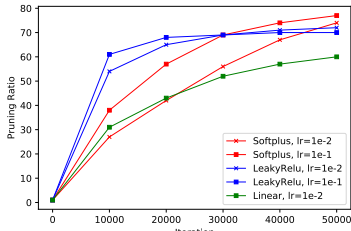


1 **Review1:** Thank you for the review and we appreciate the insightful suggestions. We will improve the paper accordingly:  
 2 In the related work section, we will first introduce structured and unstructured pruning and summarize existing papers  
 3 in these two categories, respectively. We also will include two latest papers as follows:

4 Gomez *et al.* [2019] proposed to keep neurons with high magnitude and prune only neurons with smaller magnitude in  
 5 a stochastic way. The accuracy is maintained by reducing the dependency of important neurons on unimportant neurons.  
 6 Liu *et al.* [2018] did comprehensive experiments showing that training-from-scratch on the right sparse architecture  
 7 yields better results than pruning from pre-trained models. Such a result implies that the pruning process is actually  
 8 finding the right network structure, thus bridging the gap between pruning and architecture search.

9 Our paper can be further positioned as a structure learning procedure where a neuron’s importance is learned by the  
 10 attached trainable auxiliary parameter. Unimportant neurons will be pruned if the corresponding auxiliary parameter  
 11 decreases below a threshold, and important neurons will be recovered otherwise.

12 For the experiment part, thank you for referring the AutoSlim paper. We did additional experimental comparison on  
 13 MobileNet v1 and v2 with the AutoSlim paper. Our model achieves better performance than AutoSlim under extreme  
 14 cases, where FLOPs are compressed down to 150M and 209M, respectively. We will include more neuron pruning  
 15 results and discussion on ResNet50, SuffleNet, and MNasNet to replace the LeNet5 part in the camera ready version.  
 16 We also include the performance of different STEs on VGG-like below, and will add more ablation discussion and  
 17 results. Compared to Softplus STE, LeakyRelu converge faster but with lower CR, and linear STE converge slower  
 18 with lower CR. In addition, figure 2(c) in the paper shows the comparison with and without forcing the sign.



Model	MobileNet v1			MobileNet v2		
	Base	AutoSlim	Ours	Base	AutoSlim	Ours
FLOPs	325M	325M	325M	300M	305M	300M
Top1-Err	31.6	<b>28.5</b>	28.8	28.2	<b>25.8</b>	26.1
FLOPs	150M	150M	150M	209M	207M	209M
Top1-Err	36.7	32.1	<b>31.8</b>	30.2	27	<b>26.7</b>

19 **Review2:** Thank you for the comments of the paper especially on the clarity and potential improvement. We will  
 20 mention before Eq. (9) that  $R_n$  refers to n-norm regularization. We included some theoretical analysis on the correlation  
 21 gradient direction consistency in the paper which ensures the effectiveness of gradient descent. We will also include the  
 22 necessary convergence discussion in the appendix.

23 **Review3:** We would like to thank reviewer for the valuable thoughts and detailed comments on the paper presentation.

24 **Originality:** The significant difference from other mask based pruning lies on the weak-correlation between the weights  
 25 and auxiliary parameters. We decoupled the correlation between update rule and weight magnitude which is harmful to  
 26 pruning task, as is discussed in Sec3.3 and figure 2(c). Theoretical analysis and experimental results show such weak  
 27 correlation leads to a more compact and more accurate model.

28 a) We agree that the masking formulation has been adopted by existing papers to learn sparse NNs, but optimizing the  
 29 binary masks effectively, stably and efficiently requires significant research efforts. The modified update rule results in  
 30 a better local minima of the auxiliary parameter network for the pruning task. Moreover, we have compared with [1] in  
 31 the introduction section that, "In contrast to Louizos *et al.* [2018], our method avoids inefficient and high variance single  
 32 step Monte-Carlo sampling and places no assumptions on the prior distribution". We also show in Table 1 that our model  
 33 significantly outperforms [1]. b) The update rule is carefully designed as mentioned above. The way of calculating the  
 34 partial derivations follow coarse gradients of predefined straight through estimators as discussed in Sec3.2. Specifically,  
 35 our Tensorflow implementation employs gradient overriding function for coarse gradient calculation.

36 **Quality:** 1) We agree that the high scalability and training-from-scratch ability mentioned in 4) and 5) validate the  
 37 significance of the proposed method. We will reorganize the list of contributions based on your nice suggestion. 2) The  
 38 update rule is not following an accurate gradient w.r.t.  $m$ , but a direction that forms an acute angle with the original  
 39 gradient direction. Here we kind of abuse the word "gradient" since it is the actual direction that auxiliary parameters  
 40 will be updated on. We will rephrase it as "updating rule" or "modified gradient" to avoid potential misunderstanding.  
 41 3) Thank you for pointing out the inaccurate naming of "sensitivity consistency". We will rename it "sensitivity  
 42 decoupling", since the update rule of auxiliary parameters is decoupled with the magnitude of weights. 4) Thank you  
 43 for your suggestion to improve the equations’ consistency. We will add the regularization term  $\lambda \mathcal{R}(W)$  of weight  $W$  to  
 44 Eq.(3). 5) We will add the following additional explanation on recoverable pruning to Sec3.4: Recovering a weight  $w_{ij}$   
 45 requires the gradient of  $m_{ij}$  moving toward positive direction. The weight decay will decrease the magnitude of  $w_{ij}$   
 46 and provide a negative gradient to  $m_{ij}$ , which reduces the recoverability. 6) There is a typo in the sentence. We will  
 47 rephrase it to "The goal is to find the minimum subset  $\{w : w \in \mathcal{W}\}$  that preserves the model accuracy."