

CARTL: Cooperative Adversarially-Robust Transfer Learning



Dian Chen¹ Hongxin Hu² Qian Wang¹ Yinli Li¹

Cong Wang³ Chao Shen⁴ Qi Li⁵

¹Wuhan University ²University at Buffalo ³City University of Hong Kong

⁴Xi'an Jiaotong University ⁵Tsinghua University

Overview

- We reveal that there is a trade-off between accuracy and robustness in transfer learning.
- We propose a new transfer learning strategy, CARTL, for improving the accuracy-robustness trade-off of the target model.
- We demonstrate that selectively freezing the Batch Norm layers can further boost the robustness transfer.

Training Deep Neural Networks Is Tough

Network capacity

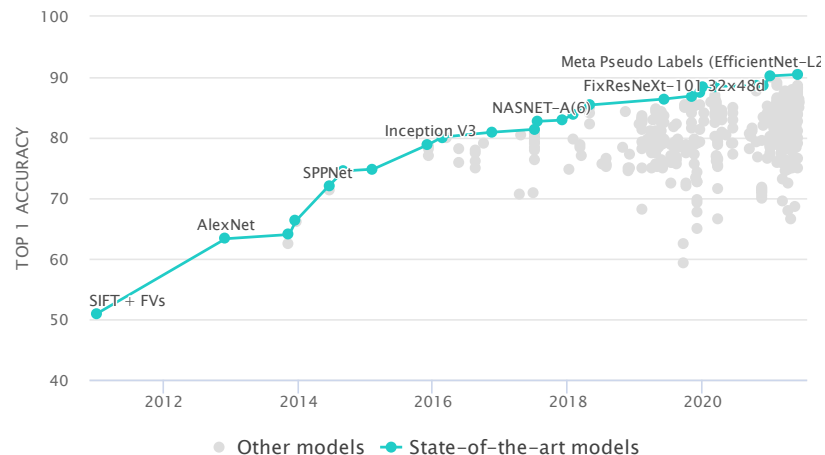
- Billions of parameters

Training data

- Extra training images except for ImageNet

Computational cost

- Thousands of core-hour



Adversarial Examples

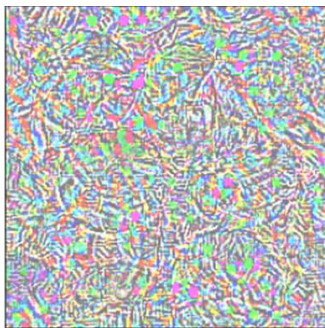
Adversarial examples are perturbed inputs that deceive DNNs into answering incorrect results.

$$\arg \max_i f_{\theta}(\mathbf{x})_i \neq \arg \max_j f_{\theta}(\mathbf{x}')_j$$

Golf ball (33.74%)



\mathbf{x}



δ



Sleeping bag (98.12%)



\mathbf{x}'

$$s. t. \|\delta\|_p \leq \epsilon, p = 0, 2, \infty$$

Adversarial Training

Adversarial training is similar to natural model training, but it takes (only) adversarial examples as the training data.

$$\min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\max_{\delta \in \mathcal{S}} \mathcal{L}(f_{\theta}(x + \delta), y) \right]$$

$AT(\theta, \mathcal{D}, \epsilon, \alpha)$:

repeat

$(x, y) \leftarrow \mathcal{D}$ // mini-batch

$\delta \leftarrow \text{rand_init}$

for i **in** $1, \dots, N$

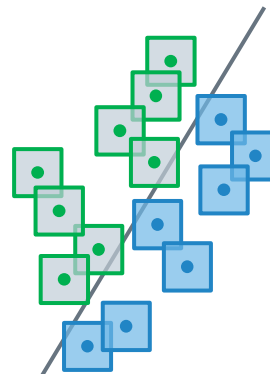
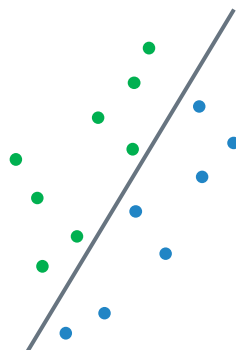
$\text{loss} = \mathcal{L}(f_{\theta}(x + \delta), y)$

$\delta = \Pi_{\epsilon}(\delta + \alpha \cdot \nabla_{\delta} \text{loss})$

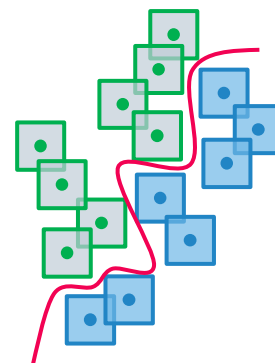
end for

$\theta = \theta + \eta \cdot \nabla_{\theta} \mathcal{L}(f_{\theta}(x + \delta), y)$

until *convergence*



Natural training



Adversarial training

Adversarial Training

Network capacity

- Adversarial robustness exhibits a strong demand on the network's capability, i.e., its depth and width.

Training data

- Adversarial training requires more training data than natural training.

Computational cost

- The training cost is N -times higher than natural training.

"Adversarial training increases the burden of model training."

- Xie et al., Intriguing Properties of Adversarial Training at Scale, ICLR 2020.
- Schmidt et al., Adversarially Robust Generalization Requires More Data, NeurIPS 2018.
- Shafahi et al., Adversarial Training for Free!, NeurIPS 2019.

Transfer Learning

Utilizing the knowledge obtained from the source domain to solve target domain tasks.

$$\min_{\bar{\theta}} \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\mathcal{L} \left(f_{\bar{\theta}}^{(L-k+1..L)} \left(f^{(L-k)}(x) \right), y \right) \right]$$

$TL(\theta_*, \mathcal{D}, k)$:

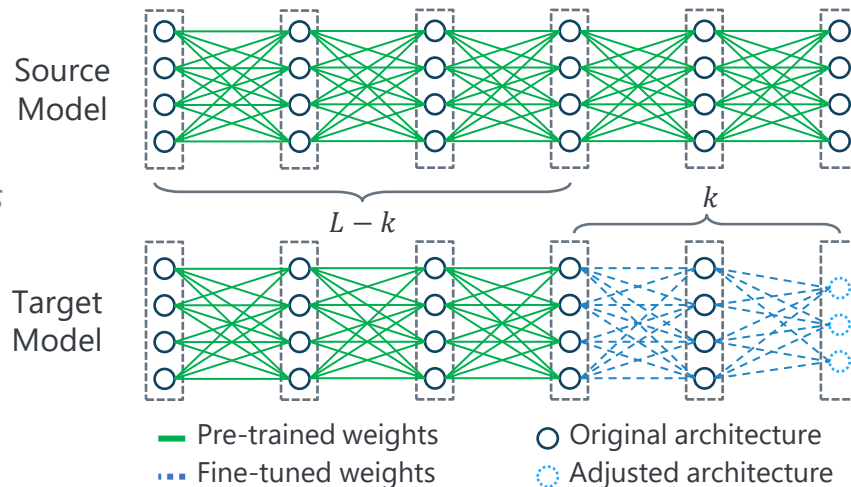
$\theta \leftarrow \theta_*$ // copy pre-trained weights

$\bar{\theta} := \{\theta_{L-k+1}, \dots, \theta_L\}$ // fine-tune last k layers

repeat

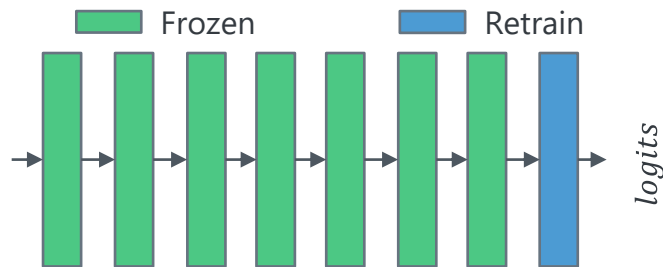
$(x, y) \leftarrow \mathcal{D}$ // mini-batch
 $loss = \mathcal{L}(f_{\bar{\theta}}(f^{(L-k)}(x)), y)$
 $\bar{\theta} = \bar{\theta} + \eta \cdot \nabla_{\bar{\theta}} loss$

until convergence



Shafahi's Work

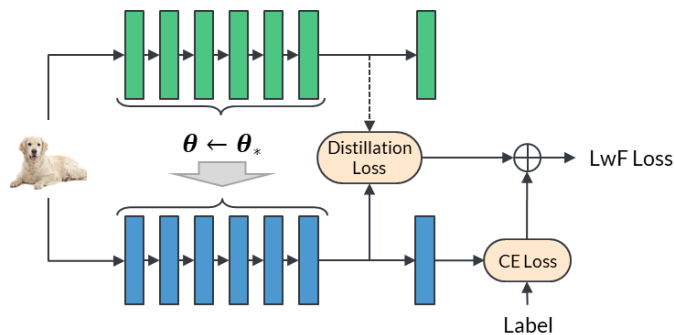
- The robustness of a hardened model is mainly due to its robust deep feature.
- Robustness does transfer when merely retraining the **last** fully-connected layer of a robust model.
- The accuracy of the target model **is poor**.



Shafahi's Work

- **End-to-end** transfer learning with a **distillation term**, called LwF

$$\min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}(f(x; \theta), y) + \lambda_d \cdot \| f^{(L-1)}(x; \theta) - f^{(L-1)}(x; \theta_*) \|_2]$$



- A trade-off between the accuracy and robustness of the target model:
 - Reduce λ_d for improving generalization, i.e., accuracy on the target domain.
 - Increase λ_d for obtaining better robustness transfer.

Problem Exploration

- How the number of fine-tuned layers affects the target model's robustness and accuracy.
- Widely-adopted architecture: wide residual network (WRN)
- Fine-tune the pre-trained robust model in the unit of the block on the target domain.

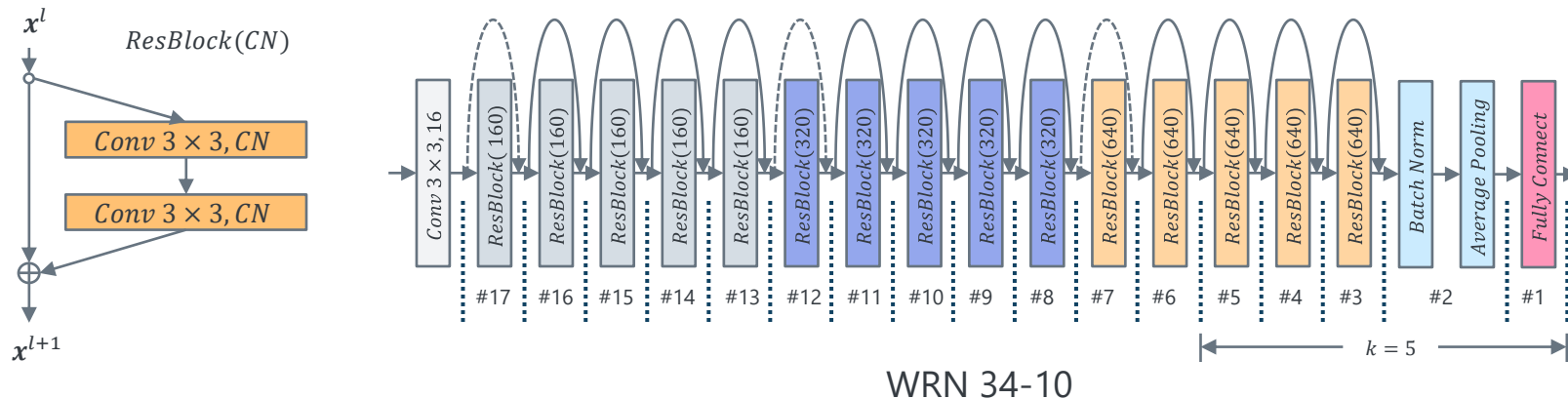
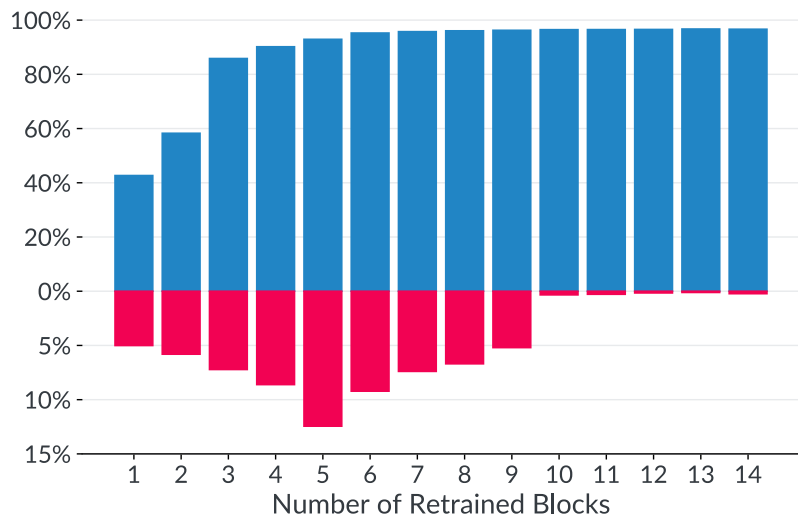
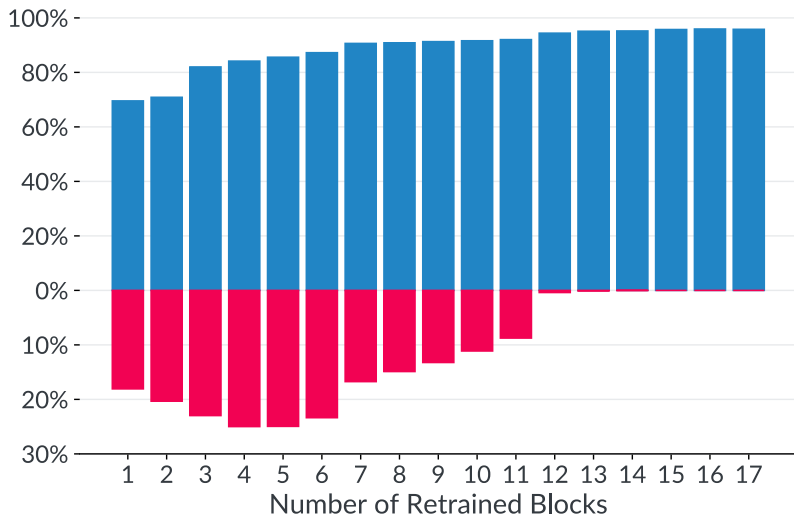


Figure from (Shafahi et al., 2020)

Problem Exploration



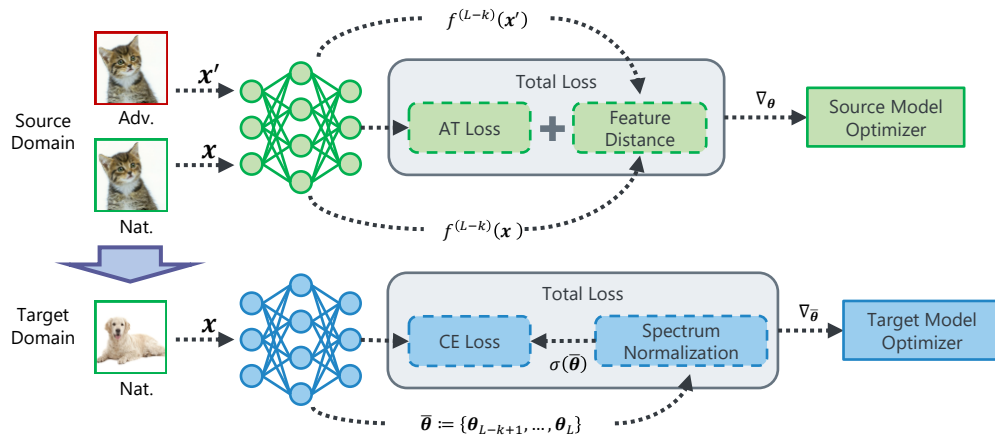
- Merely fine-tuning the last layer may not be sufficient.
- Insufficient accuracy leads to lower robustness.
- Accuracy increases together with the number of the fine-tuned layers.
- There is a **trade-off** between robustness and accuracy.



Can the target model obtain high accuracy while inheriting more robustness from the source model?

CARTL

- A cooperative approach
 - Feature distance minimization (FDM): adjusted adversarial training for the source model
 - Non-expansive fine-tuning (NEFT): constrained fine-tuning for model transfer
- Fine-tune the last k layers and freeze the first $L - k$ layers

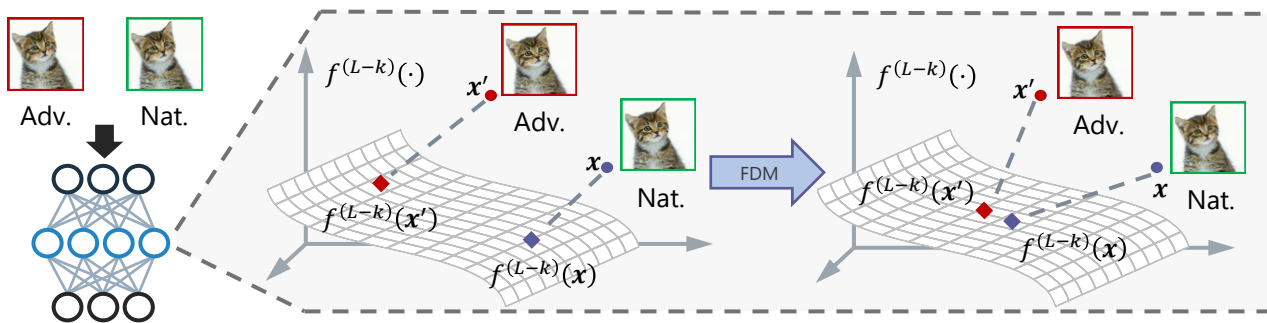


Feature Distance Minimization

The first $L - k$ layers frozen during transfer learning are taken as a feature extractor.

- Two inputs extracted similar features tend to be classified into an identical label.
- Reduce the distance between the features of adv. (x') and nat. (x) examples.

$$\mathcal{L}_{AT} + \frac{\lambda}{\sqrt{d}} \cdot \| f^{(L-k)}(x) - f^{(L-k)}(x') \|_2$$



Non-expansive Fine-tuning

- We call a function f Lipschitz continuous if

$$\|f(x) - f(x')\|_2 \leq \Lambda \cdot \|x - x'\|_2$$

- Lipschitz constant for DNN $f(\cdot; \theta) := f_{\theta_L}^L \circ f_{\theta_{L-1}}^{L-1} \circ \dots \circ f_{\theta_1}^1(\cdot)$

$$\|f(x) - f(x')\|_2 \leq \Lambda_L \cdot \Lambda_{L-1} \cdots \Lambda_1 \|x - x'\|_2$$

- A general form of the deep neural layer f^l :

$$f^l = \mathbf{W}^l \cdot \mathbf{x} + \mathbf{b}^l$$

Non-expansive Fine-tuning

- The remaining dissimilarity of features may still result in model misclassification.
- Non-expansive fine-tuning: mitigate the error caused by the dissimilarity of features.

$$\mathbf{W}_*^l := \beta \cdot \frac{\mathbf{W}^l}{\sigma(\mathbf{W}^l)}$$

- β is a hyper-parameter for further scaling down the Lipschitz constant

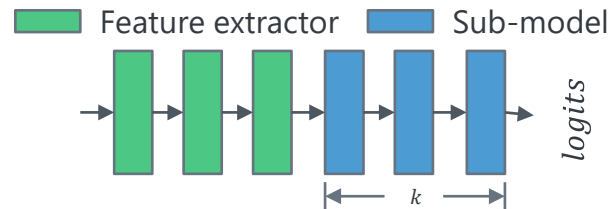
Rethinking Fine-tuning Batch Norm Layer

- An essential component for DNN: internal covariate shift, model training acceleration

$$BN(\mathbf{x}) := \mathbf{W} \cdot \frac{\mathbf{x} - \text{mean}(\mathbf{x})}{\sqrt{\text{var}(\mathbf{x}) + \varepsilon}} + \mathbf{b}$$
$$\boldsymbol{\mu} := m \cdot \boldsymbol{\mu} + (1 - m) \cdot \text{mean}(\mathbf{x})$$
$$\boldsymbol{\sigma} := m \cdot \boldsymbol{\sigma} + (1 - m) \cdot \sqrt{\text{var}(\mathbf{x})}$$

- Parameters:
 - Statistic parameters $\boldsymbol{\mu}$ and $\boldsymbol{\sigma}$: updated with a momentum (m).
 - Affine parameters \mathbf{W} and \mathbf{b} : updated through back propagation.

- Four cases:
 - Update/reuse $\boldsymbol{\mu}$ and $\boldsymbol{\sigma}$ in the feature extractor
 - Fine-tune/freeze \mathbf{W} and \mathbf{b} in the sub-model



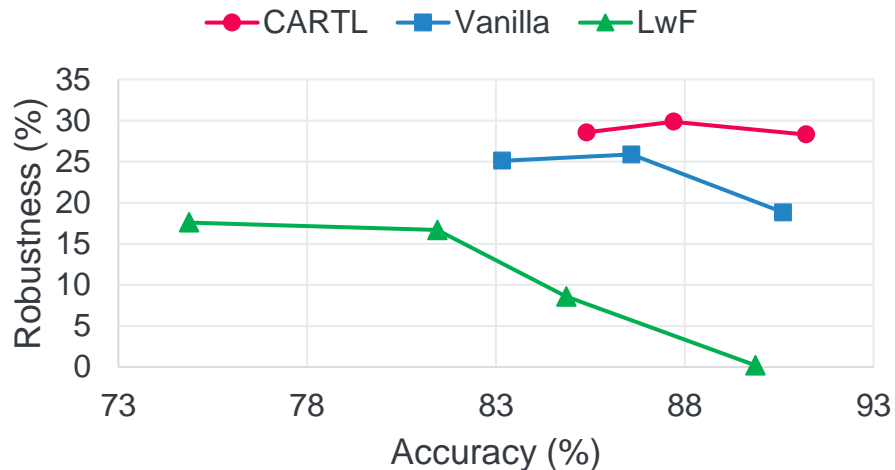
Rethinking Fine-tuning Batch Norm Layer

		W, b		μ, σ, W, b	
		Acc (%)	Rob (%)	Acc (%)	Rob (%)
CIFAR-100 → CIFAR-10 (8)	-	91.17	14.36	90.86	14.89
	W, b	90.70	17.41	90.84	18.54
CIFAR-10 → GTSRB (6)	-	93.02	30.22	89.29	32.22
	W, b	92.13	32.22	88.94	34.53
CIFAR-10 → SVHN (6)	-	95.29	3.88	95.24	9.22
	W, b	95.16	4.90	94.86	11.52
CIFAR-10 → SVHN (5)	-	93.47	4.71	92.92	12.45
	W, b	93.41	5.64	92.10	14.16

- Transferred robustness can be boosted if freezing affine parameters of the sub-model.
- Freezing statistics of the feature extractor plays a crucial role in robustness transfer.
- Reuse source domain statistics may cast negative impacts on the accuracy

Evaluations

- LwF improves the robustness but aggressively harms the accuracy and vice versa.
- Vanilla and CARTL maintain higher robustness in the case of an equivalent level of accuracy.
- CARTL further improves the accuracy-robustness trade-off compared with Vanilla.



Evaluations

- CARTL exhibits similar trends to Vanilla, achieving higher robustness at Case-6.
- A smaller λ helps robustness transfer but slightly results in lower accuracy.
- Reducing Lipschitz constants significantly improves the target models' robustness

		NEFT $\beta = 1.0$		NEFT $\beta = 0.6$		NEFT $\beta = 0.4$	
		Acc (%)	Rob (%)	Acc (%)	Rob (%)	Acc (%)	Rob (%)
Case-4	$\lambda = 0.01$	86.09	25.73	86.08	27.17	85.64	28.40
	$\lambda = 0.005$	85.41	25.75	85.47	27.14	85.51	28.47
Case-6	$\lambda = 0.01$	87.78	25.58	87.92	27.27	87.96	29.60
	$\lambda = 0.005$	87.66	25.97	88.07	27.64	87.79	30.94
Case-8	$\lambda = 0.01$	91.85	16.36	91.63	19.22	91.55	27.47
	$\lambda = 0.005$	91.71	17.62	91.10	21.60	91.30	29.34

Evaluations

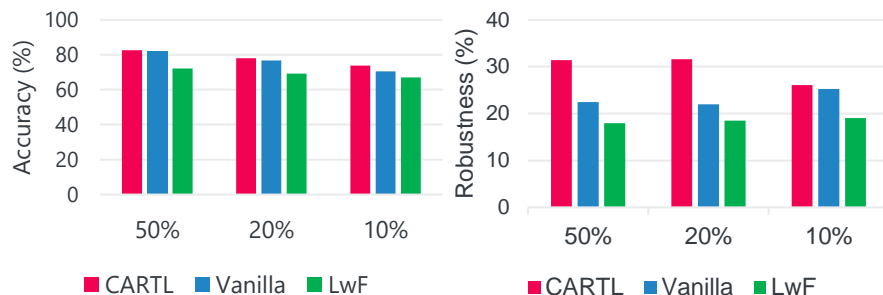
- Fine-tuning the target model with NEFT significantly increases its robustness.
- FDM further improves the robustness except for Case-8.
- By using FDM, the target model's accuracy slightly rises in all cases.

Method		Case-4		Case-6		Case-8	
Source	Transfer	Acc (%)	Rob (%)	Acc (%)	Rob (%)	Acc (%)	Rob (%)
AT	TL	83.22	25.23	86.92	25.38	90.82	18.54
AT	NEFT	83.72	26.29	86.87	27.95	90.92	29.97
AT + FDM	NEFT	85.51	28.47	87.79	30.94	91.30	29.34

CIFAR-100 → CIFAR-10

Evaluations

- CARTL outperforms both LwF and Vanilla when the data size is small.
- Target models fine-tuned with CARTL inherit superior robustness from the source model in all cases.



Source	Target	Arch.	LwF		Vanilla		CARTL	
			Acc (%)	Rob (%)	Acc (%)	Rob (%)	Acc (%)	Rob (%)
CIFAR-100	SVHN	WRN 34-10	85.90	6.67	92.83	17.64	93.96	22.21
CIFAR-100	GTSRB	WRN 34-10	70.34	15.85	80.40	30.25	83.07	47.34
CIFAR-10	SVHN	WRN 28-4	94.32	4.68	94.86	11.52	94.76	21.65
GTSRB	SVHN	WRN 28-4	81.80	1.08	93.91	6.08	94.07	15.26

Conclusion

- We conduct detailed experiments and reveal that there is a trade-off between accuracy and robustness during transfer learning.
- We propose CARTL, consisted of FDM and NEFT, for improving the accuracy-robustness trade-off of the target model.
- We demonstrate that freezing affine parameters of Batch Norm layers can further boost the robustness transfer, and Batch Norm layers' statistics play a crucial role in robustness transfer.

Thanks!

Any questions?



qianwang@whu.edu.cn



<https://github.com/NISP-official/CARTL>