

1 We thanks all the reviewers very much for the constructive comments. Below please find our response. We hope you
2 could raise your evaluation if you find that we address your concerns.

3 **General Response:** Reviewers ask why the practical gain of the local imitation method is not that significant as the
4 theory. Here we give the response.

5 **(i)** Local imitation minimizes the surrogate local discrepancy loss while global imitation directly minimizes the original
6 discrepancy loss. The surrogate loss is able to upper bound the original loss up to a constant but this constant still
7 matters in practice. As a result, despite the local imitation has faster rate, both method can not consistently outperform
8 the other one; This justifies the usefulness of comparing and choosing the better method.

9 **(ii)** Techniques such as [1] can be easily applied to our framework to reduce the constant gap between the two losses.
10 However, in this paper, we aim at giving a simple algorithm which achieves SOTA both empirically and theoretically
11 and also has future potential. We leave the further algorithm engineering for future work.

12 **(iii)** Our way of dealing with the Lipschitz constant is standard and tight. It is unavoidable to have such gap between
13 theory and practice given the complexity of deep learning. But the mathematical structure we found gives stronger
14 justification, deeper understanding of the pruning paradigm, and motivates the new algorithm, which we demonstrate
15 useful in practice.

16 **Reviewer 1:** For the ImageNet experiments, the overhead of pruning phase is generally about 0.2x of that of finetuning
17 phase. We will add more discussion on the practical computational cost of the algorithm as well as more discussion on
18 knowledge distillation and NAS in the next version.

19 **Reviewer 2: (second paragraph in ‘correctness’)** Optimizing γ (rather than fixing it) is a core component to achieve
20 faster convergence rate. However, using global reconstruction loss, does not necessarily weaken the exponential rate,
21 if we still optimize γ . The proof framework of Theorem 4 can be extended to this case (optimizing γ and minimize
22 global reconstruction loss) and gives exponential decay rate. We will add more discussion on that. Besides, Taylor
23 approximation is not used in local imitation as each iteration is already fast (line 102-107) and thus does not effect the
24 exponential rate. **(a, h)** We will move this statement to main text. **(b)** Lemma 1-3 are self contained in terms of notation.
25 We guess that you find clM and riM undefined? Their definition is in the beginning of Appendix (the def of h and \bar{h} are
26 also there). The main intuition is that local imitation actually enjoys good geometric property (line 551), which makes
27 imitating the internal layer’s output very efficient (see key inequality between line 555 and 556). **(c)** Please see general
28 response. **(d)** Yes, it is equivalent in practice. The boldsymbol denotes vector. We will improve the clarity. **(e)** Sorry
29 for the typo. $U_i = [-a_i(k)/(1 - a_i(k)), 1]$. **(f)** The time and space complexity is small. For local imitation, we only
30 add one extra parameter for each neuron and this parameter will be merged into the scale of neuron after the pruning
31 finish. Following in line 102-107 and sec 5.1 in appendix, executing the pruning algorithm only requires one forward
32 pass and the selection can be done with simple matrix multiplication using $O(\text{batch size} \times \text{num of channels})$ space
33 complexity. The Taylor approximation is only applied to global imitation, which only adds 2 extra parameters for each
34 neuron. The main time complexity for selecting one neuron is calculating the gradient of the ancillary variable, which is
35 also small. See sec 5.4 for more details. **(g, i)** Yes, b is the ancillary parameters. But a in the appendix is identical to the
36 a used in the main text (e.g. between line 61 and 62). And a is not the scale on the activation but w_1 is (see line 53). In
37 practice, we regard the weight, activation and batchnorm together as a neuron. **(j)** The bound in Theorem 1 can not be
38 adapted to Ye et al. (2020) and their bound is tight. One key difference is that we optimizing γ during selection, which
39 significantly enlarge the search space of each greedy step and thus improves the rate. **(k)** They are very different. Our
40 method is based on forward selection (starting from empty network and greedily add neurons). Taylor approximation is
41 only a speed-up technique for us. In comparison, Molchanov (2017b) eliminates neurons from full network by looking
42 at neuron importance measured by Taylor expansion.

43 **Reviewer 3: (1)** Our theoretical improvement over Ye et al. (2020) is discussed in introduction and Table 1. In terms of
44 algorithm, intuitively, our method weights each selected neurons differently instead of simply average them like Ye et
45 al. (2020), which brings the improvement. We will give more detailed discussion on difference between Ye et al. (2020)
46 and other existing papers. **(2)** There are two reasons. Firstly, the final comparison is made after finetuning and thus
47 part of the improvement on pruning phase might be smaller after finetuning. For the second point, please see general
48 response.

49 **Reviewer 4:** Please see general response (iii) on the ‘Lipschitz continuous’ comment. We will add more ablation
50 studies similar to that in sec 2.3 and sec 3.1 to compare with GFS.

51 **Reference:** [1] Zhuang, Zhuangwei, et al. "Discrimination-aware channel pruning for deep neural networks." Advances
52 in Neural Information Processing Systems. 2018.