central limit theorem.

122 Proposition B.1 can be regarded as a result of original dropout in the sense of high probability. Next,  
 123 we briefly discuss what will happen if we consider the original sampling strategy in dropout. Since

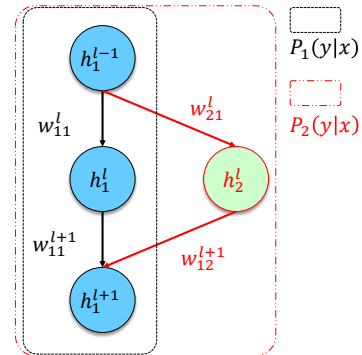


Figure 2: Example of MLP in 1).

<sup>11</sup>Please note that for self-attention,  $Q = K = V$ .

<sup>12</sup>Here, we suppose the activation function  $\sigma$  satisfies that  $w^2 \sigma(w^1 x) = \tilde{w}^2 \sigma(\tilde{w}^1 x), \forall x$  if and only if  $w^2 = \tilde{w}^2, w_1 = \tilde{w}^1$ .



Figure 3: Results of R-Drop and Transformer with double batch.

124 the original strategy contains more sub-structures, it actually brings a much stronger restriction to the  
 125 model. For example, the model has the chance to degenerate to a zero mapping.

126 1) For MLP: If we consider the original sampling strategy in dropout, the constraint of Equation (4) is  
 127 equivalent to constraining one layer of the MLP to be zero. We can consider an extreme sub-structure  
 128 that there is only one hidden node after dropout and denote its output as:

$$\mathcal{P}_1(y|x) = \text{softmax}(vh_1^L); \quad h_1^l = \sigma(w_{11}^l h_1^{l-1}), l = 1, \dots, L; \quad h_1^0 = x_1,$$

129 where  $h_1^l$  denotes the 1-th hidden node at layer  $l$  and  $w_{11}^l$  denotes the weights that connects  $h_1^{l-1}$   
 130 and  $h_1^l$ . Then we consider another sub-structure  $\mathcal{P}_2(y|x)$  which has one more hidden node at layer  $l$   
 131 (denoted as  $h_2^l$ ) than  $\mathcal{P}_1(y|x)$ , i.e.,

$$h_2^l = \sigma(w_{21}^l h_1^{l-1}); h_1^{l+1} = \sigma(w_{11}^{l+1} h_1^l + w_{12}^{l+1} h_2^l),$$

132 and other computations are the same as the calculation of  $\mathcal{P}_1(y|x)$  (see the example in Figure 2).  
 133 Therefore, the sub-structure  $\mathcal{P}_2$  has two additional parameters  $w_{12}^{l+1}$  and  $w_{21}^l$  compared with the sub-  
 134 structure  $\mathcal{P}_1$ , and they share weights on  $w_{11}^l, l = 1, \dots, L$ . Because we constrain the  $KL(\mathcal{P}_1||\mathcal{P}_2) =$   
 135  $0, \forall x$ , which is equivalent to  $\mathcal{P}_1(y|x) = \mathcal{P}_2(y|x), \forall x$ . Therefore, we need  $\sigma(w_{11}^{l+1} h_1^l + w_{12}^{l+1} h_2^l) =$   
 136  $\sigma(w_{11}^{l+1} h_1^l), \forall x$ . Then we must have  $w_{12}^{l+1} = 0$  or  $w_{21}^l = 0$ . Similarly, we can traverse all such  
 137 extreme cases to get the results.

138 2) For MLP with skip/residual connections: Using the similar technique in 1), if we consider the  
 139 MLP with skip connection  $\mathcal{P}_\xi^{res}(y|x)$  and the original sampling strategy in dropout, the constraint in  
 140 Equation (4) is equivalent to constraining all the parameters equals zero. The network will degenerate  
 141 to an identity mapping of  $x$ .

## 142 C More Studies

### 143 C.1 Batch Size Doubled Training

144 As discussed in Section 2.2, we implement the algorithm by repeating input data  $x$  once and  
 145 concatenate the  $x$  with repeated one in the same mini-batch to forward once. This is similar to  
 146 enlarge the batch size to be double size at each step. The difference is that half of the data are the  
 147 same as the other half, while directly doubling the batch size, the data in the same mini-batch are  
 148 all different. Therefore, we are still interested in the result of directly doubling the batch size to  
 149 see the performance. We conduct experiments on IWSLT14 De→En translation with Transformer,  
 150 and the batch size is enlarged from 4,096 to be 8,192. The result is 34.93 BLEU score. We can  
 151 see that though slight improvement is achieved (compared to baseline 34.64), it falls far behind our  
 152 strong performance 37.25. For the detailed training cost for each step, we represent the number here:

153 Transformer + Double Batch costs near 9ms per step, while Transformer + DR costs about 10ms  
154 per step. The additional cost is from the KL-divergence loss backward computation. We can see  
155 the cost is 1.1 times, which is a negligible cost. We also plot the valid BLEU curves along with the  
156 training for this study. The curves are shown in Figure 3. Compare to this batch size doubled training  
157 and our R-Drop, we can clearly see the advantage of R-Drop. With similar training cost, R-Drop  
158 gradually improves the performance to a much stronger one. In the figures, we also plot the curve for  
159 Transformer with the original batch size training (e.g., 4, 096) for a better comparison.

## 160 C.2 Importance of KL-divergence

161 Our method introduces a KL-divergence between the two distributions from the same sample. In this  
162 study, we specifically investigate on the importance of the KL-divergence loss. Thus, this ablation  
163 removes the  $\mathcal{L}_{KL}$  loss between the two distributions and only keeps the  $\mathcal{L}_{NLL}$  loss for training.  
164 Similar as other studies, we work on IWSLT14 De→En translation and the model is Transformer.  
165 The result is also 34.93 BLEU score (same as above enlarge batch size), it is slightly better than  
166 Transformer baseline (34.64), but far worse from our R-Drop based training result 37.25 BLEU score.  
167 This result well demonstrates the importance and effectiveness of our introduced KL-divergence loss.

## 168 References

- 169 [1] Alexei Baevski and Michael Auli. Adaptive input representations for neural language modeling.  
170 In *International Conference on Learning Representations*, 2018.
- 171 [2] Kevin Clark, Minh-Thang Luong, Quoc V Le, and Christopher D Manning. Electra: Pre-training  
172 text encoders as discriminators rather than generators. In *International Conference on Learning  
173 Representations*, 2019.
- 174 [3] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-  
175 scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern  
176 recognition*, pp. 248–255. Ieee, 2009.
- 177 [4] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,  
178 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,  
179 Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image  
180 recognition at scale. In *International Conference on Learning Representations*, 2021.
- 181 [5] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint  
182 arXiv:1412.6980*, 2014.
- 183 [6] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer  
184 Levy, Veselin Stoyanov, and Luke Zettlemoyer. Bart: Denoising sequence-to-sequence pre-  
185 training for natural language generation, translation, and comprehension. In *Proceedings of the  
186 58th Annual Meeting of the Association for Computational Linguistics*, pp. 7871–7880, 2020.
- 187 [7] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike  
188 Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining  
189 approach. *arXiv preprint arXiv:1907.11692*, 2019.
- 190 [8] Myle Ott, Sergey Edunov, David Grangier, and Michael Auli. Scaling neural machine translation.  
191 In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pp. 1–9,  
192 2018.
- 193 [9] Matt Post. A call for clarity in reporting bleu scores. In *Proceedings of the Third Conference on  
194 Machine Translation: Research Papers*. Association for Computational Linguistics, 2018.
- 195 [10] Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare  
196 words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for  
197 Computational Linguistics (Volume 1: Long Papers)*, pp. 1715–1725, 2016.
- 198 [11] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Re-  
199 thinking the inception architecture for computer vision. In *Proceedings of the IEEE conference  
200 on computer vision and pattern recognition*, pp. 2818–2826, 2016.

- 201 [12] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
202 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Proceedings of the 31st*  
203 *International Conference on Neural Information Processing Systems*, pp. 6000–6010, 2017.
- 204 [13] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman.  
205 Glue: A multi-task benchmark and analysis platform for natural language understanding. In  
206 *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural*  
207 *Networks for NLP*, pp. 353–355, 2018.
- 208 [14] Yue Wang, Yuting Liu, and Zhi-Ming Ma. The scale-invariant space for attention layer in neural  
209 network. *Neurocomputing*, 392:1–10, 2020.
- 210 [15] Jinhua Zhu, Yingce Xia, Lijun Wu, Di He, Tao Qin, Wengang Zhou, Houqiang Li, and Tieyan  
211 Liu. Incorporating bert into neural machine translation. In *International Conference on Learning*  
212 *Representations*, 2019.