1 We thank all reviewers for their efforts. Below we give detailed responses.

2 **1. "More examples and analysis of attention map" (R1&R3)** We pro-
3 vide more examples of the attention map in Figure 1. we also compute
4 the diagonal concentration for the attention map $M$ as quantitative metric.
5 It is define as $C = \frac{\sum_{|i-j|\leq 4} M_{i,j}}{\sum_{|i-j|>4} M_{i,j}}$. This indicates how much local depen-
6 dency that the attention map captures. The result in Table 1 shows that the
7 attention in BERT concentrates more on the local dependency.



Figure 1: More examples of attention maps.

8 **2. "Inference speed" (R2&R3)** We test our mixed-attention block and self-attention baseline from base-sized model
9 on Intel CPU (i7-6900K@3.20GHz). The mixed-attention has lower Flops and is much faster than self-attention, as
10 shown in Table 2. On the other hand, for the code we submitted as the supplementary material, our implementation for
11 mixed-attention on GPU and TPU is not well optimized for the efficiency yet. Thus its acceleration may not be obvious
12 when the input sequence length is short. We will work on further improvement on the low-level implementation.

13 **3. "More experiments on other NLP tasks" (R1&R2&R3)** This paper focuses on improving BERT and thus the
14 experiments are conducted in the language pre-training scenario. Thanks for reviewers' suggestions that remind us
15 span based dynamic convolution/mixed attention may be applied for other NLP tasks. Due to limited time for rebuttal,
16 we have not finished tuning the mixed attention model (of similar size as small ConvBERT), but it has shown better
17 performance than transformers (of similar size as small BERT) on language modeling task on WikiText-103, as shown
in Table 3. We will explore span-based dynamic convolution/mixed attention on other tasks in the future.
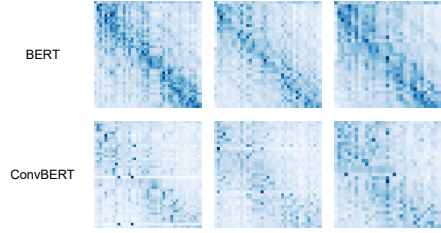
| Model | C (diagonal-concentration) |
|---|---|
| BERT | 0.941 |
| ConvBERT | 0.608 |

Table 1: Average concentration on MRPC.

| Block | Flops | Speed (ms/sample) |
|---|---|---|
| self-attention | 26.5G | 17.66 |
| mixed-attention | 19.3G | **12.94** |

Table 2: Inference speed.

| Model | Perplexity |
|---|---|
| Transformer | 34.21 |
| Ours | **32.95** |

Table 3: Result on WikiText-103.

| Model | Modification | Params | GLUE |
|---|---|---|---|
| ConvBERTmedium-small | +BNK,+GL | 14M | 81.0 |
| ConvBERTmedium-small | +BNK,+GL,+Larger | 17M | 81.1 |

Table 4: Ablation study on GLUE dev set.

19 **4. "Actual training time" (R2&R3)** We use direct implementation of our proposed
20 algorithm without dedicated low-level engineering acceleration at this stage as done in
21 ELECTRA BERT baseline. Even such, we achieve $2.67\times$ training acceleration (from 8
22 days to 3 days) on the base-sized model as shown in Table 5.

| Model | Training time |
|---|---|
| ELECTRA-small | 12h |
| ELECTRA-base | 192h |
| CONVBERT-small | 12h |
| CONVBERT-medium-small | 18h |
| CONVBERT-base | 72h |

Table 5: Training time on TPU v3-8.

23 **5. "Baseline and ablation study" (R1&R2&R4)** As suggested by R1, we add a
24 '+BNK,+GL,+Larger' baseline that increases the hidden dimension to 432. Result are
25 shown in Table 4. As expected, increasing the hidden dimension only slightly improves
26 the result (+0.1). As suggested by R2&R4, we add the experiments of only convolution based architecture (also small
27 sized) which performs poorly on downstream tasks and only achieves 64 on GLUE.

28 **6. "Definition of Grouped feed-forward" (R1)** We will add the detailed definition in the final version. The grouped
29 feed-forward module is defined as follows

$$M = \Pi_{i=0}^{g} \left[ f^i_{\frac{d}{g} \to \frac{m}{g}} \left( H_{[:, \ i-1:i\times\frac{d}{g}]} \right) \right], \quad M' = \text{GeLU}(M), \quad H' = \Pi_{i=0}^{g} \left[ f^i_{\frac{m}{g} \to \frac{d}{g}} \left( M'_{[:, \ i-1:i\times\frac{m}{g}]} \right) \right], \quad (1)$$

30 where $H, H' \in \mathbb{R}^{n\times d}$, $M, M' \in \mathbb{R}^{n\times m}$, $f_{d_1 \to d_2}(\cdot)$ indicates a fully connected layer that transforms dimension $d_1$ to
31 $d_2$, $g$ is the group number and $\Pi$ means concatenation.

32 **6. " Why use point-wise multiplication of $K_s$ and Q" (R4)** Instead of only using $K_s$, using point-wise multiplication
33 (a bi-linear operator) can merge information between a single token and its nearby tokens. Compared with concatenation
34 that makes the following fc layer occupy more parameters, it saves parameter number.

35 **7. "Comparison with Lite Transformer" (R2)** Lite transformer uses two branches of self-attention and dynamic
36 convolution. It is applied for translation, language modeling and abstractive summarization. While our ConvBERT
37 introduces span-based dynamic convolution with self-attention to form mixed-attention and is targeted for pre-training.
38 We have implemented 'Lite Transformer'-like architecture for pre-training. See results in row 4 'Dynamic' of Table 2
39 in our paper.

40 **8. "Mean and standard error on GLUE result" (R2)** The reported result on GLUE development set is the median of
41 9 runs. We list the mean and standard error of ConvBERT-small model on GLUE in Table 6.

42 **9. "Performance of base-sized model on
43 SQuAD" (R4)** As shown in Table 4 in our pa-
44 per, ConvBERT achieves similar performance
45 on SQuAD dataset with only $1/4$ training cost
46 compared to ELECTRA.

| Model | MNLI | QNLI | QQP | RTE | SST-2 | MRPC | CoLA | STS-B | Avg. |
|---|---|---|---|---|---|---|---|---|---|
| ConvBERT-small | 81.4±0.2 | 88.3±0.1 | 90.3±0.1 | 67.4±0.8 | 90.2±0.4 | 86.7±0.6 | 59.4±1.4 | 87.8±0.3 | 81.4±0.5 |

Table 6: Results on GLUE dev set.