

1 We thank the reviewers for their thoughtful comments. We will address their comments in the revised version of the  
2 paper. We believe our contributions will be interesting to the ML community for several reasons:

- 3 • General, *automatic* methods to reduce memory consumption for large deep learning models are critical. To  
4 our knowledge, techniques up to this point have been mostly focused on special case models. In practice,  
5 getting extremely large models to run on GPUs or TPUs is a time-consuming task, requiring expert knowledge  
6 of the problem.
- 7 • As Reviewer 1 pointed out, it is not always clear that the increase in running time is worth the savings in  
8 memory. But consider the case of training BERT on a TPU pod, which takes around 4 days. If that model  
9 doubled in size, it could take weeks of engineering work to get it to run. Our algorithm offers the ability  
10 to greatly reduce the memory with no human intervention, at the cost of an extra week of computation; our  
11 personal experience suggests this is a tradeoff that many would happily make. Of course, there are cases on  
12 both sides of this tradeoff, but having a reasonable starting point is valuable in itself.
- 13 • We provide a formalization of the problem with rigorous guarantees. As several reviewers noted, this is a  
14 useful contribution. In fact, while using treewidth to improve rematerialization has been suggested in an  
15 informal setting before (the tensorflow gradient checkpointing blog post), we show later in this rebuttal that the  
16 primary theoretical arguments therein are incorrect. In addition, we feel that our formalization will make the  
17 problem more attractive to the algorithmic community within NeurIPS, which would help garner additional  
18 insights. (For example, the fact that low pathwidth graphs have better theoretical guarantees is quite surprising,  
19 even if it doesn't immediately lend itself to practical algorithms.)

20 We now address a few of the specific reviewer concerns.

21 **Minimizing schedule length with memory constraint.** We have heuristics based on the TwRemat algorithm that, for  
22 a given memory limit, work by stopping the recursion early when the current set of nodes to compute can be scheduled  
23 without rematerialization under the memory limit. While minimizing the execution time of the schedule given a hard  
24 constraint on memory is an ideal version of our problem, we felt that including those heuristics would muddy the clarity  
25 of our results. However, in the revised version of this paper we will include a more thorough discussion of this. From  
26 a theoretical viewpoint, no algorithm with runtime sub-exponential in treewidth can yield a provable approximation  
27 guarantee on the schedule length (under common complexity assumptions). We will also add a discussion regarding  
28 this hardness result in the revised version.

29 **Comparisons with model parallelism.** While it is true that model parallelism is also a way to reduce the memory  
30 requirements of training, rematerialization and model parallelism are orthogonal techniques. One can simultaneously  
31 apply the two for even greater memory savings; hence our work is complementary to model parallelism. In addition, our  
32 work makes rematerialization easy to automate, whereas most model parallelism is achieved by careful hand-crafting.  
33 Since our rematerialization algorithms operate automatically on all networks, the neural network designer can obtain  
34 the benefits of memory reduction for “free” without any specialized optimization.

35 **Addressing the blog post.** Unfortunately the algorithm suggested by the TensorFlow checkpointing blog post is  
36 incorrect. That post draws on Courcelle's theorem (namely, every graph property definable in the monadic second-order  
37 logic of graphs is linear-time decidable on bounded treewidth graphs) to characterize the problems that become tractable  
38 for graphs of small treewidth. However, rematerialization is known to be PSPACE-complete [1]. Writing down a  
39 monadic second-order logic formula for rematerialization would imply that the problem is in PH (the polynomial  
40 hierarchy), and hence we would arrive at the surprising result that PSPACE collapses to PH. Thus, it is unlikely that we  
41 could apply Courcelle's theorem to rematerialization.

42 **Comparisons with previous work.** We thank the reviewers for pointing out the references that we missed and will  
43 be sure to include a comparison with those approaches in the revised version. We feel that it's more accurate to avoid  
44 the term “gradient checkpointing” here, because our technique can be applied even to inference graphs, but we will  
45 definitely include a discussion of terms.

## 46 References

- 47 [1] John R Gilbert, Thomas Lengauer, and Robert Endre Tarjan. The pebbling problem is complete in polynomial  
48 space. In *Proceedings of the 11th annual ACM Symposium on Theory of Computing*, pages 237–248. ACM, 1979.