## A. Proof of the scale-invariant property

Without loss of generality, assume a two layers neural network $f$ and $\phi$ is a ReLU-based activation function.

$$f(\boldsymbol{\theta}_f, \boldsymbol{x}) = \boldsymbol{\theta}^{(2)}\phi(\boldsymbol{\theta}^{(1)}\boldsymbol{x}). \quad (10)$$

The corresponding scaled neural network $g$ is:

$$g(\boldsymbol{\theta}_g, \boldsymbol{x}) = \frac{1}{\beta}\boldsymbol{\theta}^{(2)}\phi(\beta\boldsymbol{\theta}^{(1)}\boldsymbol{x}), \quad (11)$$

where the non-negative $\beta$ is the scaling factor.

Suppose we calculate the MRC value of the first module $\boldsymbol{\theta}^{(1)}$ and $\frac{1}{\beta}\boldsymbol{\theta}^{(1)}$.

**Theorem A.1.** The rectified function $\phi(x) = \max(x, 0)$ is a homogeneous function where

$$\forall (z, \beta) \in \mathbb{R} \times \mathbb{R}^+, \ \phi(\beta z) = \beta\phi(z). \quad (12)$$

*Proof.*

$$\phi(\beta z) = \max(\beta z, 0) = \beta\max(z, 0) = \beta\phi(z). \quad (13)$$
$\square$

**Theorem A.2.** $\forall \boldsymbol{x}, f(\boldsymbol{\theta}_f, \boldsymbol{x}) \equiv g(\boldsymbol{\theta}_g, \boldsymbol{x})$.

*Proof.*

$$g(\boldsymbol{\theta}_g, \boldsymbol{x}) = \frac{1}{\beta}\boldsymbol{\theta}^{(2)}\phi(\beta\boldsymbol{\theta}^{(1)}\boldsymbol{x}) \quad (14)$$

$$\equiv \frac{1}{\beta}\beta\boldsymbol{\theta}^{(2)}\phi(\boldsymbol{\theta}^{(1)}\boldsymbol{x}) \quad (15)$$

$$\equiv \boldsymbol{\theta}^{(2)}\phi(\boldsymbol{\theta}^{(1)}\boldsymbol{x}) \quad (16)$$

$$\equiv f(\boldsymbol{\theta}_f, \boldsymbol{x}) \quad (17)$$
$\square$

**Theorem A.3.** The robust losses of $f$ and $g$ are equal:

$$\mathcal{R}(f(\boldsymbol{\theta}_f), \mathcal{D}) \equiv \mathcal{R}(g(\boldsymbol{\theta}_g), \mathcal{D}). \quad (18)$$

*Proof.* According to Theorem A.2,

$$\forall \boldsymbol{x} + \Delta\boldsymbol{x}, f(\boldsymbol{\theta}_f, \boldsymbol{x} + \Delta\boldsymbol{x}) \equiv g(\boldsymbol{\theta}_g, \boldsymbol{x} + \Delta\boldsymbol{x}). \quad (19)$$

Thus,

$$\max_{\Delta\boldsymbol{x}\in\mathcal{S}} \ell(f(\boldsymbol{\theta}_f, \boldsymbol{x} + \Delta\boldsymbol{x}), y) \quad (20)$$

$$\equiv \max_{\Delta\boldsymbol{x}\in\mathcal{S}} \ell(g(\boldsymbol{\theta}_g, \boldsymbol{x} + \Delta\boldsymbol{x}), y). \quad (21)$$

Thus,

$$\mathcal{R}(f(\boldsymbol{\theta}_f), \mathcal{D}) = \sum_{(\boldsymbol{x},y)\in\mathcal{D}} \max_{\Delta\boldsymbol{x}\in\mathcal{S}} \ell(f(\boldsymbol{\theta}_f, \boldsymbol{x} + \Delta\boldsymbol{x}), y) \quad (22)$$

$$\equiv \sum_{(\boldsymbol{x},y)\in\mathcal{D}} \max_{\Delta\boldsymbol{x}\in\mathcal{S}} \ell(g(\boldsymbol{\theta}_g, \boldsymbol{x} + \Delta\boldsymbol{x}), y) \quad (23)$$

$$= \mathcal{R}(g(\boldsymbol{\theta}_g), \mathcal{D}) \quad (24)$$
$\square$

**Theorem A.4.** The Module Robustness Criticality (MRC) proposed in Definition 3.1 is invariant to the scaling of the parameters.

*Proof.* Let $\Delta\boldsymbol{\theta}_f = \{\Delta\boldsymbol{\theta}_f^{(1)}, \boldsymbol{0}\}, \Delta\boldsymbol{\theta}_g = \{\Delta\boldsymbol{\theta}_g^{(1)}, \boldsymbol{0}\}$ be the perturbation of the first layer for network $f$ and $g$ respectively. First, we prove

$$\max_{\Delta\boldsymbol{\theta}_f\in\mathcal{C}_{\boldsymbol{\theta}_f}} \mathcal{R}(f(\boldsymbol{\theta}_f + \Delta\boldsymbol{\theta}_f), \mathcal{D}) \quad (25)$$

$$\leq \max_{\Delta\boldsymbol{\theta}_g\in\mathcal{C}_{\boldsymbol{\theta}_g}} \mathcal{R}(g(\boldsymbol{\theta}_g + \Delta\boldsymbol{\theta}_g), \mathcal{D}). \quad (26)$$

Let

$$\Delta\boldsymbol{\theta}_f^* = \arg\max_{\Delta\boldsymbol{\theta}_f\in\mathcal{C}_{\boldsymbol{\theta}_f}} \mathcal{R}(f(\boldsymbol{\theta}_f + \Delta\boldsymbol{\theta}_f), \mathcal{D}), \quad (27)$$

$$\Delta\boldsymbol{\theta}_g^* = \arg\max_{\Delta\boldsymbol{\theta}_g\in\mathcal{C}_{\boldsymbol{\theta}_g}} \mathcal{R}(g(\boldsymbol{\theta}_g + \Delta\boldsymbol{\theta}_g), \mathcal{D}). \quad (28)$$

Consider the perturbation $\Delta\tilde{\boldsymbol{\theta}}_g = \beta\Delta\boldsymbol{\theta}_f^*$ for $g$, it is easy to show that $\Delta\tilde{\boldsymbol{\theta}}_g \in \mathcal{C}_{\boldsymbol{\theta}_g}$,

$$\mathcal{C}_{\boldsymbol{\theta}_f} = \{\Delta\boldsymbol{\theta}_f \,|\, \|\Delta\boldsymbol{\theta}_f\|_p \leq \epsilon\|\boldsymbol{\theta}_f^{(1)}\|_p\}, \quad (29)$$

$$\mathcal{C}_{\boldsymbol{\theta}_g} = \{\Delta\boldsymbol{\theta}_g \,|\, \|\Delta\boldsymbol{\theta}_g\|_p \leq \epsilon\|\boldsymbol{\theta}_g^{(1)}\|_p\} \quad (30)$$

$$= \{\Delta\boldsymbol{\theta}_g \,|\, \Delta\boldsymbol{\theta}_g = \beta\|\Delta\boldsymbol{\theta}_f\|_p \leq \epsilon\beta\|\boldsymbol{\theta}_f^{(1)}\|_p\}. \quad (31)$$

Therefore,

$$\mathcal{R}(g(\boldsymbol{\theta}_g + \Delta\tilde{\boldsymbol{\theta}}_g), \mathcal{D}) \leq \max_{\Delta\boldsymbol{\theta}_g\in\mathcal{C}_{\boldsymbol{\theta}_g}} \mathcal{R}(g(\boldsymbol{\theta}_g + \Delta\boldsymbol{\theta}_g), \mathcal{D}) \quad (32)$$

$$= \mathcal{R}(g(\boldsymbol{\theta}_g + \Delta\boldsymbol{\theta}_g^*), \mathcal{D}). \quad (33)$$

Repeat the same analysis as presented in Theorem A.2,

$$g(\boldsymbol{\theta}_g + \Delta\tilde{\boldsymbol{\theta}}_g) \quad (34)$$

$$= \frac{1}{\beta}\boldsymbol{\theta}^{(2)}\phi((\beta\boldsymbol{\theta}^{(1)} + \beta\Delta\boldsymbol{\theta}_f^*)\boldsymbol{x}) \quad (35)$$

$$= \boldsymbol{\theta}^{(2)}\phi((\boldsymbol{\theta}^{(1)} + \Delta\boldsymbol{\theta}_f^*)\boldsymbol{x}) \quad (36)$$

$$\equiv f(\boldsymbol{\theta}_f + \Delta\boldsymbol{\theta}_f^*). \quad (37)$$

According to Theorem A.3,

$$\mathcal{R}(f(\boldsymbol{\theta}_f + \Delta\boldsymbol{\theta}_f^*), \mathcal{D}) \equiv \mathcal{R}(g(\boldsymbol{\theta}_g + \Delta\tilde{\boldsymbol{\theta}}_g), \mathcal{D}) \quad (38)$$

$$\leq \mathcal{R}(g(\boldsymbol{\theta}_g + \Delta\boldsymbol{\theta}_g^*), \mathcal{D}). \quad (39)$$

Similarly, we can prove

$$\max_{\Delta\boldsymbol{\theta}_g\in\mathcal{C}_{\boldsymbol{\theta}_g}} \mathcal{R}(g(\boldsymbol{\theta}_g + \Delta\boldsymbol{\theta}_g), \mathcal{D}) \quad (40)$$

$$\leq \max_{\Delta\boldsymbol{\theta}_f\in\mathcal{C}_{\boldsymbol{\theta}_f}} \mathcal{R}(f(\boldsymbol{\theta}_f + \Delta\boldsymbol{\theta}_f), \mathcal{D}). \quad (41)$$

Thus,

$$\max_{\Delta\boldsymbol{\theta}_f\in\mathcal{C}_{\boldsymbol{\theta}_f}} \mathcal{R}(f(\theta_f + \Delta\boldsymbol{\theta}_f), \mathcal{D}) \quad (42)$$

$$= \max_{\Delta\boldsymbol{\theta}_g\in C_{\boldsymbol{\theta}_g}} \mathcal{R}(g(\theta_g + \Delta\boldsymbol{\theta}_g), \mathcal{D}). \quad (43)$$

Such that the proof ends. $\square$

# B. Algorithm of RiFT

The complete algorithm of RiFT is presented in Algorithm 2.

---

**Algorithm 2** Robust Critical Fine-Tuning

---

**Input:** adversarially trained model weights $\boldsymbol{\theta}_{AT}$, standard dataset $\mathcal{D}_{std}$, weight perturbation scaling factor $\alpha$, fine-tuning optimization iteration steps $T$ and learning rate $\gamma$, weight decay facotr $\lambda$.

**Output:** The fine-tuned model weights $\boldsymbol{\theta}_{AT}^*$.

1: **Step 1**: Calculate MRC for each module
2: **for** Module weight $\boldsymbol{\theta}^{(j)}$ **do**
3:     Calculate MRC value of $\boldsymbol{\theta}^{(j)}$ using Algorithm 1.
4: **end for**
5: Select the module with lowest MRC value, denote as non-robust critical module $\boldsymbol{\theta}^{(i)}$
6: **Step 2**: Fine-tuning on Non-robust critical module
7: $\boldsymbol{\theta}_1 = \boldsymbol{\theta}_{AT}$
8: **for** $t = 1, \dots, T$ **do**       ▷ Fine-tuning $T$ epochs
9:     **for** Batch $\mathcal{B}_k \in \mathcal{D}_{std}$ **do**
10:         Calculate loss: $\mathcal{L}(f(\boldsymbol{\theta}_t), \mathcal{B}_k))$
11:         $\boldsymbol{\theta}_{t+1}^{(i)} = \boldsymbol{\theta}_{t+1}^{(i)} - \gamma\nabla_{\boldsymbol{\theta}_t}(\mathcal{L})$    ▷ Gradient Descent
12:     **end for**
13:     $\boldsymbol{\theta}_{FT} = \boldsymbol{\theta}_t$ if $\boldsymbol{\theta}_t$ obtain highest std test acc.
14: **end for**
15: **Step 3**: Interpolation
16: **for** $\alpha \in (0, 1, 0.05)$ **do**
17:     $\boldsymbol{\theta}_\alpha = (1-\alpha)\boldsymbol{\theta}_{AT} + \alpha\boldsymbol{\theta}_{FT}$
18:     $\boldsymbol{\theta}_{FT}^* = \boldsymbol{\theta}_\alpha$ if it reaches best standard test acc while preserve the robustness as $\boldsymbol{\theta}_{AT}$.
19: **end for**
20: **Return** Fine-tuned model weights $\boldsymbol{\theta}_{FT}^*$

---

# C. Training Details

## C.1. Experiment Environment

All experiments are conducted on a workstation equipped with an NVIDIA GeForce RTX 3090 GPU with 24GB memory and NVIDIA A100 with 80GB memory. The PyTorch version is 1.11.0.

## C.2. Adversarial Training Details

For vanilla adversarial training, We set the initial learning rate as 0.1, which decays at 100 and 105 epochs with factor 10. When generating adversarial examples, we set BN as train mode since it usually achieves higher robustness.

When incorporating RiFT with other adversarial training methods, the SCORE method is incorporated with TRADES. For the CIFAR100 training, we ran with three different learning rate and select the best model weights as the one with highest robust accuracy. The hyper-parameter settings are either based on their original paper or same as the vanilla AT, depends on which method achieves better robust accuracy.

## C.3. Fine-tuning Details

The hyper-parameter that most affects fine-tuning is the initial learning rate. According to our experience, we find a small learning rate usually performs better. If the adversarial robustness of the final fine-tuned weights is still higher than the robustness of the initial adversarial training, we then increase the learning rate.

## C.4. The MRC value of ResNet34 and WRN34-10

Figure C.1 and Figure C.2 shows the Module Robust Criticality (MRC) value of each module in ResNet34 trained on CIFAR100 and WideResNet34 trained on Tiny-ImageNet, respectively. It can be observed that both models exhibit redundant capacity. Additionally, Figure C.3 and Figure C.4 shows the MRC value of each module in ResNet18 trained on CIFAR100 and Tiny-ImageNet, respectively. As we discussed in Section 5.3 and Section 5.5, ResNet18 has a lower redundant capacity compared to ResNet34 and WideResNet34, and the redundant capacity decreases as the classification task becomes more complex.

## C.5. More interpolation results

Figure C.5 shows the interpolation results of different modules of ResNet18 trained on CIFAR100 dataset. It can be observed that fine-tuning on robust-critical module can also help improve generalization and robustness. This does not mean that our MRC is wrong, as we claimed in Section 5.2, fine-tuning on robust-critical module does not necessarily hurt robustness. The MRC provides guidance on which module to fine-tune for optimal results, and still, fine-tuning on non-robust-critical module achieves the highest test accuracy while preserving robustness.

# D. Analysis of the complexity of MRC algorithm

When identifying the most non-robust-critical module, it is required to iterate all modules of the model. Suppose a model with $n$ modules, for each module, the calculation complexity depends on the iteration steps in Algorithm 1. Considering the different overheads for each iterative computation of the modules at different locations, for example, when calculating the last module's MRC value, it only requires forward-backward iteration of the last layer of parameters. Thus, the average total forward-backward iteration of each module is $n/2$. In our experiments, we set the learning rate as 1 and the iteration step as 10. Thus, in our experiments, the complexity of MRC algorithm cost $5n$ total forward-backward propagation.
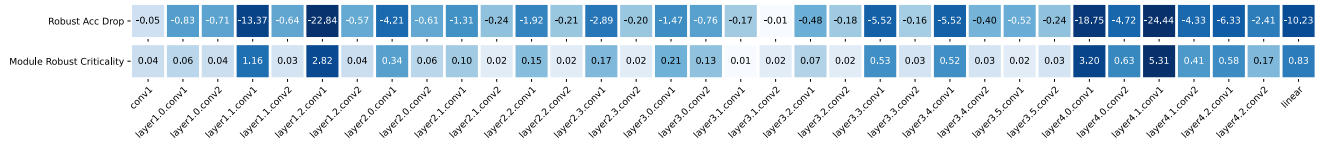
Figure C.1. Example of module robust criticality (MRC) and its corresponding robust accuracy drop of ResNet34 trained on CIFAR100.
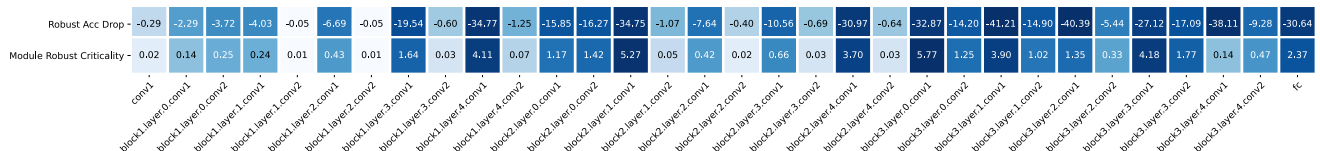
| | conv1 | layer1.0.conv1 | layer1.0.conv2 | layer1.1.conv1 | layer1.1.conv2 | layer1.2.conv1 | layer1.2.conv2 | layer2.0.conv1 | layer2.0.conv2 | layer2.1.conv1 | layer2.1.conv2 | layer2.2.conv1 | layer2.2.conv2 | layer2.3.conv1 | layer2.3.conv2 | layer3.0.conv1 | layer3.0.conv2 | layer3.1.conv1 | layer3.1.conv2 | layer3.2.conv1 | layer3.2.conv2 | layer3.3.conv1 | layer3.3.conv2 | layer3.4.conv1 | layer3.4.conv2 | layer3.5.conv1 | layer3.5.conv2 | layer4.0.conv1 | layer4.0.conv2 | layer4.1.conv1 | layer4.1.conv2 | layer4.2.conv1 | layer4.2.conv2 | linear |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Robust Acc Drop | -0.05 | -0.83 | -0.71 | -13.37 | -0.64 | -22.84 | -0.57 | -4.21 | -0.61 | -1.31 | -0.24 | -1.92 | -0.21 | -2.89 | -0.20 | -1.47 | -0.76 | -0.17 | -0.01 | -0.48 | -0.18 | -5.52 | -0.16 | -5.52 | -0.40 | -0.52 | -0.24 | -18.75 | -4.72 | -24.44 | -4.33 | -6.33 | -2.41 | -10.23 |
| Module Robust Criticality | 0.04 | 0.06 | 0.04 | 1.16 | 0.03 | 2.82 | 0.04 | 0.34 | 0.06 | 0.10 | 0.02 | 0.15 | 0.02 | 0.17 | 0.02 | 0.21 | 0.13 | 0.01 | 0.02 | 0.07 | 0.02 | 0.53 | 0.03 | 0.52 | 0.03 | 0.02 | 0.03 | 3.20 | 0.63 | 5.31 | 0.41 | 0.58 | 0.17 | 0.83 |



Figure C.2. Example of module robust criticality (MRC) and its corresponding robust accuracy drop of WideResNet34 trained on Tiny-ImageNet.

| | conv1 | block1.layer0.conv1 | block1.layer0.conv2 | block1.layer1.conv1 | block1.layer1.conv2 | block1.layer2.conv1 | block1.layer2.conv2 | block1.layer3.conv1 | block1.layer3.conv2 | block1.layer4.conv1 | block1.layer4.conv2 | block2.layer0.conv1 | block2.layer0.conv2 | block2.layer1.conv1 | block2.layer1.conv2 | block2.layer2.conv1 | block2.layer2.conv2 | block2.layer3.conv1 | block2.layer3.conv2 | block2.layer4.conv1 | block2.layer4.conv2 | block3.layer0.conv1 | block3.layer0.conv2 | block3.layer1.conv1 | block3.layer1.conv2 | block3.layer2.conv1 | block3.layer2.conv2 | block3.layer3.conv1 | block3.layer3.conv2 | block3.layer4.conv1 | block3.layer4.conv2 | fc |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Robust Acc Drop | -0.29 | -2.29 | -3.72 | -4.03 | -0.05 | -6.69 | -0.05 | -19.54 | -0.60 | -34.77 | -1.25 | -15.85 | -16.27 | -34.75 | -1.07 | -7.64 | -0.40 | -10.56 | -0.69 | -30.97 | -0.64 | -32.87 | -14.20 | -41.21 | -14.90 | -40.39 | -5.44 | -27.12 | -17.09 | -38.11 | -9.28 | -30.64 |
| Module Robust Criticality | 0.02 | 0.14 | 0.25 | 0.24 | 0.01 | 0.43 | 0.01 | 1.64 | 0.03 | 4.11 | 0.07 | 1.17 | 1.42 | 5.27 | 0.05 | 0.42 | 0.02 | 0.66 | 0.03 | 3.70 | 0.03 | 5.77 | 1.25 | 3.90 | 1.02 | 1.35 | 0.33 | 4.18 | 1.77 | 0.14 | 0.47 | 2.37 |



Figure C.3. Example of module robust criticality (MRC) and its corresponding robust accuracy drop of ResNet18 trained on CIFAR100.

| | conv1 | layer1.0.conv1 | layer1.0.conv2 | layer1.1.conv1 | layer1.1.conv2 | layer2.0.conv1 | layer2.0.conv2 | layer2.1.conv1 | layer2.1.conv2 | layer3.0.conv1 | layer3.0.conv2 | layer3.1.conv1 | layer3.1.conv2 | layer4.0.conv1 | layer4.0.conv2 | layer4.1.conv1 | layer4.1.conv2 | linear |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Robust Acc Drop | -1.13 | -23.00 | -14.87 | -29.10 | -23.87 | -27.67 | -24.41 | -29.09 | -17.45 | -28.35 | -27.03 | -29.09 | -12.01 | -29.02 | -28.94 | -29.09 | -25.51 | -26.45 |
| Module Robust Criticality | 0.09 | 2.31 | 1.30 | 6.53 | 2.92 | 6.06 | 3.66 | 10.54 | 1.61 | 6.53 | 13.06 | 21.28 | 1.06 | 21.33 | 24.24 | 41.72 | 4.99 | 3.26 |



Figure C.4. Example of module robust criticality (MRC) and its corresponding robust accuracy drop of ResNet18 trained on Tiny-ImageNet.
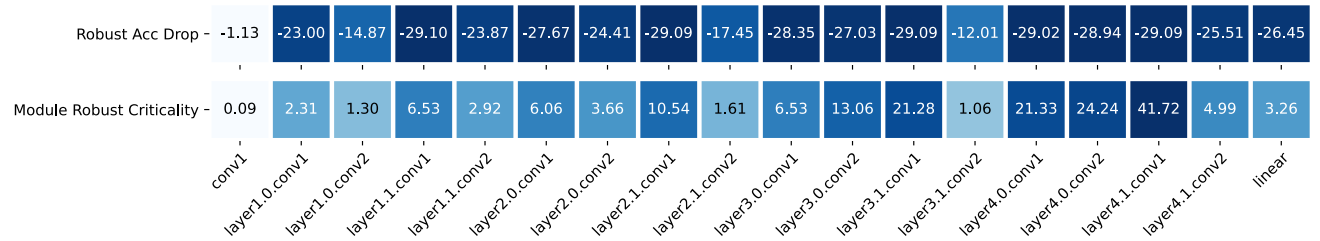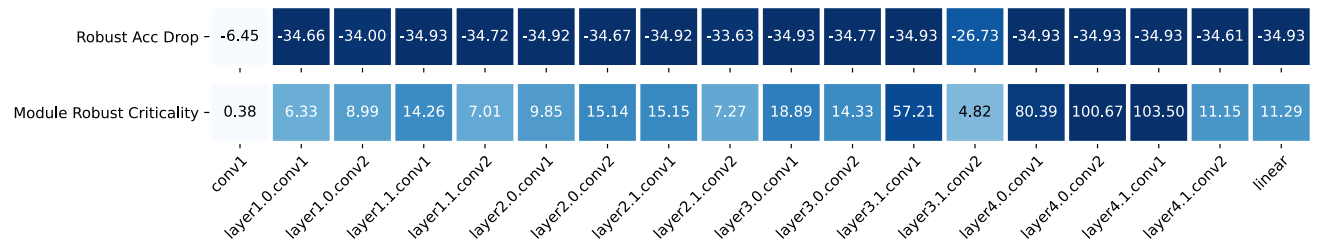
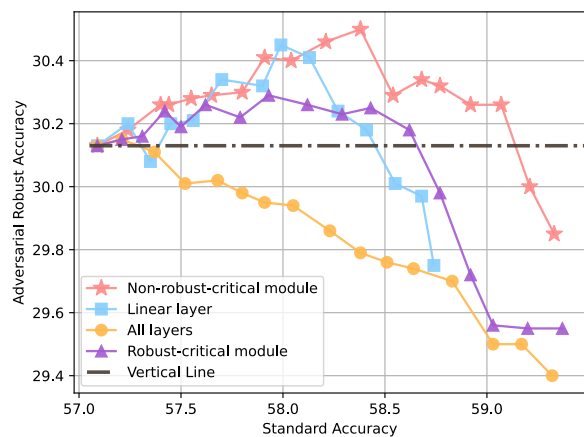| | conv1 | layer1.0.conv1 | layer1.0.conv2 | layer1.1.conv1 | layer1.1.conv2 | layer2.0.conv1 | layer2.0.conv2 | layer2.1.conv1 | layer2.1.conv2 | layer3.0.conv1 | layer3.0.conv2 | layer3.1.conv1 | layer3.1.conv2 | layer4.0.conv1 | layer4.0.conv2 | layer4.1.conv1 | layer4.1.conv2 | linear |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Robust Acc Drop | -6.45 | -34.66 | -34.00 | -34.93 | -34.72 | -34.92 | -34.67 | -34.92 | -33.63 | -34.93 | -34.77 | -34.93 | -26.73 | -34.93 | -34.93 | -34.93 | -34.61 | -34.93 |
| Module Robust Criticality | 0.38 | 6.33 | 8.99 | 14.26 | 7.01 | 9.85 | 15.14 | 15.15 | 7.27 | 18.89 | 14.33 | 57.21 | 4.82 | 80.39 | 100.67 | 103.50 | 11.15 | 11.29 |

Figure C.5. Interpolation results of fine-tuning on different modules of ResNet18 on CIFAR100 dataset. Dots denote different interpolation points between the final fine-tuned weights of RiFT and the initial adversarially trained weights.