

# Supplementary of FLOAT: Fast Learnable Once-for-All Adversarial Training for Tunable Trade-off between Accuracy and Robustness

Souvik Kundu<sup>1,2\*</sup>, Sairam Sundaresan<sup>1</sup>, Massoud Pedram<sup>2</sup>, Peter A. Beerel<sup>2</sup>

<sup>1</sup>Intel Labs, USA <sup>2</sup>University of Southern California, Los Angeles, USA

{souvik.kundu, sairam.sundaresan}@intel.com {pedram, pabeerel}@usc.edu

## 1. Training Details

### 1.1. Hyperparameters

We trained the models for 200 epochs with a starting learning rate (LR) of 0.1 and a cosine annealing LR scheduler [5]. We used the SGD optimizer with a weight decay and momentum value of  $5e - 4$  and 0.9, respectively. For CIFAR-10, CIFAR-100, and SVHN we used a batch-size of 128 and for STL-10 which has a comparatively larger resolution ( $96 \times 96$  as image height/width compared to  $32 \times 32$  of others), we used a batch-size of 64. We used half of each image-batch as clean and the remaining half as perturbed (using PGD-7 attack). For test-time random autoattack [1] generation we used  $L_\infty$  norm with  $\epsilon$  as  $8/255$ . To reproduce the results with OAT and OATS we used the official repository of [6] and used their recommended  $\lambda$  distribution, encoding dimension, sampling scheme as well as  $\sigma$  value for attack.

### 1.2. ImageNet Training

To train on ImageNet dataset on limited compute resources, we first sampled 100 classes from it, with 500 training and 50 test samples from each class, each with resolution of  $224 \times 224 \times 3$ . We then trained a ResNet18 on the sampled subset of ImageNet with FLOAT for 60 epochs. We used multi step LR scheduler that decayed by a factor of 0.1 respectively at epoch 30 and 45. We used a training batch size of 50 and partitioned each batch equally to create clean and adversarial samples. To create train and test adversarial samples we used similar adversarial strengths as other datasets.

### 1.3. Models

We used ResNet34 to evaluate CIFAR-10 and CIFAR-100, WRN16-8 to evaluate SVHN, and WRN40-2 to evaluate STL10. Table 1 shows the models' parameter and FLOP count (contributed by the convolutions and FC layers as they attribute to majority of the FLOPs and parameters).

Table 1. Model parameters and FLOPs details.

Model type	Parameters (M)		FLOPs (G)	
	OAT	FLOAT	OAT	FLOAT
ResNet34	31.4	21.28	1.18	1.165
WRN40-2	3.22	2.25	3	2.96
WRN16-8	15.43	10.97	2.1	2.01

### 1.4. Trained Scaling Factors

The trained noise scaling-factor  $\alpha^l$  for each layer of a FLOAT model is shown in Fig. 1. The  $\alpha^l$  values are generally high at the initial layers of the models while and reduce to near zero at later layers. [3] has also observed a similar trend in trained alpha values while training targeting only robustness.

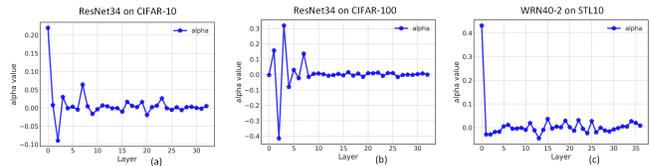


Figure 1. Trained noise scaling factor value (layer-wise) for (a) ResNet34 on CIFAR-10, (b) ResNet34 on CIFAR-100, and (c) WRN40-2 on STL10.

### 1.5. Alternate Results Plots in Terms of $\lambda$ vs. Accuracy

To bolster Fig. 7 in the original manuscript, Fig. 2 show an alternate view of the CA-RA comparisons between FLOAT, OAT, and PGD-AT.

### 1.6. Training Time as a Function of Size of $S_\lambda$ used for training.

The comparison of training times of OAT and FLOAT as a function of different numbers of  $\lambda$ s is shown in Fig. 3. Because FLOAT always has two possible training  $\lambda$ s, the time is constant. However, for OAT, the training time increases when supporting more training  $\lambda$ s. Note, in the

\*Part of the work was done when the first author was with USC.

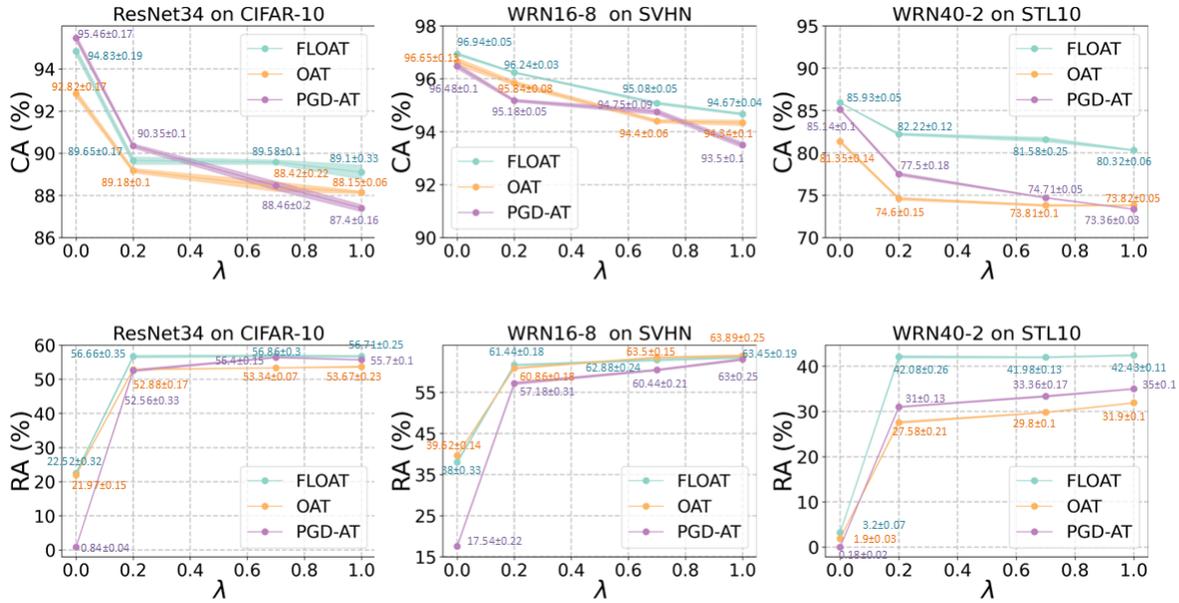


Figure 2. Comparison of trade-off between accuracy and robustness of FLOAT, OAT, and PGD-AT. (a)-(c) and (d)-(f) show CA and RA plots vs different  $\lambda$  values, respectively. The values corresponds to the mean  $\pm$  std error for each point.

main manuscript, we report the training time comparisons for OAT training with four training  $\lambda$ s (0.0, 0.2, 0.7, 1.0).

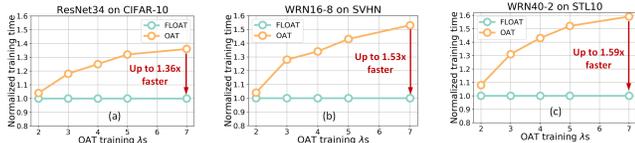


Figure 3. Comparison of per epoch normalized training time between FLOAT and OAT for different number of OAT training  $\lambda$ s, on (a) CIFAR-10, (b) SVHN, and (c) STL10. Note, lower normalized value implies faster training time. meaning better training efficiency.

## 2. More Analysis on Pruned Models

We now show per layer parameter density for a specific target global density  $d$ . In particular, Fig. 4 shows the layer-wise density of non-zero weight. As we can see in the Fig, the earlier layers generally have higher parameter density (for both the densities). This trend is similar to that observed in [4]. Also, earlier literature have used the parameter density as a proxy of the layer *sensitivity* [4, 2]. Here sensitivity of a layer can be measured by the accuracy reduction caused by pruning a certain ratio of parameters from that layer. Thus, we can use FLOATS as an automated to evaluate layer-sensitivity for a target  $d$  eliminating the need for heavy human works to tune the sparsity ratios as hyper-parameters.

Fig. 5 shows the density of non-zero channel present per layer for a model pruned via irregular and channel pruned

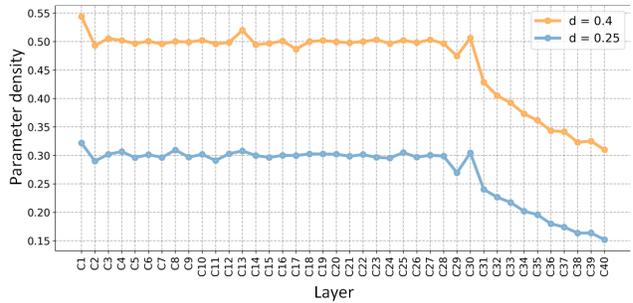


Figure 4. Per layer parameter density plot for two different target parameter density ( $d$ ) with irregular pruning. We used WRN40-2 on STL10 for training.

ing. Interestingly, as shown in the figure, irregular pruning keeps the channel density close to 1.0 for all the layers, while channel pruning can reduce the density to  $\sim 0.1$  for some layers. We attribute this large reduction of channels for some layers to the non-significant accuracy drop of the channel pruned models compared to the baseline.

### 2.1. FLOAT and FLOATS Slim Training Algorithms

Algorithm 1 and 2 describes the training algorithms of FLOAT and FLOATS slim respectively. In particular, FLOATS slim has an additional loop running over the set of SFs that the network can dynamically slim down to during inference.

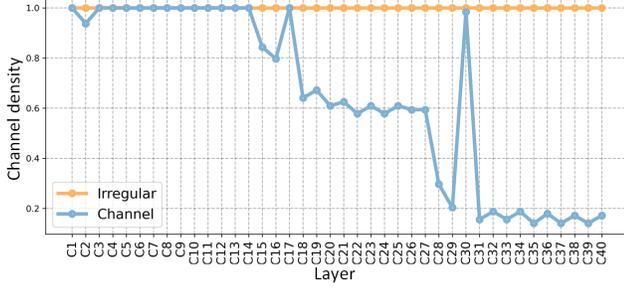


Figure 5. Per layer channel density for irregular vs channel pruning for the same target parameter density ( $d$ ) of 0.3. We used WRN40-2 on STL10 for training.

---

### Algorithm 1: FLOAT Algorithm

---

**Data:** Training set  $\mathbf{X} \sim D$ , model parameters  $\Theta$ , trainable noise scaling factor  $\alpha$ , binary conditioning parameter  $\lambda$ , mini-batch size  $\mathcal{B}$ .

- 1 **Output:** trained model parameters  $\Theta$ ,  $\alpha$ .
- 2 **for**  $i \leftarrow 0$  **to**  $ep$  **do**
- 3     **for**  $j \leftarrow 0$  **to**  $n_B$  **do**
- 4         Sample clean image-batch of size  $\mathcal{B}/2$  ( $\mathbf{X}_{0:\mathcal{B}/2}$ ,  $\mathbf{Y}_{0:\mathcal{B}/2}$ )  $\sim D$
- 5          $\mathcal{L}_C \leftarrow \text{computeLoss}(\mathbf{X}_{0:\mathcal{B}/2}, \Theta, \lambda = 0, \alpha; \mathbf{Y}_{0:\mathcal{B}/2})$  // condition to use weights w/o noise
- 6          $\hat{\mathbf{X}}_{\mathcal{B}/2:\mathcal{B}} \leftarrow \text{createAdv}(\mathbf{X}_{\mathcal{B}/2:\mathcal{B}}, \mathbf{Y}_{\mathcal{B}/2:\mathcal{B}})$  // adversarial image
- 7         // creation
- 8          $\mathcal{L}_A \leftarrow \text{computeLoss}(\hat{\mathbf{X}}_{\mathcal{B}/2:\mathcal{B}}, \Theta, \lambda = 1, \alpha; \mathbf{Y}_{\mathcal{B}/2:\mathcal{B}})$  // condition to use transformed weights
- 9          $\mathcal{L} \leftarrow 0.5 * \mathcal{L}_C + 0.5 * \mathcal{L}_A$
- 10          $\text{updateParam}(\Theta, \alpha, \nabla_{\mathcal{L}})$
- 11     **end**
- 12 **end**

---

### 3. Reproducibility Statement

We have detailed our training algorithms for both FLOAT and FLOATS slim in Supplementary Algorithm 1 and 2, respectively. We have further provided detailed hyperparameter and model information in the main manuscript and further additional details in the Supplementary section 1.1. Finally, the code for this project will be made publicly available upon acceptance of the paper.

### References

- [1] Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In *ICML*, 2020.
- [2] Xiaohan Ding, Xiangxin Zhou, Yuchen Guo, Jungong Han, Ji Liu, et al. Global sparse momentum sgd for pruning very deep

---

### Algorithm 2: FLOATS slim Algorithm

---

**Data:** Training set  $\mathbf{X} \sim D$ , model parameters  $\Theta$ , trainable noise scaling factor  $\alpha$ , binary conditioning parameter  $\lambda$ , mini-batch size  $\mathcal{B}$ , global parameter density  $d$ , initial mask  $\Pi$ , prune type (irregular/channel)  $t_p$ , slimming-factor set  $S_f$ .

- 1 **Output:** trained model parameters with density  $d$ .
- 2  $\Theta \leftarrow \text{applyMask}(\Theta, \Pi)$
- 3 **for**  $i \leftarrow 0$  **to**  $ep$  **do**
- 4     **for**  $j \leftarrow 0$  **to**  $n_B$  **do**
- 5          $\mathcal{B}/2$  ( $\mathbf{X}_{0:\mathcal{B}/2}, \mathbf{Y}_{0:\mathcal{B}/2}$ )  $\sim D$
- 6         **for** slimming-factor  $s_f$  in sorted  $S_f$  **do**
- 7              $\text{sampleSBN}(s_f)$  //sample the  $BN_C$  and  $BN_A$
- 8             // corresponding to  $w$
- 9              $\mathcal{L}_C \leftarrow \text{computeLoss}(\mathbf{X}_{0:\mathcal{B}/2}, \Theta, \lambda = 0, \alpha; \mathbf{Y}_{0:\mathcal{B}/2})$
- 10              $\hat{\mathbf{X}}_{\mathcal{B}/2:\mathcal{B}} \leftarrow \text{createAdv}(\mathbf{X}_{\mathcal{B}/2:\mathcal{B}}, \mathbf{Y}_{\mathcal{B}/2:\mathcal{B}})$
- 11              $\mathcal{L}_A \leftarrow \text{computeLoss}(\hat{\mathbf{X}}_{\mathcal{B}/2:\mathcal{B}}, \Theta, \lambda = 1, \alpha; \mathbf{Y}_{\mathcal{B}/2:\mathcal{B}})$
- 12              $\mathcal{L} \leftarrow 0.5 * \mathcal{L}_C + 0.5 * \mathcal{L}_A$
- 13              $\text{accumulateGrad}(\mathcal{L})$
- 14             **end**
- 15              $\text{updateParam}(\Theta, \alpha, \nabla_{\mathcal{L}}, \Pi)$
- 16         **end**
- 17          $\text{updateLayerMomentum}(\mu)$
- 18          $\text{pruneRegrow}(\Theta, \Pi, \mu, d)$  // Prune fixed % of active and
- 19         regrow fraction of inactive weights
- 20          $\Pi \leftarrow \text{updateMask}(\Pi, t_p, \mu)$
- 21 **end**

---

neural networks. *Advances in Neural Information Processing Systems*, 32, 2019.

- [3] Zhezhi He, Adnan Siraj Rakin, and Deliang Fan. Parametric noise injection: Trainable randomness to improve deep neural network robustness against adversarial attack. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 588–597, 2019.
- [4] Souvik Kundu, Mahdi Nazemi, Peter A Beerel, and Massoud Pedram. Dnr: A tunable robust pruning framework through dynamic network rewiring of dnns. In *Proceedings of the 26th Asia and South Pacific Design Automation Conference*, pages 344–350, 2021.
- [5] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016.
- [6] Haotao Wang, Tianlong Chen, Shupeng Gui, Ting-Kuei Hu, Ji Liu, and Zhangyang Wang. Once-for-all adversarial training: In-situ tradeoff between robustness and accuracy for free. *arXiv preprint arXiv:2010.11828*, 2020.