

1 We thank all the reviewers for their positive and constructive comments.

2

3 To Reviewer #1

4 **Q1:** “My main question is whether authors have tried training much larger Transformer models that don’t fit into one  
5 GPU using their algorithm.”

6 **A1:** Yes, in our experiment, training 12-layer vanilla Transformer and Transformer-XL models with batch size 22  
7 cannot fit into one GPU. So we need to split them into multiple modules with each module placed on one GPU. Our  
8 method also works for models which are so large that even training with batch size 1 cannot fit into one GPU.

9

10 To Reviewer #2

11 **Q1:** “I strongly encourage authors to share the code used in experiments for researchers and practitioners in the future  
12 as soon as possible as noted in the paper.”

13 **A1:** Yes, we will release our source code on GitHub to encourage further research.

14

15 To Reviewer #3

16 **Q1:** “74: the notation  $\text{grad } f_{l, x_{i(t)}}$  is not used in equation (6). It would also be useful to remind what this notation  
17 means next to equation (8).”

18 **A1:** The notation  $\text{grad } f_{l, x_{i(t)}}$  is trying to clarify the relation between equations (5) and (6). Thanks for your good  
19 suggestion and we will also add it close to the equation (8).

20

21 **Q2:** “84: should be ‘any ... method’.”

22 **A2:** We will correct this typo in the future version.

23

24 **Q3:** “116 or 127: it would be good to say that proofs are in the supplementary material, rather than leaving  
25 it unstated whether the authors have proven the key theorems.”

26 **A3:** Thank you for the useful suggestion and we will revise it in the updated version.

27

28 **Q4:** “Could the authors please comment on why speed gains remain below  $2\times$ , even as up to  $4\times$  the num-  
29 ber of GPUs is used? The speed-ups should break down into two parts: the ‘time per step’ decreases as more GPUs  
30 are added, but the ‘number of steps to convergence’ increases as a function of  $K$ . The first of these could potentially  
31 improve with new hardware and systems software, while the second is inherent to the proposed method. What is the  
32 breakdown between these two components?”

33 **A4:** It is an ideal case to obtain linear speedup, using  $K\times$  machines to achieve  $K\times$  speedup regarding time. However, it  
34 is impossible to achieve even for data parallelism. It is also hard to define the concept of “step” in our model parallelism,  
35 because devices compute gradients from different steps in parallel. The goal of our method is to guarantee that there is no  
36 idle machines during the training and fully utilize all computing resources. On the contrary, the vanilla backpropagation  
37 algorithm is a sequential process and the other devices are idle when one device is processing. Regarding the effect of  
38  $K$  on the number of steps to converge, it is true that increasing  $K$  may require more steps to converge as the Theo-  
39 rem 1 suggests. A better algorithm which can mitigate the effect of  $K$  is one of the potential future direction of this paper.

40

41 **Q5:** “Relatedly, I also would be interesting in being able to better extrapolate how the algorithm might be-  
42 have in the following regimes: (a) A large number of layers (64+) is split across a relatively small number of devices  
43 ( $\sim 4$ ); and (b) the same 12-layer model used in the experiments is split across 11 devices.”

44 **A5:** For the regime (a), we think the speedup has little relation with the number of layers. Thus, it could get about the  
45 same  $2\times$  speedup as the 12-layer Transformer. As for the regime (b), due to the limited resources, we validate our  
46 algorithm by varying  $K$  from 3 to 5 and we didn’t test the performance of our method when there are more (e.g., 11)  
47 devices. We guess there exists a number that increasing the number of devices will not get further speedup. Finally, we  
48 want to emphasize that we can combine data parallelism with our model parallelism to obtain further speedup, which is  
49 nontrivial for practitioners.