

# Multi-head or Single-head? An Empirical Comparison for Transformer Training

Anonymous ACL submission

## Abstract

Multi-head attention plays a crucial role in the recent success of Transformer, which leads to consistent performance improvements over conventional attention in various applications. The popular belief is that its effectiveness stems from attending to information from multiple representation subspaces. In this paper, we first demonstrate that using multiple subspaces is not a unique feature of *multi-head* attention, as *multi-layer single-head* attention also leverages multiple subspaces. Then, we suggest the main advantage of the multi-head attention is the training stability, since it has fewer layers than the single-head attention when using the same number of subspaces. For example, *24-layer 16-head* Transformer (BERT-large) and *384-layer single-head* Transformer have roughly the same model size and employ the same total subspace number (attention head number), while the multi-head one is significantly shallower. Meanwhile, we show that, with recent advances in deep learning, we can successfully stabilize the training of the deep single-head Transformer. As the training difficulty is no longer a bottleneck, substantially deeper single-head Transformers achieve consistent performance improvements<sup>1</sup>.

## 1 Introduction

Transformers (Vaswani et al., 2017) have led to a series of breakthroughs in various deep learning tasks (Devlin et al., 2019; Velickovic et al., 2018b). One crucial component of Transformer is the multi-head attention, which has been observed to be one major reason behind the success of the Transformer. For example, on machine translation benchmarks, Recurrent Neural Networks (RNNs) can outperform Transformers when both are using the multi-head encoder-decoder attention and would underperform without using the multi-head

<sup>1</sup>Our model implementations and data preparation scripts will be made publicly available.

attention (Chen et al., 2018). Besides Transformer, multi-head attention has also been incorporated into other models (Chen et al., 2018; Velickovic et al., 2018a; Fang et al., 2019). More discussions on related work is available at Appendix A.

Multi-head attention projects the inputs into multiple different subspaces and attend to information from them, while one conventional attention can only attend to information from one subspace,

**Our Contributions.** Our point of start is demonstrating that leveraging multiple subspaces is not a unique feature of multi-head attention. In fact, stacking multiple conventional attention modules also leverage multiple subspaces.

As in Figure 1, a multi-head attention module can be viewed as an ensemble model, which combines multiple single-head attention modules by calculating their average (more elaborations are included in Appendix C). Thus, by integrating these modules differently, we can reconstruct a Transformer to be single-head<sup>2</sup> and substantially deeper, without changing the number of subspaces or the inference computation complexity.

In our experiments, compared to the shallower multi-head Transformer, the deeper single-head Transformer performs better but is harder to train. It matches the common wisdom that model depth can increase model capacity at the cost of training difficulty. We also observe that, benefited from the recent advance of deep learning (Liu et al., 2020b), the training difficulty is no longer an obstacle.

## 2 Experiment Overview

Here, we discuss the experiment setup (more in Appendix D). Then, we compare the shallow multi-head Transformer and deep single-head Transformer from three aspects, i.e., stability (Sec. 3), performance (Sec. 4), and efficiency (Sec. 5).

<sup>2</sup>We use single-head/multi-head Transformer to refer Transformer with single-head/multi-head attention.

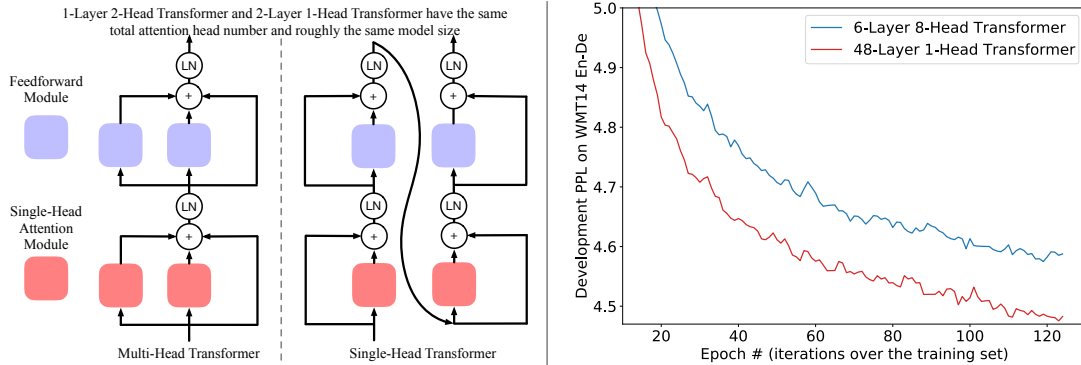


Figure 1: Left: both multi-head and single-head Transformer can attend to information from multiple subspaces. Right: comparing to the shallow multi-head Transformer, the deep single-head Transformer can achieve a lower PPL score, while its training is more challenging (without Admin, the 48-layer 1-head Transformer training failed).

**Tasks.** We conduct experiments on language model pre-training and translation. For translation, we use the WMT’14 English-German (En-De) as the benchmark. For language model pre-training, we evaluate models on SQuAD 2.0 (Rajpurkar et al., 2018) and GLUE (Wang et al., 2018).

**Model Specificity.** For machine translation, the original Transformer-base model is 8H-6L-6L encoder-decoder (Vaswani et al., 2017)<sup>3</sup>. Here, we compare it with 1H-48L-48L encoder-decoder, and both models have 512-dimension word embedding, 64-dimension per-head attention output, and  $256 \cdot \gamma$ -dimension feedforward network ( $\gamma$  is the number of heads). For language model pre-training, we compare BERT-base model is 12H-12L encoder and BERT-large model is 16H-24L encoder. Here, we compare them with deep single-head BERT-base model (1H-144L) and deep single-head BERT-large model (1H-384L). All language models have 768-dimension word embedding, 64-dimension per-head attention output, and  $256 \cdot \gamma$ -dimension word embedding ( $\gamma$  is the number of heads). Moreover, we employed the Admin initialization (Liu et al., 2020b) to stabilize 1H-48L-48L Transformer-base and 1H-384L BERT-large.

### 3 Stability Comparison

As in Table 1, after changing the shallow multi-head Transformer to the deep single-head Transformer, the training fails to converge well for 2 out of 3 models. Note that, although the 1H-144L BERT-base model converges successfully, the model is sensitive to the choice of initialization. Specifically, the BERT-base model and BERT-

<sup>3</sup>“ $\gamma$ H- $\alpha$ L(- $\beta$ L)” indicates a Transformer model has  $\gamma$ -head  $\alpha$ -layer encoder (and  $\gamma$ -head  $\beta$ -layer decoder).

large model are initialized with truncated normal distribution with 0.02 variance, instead of following the common practice (e.g., using the Kaiming initialization (He et al., 2015) or the Xavier initialization (Glorot and Bengio, 2010)). We observe that after changing the variance of the initialization, or following the common practice, the training of the 1H-144L BERT-base model would also fail.

Meanwhile, we show that, with recent deep learning advances, we can successfully stabilize Transformer training. After employing the Admin initialization (Liu et al., 2020b), all deep single-head Transformer models are trained successfully, without changing any hyper-parameters. This shows that, although the deep single-head Transformer is harder to train, the training difficulty is no longer an obstacle.

### 4 Performance Comparison

For the machine translation task, we summarize the results in Table 2. The deep single-head Transformer (1H-48L-48L) achieves a 0.5 BLEU improvements over the shallow multi-head Transformer. Also, the deep single-head Transformer achieves the same performance with the architecture search algorithm (Evolved Transformer (So et al., 2019) and DARTSformer (Zhao et al., 2021)), with slightly less parameters. Specifically, Evolved Transformer and DARTSformer conducts neural architecture search on Transformer, and treat the multi-head attention as the basic module (i.e., the deep single-head Transformer is not in their search space). Deep single-head Transformer achieves comparable performance without hyper-parameter tuning, which further verifies its effectiveness.

For language model pre-training, the deep single-head Transformer also achieves consistent perfor-

	Transformer-base		BERT-base		BERT-large	
	8H-6L-6L	1H-48L-48L	12H-12L	1H-144L	16H-24L	1H-384L
Training	✓	×/✓(w. Admin)	✓	✓	✓	×/✓(w. Admin)

Table 1: Deep single-head Transformers are harder to train than shallow multi-head Transformers.

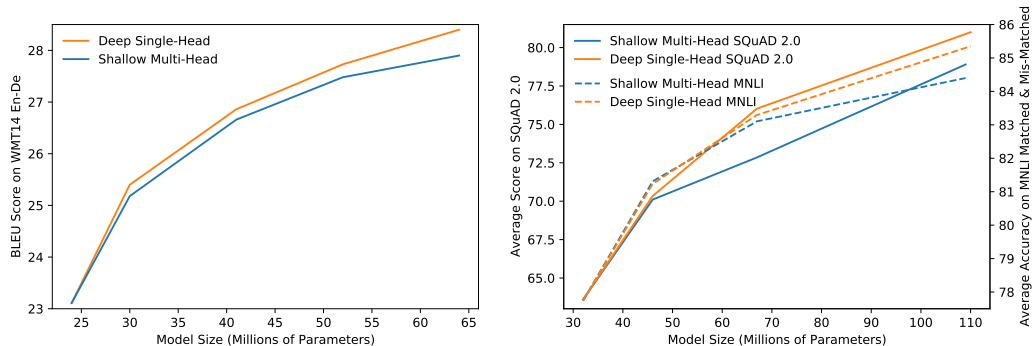


Figure 2: Performance with Different Model Size. Left: the performance of  $\alpha$ H-6L-6L ( $\alpha=1, 2, 4, 6, 8$ ) and 1H- $\beta$ L- $\beta$ L ( $\beta=6, 12, 24, 36, 48$ ), whose per-head dimension is the same with Transformer-base. Right: the performance of  $\alpha$ H-12L ( $\alpha=1, 3, 6, 12$ ) and 1H- $\beta$ L ( $\beta=12, 36, 72, 144$ ), whose per-head dimension is the same with BERT-base.

Model	BLEU	Param.
8H-6L-6L	27.90	63.2M
1H-48L-48L	28.40	63.6M
2D-CSANs <sup>†</sup> (Yang et al., 2019)	28.18	88.0M
Evolved*(So et al., 2019)	28.4	64.1M
DARTSformer <sup>†</sup> *(Zhao et al., 2021)	28.4	65.2M

Table 2: Performance on the WMT’14 En-De dataset. \* indicates neural architecture search methods. <sup>†</sup> indicates the results may not be directly comparable to others, due to the difference on pre-processing and evaluation.

146 mance improvements over the original shallow  
 147 multi-head Transformer (as in Table 3). Table 4  
 148 shows the test performance on the GLUE bench-  
 149 mark. The deep single-head Transformer outper-  
 150 forms the shallow multi-head Transformer on 7 out  
 151 of 9 tasks, and improves the average score (GLUE)  
 152 by roughly 1 point. In the mean time, it is worth  
 153 mentioning that, on 2 out of 3 sentence similar-  
 154 ity/paraphrase tasks, the shallow multi-head Trans-  
 155 former achieves better performance. This indicates  
 156 the deep single-head Transformer can be further  
 157 improved, and we will further explore this in the  
 158 future work. These observations verified that the  
 159 deep single-head Transformer could perform better  
 160 than the shallow multi-head Transformer.

161 **Impact of Model Initialization.** Here, we aim to  
 162 understand the impact of model initialization on  
 163 model performance. As the 1H-144L BERT-base  
 164 model converges well with both the vanilla initial-

165 ization and the Admin initialization, we not only  
 166 conduct training with the Admin initialization, but  
 167 also the vanilla initialization. As summarized in  
 168 Table 3, the default initialization and the Admin  
 169 initialization achieve similar performance. This  
 170 observation supports our intuition that the major  
 171 benefit of the Admin initialization is on training sta-  
 172 bility, and the performance improvements mostly  
 173 come from the change from shallow multi-head  
 174 Transformer to deep single-head Transformer.

175 **Impact of Head Number.** Intuitively, the differ-  
 176 ence between deep single-head Transformers and  
 177 shallow multi-head Transformers is proportional to  
 178 the model size/head number (e.g., the difference  
 179 between 2H-6L and 1H-12L should be smaller than  
 180 the difference between 4H-6L and 1H-24L). We  
 181 conduct experiments on Transformers with differ-  
 182 ent head numbers, and visualize the results in Fig-  
 183 ure 2. It shows that when the architecture differ-  
 184 ence is larger (i.e., with more number of heads),  
 185 the performance improvement is also larger.

## 186 5 Efficiency Comparison

187 **Inference Speed.** The shallow multi-head Trans-  
 188 former and the deep single-head Transformer have  
 189 roughly the same model size and computation com-  
 190 plexity. Here, we calculated the average inference  
 191 speed on an idle RTX 3060 GPU<sup>4</sup>. We find that,  
 192 with an optimized implementation, the inference

<sup>4</sup>We used the FasterTransformer (version 4.0) as in <https://github.com/NVIDIA/FasterTransformer>

	FLOPs #	Param. #	MNLI Acc.		SQuAD v2.0	
			match	mis-match	exact match	F1
<b>12H-12L</b> BERT <sub>BASE</sub>	46.3B	109.5M	84.4	84.4	77.4	80.4
<b>1H-144L</b> BERT <sub>BASE</sub> default	46.9B	110.0M	85.6	85.1	79.6	82.4
<b>1H-144L</b> BERT <sub>BASE</sub> Admin	46.9B	110.0M	85.2	85.4	79.2	82.5
<b>16H-24L</b> BERT <sub>LARGE</sub>	161.8B	335.1M	86.3	86.4	81.0	84.3
<b>1H-384L</b> BERT <sub>LARGE</sub> Admin	164.1B	337.4M	87.7	87.5	82.6	85.7

Table 3: The model performance on dev sets of MNLI and SQuAD 2.0. The FLOPs are calculated for the inference computation of one 512-length input sequence.

	GLUE	CoLA	SST-2	MRPC	SST-B	QQP	MNLI-m/mm	QNLI	RTE	WNLI
<b>12H-12L</b>	78.3	52.1	93.5	88.9/84.8	<b>87.1/85.8</b>	<b>71.2/89.2</b>	84.6/83.4	90.5	66.4	<b>65.1</b>
<b>1H-144L</b>	79.4	<b>59.2</b>	<b>94.2</b>	<b>89.3/85.4</b>	84.3/83.5	70.9/88.9	<b>85.1/84.3</b>	<b>91.0</b>	<b>69.0</b>	<b>65.1</b>
<b>16H-24L</b>	80.5	60.5	94.9	89.3/85.4	<b>87.6/86.5</b>	<b>72.1/89.3</b>	86.7/85.9	92.7	70.1	<b>65.1</b>
<b>1H-384L</b>	81.3	<b>62.7</b>	<b>95.1</b>	<b>90.5/87.2</b>	86.9/86.3	71.3/89.1	<b>87.4/86.5</b>	<b>93.3</b>	<b>72.7</b>	<b>65.1</b>

Table 4: The test performance on the GLUE benchmark with metrics described in Table 5.

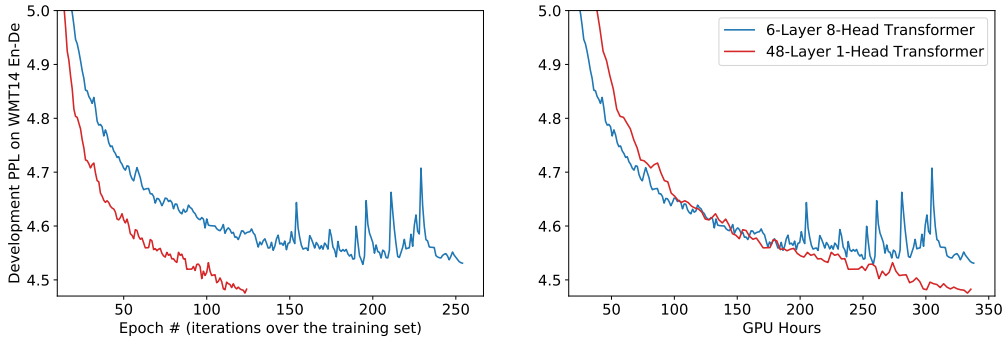


Figure 3: Transformer Training Efficiency (GPU Hours are calculated on an idle RTX 3060).

193 efficiency of the shallow multi-head Transformer  
 194 and the deep single-head Transformer are roughly  
 195 the same (visualized in Appendix, Figure 4).

196 **Training Speed.** As in Figure 3, we can find that  
 197 the training computation speed of the 1H-48L-48L  
 198 Transformer is about two times slower than the  
 199 8H-6L-6L Transformer. Meanwhile, the 8H-6L-6L  
 200 Transformer converges faster with regard to epoch  
 201 number, or GPU hours. This phenomenon verifies  
 202 our intuition that the network depth of the 6-Layer  
 203 Transformer has become a bottleneck of the model  
 204 capacity, which restricts the model performance.  
 205 Also, it indicates an limitation of the deep single-  
 206 head network, i.e., the computation time per update  
 207 is longer than the shallow multi-head network.

## 208 6 Conclusion

209 Here, we aim to understand the benefits of the  
 210 multi-head Transformer. We first show that deep  
 211 single-head Transformer also leverages multiple

212 representation subspaces and performs better than  
 213 the popular shallow multi-head Transformer. Then,  
 214 we suggest the main advantage of multi-head atten-  
 215 tion is the training stability since it has fewer layers  
 216 than the single-head attention when using the same  
 217 number of subspaces (number of attention heads).  
 218 We also show that, with recent advances in deep  
 219 learning, the training stability is no longer an ob-  
 220 stacle and it can lead to consistent performance im-  
 221 provements by turning shallow single-head Trans-  
 222 former into deep multi-head Transformer.

223 Our work opens up new possibilities to not only  
 224 further push the state-of-the-art but understand the  
 225 effectiveness of Transformer better. It leads to var-  
 226 ious interesting future work. For example, intu-  
 227 itively, both shallow multi-head Transformer and  
 228 deep single-head Transformer should not be the  
 229 optimal architecture, and neural architecture search  
 230 can be employed to find a good balance between  
 231 multi-head and single-head.

232  
233  
234  
  
235  
236  
237  
  
238  
239  
240  
  
241  
242  
243  
244  
245  
246  
  
247  
248  
249  
250  
  
251  
252  
253  
254  
  
255  
256  
257  
  
258  
259  
  
260  
261  
262  
263  
264  
  
265  
266  
  
267  
268  
269  
270  
271  
  
272  
273  
274  
  
275  
276  
277  
278  
  
279  
280  
281  
282

## References

Jimmy Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. 2016. Layer normalization. *ArXiv*, abs/1607.06450.

Alexei Baevski and Michael Auli. 2019. Adaptive input representations for neural language modeling. In *ICLR*.

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. *CoRR*, abs/1409.0473.

Mia Xu Chen, Orhan Firat, Ankur Bapna, Melvin Johnson, Wolfgang Macherey, George Foster, Llion Jones, Niki Parmar, Michael Schuster, Zhi-Feng Chen, Yonghui Wu, and Macduff Hughes. 2018. The best of both worlds: Combining recent advances in neural machine translation. In *ACL*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*.

Yong Fang, J. Gao, C. Huang, H. Peng, and R. Wu. 2019. Self multi-head attention-based convolutional neural networks for fake news detection. *PLoS ONE*, 14.

Xavier Glorot and Yoshua Bengio. 2010. Understanding the difficulty of training deep feedforward neural networks. In *AISTATS*.

A. Graves, Greg Wayne, and Ivo Danihelka. 2014. Neural Turing machines. *ArXiv*, abs/1410.5401.

Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. 2015. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 1026–1034.

Dan Hendrycks and Kevin Gimpel. 2016. Gaussian error linear units (gelus). *arXiv: Learning*.

Cheng-Zhi Anna Huang, Ashish Vaswani, Jakob Uszkoreit, Ian Simon, Curtis Hawthorne, Noam Shazeer, Andrew M. Dai, Matthew D. Hoffman, Monica Dinulescu, and Douglas Eck. 2019. Music transformer: Generating music with long-term structure. In *ICLR*.

Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. 2020. Reformer: The efficient transformer. *ArXiv*, abs/2001.04451.

Zhouhan Lin, Minwei Feng, Cícero Nogueira dos Santos, Mo Yu, Bing Xiang, Bowen Zhou, and Yoshua Bengio. 2017. A structured self-attentive sentence embedding. *ArXiv*, abs/1703.03130.

Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Jiawei Han. 2020a. On the variance of the adaptive learning rate and beyond. In *ICLR*.

Liyuan Liu, X. Liu, Jianfeng Gao, Weizhu Chen, and J. Han. 2020b. Understanding the difficulty of training transformers. In *EMNLP*. 283  
284  
285

Thang Luong, Hieu Pham, and Christopher D. Manning. 2015. Effective approaches to attention-based neural machine translation. *ArXiv*, abs/1508.04025. 286  
287  
288

Paul Michel, Omer Levy, and Graham Neubig. 2019. Are sixteen heads really better than one? In *NeurIPS*. 289  
290

Toan Q. Nguyen and Julian Salazar. 2019. Transformers without tears: Improving the normalization of self-attention. In *IWSLT*. 291  
292  
293

Niki Parmar, Ashish Vaswani, Jakob Uszkoreit, Lukasz Kaiser, Noam Shazeer, Alexander Ku, and Dustin Tran. 2018. Image transformer. In *ICML*. 294  
295  
296

Hao Peng, Roy Schwartz, Dianqi Li, and Noah A. Smith. 2020. A mixture of h - 1 heads is better than h heads. In *ACL*. 297  
298  
299

Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. Know what you don’t know: Unanswerable questions for squad. In *ACL*. 300  
301  
302

Prajit Ramachandran, Niki Parmar, Ashish Vaswani, Irwan Bello, Anselm Levskaya, and Jonathon Shlens. 2019. Stand-alone self-attention in vision models. In *NeurIPS*. 303  
304  
305  
306

Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. 2018. Self-attention with relative position representations. In *NAACL-HLT*. 307  
308  
309

David So, Quoc Le, and Chen Liang. 2019. The evolved transformer. In *International Conference on Machine Learning*, pages 5877–5886. PMLR. 310  
311  
312

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 *NIPS*. 313  
314  
315  
316

Petar Velickovic, Guillem Cucurull, A. Casanova, Adriana Romero, P. Lio’, and Yoshua Bengio. 2018a. Graph attention networks. 317  
318  
319

Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018b. Graph attention networks. In *ICLR*. 320  
321  
322

Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. 2018. Glue: A multi-task benchmark and analysis platform for natural language understanding. In *ICLR*. 323  
324  
325  
326

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 *ACL*. 327  
328  
329  
330

Ruibin Xiong, Yunchang Yang, Di He, Kai Zheng, Shu xin Zheng, Chen Xing, Huishuai Zhang, Yanyan Lan, Li-Wei Wang, and Tie-Yan Liu. 2019. On layer normalization in the transformer architecture. *ArXiv*, abs/2002.04745. 331  
332  
333  
334  
335

336 Baosong Yang, Longyue Wang, Derek Wong, Lidia S  
337 Chao, and Zhaopeng Tu. 2019. Convolutional  
338 self-attention networks. *arXiv preprint*  
339 *arXiv:1904.03107*.

340 Yuekai Zhao, Li Dong, Yelong Shen, Zhihua Zhang,  
341 Furu Wei, and Weizhu Chen. 2021. Memory-efficient  
342 differentiable transformer architecture search. *ArXiv*,  
343 abs/2105.14669.

## 344 A Related Work

345 There exist two aspects of related work regarding  
346 the topic here, *i.e.*, Attention and Transformer.

347 **Attention and Multi-Head Structure.** Attention  
348 modules are first proposed to capture the long-  
349 term dependency in sequence-to-sequence mod-  
350 els (Graves et al., 2014; Bahdanau et al., 2015).  
351 To calculate the output for a token in the target  
352 sequence, the attention module would calculate a  
353 weighted average of source token representations,  
354 while the weight is calculated by applying softmax  
355 on attention scores. Different variants of atten-  
356 tion modules calculate attention scores differently,  
357 *e.g.*, Graves et al. (2014) uses the cosine similar-  
358 ity, Bahdanau et al. (2015) uses the perception  
359 network, and Luong et al. (2015) uses dot prod-  
360 uct. While these modules only employ one sub-  
361 space, attempts like multi-head attention have been  
362 made to jointly attend to information from multiple  
363 subspaces (Lin et al., 2017; Vaswani et al., 2017),  
364 which is identified as one major reason behind the  
365 success of Transformer (Chen et al., 2018). Also,  
366 it has inspired several follow-up studies to analyze  
367 the multi-head structure (Michel et al., 2019; Peng  
368 et al., 2020). Specifically, Michel et al. (2019) ob-  
369 serves single-head Transformer performing better  
370 than multi-head Transformer for model pruning.  
371 Still, no study has been done on deep single-head  
372 Transformer training, due to its training difficulty.

373 **Transformer.** Transformer (Vaswani et al., 2017)  
374 has led to a series of breakthroughs in various do-  
375 mains (Devlin et al., 2019; Velickovic et al., 2018b;  
376 Huang et al., 2019; Parmar et al., 2018; Ramachan-  
377 dran et al., 2019). Meanwhile, Transformer train-  
378 ing has been found to be more challenging and  
379 attracted lots of attention to analyze why Trans-  
380 former is harder to train and how to stabilize Trans-  
381 former training (Liu et al., 2020a; Baevski and Auli,  
382 2019; Nguyen and Salazar, 2019; Wang et al., 2019;  
383 Xiong et al., 2019; Liu et al., 2020b). Many efforts  
384 have been made to improve Transformer, *e.g.*, rela-  
385 tive position encoding (Shaw et al., 2018) or replac-

386 ing dot-product attention with locality-sensitive  
387 hashing (Kitaev et al., 2020). Here, we choose to  
388 focus our study on the original Transformer model  
389 as proposed in Vaswani et al. (2017), and uses the  
390 initialization technique Admin to stabilize model  
391 training (Liu et al., 2020b), since this method does  
392 not include any additional hyper-parameters and  
393 its final model is equivalent to the original Trans-  
394 former.

## 395 B Transformer Architecture

396 The Transformer architecture contains two types  
397 of sub-layers, *i.e.*, Attention sub-layers and Feed-  
398 forward sub-layers. Each sub-layer is constructed  
399 with the shortcut connection and the Layer Norm.  
400 Specifically, it calculates the output as  $\mathbf{x}_{i+1} =$   
401  $f_{\text{LN}}(\mathbf{x}_i + f(\mathbf{x}_i))$ , where  $\mathbf{x}_i$  is the input of layer  $i$   
402 and the output of layer  $i - 1$  (top layers have larger  
403 indexes),  $f_{\text{LN}}$  is the Layer Norm, and  $f(\cdot)$  is multi-  
404 head attention  $f_{\text{ATT}}(\cdot)$  or feedforward  $f_{\text{FFN}}(\cdot)$  for  
405 Attention sub-layers and Feedforward sub-layers  
406 respectively.

407 **Layer Norm.** Layer norm (Ba et al., 2016) plays  
408 a vital role in the Transformer architecture. It is  
409 defined as  $f_{\text{LN}}(\mathbf{x}) = \gamma \frac{\mathbf{x} - \mu}{\sigma} + \nu$ , where  $\mu$  and  $\sigma$  are  
410 the mean and standard deviation of  $\mathbf{x}$ ,  $\gamma$  and  $\nu$  are  
411 learnable parameters.

412 **Feedforward.** Transformers use two-layer per-  
413 ceptrons as feedforward networks, *i.e.*,  $f_{\text{FFN}}(\mathbf{x}) =$   
414  $\phi(\mathbf{x}W^{(1)})W^{(2)}$ , where  $W^{(\cdot)}$  are parameters, and  
415  $\phi(\cdot)$  is the non-linear function. Specifically, the  
416 original Transformer ReLU as the activation func-  
417 tion, while later study uses other types of activation  
418 function, *e.g.*, BERT uses GELU as the activation  
419 function (Hendrycks and Gimpel, 2016).

420 **Attention.** Transformers use the multi-head atten-  
421 tion to capture the dependency among input to-  
422 kens, which is based on the scaled dot-product  
423 attention. Scaled dot-product attention tries to  
424 query information from the source sequence that  
425 is relevant to the target sequence. Specifically,  
426 assuming the length of the source sequence and  
427 the target sequence to be  $n$  and hidden dimen-  
428 sion to be  $m$ , the target sequence would be en-  
429 coded as  $Q \in R^{n \times m}$ , source sequence would  
430 be encoded as  $K \in R^{n \times m}$  and  $V \in R^{n \times m}$ .  
431 The scaled dot-product attention would calculate  
432 the output as  $f_{\text{Scaled Dot-Product Attention}}(Q, K, V) =$   
433  $\text{softmax}(\frac{QK^T}{\sqrt{m}})V$ , where  $\text{softmax}(\cdot)$  is the row-  
434 wise softmax.

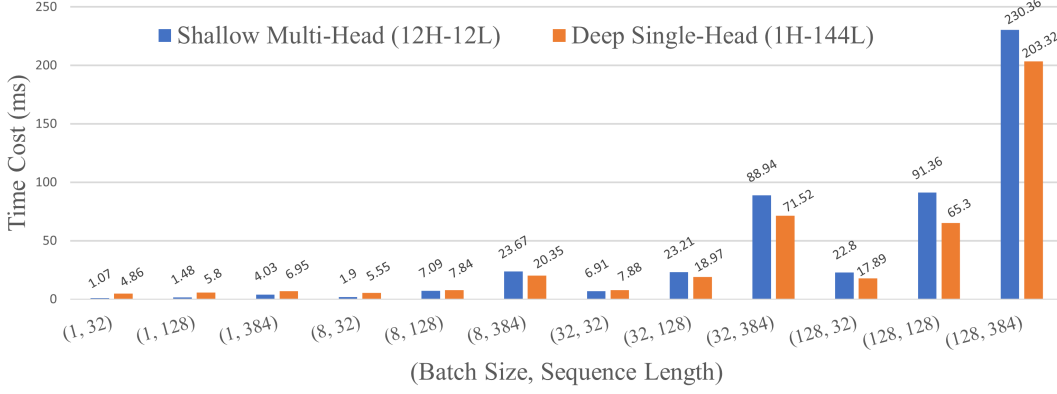


Figure 4: The Inference Speed of BERT-base with Different Batch Size and Sequence Length.

One scaled dot-product attention is believed to attend only one position in each row (for each target token), since the output of softmax typically would have one dimension significantly larger than other dimensions in each row. Multi-head attention was proposed to jointly attend to information from multiple representation subspaces, which employs multiple scaled dot-product attention in parallel. Specifically, it calculates the output as  $f_{\text{ATT}}(Q, K, V) = [\text{head}_1; \dots; \text{head}_h]W^{(O)}$ , where  $\text{head}_i = f_{\text{Scaled Dot-Product Attention}}(QW_i^{(Q)}, KW_i^{(K)}, VW_i^{(V)})$ ,  $W^{(\cdot)}$  are learnable parameters, and  $h$  is the number of heads.

**Transformer.** Transformer has two types of layer configurations when serving as the encoder and the decoder respectively. Here, we use  $\mathbf{x}_i$  to mark the input of sub-layer  $i$ . Each Transformer encoder layer contains two sub-layers, i.e., one attention sub-layer in a self-attention manner and one feedforward sublayer. Specifically, the attention sub-layer calculates outputs as  $\mathbf{x}_{2i+1} = f_{\text{LN}}(\mathbf{x}_{2i} + f_{\text{ATT}}(\mathbf{x}_{2i}, \mathbf{x}_{2i}, \mathbf{x}_{2i}))$  and the feedforward sub-layer calculates outputs as  $\mathbf{x}_{2i+2} = f_{\text{LN}}(\mathbf{x}_{2i+1} + f_{\text{FFN}}(\mathbf{x}_{2i+1}))$ . Notice that the attention sub-layer sets  $Q$ ,  $K$ , and  $V$  as the same value  $\mathbf{x}_{2i}$ , capturing the dependency among tokens within the same sequence, which is referred to as self-attention.

Each Transformer decoder layer contains three sub-layers, besides the self-attention sublayer and the feedforward sublayer, it also includes an encoder-decoder attention sub-layer between them. Specifically, the encoder-decoder attention sub-layer calculates outputs as  $\mathbf{x}_{3i+2} = f_{\text{LN}}(\mathbf{x}_{3i+1} + f_{\text{ATT}}(\mathbf{x}_{3i+1}, \mathbf{h}, \mathbf{h}))$ , where  $K$  and  $V$  are set to the encoder output  $\mathbf{h}$ .

### C Implicit Ensemble Structure

As in Figure 1, multi-head attention sub-layers and feedforward sub-layers have the implicit ensemble structure, i.e., each of these sub-layers can be viewed as an ensemble of smaller models. Now let us proceed to introduce those parallel structures in detail. Notations are introduced in Section B.

**Attention.** We split the weight matrix  $W^{(O)}$  into  $h$  parts by rows, i.e., we mark  $W^{(O)} = [W_1^{(O)T}; \dots; W_h^{(O)T}]^T$ . Then, the multi-head attention calculates outputs as:

$$f_{\text{ATT}}(Q, K, V) = [\text{head}_1; \dots; \text{head}_h]W^{(O)} = \sum_{i=1}^h \text{head}_i W_i^{(O)}$$

$$= \sum_{i=1}^h \text{softmax}\left(\frac{QW_i^{(Q)}W_j^{(K)T}K^T}{\sqrt{m}}\right)VW_i^{(V)}W_i^{(O)}$$

Note that each head can be viewed as a low-rank version of the general attention (Luong et al., 2015).

Thus, the multi-head attention can be viewed as jointly attending multiple places by ensembling multiple conventional attention modules. Specifically, the general attention module (Luong et al., 2015) calculates outputs as:

$$f_{\text{General Attention}}(Q, K, V) = \text{softmax}(QW_1K^T)VW_2$$

Comparing  $f_{\text{ATT}}$  and  $f_{\text{General Attention}}$ , we can find their major difference is that the multi-head attention decomposes the  $m \times m$  matrix  $W_1$  and  $W_2$  into  $\frac{W_i^{(Q)}W_i^{(K)T}}{\sqrt{m}}$  and  $W_i^{(V)}W_i^{(O)}$ , where  $W_i^{(Q)}, W_i^{(K)}, W_i^{(V)}, W_i^{(O)T} \in R^{m \times \frac{m}{h}}$ . With this low-rank decomposition, the parameter number

and computation complexity of the multi-head attention module would stay the same no matter what the value of  $h$  is (i.e., how many heads one layer has).

**Feedforward.** Similar to the Attention module, we can also rewrite the Feedforward sub-layer as an ensemble of  $h$  modules.<sup>5</sup> Specifically, we split the weight matrix  $W^{(1)}$  into  $h$  parts by rows and  $W^{(2)}$  into  $h$  parts by columns, i.e., we mark  $W^{(1)} = [W_1^{(1)}; \dots; W_h^{(1)}]$  and  $W^{(2)} = [W_1^{(2)T}; \dots; W_h^{(2)T}]^T$ . Then, the feedforward sub-layer calculates outputs can be rewrote as:

$$f_{\text{FFN}}(\mathbf{x}) = \phi(\mathbf{x}W^{(1)})W^{(2)} = \sum_{i=1}^h \phi(\mathbf{x}W_i^{(1)})W_i^{(2)}$$

Thus, the Feedforward sub-layer can be viewed as an ensemble of  $h$  sub-modules. Note that since the sum of the  $h$  sub-modules would be normalized by Layer Norm, their outputs are integrated in an averaging manner.

**Average Ensemble.** Each Transformer sub-layer calculates outputs as  $f_{\text{LN}}(\mathbf{x} + f(\mathbf{x}))$ , where  $f(\cdot)$  could be  $f_{\text{FFN}}(\cdot)$  and  $f_{\text{ATT}}(\cdot)$ . Thus, the sum calculated in the computation of  $f_{\text{ATT}}$  and  $f_{\text{FFN}}$  would be normalized by  $\text{Var}[\mathbf{x} + f(\mathbf{x})]$ . In this way, the joint effect of layer norm and the sum would be similar to combining these modules in an average ensemble manner.

## D Experiment Details

### D.1 Experiment Setup

In our experiments, we adopt hyper-parameter settings from previous work (Liu et al., 2020b; Devlin et al., 2019). More experiment details can be found as below.

**Transformer Model Configurations.** We conduct experiments with three Transformer models, i.e., Transformer-base for the WMT’14 English-German (EN-DE) translation task, BERT-base, and BERT-large for the language model pre-training. Specifically, the original Transformer-base model is 8H-6L-6L, and we compare it with 1H-48L-48L. The original BERT-base and BERT-large models are 12H-12L and 16H-24L, and we compare them with 1H-144L and 1H-384L. We use the Admin initialization (Liu et al., 2020b) to stabilize 1H-48L-48L Transformer-base and 1H-384L BERT-large.

<sup>5</sup>Note  $h$  here is decided to be consistent with the Multi-Head Attention sub-layers.

More detailed configurations are included in the appendix.

**Translation.** Here, we conduct experiments on WMT’14 EN-DE and evaluate model performance based on their BLEU score on the test set and perplexity score on the development set<sup>6</sup>.

**BERT.** Here, we follow the training setting from Devlin et al. (2019) and evaluate pre-trained language models on the SQuAD 2.0 (Rajpurkar et al., 2018) datasets for question answering, and the GLUE benchmark (Wang et al., 2018), which includes 9 subtasks (as in Table 5).

### D.2 Model Specificity

For machine translation, the original Transformer-base model is 8H-6L-6L Transformer encoder-decoder with 512-dimension word embedding, 64-dimension per-head attention output, and 2048-dimension feedforward network (Vaswani et al., 2017). Here, we compare it with 1H-48L-48L Transformer encoder-decoder with 512-dimension word embedding, 64-dimension per-head attention output, and 256-dimension feedforward network. For language model pre-training, BERT-base model is 12H-12L Transformer encoder with 768-dimension word embedding, 64-dimension per-head attention output, and 3072-dimension feedforward network; BERT-large model is 16H-24L Transformer encoder with 1024-dimension word embedding, 64-dimension per-head attention output, and 4096-dimension feedforward network (Devlin et al., 2019). Here, we compare them with deep single-head BERT-base model (1H-144L Transformer encoder with 768-dimension word embedding, single-head 64-dimension per-head attention output, and 256-dimension word embedding) and deep single-head BERT-large model (1H-384L Transformer encoder with 768-dimension word embedding, 64-dimension per-head attention output, and 256-dimension word embedding). To stabilize 1H-48L-48L Transformer-base and 1H-384L BERT-large, we use the Admin initialization (Liu et al., 2020b).

### D.3 Implementation Detail

Besides the layer number and head number, we adopted all hyper-parameters from previous work. Specifically, we followed (Liu et al., 2020b) for

<sup>6</sup>We mimicked the pre-processing setting from So et al. (2019). BLEU score is calculated by the BLEU implementation of fairseq (0.8.0).



Corpus	Train	Label	Task	Metric(s)	Domain
Single-Sentence Classification					
CoLA	8.5k	2	acceptability	Matthews corr.	misc.
SST-2	67k	2	sentiment	accuracy	movie reviews
Sentence Similarity/Paraphrase					
MRPC	3.7k	2	paraphrase	accuracy/F1	news
STS-B	5.7k	-	similarity	Pearson/Spearman corr.	misc.
QQP	364k	2	similarity	accuracy/F1	social QA questions
Natural Language Inference (NLI)					
MNLI	393k	3	NLI	(mis)matched acc.	misc.
QNLI	108k	2	QA/NLI	accuracy	Wikipedia
RTE	2.5k	2	NLI	accuracy	misc.
WNLI	634	2	coreference/NLI	accuracy	fiction books

Table 5: GLUE task descriptions and statistics. The second and fourth column denotes the number of training examples and the number of classes. Note that STS-B is a regression task.

591 machine translation experiments and (Devlin et al.,  
592 2019) for language model pre-training experiments.  
593 It is worth mentioning that, in (Liu et al., 2020b),  
594 the default initialization method is the Xavier initial-  
595 ization (Glorot and Bengio, 2010), which de-  
596 pends on the size of the weight matrix. Here, to  
597 control variables, we fix the initialization scale  
598 to be the same with original multi-head shallow  
599 Transformer. Meanwhile, for language model pre-  
600 training, since (Devlin et al., 2019) fixes the initial-  
601 ization scale for all models, we directly adopt the  
602 initialization strategy without modification.

#### 603 D.4 Training Detail

604 For machine translation experiments, we followed  
605 (Liu et al., 2020b) to conduct data pre-processing,  
606 conduct model training on Nvidia GPUs (includ-  
607 ing Quadro RTX 8000, GeForce RTX 3060, and  
608 Quadro RTX A6000). As to language model pre-  
609 training experiments, we followed (Devlin et al.,  
610 2019) to conduct data pre-processing, conduct  
611 model training with Google TPU v3.