

1 We thank the reviewers for their detailed comments. Please see our response below.

2 **R1:** ○ “... common implementation of weight decay [1] will usually multiply the amount of weight decay by the learning  
3 rate.” The same holds in our setup: We have an  $L_2$  regularization term in the loss. In the gradient descent update  
4 equations, this results in a weight decay term where the  $L_2$  coefficient is multiplied by the learning rate, as in [1].  
5 Perhaps this is not immediately obvious in Theorem 2 where we consider gradient flow rather than gradient descent, but  
6 the statement holds in that case as well.

7 ○ “How do different learning rate schedules affect the conclusion?”: We address LR schedule questions below.

8 ○ “It would be great if the authors can provide more experiments on ... AUTOL2” We ran additional experiments  
9 training Wide ResNets and ConvNets on CIFAR-100 and SVHN, with similar conclusions: AutoL2 either beats or  
10 matches the performance and training speed of a tuned, constant  $L_2$  parameter. As an example, for CIFAR-100 with a  
11 fixed lr and evolved for 500 epochs, the optimal  $L_2$  parameter gives a test accuracy of 0.76 while AutoL2 gives 0.79.

12 **R2:** ○ “((1)) If I could have access to the test set...”. We reject the claim that our submission “violates the ethics of  
13 machine learning research”. The method we propose for finding the optimal  $L_2$  parameter does not rely on a validation  
14 or test set — it can be implemented using training data alone! Our baseline is to compare against tuning on the test set,  
15 as is commonly done in the deep learning literature. If our setup introduces any bias (compared with having a separate  
16 validation set), it is in favor of the baseline.

17 ○ “((2)) I have concerns on comparing AutoL2...”. Indeed, the experiment of Fig. 1c does not include lr decay (as  
18 discussed in the text). Experiments with lr decay and AutoL2 are presented in the SM. Please see below for additional  
19 discussion.

20 ○ “((3)) The practicality of the proposed work... The most interesting plots might be Test Acc vs Lambda in Fig. 2c and  
21 2f, but the authors didn’t report multiple runs of the same configuration...” Please see Fig. 1a, 2b, and 2e for many  
22 additional runs illustrating the same effect. We made an effort in the paper to understand the settings under which our  
23 conclusions hold; these are summarized in the Discussion and in SM.D.

24 **R3:** ○ “... more insights on the relation between learning rate scheduler and AutoL2...” We address this point in the  
25 learning decay discussion at the bottom of the page.

26 ○ “... the lambda update refractory period is not detailed ...” The refractory period lasts for  $\frac{\lambda(t=0)}{\lambda(t)}$  steps. This is  
27 explained in the SM, and we will clarify it in the main text. The results are not sensitive to the choice of  $\lambda(t=0)$  and  
28 we pick 0.1 in our experiments.

29 ○ “It would be interesting to see on the same graph, training with learning rate scheduler...” In the SM we have the  
30 training curves for the models trained with a learning rate schedule (see figure S10).

31 ○ “In Figure 1a and 1b, how is the best test accuracy determined?...” In Figs. 1a,2,3, the model is trained for a specified  
32 number of epochs and we report the best test accuracy during training. In Fig. 1b, we have multiple runs with different  
33  $L_2$  and the ‘optimal  $L_2$ ’ corresponds to the value which achieved the best test accuracy during training.

34 **R4:** ○ “Default test accuracy of WRN of 0.92 seems a low ...” The accuracies of Figs. 1a,1b exceed 0.96 for some  
35 choices of the  $L_2$  parameter. For figure 1c, we limited the training time to 200 epochs and used a fixed lr, which explains  
36 why the accuracy is lower.

37 ○ “It would be great to have a better description on the relationship of learning rate schedules...” Please see LR  
38 schedule comments below.

39 ○ “An imagenet training with an overparameterized (resnet50) would boost...” We agree. We have run several additional  
40 experiments on CIFAR-100 and SVHN and the results show consistent improvement when using AutoL2. It is not  
41 clear whether ResNet50 on ImageNet (when trained with data augmentation) is sufficiently overparameterized for our  
42 purposes — we comment on this in SM.D.

43 **Comments on learning rate schedules.** Here we address reviewers’ questions regarding learning rate schedules, and  
44 we will revise the paper to include these points. Our empirical observations apply in the presence of learning rate  
45 schedules, as illustrated in the paper. Figure 1a shows that the test accuracy remains roughly the same if we scale the  
46 training time as  $1/\lambda$ . As to our proposed algorithms, we found that our predicted optimal  $L_2$  value is within an order  
47 of magnitude of the tuned value when using learning rate schedules, when accounting for the schedule by weighing  
48 each step according to the learning rate. Therefore, as in the case of fixed learning rate, our prediction is beneficial  
49 for hyperparameter tuning. For the AutoL2 scheduler, in the presence of learning rate decay we find that it performs  
50 similarly to a tuned, fixed  $L_2$  parameter but does not exceed it. This is still beneficial, in that it saves the need to tune  
51 the  $L_2$  parameter.