

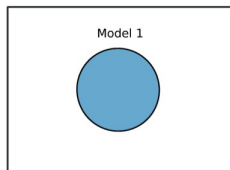
# LOTOS: LAYER-WISE ORTHOGONALIZATION FOR TRAINING ROBUST ENSEMBLES

**A. Ebrahimpour-Boroojeny**, H. Sundaram, and V. Chandrasekaran

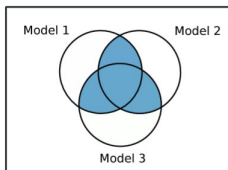
University of Illinois at Urbana-Champaign



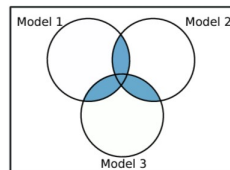
- Ensembles of models is known to be more robust to adversarial attacks.
  - It is harder to fool multiple models!



(a) Single Model



(b) Ensemble of 3 Models



(c) Diverse Ensemble

- The adversarial examples are transferable!
  - ▶ this reduces the effectiveness of using multiple models
  - ▶ introducing diversity among the models helps to counteract this effect
  - ▶ this is the basis of methods for training robust ensembles

---

<sup>1</sup> Kariyappa, S., & Qureshi, M. K. (2019). Improving adversarial robustness of ensembles with diversity training. arXiv preprint arXiv:1901.09981.

- For model  $\mathcal{F} : \mathcal{X} \rightarrow \mathcal{Y}$ , untargeted attack  $\mathcal{A}(x) = x + \delta_x$ , and a sample  $(x, y)$ ,  $\mathcal{A}$  maximizes  $\ell_{\mathcal{F}}(x + \delta_x, y)$ , s.t.  $\|\delta_x\|_2 \leq \epsilon$ .

- For model  $\mathcal{F} : \mathcal{X} \rightarrow \mathcal{Y}$ , untargeted attack  $\mathcal{A}(x) = x + \delta_x$ , and a sample  $(x, y)$ ,  $\mathcal{A}$  maximizes  $\ell_{\mathcal{F}}(x + \delta_x, y)$ , s.t.  $\|\delta_x\|_2 \leq \epsilon$ .
- Given an untargeted attack  $\mathcal{A}(x)$  on a surrogate model  $\mathcal{F}$ , transferability between  $\mathcal{F}$  and the target model  $\mathcal{G}$  is defined as:
  - Consider the samples for which:
    - 1  $\mathcal{F}$  makes correct prediction,
    - 2  $\mathcal{G}$  makes correct prediction,
    - 3  $\mathcal{F}$  makes wrong prediction on  $\mathcal{A}(x)$
  - And compute the ratio of them on which  $\mathcal{G}$  makes wrong prediction on  $\mathcal{A}(x)$  as well.

- For model  $\mathcal{F} : \mathcal{X} \rightarrow \mathcal{Y}$ , untargeted attack  $\mathcal{A}(x) = x + \delta_x$ , and a sample  $(x, y)$ ,  $\mathcal{A}$  maximizes  $\ell_{\mathcal{F}}(x + \delta_x, y)$ , s.t.  $\|\delta_x\|_2 \leq \epsilon$ .
- Given an untargeted attack  $\mathcal{A}(x)$  on a surrogate model  $\mathcal{F}$ , transferability between  $\mathcal{F}$  and the target model  $\mathcal{G}$  is defined as:
  - ▶ Consider the samples for which:
    - 1  $\mathcal{F}$  makes correct prediction,
    - 2  $\mathcal{G}$  makes correct prediction,
    - 3  $\mathcal{F}$  makes wrong prediction on  $\mathcal{A}(x)$
  - ▶ And compute the ratio of them on which  $\mathcal{G}$  makes wrong prediction on  $\mathcal{A}(x)$  as well.
- More specifically:

$$T_{rate}(\mathcal{A}_{\mathcal{F}}, \mathcal{F}, \mathcal{G}) = \mathbb{P}_{(x,y) \in \mathcal{X} \times \mathcal{Y}} [\mathcal{G}(\mathcal{A}_{\mathcal{F}}(x)) \neq y \mid \mathcal{F}(x) = \mathcal{G}(x) = y \\ \wedge \mathcal{F}(\mathcal{A}_{\mathcal{F}}(x)) \neq y],$$

- Function  $f$  is  $L$ -Lipschitz if:

$$\|f(x_1) - f(x_2)\|_2 \leq L\|x_1 - x_2\|_2.$$

- Consider a model

$$f(x, \theta) = W^{L+1} \sigma_L(W^L(\sigma_{L-1}(W^{L-1}(\dots \sigma_1(W^1 x) \dots))),$$

- ▶ Then:

$$\|f\|_{\text{Lip}} \leq \|(W^{L+1})\|_{\text{Lip}} \cdot \|\sigma_L\|_{\text{Lip}} \cdot \|(W^L)\|_{\text{Lip}} \cdots \|\sigma_1\|_{\text{Lip}} \cdot \|(W^1)\|_{\text{Lip}}$$

Controlling the Lipschitz factor of a model makes it more robust. <sup>[2]</sup>

We investigated how it affects the transferability rate among models!

---

<sup>2</sup> Ebrahimpour Boroojeny, A., Telgarsky, M., & Sundaram, H. (2024, April). Spectrum extraction and clipping for implicitly linear layers. In International Conference on Artificial Intelligence and Statistics (pp. 2971-2979). PMLR.

Controlling the Lipschitz factor of a model makes it more robust. <sup>[2]</sup>

We investigated how it affects the transferability rate among models!

## Proposition

*If  $\ell_{\mathcal{F}}$  and  $\ell_{\mathcal{G}}$  are  $L$ -Lipschitz w.r.t the inputs, and  $\|\delta\|_2 \leq r$ , for attack  $\mathcal{A}$ :*

$$|R_{\mathcal{F}}(A(x), y) - R_{\mathcal{G}}(A(x), y)| \leq 2Lr + |R_{\mathcal{F}}(x, y) - R_{\mathcal{G}}(x, y)|.$$

---

<sup>2</sup> Ebrahimpour Boroojeny, A., Telgarsky, M., & Sundaram, H. (2024, April). Spectrum extraction and clipping for implicitly linear layers. In International Conference on Artificial Intelligence and Statistics (pp. 2971-2979). PMLR.



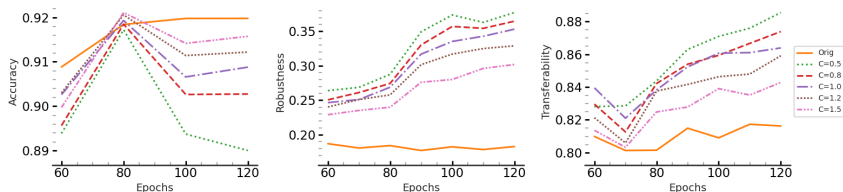
Controlling the Lipschitz factor of a model makes it more robust. [2]

We investigated how it affects the transferability rate among models!

## Proposition

If  $\ell_{\mathcal{F}}$  and  $\ell_{\mathcal{G}}$  are  $L$ -Lipschitz w.r.t the inputs, and  $\|\delta\|_2 \leq r$ , for attack  $\mathcal{A}$ :

$$|R_{\mathcal{F}}(A(x), y) - R_{\mathcal{G}}(A(x), y)| \leq 2Lr + |R_{\mathcal{F}}(x, y) - R_{\mathcal{G}}(x, y)|.$$

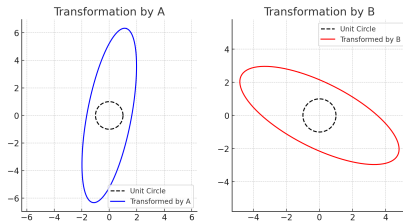


<sup>2</sup> Ebrahimpour Boroojeny, A., Telgarsky, M., & Sundaram, H. (2024, April). Spectrum extraction and clipping for implicitly linear layers. In International Conference on Artificial Intelligence and Statistics (pp. 2971-2979). PMLR.



- How to preserve the robustness of individual models *and* increase the diversity?
- We consider the layers-wise transformations of linear layers!

- How to preserve the robustness of individual models *and* increase the diversity?
- We consider the layers-wise transformations of linear layers!
  - ▶ Consider two corresponding layers from the two models in the ensemble:
  - ▶ The direction along the largest singular vector of each one should cause small changes in the other!



- We define a new notion of similarity between the  $j$ th layers of models  $f$  and  $g$ , based on their top- $k$  sub-spaces:

$$S_k^{(j)}(f^{(j)}, g^{(j)}, \text{mal}) := \sum_{i=1}^k w_i (\text{ReLU}(\|f^{(j)}(v'_i)\|_2 - \text{mal}) + \text{ReLU}(\|g^{(j)}(v_i)\|_2 - \text{mal})),$$

- where `mal` is a threshold hyperparameter

- We define a new notion of similarity between the  $j$ th layers of models  $f$  and  $g$ , based on their top- $k$  sub-spaces:

$$S_k^{(j)}(f^{(j)}, g^{(j)}, \text{mal}) := \sum_{i=1}^k w_i (\text{ReLU}(\|f^{(j)}(v'_i)\|_2 - \text{mal}) + \text{ReLU}(\|g^{(j)}(v_i)\|_2 - \text{mal})),$$

- where `mal` is a threshold hyperparameter
- LOTOS uses a new component in the loss that penalizes the similarity for all pairs of corresponding layers for the models of the ensemble:

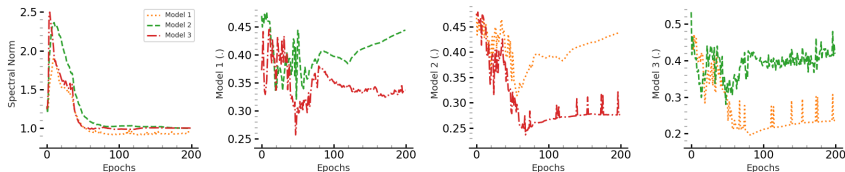
$$\mathcal{L}_{\text{train}} = \frac{1}{N} \sum_{i=1}^N \mathcal{L}_{\text{CE}}(\mathcal{F}_i(x), y) + \frac{\lambda}{M N (N-1)} \sum_{z=1}^{N-1} \sum_{j=z+1}^N \sum_{l=1}^M S_k^{(l)}(f_z^{(l)}, f_j^{(l)}, \text{mal})$$

# Does the new loss function work?



The effect of orthogonalization with different values of `mal`:

- When `mal` = 0.5:

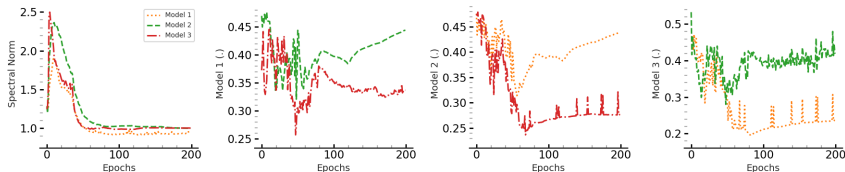


# Does the new loss function work?

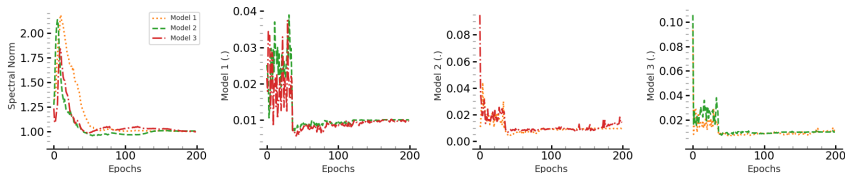


The effect of orthogonalization with different values of `mal`:

- When `mal = 0.5`:

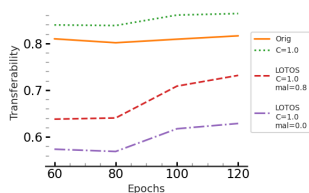
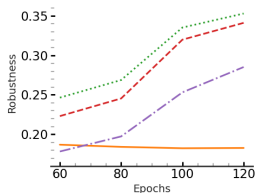
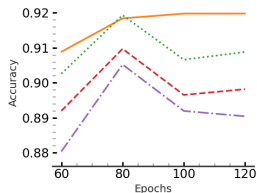


- When `mal = 0.01`:



Experimental setup: ensembles of three ResNet-18 models

- LOTOS effectively *decreases the transferability rate* among pairs of models
- LOTOS still benefits from *robustness of individual models*
- By decreasing *the value of  $mal$*  transferability rate decreases, but individual accuracy and robust accuracy decreases, and therefore, *introduces a trade-off*





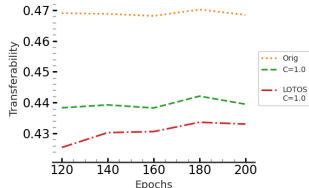
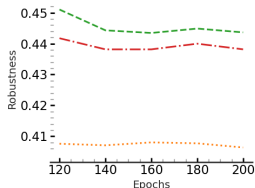
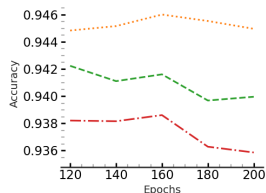
- Comparing the effectiveness in overall robustness of the ensembles:
  - Two different architectures (ResNet-18 and DLA)
  - Two different datasets (CIFAR-10 and CIFAR-100)

	CIFAR-10			CIFAR-100		
	ORIG	$C = 1.0$	LOTOS	ORIG	$C = 1.0$	LOTOS
ENSEMBLES OF RESNET-18 MODELS						
TEST ACC	<b>95.3 <math>\pm</math> 0.06</b>	94.7 $\pm$ 0.24	94.6 $\pm$ 0.19	<b>77.2 <math>\pm</math> 0.17</b>	76.6 $\pm$ 0.01	76.6 $\pm$ 0.10
ROBUST ACC	30.3 $\pm$ 1.63	35.2 $\pm$ 0.72	<b>36.3 <math>\pm</math> 0.88</b>	15.2 $\pm$ 0.45	18.9 $\pm$ 0.40	<b>20.2 <math>\pm</math> 0.47</b>
ENSEMBLES OF DLA MODELS						
TEST ACC	<b>95.4 <math>\pm</math> 0.12</b>	95.2 $\pm$ 0.05	95.05 $\pm$ 0.09	77.1 $\pm$ 0.09	<b>78.8 <math>\pm</math> 0.31</b>	78.3 $\pm$ 0.38
ROBUST ACC	26.7 $\pm$ 0.58	32.8 $\pm$ 1.28	<b>34.5 <math>\pm</math> 0.63</b>	16.5 $\pm$ 0.78	19.4 $\pm$ 0.32	<b>21.0 <math>\pm</math> 0.39</b>

# What about heterogeneous ensembles?



- The singular vectors of the first affine layer have the same dimension, which is equal to the input dimension!
  - They can be made orthogonal using LOTOS
- Ensembles of ResNet-18, ResNet-34, and DLA trained on CIFAR-10:
  - LOTOS still effective, but not as much as homogeneous ensembles



- Comparing the effectiveness of LOTOS when used along with TRS<sup>[3]</sup>, which is one of the prior SOTA in ensemble robust training
  - Two different architectures (ResNet-18 and DLA)
  - Two different datasets (CIFAR-10 and CIFAR-100)

	CIFAR-10			CIFAR-100		
	TRS	TRS + $C = 1$	TRS + LOTOS	TRS	TRS + $C = 1$	TRS + LOTOS
ENSEMBLES OF RESNET-18 MODELS						
TEST ACC	<b>94.4 ± 0.05</b>	94.1 ± 0.17	92.7 ± 0.09	<b>73.28 ± 0.46</b>	72.94 ± 0.29	67.23 ± 1.22
ROBUST ACC	30.8 ± 0.65	35.9 ± 1.35	<b>41.5 ± 1.04</b>	12.3 ± 0.53	16.3 ± 0.57	<b>20.7 ± 0.99</b>
ENSEMBLES OF DLA MODELS						
TEST ACC	<b>94.72 ± 0.06</b>	92.79 ± 0.13	93.18 ± 0.14	<b>72.6 ± 0.54</b>	63.3 ± 1.20	66.8 ± 1.26
ROBUST ACC	31.2 ± 0.80	32.9 ± 0.77	<b>35.3 ± 0.39</b>	23.2 ± 0.41	23.7 ± 2.36	<b>24.3 ± 1.67</b>

- Similar results with DVERGE<sup>[4]</sup>, another prior SOTA

<sup>3</sup> Yang, Z., Li, L., Xu, X., Zuo, S., Chen, Q., Zhou, P., ... Li, B. (2021). Trs: Transferability reduced ensemble via promoting gradient diversity and model smoothness. Advances in Neural Information Processing Systems, 34, 17642-17655.

<sup>4</sup> Yang, H., Zhang, J., Dong, H., Inkawhich, N., Gardner, A., Touchet, A., ... Li, H. (2020). Dverge: diversifying vulnerabilities for enhanced robust generation of ensembles. Advances in Neural Information Processing Systems, 33, 5505-5515.

- LOTOS promotes orthogonality among the top- $k$  sub-spaces.
- The  $(k + 1)$ st singular vectors can still be highly correlated!
- The  $(k + 1)$ st singular value can be as large as the 1st one!

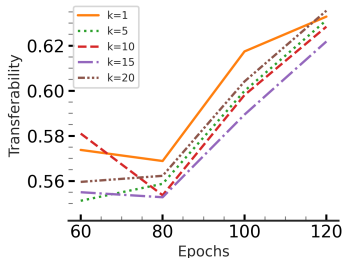
## Theorem

*Given two convolutional layers,  $M_1$  and  $M_2$  with a single input and output channel and circular padding for which  $\mathbf{f}$  is the vectorized form of the filter with a length of  $T$ , and considering  $n$  to be the length of the vectorized input, if  $\|Av'_1\|_2 \leq \epsilon$ , then:*

$$\|Av'_p\|_2 \leq \sqrt{\epsilon^2 + \pi \|\mathbf{f}\|_2^2 T^2 \frac{p}{n}},$$

*where  $A$  is the corresponding linear transformation of  $M_1$  and  $v'_p$  is singular vector of  $M_2$  corresponding to its  $p$ -th largest singular value.*

- Observing the effect of  $k$  in transferability rate:



- Increasing  $k$  slightly decreases the transferability rate, but up to  $k = 15$
- $k$  larger than 15 does not show further improvements due to further optimization constraints and restrictions of convolutional layer's spectrum

# Thank You!



Our Paper

- Observing the effect of  $k$  in transferability rate:

	ORIG	ORIG+ $C = 1$	ORIG + LOTOS
TRS	$10.3 \pm 0.33$	$13.2 \pm 1.12$	<b><math>14.5 \pm 0.81</math></b>
DVERGE	$19.7 \pm 2.34$	$26.8 \pm 0.75$	<b><math>29.2 \pm 0.56</math></b>

- Comparing the running times:

	ORIG	$C = 1$	LOTOS	TRS	TRS + $C = 1$	TRS + LOTOS	Adv	Adv + $C = 1$	Adv + LOTOS
RESNET-18	33.2	74.9 $\times 2.3$	79.3 $\times 2.4$	158.2	224.4 $\times 1.4$	227.3 $\times 1.4$	312.6	479.2 $\times 1.5$	485.2 $\times 1.5$
DLA	63.1	155.4 $\times 2.5$	165.4 $\times 2.6$	326.2	466.1 $\times 1.4$	477.4 $\times 1.5$	758.6.2	942.5 $\times 1.2$	949.2 $\times 1.2$

- For model  $\mathcal{F} : \mathcal{X} \rightarrow \mathcal{Y}$  and, attack  $\mathcal{A}(x) = x + \delta_x$ , and a sample  $(x, y)$ :
  - for untargeted attack:  $\mathcal{A}$  maximizes  $\ell_{\mathcal{F}}(x + \delta_x, y)$ , s.t.  $\|\delta_x\|_2 \leq \epsilon$ .
  - and targeted attack for  $y_t$ :  $\mathcal{A}_t(x) = \min_{\delta} \ell_{\mathcal{F}}(x + \delta, y_t)$
- Given an untargeted attack  $\mathcal{A}(x)$  on a surrogate model  $\mathcal{F}$ , transferability between  $\mathcal{F}$  and the target model  $\mathcal{G}$  is defined as:

$$T_{rate}(\mathcal{A}_{\mathcal{F}}, \mathcal{F}, \mathcal{G}) = \mathbb{P}_{(x,y) \in \mathcal{X} \times \mathcal{Y}} [\mathcal{G}(\mathcal{A}_{\mathcal{F}}(x)) \neq y \mid \mathcal{F}(x) = \mathcal{G}(x) = y \\ \wedge \mathcal{F}(\mathcal{A}_{\mathcal{F}}(x)) \neq y],$$



- For model  $\mathcal{F} : \mathcal{X} \rightarrow \mathcal{Y}$  and, attack  $\mathcal{A}(x) = x + \delta_x$ , and a sample  $(x, y)$ :
  - for untargeted attack:  $\mathcal{A}$  maximizes  $\ell_{\mathcal{F}}(x + \delta_x, y)$ , s.t.  $\|\delta_x\|_2 \leq \epsilon$ .
  - and targeted attack for  $y_t$ :  $\mathcal{A}_t(x) = \min_{\delta_x} [\ell_{\mathcal{F}}(x + \delta_x, y_t)]$
- Given an untargeted attack  $\mathcal{A}(x)$  on a surrogate model  $\mathcal{F}$ , transferability between  $\mathcal{F}$  and the target model  $\mathcal{G}$  is defined as:

$$T_{rate}(\mathcal{A}_{\mathcal{F}}, \mathcal{F}, \mathcal{G}) = \mathbb{P}_{(x,y) \in \mathcal{X} \times \mathcal{Y}} [\mathcal{G}(\mathcal{A}_{\mathcal{F}}(x)) \neq y \mid \mathcal{F}(x) = \mathcal{G}(x) = y \\ \wedge \mathcal{F}(\mathcal{A}_{\mathcal{F}}(x)) \neq y],$$

- And for targeted attack  $\mathcal{A}_t(x)$ , it is defined as:

$$T_{rate}(\mathcal{A}_{\mathcal{F}}^{(t)}, \mathcal{F}, \mathcal{G}) = \mathbb{P}_{(x,y) \in \mathcal{X} \times \mathcal{Y}} [\mathcal{G}(\mathcal{A}_{\mathcal{F}}^{(t)}(x)) = y_t \mid \mathcal{F}(x) = \mathcal{G}(x) = y \\ \wedge \mathcal{F}(\mathcal{A}_{\mathcal{F}}^{(t)}(x)) = y_t].$$