

# The Implicit Bias of Gradient Descent toward Collaboration between Layers: A Dynamic Analysis of Multilayer Perceptions

*Zheng Wang<sup>1</sup>, Geyong Min<sup>1</sup>, Wenjie Ruan<sup>2</sup>*

*<sup>1</sup>University of Exeter    <sup>2</sup>University of Science and Technology of China*



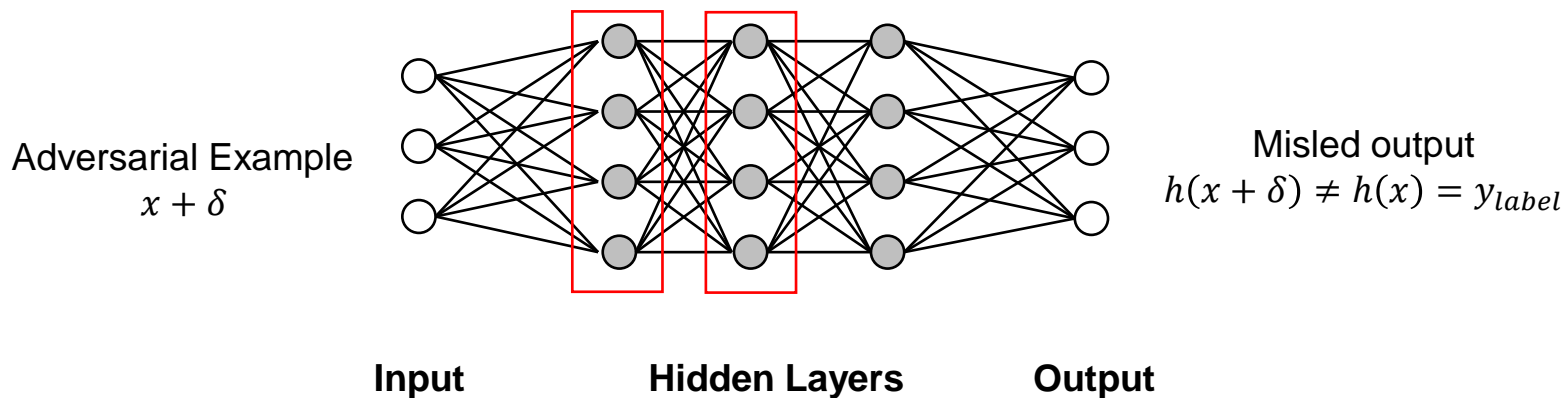
University  
of Exeter



*[github.com/squarewang2077/co-correlation](https://github.com/squarewang2077/co-correlation)*

# Motivation

***Whether layers in neural networks collaborate to strengthen adversarial robustness during gradient descent?***



# Problem Setting

- **Binary classification**
- **Training Set**  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$

## Neural Networks

$$f_{linear}(\mathbf{x}, W) = \mathbf{a}^T(W\mathbf{x})$$

$$f_{mlp}(\mathbf{x}, W) = \mathbf{a}^T \sigma(W\mathbf{x})$$

Activation function

Sigmoid function

$$\text{sig}(x) = 1/(1 + e^{-x})$$

Outputs

$$u_i = f(x_i, W), i \in [n]$$

Prediction

$$y_{pred} = \begin{cases} 1, & \text{sig}(f_W(\mathbf{x})) > 0.5 \\ 0, & \text{sig}(f_W(\mathbf{x})) \leq 0.5 \end{cases}$$

Loss function

$$L(f, y) = -\frac{1}{n} \sum_{i=1}^n \left[ y_i \log(\text{sig}(u_i)) + (1 - y_i) \log(1 - \text{sig}(u_i)) \right]$$

# Measure Adversarial Risk by Dirichlet Energy

**Theorem 4.1.** Given data points  $(\mathbf{x}, y) \sim P$  and  $\mathbf{x} \sim P_{\mathbf{x}}$ , the relationship between adversarial risk and Dirichlet energy for classifier  $f$  with differentiable loss function  $L$  is shown as

$$R^{rob}(f, r) \lesssim R(f) + r \mathfrak{S}(L(f)), \quad (8)$$

where  $r > 0$  is the largest perturbation budget and  $\mathfrak{S}(L(f)) = \sqrt{\mathbb{E}_{\mathbf{x} \sim P_{\mathbf{x}}} [\|\nabla_{\mathbf{f}} L^T \cdot J_{\mathbf{f}}(\mathbf{x})\|_2^2]}$  indicating the Dirichlet energy of the classifier on loss  $L$ .

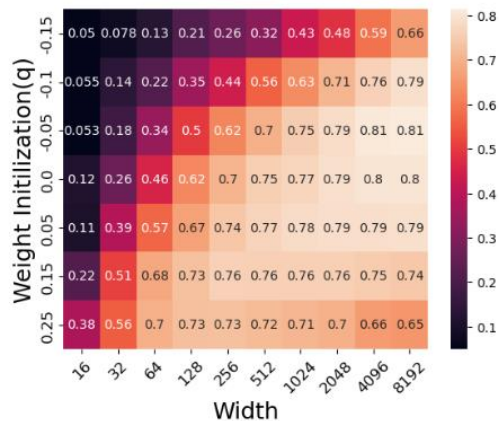
$$R^{rob}(f, r) = \mathbb{E}_{(\mathbf{x}, y) \sim P} \left[ \sup_{\boldsymbol{\varepsilon} \in B_r} L(f(\mathbf{x} + \boldsymbol{\varepsilon}), y) \right]$$

$$R(f) = \mathbb{E}_{(\mathbf{x}, y) \sim P} [L(f(\mathbf{x}), y)]$$

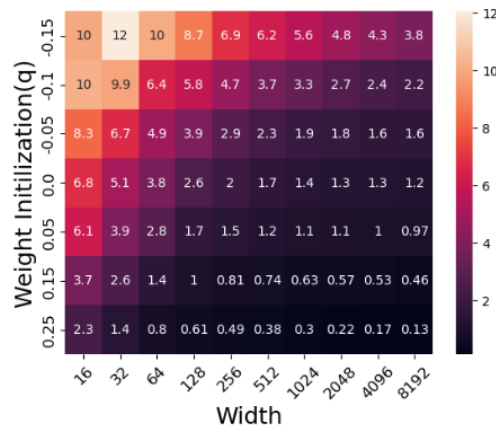
$$B_r = \{\|\boldsymbol{\varepsilon}\|_2 \leq r\}$$

▪ The proof is based on 1<sup>st</sup> order Taylor's expansion

# Measure Adversarial Risk by Dirichlet Energy



(a) Robust Acc.



(b)  $\mathcal{S}(f)$

- 2-Layer MLPs with width from  $2^4$  to  $2^{13}$
- Initialize the weight matrix  $w_{\{i,j\}} \sim N\left(0, \frac{1}{m^{1+2q}}\right)$
- Let  $q$  change from  $-0.15$  to  $0.25$
- Dirichlet Energy of  $f$  can be a good representation
- It can measure individual layers therefore the correlations

**Theorem 4.5** (Robustness Decomposition). *Given the same assumption in Definition 4.2, the measurement for overall adversarial robustness can be decomposed as*

$$\begin{aligned} \mathfrak{S}(\phi \circ \varphi) &= \left( \mathbb{E}_{\mathbf{x} \sim P} [\|J_{\phi \circ \varphi}(\mathbf{x})\|_2^2] \right)^{\frac{1}{2}} \\ &= \varrho_{\phi, \varphi} \left( 1 + \frac{\text{var}_{\phi, \varphi}}{\mu_{\phi, \varphi}^2} \right)^{\frac{1}{2}} \rho_{\phi, \varphi} \mathfrak{S}(\phi) \mathfrak{S}(\varphi) \end{aligned} \quad (15)$$

$$f_{\text{linear}}(\mathbf{x}, W) = \mathbf{a}^T(W\mathbf{x})$$

$$f_{\text{mlp}}(\mathbf{x}, W) = \mathbf{a}^T(\sigma(W\mathbf{x}))$$

Co-correlation

$$\varrho_{\phi, \varphi} \triangleq \frac{\left( \mathbb{E}_{\mathbf{x} \sim P_{\mathbf{x}}} [\|J_{\phi \circ \varphi}(\mathbf{x})\|_2^2] \right)^{\frac{1}{2}}}{\left( \mathbb{E}_{\mathbf{x} \sim P_{\mathbf{x}}} [\|J_{\phi}(\varphi)\|_2^2 \cdot \|J_{\varphi}(\mathbf{x})\|_2^2] \right)^{\frac{1}{2}}}$$

$$\text{var}_{\phi, \varphi} \triangleq \text{Var}_{\mathbf{x} \sim P_{\mathbf{x}}} [\|J_{\phi}(\varphi)\|_2 \cdot \|J_{\varphi}(\mathbf{x})\|_2]$$

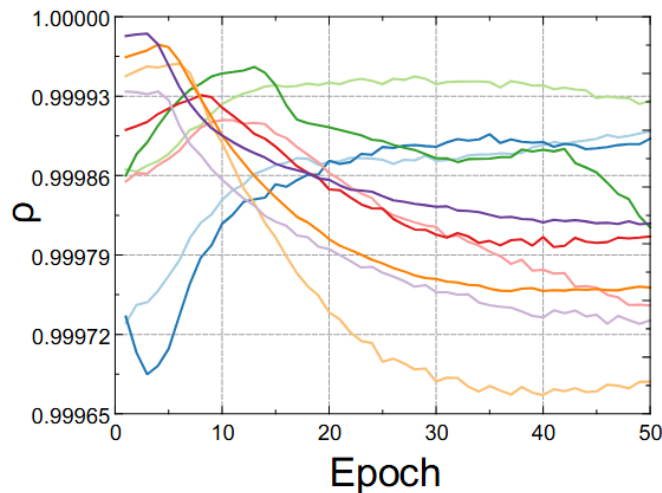
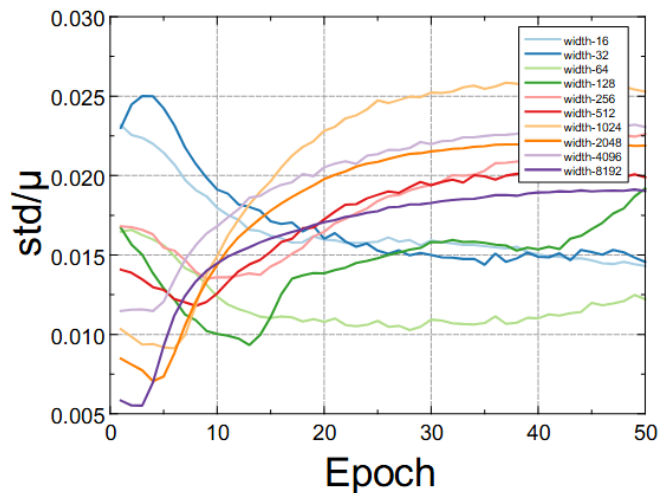
$$\rho_{\phi, \varphi} \triangleq \frac{\mathbb{E}_{\mathbf{x} \sim P_{\mathbf{x}}} [\|J_{\phi}(\varphi)\|_2 \cdot \|J_{\varphi}(\mathbf{x})\|_2]}{\left( \mathbb{E}_{\varphi \sim \varphi(\mathbf{x})} [\|J_{\phi}(\varphi)\|_2^2] \right)^{\frac{1}{2}} \left( \mathbb{E}_{\mathbf{x} \sim P_{\mathbf{x}}} [\|J_{\varphi}(\mathbf{x})\|_2^2] \right)^{\frac{1}{2}}}$$

$$\mu_{\phi, \varphi} \triangleq \mathbb{E}_{\mathbf{x} \sim P_{\mathbf{x}}} [\|J_{\phi}(\varphi)\|_2 \cdot \|J_{\varphi}(\mathbf{x})\|_2]$$

- Regard the *neural network* as **function composition**
- The proof is straight forward

# On Dynamics of Co-Correlation

## *Robustness Decomposition*



- Empirically, co-correlation is more influential
- We focus on the **co-correlation**

# On Dynamics of Co-Correlation

## Linear Model

**Assumption 5.1.** We assume that each element  $w_{i,j}$  in the weight matrix  $W(0) \in \mathbb{R}^{m \times d}$  at initialization follows the Gaussian distribution  $N(0, \frac{1}{m^{1+2q}})$ , with  $q > 0$ . Additionally, each element  $a_r, r \in [m]$  in  $\mathbf{a}$  is randomly selected from the set  $\{-\frac{1}{\sqrt{m}}, \frac{1}{\sqrt{m}}\}$ , and fixed during training.

**Assumption 5.2.** We assume that for each  $(\mathbf{x}_i, y_i) \in D, i \in [n]$ ,  $\mathbf{x}_i$  is  $L_2$  norm bounded such that  $\|\mathbf{x}_i\|_2 = 1$  for all  $i \in [n]$ .

**Theorem 5.3** (Dynamics of the Co-correlation for Linear Model). *Given the linear model defined in Equation (3a) and training dataset  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ . Assume that assumptions 5.1 and 5.2 hold for  $W$  and  $\mathbf{a}$ . The gradient descent applied to the weights results in the dynamics of the co-correlation being expressed as:*

$$\dot{\varrho}_{\mathbf{a}, W}(t) = \eta C(t) \varrho_{\mathbf{a}, W}, \quad (18)$$

and with high probability,

$$C(t) \geq \frac{\sum_{\tau=1}^t \tilde{\mathbf{x}}(\tau)^T \tilde{\mathbf{x}}(t)}{\|W(t)\|_2^2} \cdot \left(1 - (\mathbf{v}(t)^T \mathbf{a})^2\right) + \mathcal{O}\left(\frac{1}{m^q}\right) \quad (19)$$

where the  $\mathbf{v}(t)$  is the dominate eigenvector for  $W(t)W(t)^T$ .

When  $m$  is sufficiently large, and during the initial steps of the optimization process,  $\tilde{\mathbf{x}}(\tau), \tau \in [t]$  are quite similar to each other in terms of cosine similarity, implying an acute angle to each other, which leads to  $\sum_{\tau=1}^t \tilde{\mathbf{x}}(\tau)^T \tilde{\mathbf{x}}(t) \geq 0$ . As a result, we can conclude that  $C(t) \geq 0$ .

- $\varrho_{\mathbf{a}, W}$  increase during the initial stages and become saturated to its later stages.
- The speed of the accumulation of  $\varrho_{\mathbf{a}, W}$  is inversely related to  $\|W(t)\|_2$

$$\tilde{\mathbf{x}}(t) = \frac{1}{n} \sum_{i=1}^n \left[ y_i - \text{sig}(u_i(t)) \right] \mathbf{x}_i$$



# On Dynamics of Co-Correlation

## 2-Layer MLP

**Assumption 5.4.** The derivative of the activation function  $\sigma'(x)$  in non-linear neural networks is bounded by  $M$ . In other words, we have  $|\sigma'(x)| \leq M$ .

**Theorem 5.5.** (Dynamics of the Co-correlation for MLP) Given the MLP defined in Equation (3) with training dataset  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ ,  $\mathbf{x} \in \mathcal{X}$  such that  $\mathbf{x} \sim P_{\mathbf{x}}$ . Assume that Assumption 5.1 and 5.2 hold for  $W$  and  $\mathbf{a}$ , and Assumption 5.4 holds for the activation function. we have

$$\dot{\varrho}_{\mathbf{a}, \sigma \circ W}(t) = \eta C(t) \varrho_{\mathbf{a}, \sigma \circ W}(t).$$

With high probability,

$$C(t) \geq \frac{\sum_{\tau=1}^t (1 - \mathbf{a}^T \mathbf{v}(\tau) \mathbf{a}^T \mathbf{v}(t)) \mathbb{E}_{\mathbf{x} \sim P_{\mathbf{x}}} [\tilde{\mathbf{x}}_*^T(\tau) \tilde{\mathbf{x}}_*(t)]}{\mathbb{E}_{\mathbf{x} \sim P_{\mathbf{x}}} \|D(t)W(t)\|_2^2} + \max \left\{ \mathcal{O}\left(\frac{1}{\sqrt{m}}\right), \mathcal{O}\left(\frac{1}{m^q}\right) \right\},$$

where

$$D(t) = \text{diag}(\sigma'(\mathbf{w}_1(t)^T \mathbf{x}), \dots, \sigma'(\mathbf{w}_m(t)^T \mathbf{x})),$$

and  $\mathbf{v}(t)$  denotes the dominant eigenvector for  $W(t)W(t)^T$ , with  $\tilde{\mathbf{x}}_*^T$  is defined in Equation (21). Similar to the Theorem 5.3 when  $m$  is sufficiently large, and during the initial steps of the optimization where the error-weighted inputs  $\tilde{\mathbf{x}}_*^T(\tau)$ ,  $\tau \in [t]$  do not significantly fluctuate, we have that  $C(t) \geq 0$ .

- The dynamics for  $\varrho_{\mathbf{a}, \sigma \circ W}$  is the same to  $\varrho_{\mathbf{a}, W}$
- The speed of the accumulation of  $\varrho_{\mathbf{a}, \sigma \circ W}$  is inversely related to  $\|D(t)W(t)\|_2$

Serve as similar purpose of  $\tilde{\mathbf{x}}(t)$

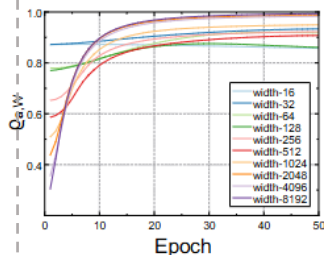
$$\alpha_i(t, \mathbf{x}) \triangleq \mathbb{E}_{W(0)} \left[ \sigma'(\mathbf{w}(t)^T \mathbf{x}) \sigma'(\mathbf{w}(t)^T \mathbf{x}_i) \right]$$

$$\tilde{\mathbf{x}}_*(t) \triangleq \frac{1}{n} \sum_{i=1}^n \alpha_i(t, \mathbf{x}) (y_i - \text{sig}(u_i(t))) \mathbf{x}_i$$

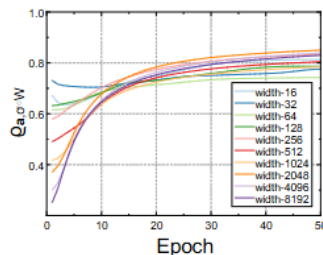
# Experiments

## 2-Layer MLP

### Different width



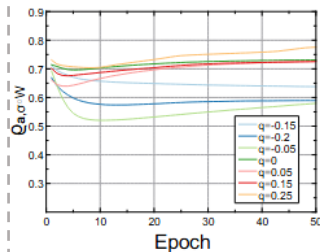
(c) Linear Model



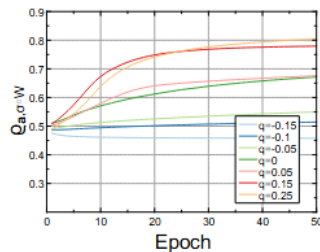
(d) MLP

- The upward trend is true for all 2-Layer MLPs
- The theorem is quite tight on  $q$

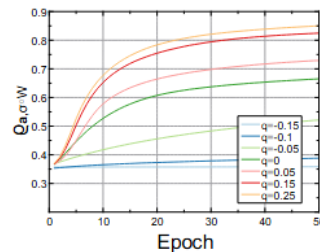
### Different weight initialization parameter $q$



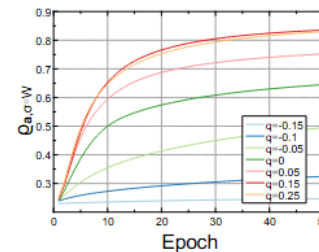
(a) width-32



(b) width-512



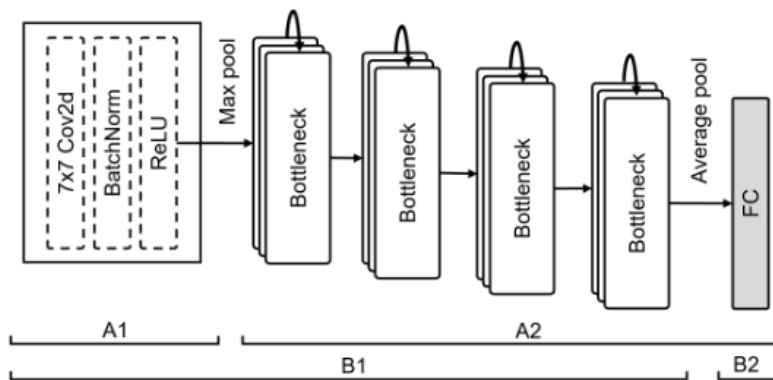
(c) width-2048



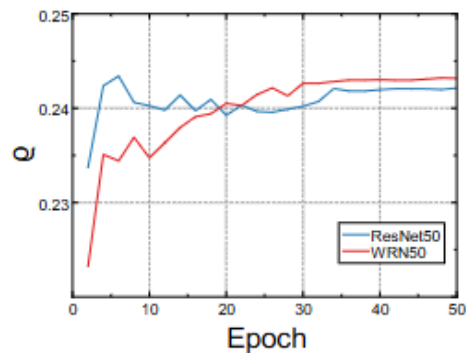
(d) width-8192

# Experiments

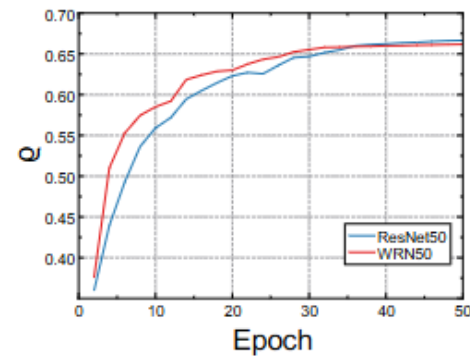
## ResNets



(a) ResNets



(b) A1-A2

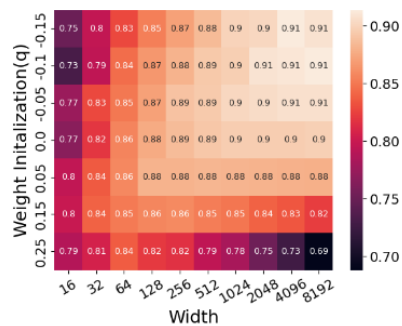


(c) B1-B2

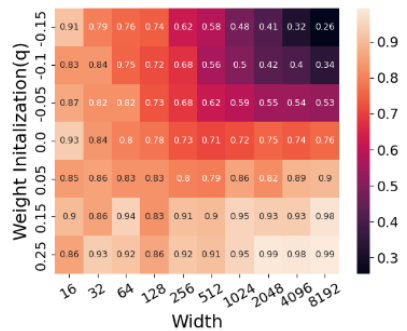
- Divide the ResNet50 and WRN in 2 ways
- w/o specific weight initialization
- On CIFAR10 with Adam optimizer

# Experiments

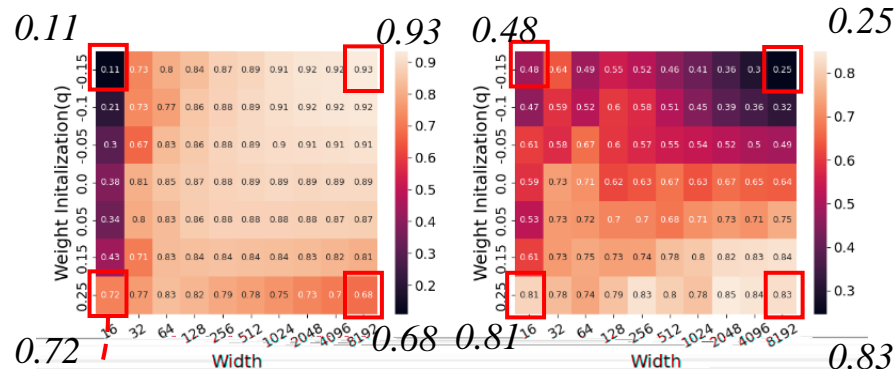
*Different behavior for wide and narrow MLPs*



(a) Linear Model/Acc.



(b) Linear Model/ $q$



(c) MLP ReLU/Acc.

(d) MLP ReLU/ $q$

- *Difference in behavior between wider and narrower neural networks*

# Conclusion

- ❑ By quantifying the interactions between layers, we found that it not only fails to collaborate against adversarial perturbations but may even hinder resistance to them during gradient descent.
- ❑ Wider MLPs exhibit more resistance to increased co-correlation and, therefore, are more adversarial robust.
- ❑ Future research can expand upon this by examining the effects of increased network depth and more sophisticated structures on the observed phenomena.