

	<b>Models</b>	En-Vi	De-En
	NPMT+LM (Huang et al., 2017)	28.1	30.1
	Risk (Edunov et al., 2018)	-	32.8
	Var-Attn (Deng et al., 2018)	-	33.7
	Transformer	30.1	34.2
	+KD (Tan et al., 2018)	28.7	34.0
	+Fixup (Zhang et al., 2019)	-	34.5
	+AdaNorm	<b>30.7</b>	<b>35.0</b>

Table 1: Results on the IWSLT15 English-to-Vietnamese translation test set and IWSLT14 German-to-English test set.

1 **To all reviewers:**  
2 We thank the reviewers for their detailed comments. The  
3 followings are our responses.

4 **To reviewer 1:**  
5 [1] Are the empirical improvements strong?  
6 We have strong confidence in its empirical improvements.  
7 Take the results on translation tasks as an example (shown  
8 in Table 1). AdaNorm has brought improvements of 0.6  
9 BLEU on En-Vi and 0.8 BLEU on De-En, much higher  
10 than other techniques do like Fixup and KD.

11 [2] More intuitions for which task AdaNorm can really  
12 improve, and why.  
13 AdaNorm works better for tasks requiring complex model  
14 structures. The reason is that deeper models usually have  
15 the tendency to over-fit training data, and AdaNorm alleviates the over-fitting by adaptively controlling scaling  
16 weights towards different inputs on affine transformation. Comparing to LayerNorm that ignores the input  
17 distribution when testing, our proposed AdaNorm has achieved better empirical improvements.

18 **To reviewers 2:**  
19 Thanks for your comments and suggestions.

20 **To reviewer 3:**  
21 [1] Equation for variance in (1) seems wrong.  
22 The equation for variance in (1) is correct. It is a variant of the traditional variance equation. The followings are  
23 the derivation process. If  $\sigma^2$  is the variance of  $X$ , then

$$\sigma^2 = E[(X - E[X])^2] = E[X^2 - 2X E[X] + E[X]^2] = E[X^2] - 2E[X]E[X] + E[X]^2 = E[X^2] - E[X]^2 \quad (1)$$

24 where  $E$  is a mean function. In this paper,  $X = x_1, x_2, \dots, x_H$  and the variance can be written as  $\sigma^2 = \frac{1}{n}(\sum_{i=1}^n x_i^2 - n\mu^2)$ .  
25 [2] In DetachNorm, the gradient is simply wrong (due to parts of the gradient being detached and essentially  
26 random noise is added into the model through the special copy function).  
27 Here we illustrate its correctness by analyzing the two mentioned operations. First the detaching operation  
28 simulates the situation of constant variance and mean that have zero gradient to the input. Comparing to  
29 LayerNorm, they are two settings to evaluate the effect of variance and mean on gradients. The gradients in  
30 these two settings are different, but they are both right. Second, the special copy function is a simple assignment  
31 operation. It has extremely weak effect on model performance considering the huge amount of assignment  
32 operations in neural networks.  
33 [3] The proposed AdaNorm does not really directly address the items discussed in the first part of the paper.  
34 As described in lines 193-197, AdaNorm is proposed to address the over-fitting problem discussed in the  
35 first part. The first part analyzes which parts in LayerNorm work and which parts do not. Empirical results  
36 show that "bias and gain", parameters of LayerNorm, are not always beneficial because they increase the risk  
37 of over-fitting. Motivated by this fact, we propose a new normalization approach, AdaNorm, to address the  
38 over-fitting problem. Experiment results demonstrate that AdaNorm outperforms LayerNorm on seven datasets  
39 with better convergence.  
40 [4] In Theorem 2 and above, should the absolute value be only around  $z_i$  and not the entire sum. What if  $z_i$  are  
41 large but they cancel each other out?  
42 Thanks for your suggestions. We will consider replacing  $|\sum_{i=1}^H z_i|/H < M$  with  $\sum_{i=1}^H |z_i|/H < M$ . For the  
43 proof of the theorem, we only need  $|\sum_{i=1}^H z_i|/H < M$ . Since  $\sum_{i=1}^H |z_i|/H < M$  is a stronger constraint, it  
44 does not affect the proof.  
45 [5] "To prevent ... dismissing the feature of gradient", what does this even mean?  
46 It means that LayerNorm has an advantage of re-centering and re-scaling gradients. The proposed AdaNorm still  
47 keeps this advantage when avoiding the over-fitting problem.