

1 We thank the reviewers for the time and effort they put into the reviews and address their questions and comments
2 below. In particular R3 and R4 asked for more information about baselines and we point them below to the relevant
3 baselines in the paper (we compare to 6 different methods) and try our best to address their concerns about the method.

4 **re CIFAR/other architectures and code (R1/R2):** We agree that more datasets/architectures would always be useful
5 however in practice CIFAR is easier to sparsify and train sparse than ImageNet, thus we believe that ImageNet is
6 a better demonstration of the capabilities of our algorithm. We would argue that CIFAR should not be used for
7 comparisons of these techniques as findings do not always translate to larger datasets. We provide results from two
8 completely different/popular families of models (resnet and transformer) and two different tasks (language modelling
9 and classification) on large scale datasets and believe this demonstrates the efficacy of the method. We will add
10 pseudocode for how to implement this method. In general this can be implemented with a custom getter in TensorFlow,
11 such that the weight matrix is sparsified before any computation occurs.

12 **re tuning backwards sparsity (R2):** In practice the correlation between backward sparsity and performance is
13 straightforward: lower backwards sparsity generally allows us to discover better sparsity masks. A benefit of this aspect
14 of our method’s design is that backwards sparsity can thus be chosen based on the FLOP budget available and the amount
15 of memory available – in general the recommendation would be to use as dense a backward pass as the budget allows.

16 **Re baselines (R3/R4):** Both reviewers ask for a comparison to dense-to-sparse methods. We point the reviewers
17 to Figure 2 in the paper **where we compare with six different baselines, of which two (Pruning and SNIP) are**
18 **dense-to-sparse methods.** Moreover we outperform most of them on a per FLOP basis while staying entirely
19 sparse. Reviewer 4 asks for comparisons to SET and DSR. We compare to RigL which consistently outperforms
20 both SET and DSR so we did not see the value in including this baseline, although we do also compare to SET
21 in Figure 2 and are happy to add DSR.

22 **Re Novelty and Tuning (R3):** We disagree that the model is a simple tweak of pruning given we maintain constant
23 sparsity throughout training. Moreover we match state-of-the-art performance on Imagenet and sparsify TransformerXLs
24 – something which has thus far not been done in a sparse-to-sparse manner. On the point about tuning hyperparameters:
25 we in fact reduce the hyperparamtere compared to pruning (which requires a hand-tuned schedule). Additionally our
26 results outperform or are similar to those published by other sparse-to-sparse methods and are also the first method
27 that allows differing backward and forward sparsities.

28 **Re Wiki103 (R3):** There is no published method that has been applied to Wiki103 in a completely sparse-to-sparse
29 manner. We also disagree that there are huge performance drops - as we show **we match the performance of a 97M**
30 **parameter dense model with 57M and 60% backward sparsity.** On the point about comparing to pruning, we
31 do this on a smaller model (as we could not fit the pruning masks for bigger models in memory) and compare
32 in Table 5 in the appendix.

33 **Re theoretical flops (R4):** Every single published method we compare to also relies on theoretical FLOP reduction.
34 Current industry trends strongly suggest that future advancements in hardware/software kernels will allow us to
35 implement sparse methods with native-sparse support. While valid, the above criticism can be equally levied against
36 every previously published sparse-to-sparse method and would devalue every single published paper in the field.

37 **re theoretical proofs/limited practical value(R4):** We do not have a proof that using top-k is a best/better choice
38 and is unclear whether this is tractable theoretically. We do however point the reviewer to section 2.1 where we
39 explain with a taylor approximation why this approximation holds and that the algorithm will converge to some local
40 minima. On structured sparsity: our method is equally applicable to block sparsity given a function such as max-or
41 sum-of-absolute-values that reduces a block to a scalar value.

42 **re regularization loss (R4):** this is used in all experiments and is an integral part of Top-KAST (not a potential
43 add-on). We will make this clearer in the paper.

44 **re correctness and other experiments (R2):** We thank the reviewer for their comments and suggestions – on random
45 vs topk, we will explore this further. One added benefit of top-k over random is also the stability of the mask and not
46 having to constantly resample it. On corner cases, bwd sparsity=fwd would be similar to SET with a difference in how
47 we sample new parameters and bwd=0 would be like DNW but with the special regularization. On the comparison
48 with RigL and some of the subjective words, we will further clarify our language in the paper as our intention was
49 not to overclaim. In general for lower values of sparsity, RigL outperforms Top-KAST, and Top-KAST outperforms
50 for higher. We view the overall results as comparable (specifically both methods significantly outperform previous
51 baselines like SET) and note in the paper the main advantages of Top-KAST over RigL in section 4.

52 **Typos and minor comments and ablations:** We thank all reviewers for their notes on these and will correct this
53 in future versions. We thank R1,R2 and R4 for their suggestions for ablations. While unfortunately out of scope
54 for this document, we will strive to include this in the paper.