

## A. Appendix

### A.1. Second order term from bilevel optimization

For the second order term for the optimization of augmentation parameters, we follow the formulation in [23], which we summarize below. We treat the optimization of augmentation parameters and weights of the neural network as a bilevel optimization problem, where  $\alpha$  are the augmentation parameters and  $w$  are the weights of the neural network. Then the goal is to find the optimal augmentation parameters  $\alpha$  such that when weights are optimized on the training set using data augmentation given by  $\alpha$  parameters, the validation loss is minimized. In other words:

$$\min_{\alpha} \mathcal{L}_{val}(w^*(\alpha), \alpha) \text{ s.t. } w^*(\alpha) = \operatorname{argmin}_w \mathcal{L}_{train}(w, \alpha). \quad (1)$$

Then, again following [23], we approximate this bilevel optimization by a single virtual training step,

$$\nabla_{\alpha} \mathcal{L}_{val}(w^*(\alpha), \alpha) \approx \nabla_{\alpha} \mathcal{L}_{val}(w - \xi \nabla_w \mathcal{L}_{train}(w, \alpha), \alpha), \quad (2)$$

where  $\xi$  is the virtual learning rate. Eq. 2 can be expanded as

$$\nabla_{\alpha} \mathcal{L}_{val}(w^*(\alpha), \alpha) \approx \nabla_{\alpha} \mathcal{L}_{val}(w - \xi \nabla_w \mathcal{L}_{train}(w, \alpha), \alpha) - \xi \nabla_{\alpha, w}^2 \mathcal{L}_{train}(w, \alpha) \nabla_w \mathcal{L}_{val}(w', \alpha), \quad (3)$$

where  $w' = w - \xi \nabla_w \mathcal{L}_{train}(w, \alpha)$ . In the case where the virtual learning rate,  $\xi$ , is zero, the second term disappears and the first term becomes  $\nabla \mathcal{L}_{val}(w, \alpha)$ , which was called the first-order approximation [23]. This first-order approximation was found to be highly significant for architecture search, where most of the improvement (0.3% out of 0.5%) could be achieved using this approximation in a more efficient manner (1.5 days as opposed to 4 days). Unfortunately, when  $\alpha$  represents augmentation parameters, first-order approximation is irrelevant since the predictions of a model on the clean validation images do not depend on the augmentation parameters  $\alpha$ . Then we are left with just the second order approximation, where  $\xi > 0$ , which we approximate via finite difference approximation as

$$\nabla_{\alpha, w}^2 \mathcal{L}_{train}(w, \alpha) \nabla_w \mathcal{L}_{val}(w', \alpha) \approx \frac{\nabla_{\alpha} \mathcal{L}_{train}(w^+, \alpha) - \nabla_{\alpha} \mathcal{L}_{train}(w^-, \alpha)}{2\epsilon}, \quad (4)$$

where  $w^{\pm} = w \pm \epsilon \nabla_w \mathcal{L}_{val}(w', \alpha)$  and  $\epsilon$  is a small number.

#### A.1.1 Magnitude methods

A random magnitude uniformly randomly samples the distortion magnitude between two values. A constant magnitude sets the distortion magnitude to a constant number

Magnitude Method	Accuracy
Random Magnitude	97.3
Constant Magnitude	97.2
Linearly Increasing Magnitude	97.2
Random Magnitude with Increasing Upper Bound	97.3

Table 7. **Results for different ways of setting the global magnitude parameter  $M$ .** All magnitude methods were run on CIFAR-10 with Wide-ResNet-28-10 for 200 epochs. The reported accuracy is the average of 10 runs on the validation set for the best hyperparameter setting for that magnitude method. All magnitude methods searched over had 48 different hyperparameter settings tried.

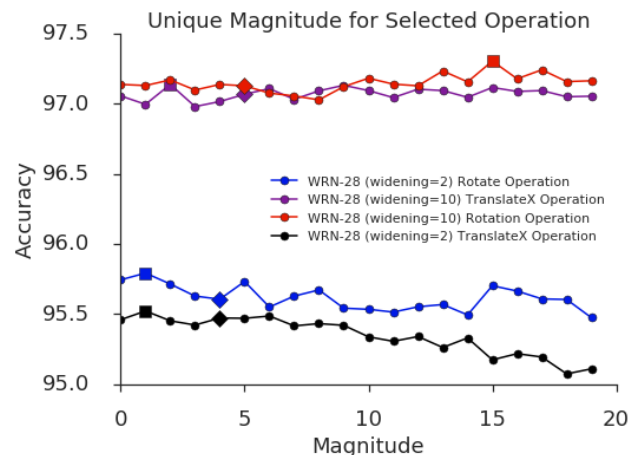


Figure 5. **Performance when magnitude is changed for one image transformation.** This plot uses a shared magnitude for all image transformations and then changes the magnitude of only one operation while keeping the others fixed. Two different architectures were tried (WRN-28-2 and WRN-28-10) and two different image transformations were changed (Rotate and TranslateX), which results in the 4 lines shown. Twenty different magnitudes were tried for the selected transformation ([0 – 19]). The squares indicate the optimal magnitude found and the diamonds indicate the magnitude used for all other transformations (4 for WRN-28-2 and 5 for WRN-28-10).

during the course of training. A linearly increasing magnitude interpolates the distortion magnitude during training between two values. A random magnitude with increasing upper bound is similar to a random magnitude, but the upper bound is increased linearly during training. In preliminary experiments, we found that all strategies worked equally well. Thus, we selected a constant magnitude because this strategy includes only a single hyper-parameter, and we employ this for the rest of the work. The results from our experiment on trying the different magnitude strategies can be seen in Table 7.

### A.1.2 Optimizing individual transformation magnitudes

Figure 5 demonstrates that changing the magnitude for one transformation, when keeping the rest fixed results in a very minor accuracy change. This suggests that tying all magnitudes together into a single value  $M$  is not greatly hurting the model performance. Across all for settings in Figure 5 the difference in accuracy of the tied magnitude vs the optimal one found was 0.19% 0.18% for the rotation operation experiments and 0.07% 0.05% for the TranslateX experiments. Changing one transformation does not have a huge impact on performance, which leads us to think that tying all magnitude parameters together is a sensible approach that drastically reduces the size of the search-space.

## A.2. Experimental Details

### A.2.1 CIFAR

The Wide-ResNet models were trained for 200 epochs with a learning rate of 0.1, batch size of 128, weight decay of  $5e-4$ , and cosine learning rate decay. Shake-Shake [10] model was trained for 1800 epochs with a learning rate of 0.01, batch size of 128, weight decay of  $1e-3$ , and cosine learning rate decay. ShakeDrop [45] models were trained for 1800 epochs with a learning rate of 0.05, batch size of 64 (as 128 did not fit on a single GPU), weight decay of  $5e-5$ , and cosine learning rate decay.

On CIFAR-10, we used 3 for the number of operations applied ( $N$ ) and tried 4, 5, 7, 9, and 11 for magnitude. For Wide-ResNet-2 and Wide-ResNet-10, we find that the optimal magnitude is 4 and 5, respectively. For Shake-Shake (26 2x96d) and PyramidNet + ShakeDrop models, the optimal magnitude was 9 and 7, respectively.

### A.2.2 SVHN

For both SVHN datasets, we applied cutout after RandAugment as was done for AutoAugment and related methods. On core SVHN, for both Wide-ResNet-28-2 and Wide-ResNet-28-10, we used a learning rate of  $5e-3$ , weight decay of  $5e-3$ , and cosine learning rate decay for 200 epochs. We set  $N = 3$  and tried 5, 7, 9, and 11 for magnitude. For both Wide-ResNet-28-2 and Wide-ResNet-28-10, we find the optimal magnitude to be 9.

On full SVHN, for both Wide-ResNet-28-2 and Wide-ResNet-28-10, we used a learning rate of  $5e-3$ , weight decay of  $1e-3$ , and cosine learning rate decay for 160 epochs. We set  $N = 3$  and tried 5, 7, 9, and 11 for magnitude. For Wide-ResNet-28-2, we find the optimal magnitude to be 5; whereas for Wide-ResNet-28-10, we find the optimal magnitude to be 7.

### A.2.3 ImageNet

The ResNet models were trained for 180 epochs using the standard ResNet-50 training hyperparameters. The image size was 224 by 224, the weight decay was 0.0001 and the momentum optimizer with a momentum parameter of 0.9 was used. The learning rate was 0.1, which gets scaled by the batch size divided by 256. A global batch size of 4096 was used, split across 32 workers. For ResNet-50 the optimal distortion magnitude was 9 and ( $N = 2$ ). The distortion magnitudes we tried were 5, 7, 9, 11, 13, 15 and the values of  $N$  that were tried were 1, 2 and 3.

The EfficientNet experiments used the default hyperparameters and training schedule, which can be found in [41]. We trained for 350 epochs, used a batch size of 4096 split across 256 replicas. The learning rate was 0.016, which gets scaled by the batch size divided by 256. We used the RMSProp optimizer with a momentum rate of 0.9, epsilon of 0.001 and a decay of 0.9. The weight decay used was  $1e-5$ . For EfficientNet B5 the image size was 456 by 456 and for EfficientNet B7 it was 600 by 600. For EfficientNet B5 we tried  $N = 2$  and  $N = 3$  and found them to perform about the same. We found the optimal distortion magnitude for B5 to be 17. The different magnitudes we tried were 8, 11, 14, 17, 21. For EfficientNet B7 we used  $N = 2$  and found the optimal distortion magnitude to be 28. The magnitudes tried were 17, 25, 28, 31.

The default augmentation of horizontal flipping and random crops were used on ImageNet, applied before RandAugment. The standard training and validation splits were employed for training and evaluation.

### A.3. COCO

We applied horizontal flipping and scale jitters in addition to RandAugment. We used the same list of data augmentation transformations as we did in all other classification tasks. Geometric operations transformed the bounding boxes the way it was defined in Ref. [51]. We used a learning rate of 0.08 and a weight decay of  $1e-4$ . The focal loss parameters are set to be  $\alpha = 0.25$  and  $\gamma = 1.5$ . We set  $N = 1$  and tried distortion magnitudes between 4 and 9. We found the optimal distortion magnitude for ResNet-101 and ResNet-200 to be 5 and 6, respectively.