

Institute for the Wireless Internet of Things at Northeastern University

HyperAdv: Dynamic Defense Against Adversarial Radio Frequency Machine Learning Systems

Milin Zhang, Michael De Lucia, Ananthram Swami, Jonathan Ashdown, Kurt Turck and Francesco Restuccia

> Presenting: **Milin Zhang** Northeastern University, United States Email: <u>frestuc@northeastern.edu</u> Website: http://mentislab.info

Adversarial Machine Learning





Adversarial Machine Learning





[1] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards Deep Learning Models Resistant to Adversarial Attacks," in Proceedings of International Conference on Learning Representations (ICLR), 2018.



Min-max optimization in adversarial training (AT) [1]:



Natural Training (NT) and AT on CIFAR 10 (data reported in original paper [1])

Adversarial training sacrifice performance on clean data

[1] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards Deep Learning Models Resistant to Adversarial Attacks," in Proceedings of International Conference on Learning Representations (ICLR), 2018.

Dynamic Adversarial Defense: HyperAdv



Key idea: dynamically change the DNN at each inference step so that the adversarial information at step t is not effective at step t+1

Institute for the Wireless

Internet of Things

at Northeastern





HyperNet [1]: a neural network to generate weight Wc for CNN classifier dynamically

CNN: a target network to perform RF signal classification

Attacker: Adversary to create noise δ based on the gradient $\nabla L(f(Wc, x), y')$

Context: Random vector to trigger the weight generation as well as calibrate the output y to y'

[1] D. Ha, A. M. Dai, and Q.V. Le, "HyperNetworks," in International Conference on Learning Representations (ICLR), 2017.



Training:

- 1. a set of context $C_i = [C_i^1, C_i^2]$ is randomly sampled from U(-1, 1)
- 2. HyperNet $H(\cdot): C_i \mapsto W_c^i$ generate weights for target CNN
- 3. CNN $f_{W_c^i}(\cdot): x \mapsto y$ maps input to label space
- 4. logit calibration $y'_i = C_i^1 * y + C_i^2$ to map output to diversified space
- 5. Ensemble learning

$$\mathcal{L} = \frac{1}{n} \sum_{i} \mathcal{L}_i(f_{W_c^i}(x), y_i')$$

can combine other defensive training such as AT and TRADES [I] to further improve robustness

[1] H. Zhang, Y. Yu, J. Jiao, E. Xing, L. El Ghaoui, and M. Jordan, "Theoretically Principled Trade-off between Robustness and Accuracy," in Proceedings of International Conference on Machine Learning (ICML), pp. 7472–7482, PMLR, 2019.





The size of Hypernetwork is **~M times** larger (e.g. 256) than a single target CNN



Effectively reduce the size by ~32 times





Robustness-Performance Trade-off:

respectively.

AT and TRADES reduce accuracy on clean data by 27.92% and 6.65%

Institute for the Wireless Internet of Things at Northeastern

TABLE I: Accuracy of HyperAdv and its static counterparts

	NT		AT		TRADES	
	Clean	PGD	Clean	PGD	Clean	PGD
CNN	94.20	3.72	66.28	50.42	87.55	43.16
HyperAdv -R	95.20	45.02	74.01	63.98	90.37	60.92
HyperAdv- ${f S}$	95.11	23.70	72.74	50.10	91.27	43.08
HyperAdv - E	96.34	27.88	78.36	61.34	92.79	49.50

Ablation Study:

- The dynamic HyperAdv has best robust accuracy
- The **ensemble** HyperAdv

performs best on **clean** data

Computational Cost of Adversary:

NT: $3.31 \rightarrow 19.83$ iterations AT: $19.83 \rightarrow 22.66$ iterations TRADES: $17.64 \rightarrow 21.01$ iterations





- More scalable HyperNetwork design
- Latency and energy consumption
- Real-world implementation





Code and pretrained models are available: https://github.com/Restuccia-Group/HyperAdv