

1 We thank the reviewers for the generous comments and invaluable feedback on the manuscript. Please find below our
2 point-by-point response to reviewers' concerns.

3 **Reviewer #1**

- 4 • *“Analysis of wall-clock time breakdown”*:

5 We will move the wall-clock time breakdown analysis into the main text of the paper.

- 6 • *“Literature of architectures that are explicitly easy to distribute”*:

7 Thank you for pointing out the related work, which we will cite in the camera-ready version. We would also
8 like to emphasize that although mixture-of-experts-type architectures are relatively easy for device placement,
9 the dispatching of examples and gathering of experts outputs is often non-trivial, and comes at the cost of
10 expert load balancing, expert utilization, batch-size tuning etc. We will add a detailed discussion regarding this
11 in the camera-ready version.

- 12 • *“Overview the architectures for AmoebaNet and Transformer”*:

13 We will include a brief overview and figures for the architectures considered here in the appendix.

- 14 • *“Suggestions on improvements”*:

15 Thank you very much for the suggestions. We have been running experiments with very deep transformer
16 language models and are currently exploring approaches to stabilize training for these models, beyond those
17 discussed in the paper. We will try to share our findings in the camera-ready version. At the same time, we
18 will soon open source the code to train 128+ block transformers on Github to allow researchers to train these
19 models with their own data.

20 **Reviewer #2**

- 21 • *“Batch size experiments feel a bit out of place”*:

22 Thank you for the suggestion. We added this discussion in order to provide more details for our machine
23 translation experiments. We will move this set of results into the supplementary material and move the
24 wall-clock time breakdown analysis into the main paper as suggested by reviewer #1.

25 **Reviewer #3**

- 26 • *“Elaborate the scheduling algorithm”*:

27 We agree with the reviewer that we should have done a better job at explaining our algorithms in greater
28 details. Currently, we allow the user to specify the per-layer computation cost for each layer (c_i) as measured
29 in flops, as briefly mentioned in Section 2.1. For distributing layers across accelerators we try to balance the
30 per-accelerator computation cost by minimizing the variance in the estimated cost of all accelerators. We have
31 open sourced our heuristic algorithm in our open-source repository. Since actual computational time might
32 differ from the flops estimation c_i , depending on different accelerators and architectures, we also allow the
33 user to manually specify layer placement. We are also planning to further improve our scheduling algorithm
34 to balance both memory utilization and computational time in order to maximize batch sizes and training
35 efficiency. We will describe the scheduling algorithm more clearly in Section 2 in the camera-ready version.

- 36 • *“How possible node failures are taken into account”*:

37 Our system uses synchronous training so any one of the machine failures would force the whole system to
38 restart from the previous checkpoint, as with standard non-model parallelism training.

- 39 • *“The micro-batching algorithm details are not presented in the paper.”*:

40 The overview of the micro-batch algorithm is depicted in Figure 2(c). We first split the tensor of mini-batch
41 input along the batch dimension into micro-batches and apply control-flow ops to loop over each micro-batch
42 to compute the gradients. Accumulated gradients over all micro-batches are applied to variables at the end of
43 each global step. We will describe the micro-batching algorithm more clearly in Section 2 of camera-ready.

- 44 • *“The lack of empirical benchmark system.”*:

45 We agree that providing a unified benchmark system for scaling giant neural network would be beneficial to
46 the research community. However, we are afraid it is a bit beyond the scope of our paper to define such a
47 benchmark system to evaluate all model parallelism libraries. For example, Mesh TensorFlow does not support
48 convolution operations, making it infeasible for us to compare our image model performance with theirs.

- 49 • *“Missing reference in the supplementary material”*:

50 Thank you for pointing this out. We will fix this in the camera-ready version.