

Robustness to Adversarial Examples through an Ensemble of Specialists

Mahdieh Abbasi Christian Gagné

1 PROBLEM

Adding **small** but **smart** perturbations to an input image generates another image, called *adversarial examples*, that is visually similar to the original one.

While a CNN can correctly classify a clean sample, it can confidently misclassify its corresponding adversaries.

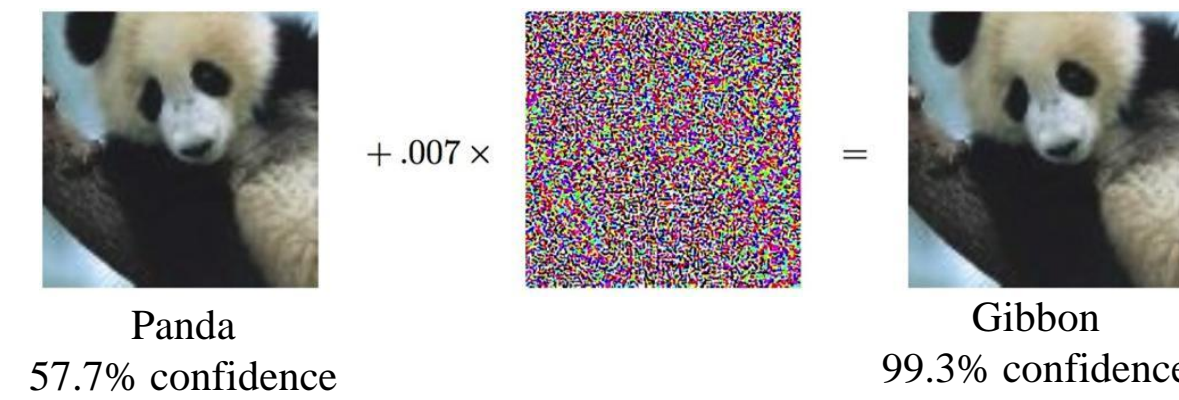


Figure 1: An adversarial example generated by Fast Gradient Sign (FGS) [Goodfellow et al. 2014].

Due to the cross-model generalization property of adversaries, an attacker can easily attack a CNN based system by generating some adversarial examples with another CNN.

2 RELATED WORK

There are two general trends for robustifying CNNs:

- Training CNNs on adversaries
 - [Goodfellow et al., 2014, Huang et al., 2016]: training on Fast Gradient Sign (FGS) adversaries
 - [Moosavi-Dezfooli et al., 2016]: training on DeepFool (DF) adversaries
 - [Rozsa et al., 2016]: training on diverse types of adversaries
- Identifying and rejecting adversaries as unknown
 - [Bendale & Boulton, 2016]: adapting CNNs for recognizing unknown samples as coming from unknown classes or from fooling examples

3 MOTIVATION

- Instead of training a CNN on all possible types of adversaries, developing a generic framework that can identify and reject adversarial examples.
- An ensemble of diverse CNNs can provide the following properties:
 - In presence of adversaries**, disagreement (i.e. high entropy) in the ensemble leads to identifying and rejecting them.
 - In presence of clean samples**, the ensemble can correctly and confidently classify them.

4 OBSERVATION

- The confusion matrices of FGS adversaries reveal some interesting patterns among labels
 - Samples from each class have a high tendency to be fooled toward a limited number of classes

True Labels	0	1	2	3	4	5	6	7	8	9	Airplane	Auto	Bird	Cat	Deer	Dog	Frog	Horse	Ship	Truck
0	0.00	1.40	6.80	20.20	1.00	17.00	7.60	6.20	13.20	26.61	0.00	1.60	22.00	3.88	15.20	0.22	18.46	1.70	17.70	19.24
1	0.00	0.00	5.80	17.80	7.20	1.80	0.80	37.59	29.00	0.00	2.70	0.00	7.20	1.94	5.46	0.06	7.86	0.90	8.50	65.38
2	3.80	12.41	0.00	24.19	2.60	2.40	1.80	28.20	23.80	0.80	3.70	0.58	0.00	13.34	32.24	0.98	33.02	2.78	3.28	10.98
3	0.20	0.60	10.20	0.00	0.00	58.00	0.00	5.00	23.80	2.20	0.82	0.30	12.58	0.00	20.92	6.86	43.38	6.24	1.46	7.44
4	0.00	0.40	6.80	10.00	0.00	2.40	11.40	17.39	11.80	39.78	0.90	0.24	21.00	7.30	0.00	0.58	51.94	8.94	0.94	8.16
5	0.20	0.20	4.40	22.20	1.20	0.00	3.20	1.40	35.59	31.59	0.66	0.26	16.94	14.92	15.88	0.00	36.98	5.86	1.74	6.78
6	32.00	2.80	5.00	1.40	7.40	20.80	0.00	0.20	29.39	1.00	0.60	0.76	28.02	16.94	41.14	0.52	0.00	3.48	0.74	7.80
7	0.00	7.80	8.40	41.00	2.80	3.00	0.20	0.00	7.00	29.81	0.92	0.36	6.40	5.66	32.44	1.62	47.16	0.00	0.98	4.46
8	0.00	1.00	40.78	21.41	0.40	10.60	8.80	4.00	0.00	13.00	11.10	4.12	20.94	4.42	7.12	0.10	14.28	0.84	0.00	37.08
9	0.80	0.40	2.20	22.20	41.00	4.20	0.00	10.00	19.20	0.00	3.88	23.92	7.24	2.64	13.44	0.22	42.06	1.74	4.86	0.00

Figure 2: The confusion matrices of adversaries for MNIST (left) and CIFAR10 (right). Each number in row i and column j presents the percentage of the sample from class i that is being fooled as class j .

5 SPECIALISTS+1 ENSEMBLE

5.1 Definition of Expertise Domains

The expertise domains are defined based on some **subsets of classes** for a classification problem with K classes, $C = \{c_1, c_2, \dots, c_K\}$.

- For each class c_i , two subsets are identified according to its corresponding row from the adversaries confusion matrix:
 - The confusing target subset (U_i): built by adding classes sequentially in decreasing c_i -related confusion values order until at least 80% of confusions are covered
 - The less-confusing target subset: $U_{i+K} = C \setminus U_i$

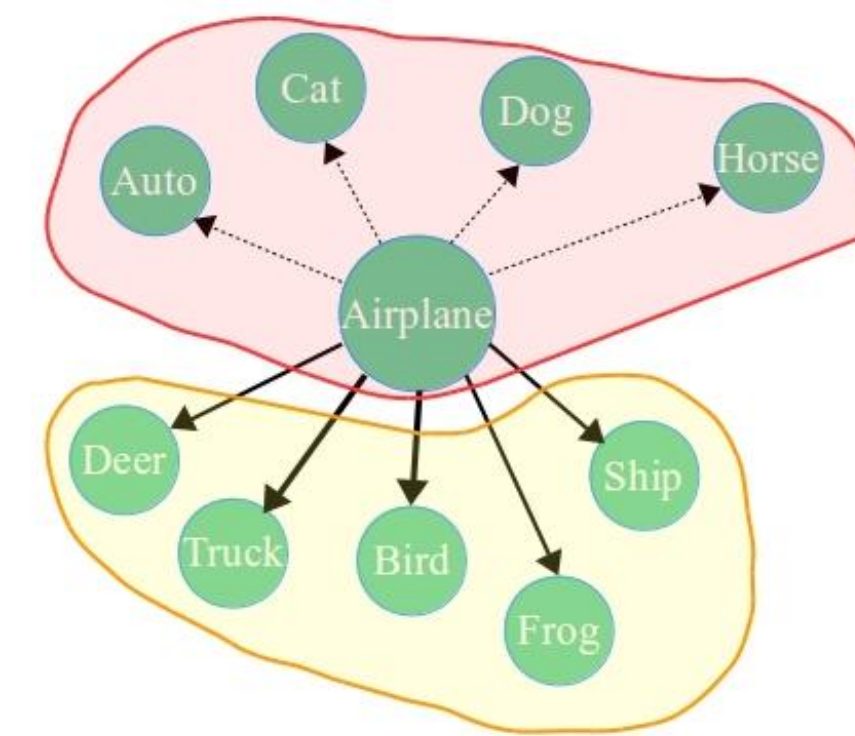


Figure 3: Schematic illustration of the expertise domains for class "Airplane". From the confusion matrix depicted in Fig. 2 (right), "Airplane" samples mostly get fooled toward the classes in yellow zone, while these samples get less fooled toward the classes in red zone.

5.2 Ensemble Creation

- An ensemble of specialist CNNs generated by training a CNN for each expertise domain, i.e. label subset.
- The ensemble also includes a generalist CNN trained on the whole set of classes.

6 VOTING MECHANISM-ALGORITHM

Input:

- Given ensemble $\mathcal{H} = \{h^1, \dots, h^M\}$ with $h^j \in \mathbb{R}^K$
- Given label subsets (expertise domains) $\mathcal{U} = \{U_1, \dots, U_M\}$
- The maximum expected number of votes to class c_k , $V_k = K + 1$

Output:

- Final prediction $\bar{h}(x) \in \mathbb{R}^K$

Indicating the winner class:

Given an input x ,

- Computing the number of votes for each class, c_k

$$v_k(x) \leftarrow \sum_{j=1}^M \mathbb{I}[c_k = \operatorname{argmax}_{i=1}^K h_i^j(x)]$$

- Indicating the winner class (k^*): the class with maximum number of votes

Computing final prediction:

- If $v_{k^*}(x) = V_{k^*}$, activate the CNNs that vote to k^* :

$$S \leftarrow \{h^i \in \mathcal{H} | c_{k^*} \in U_i\}$$

$$\bar{h}(x) \leftarrow \frac{1}{K+1} \sum_{h^i \in S} h^i(x)$$

- Otherwise, activate all of the CNNs:

$$\bar{h}(x) \leftarrow \frac{1}{M} \sum_{h^i \in \mathcal{H}} h^i(x)$$

7 EVALUATION METRICS

Consider $h(x) = [h_1(x), \dots, h_K(x)]$ as a multi-classification system:

- Rejecting instances with confidence lower than a threshold (τ) to a "reject class" c_{K+1}

Two types of error should be considered:

- Error E_D on the clean set**, counts both clean samples that are misclassified and correctly classified rejected clean samples
- Error E_A on the adversaries set**, considers misclassified adversarial instances that are not rejected

8 EXPERIMENTAL RESULTS

Specialists+1 ensemble is compared with

- Ensemble of 5 generalists, i.e. pure ensemble,
- Naïve CNN*

Tested on three types of adversaries: FGS, DeepFool (DF), [Szegedy et al., (2013)]

- Distribution of confidence**

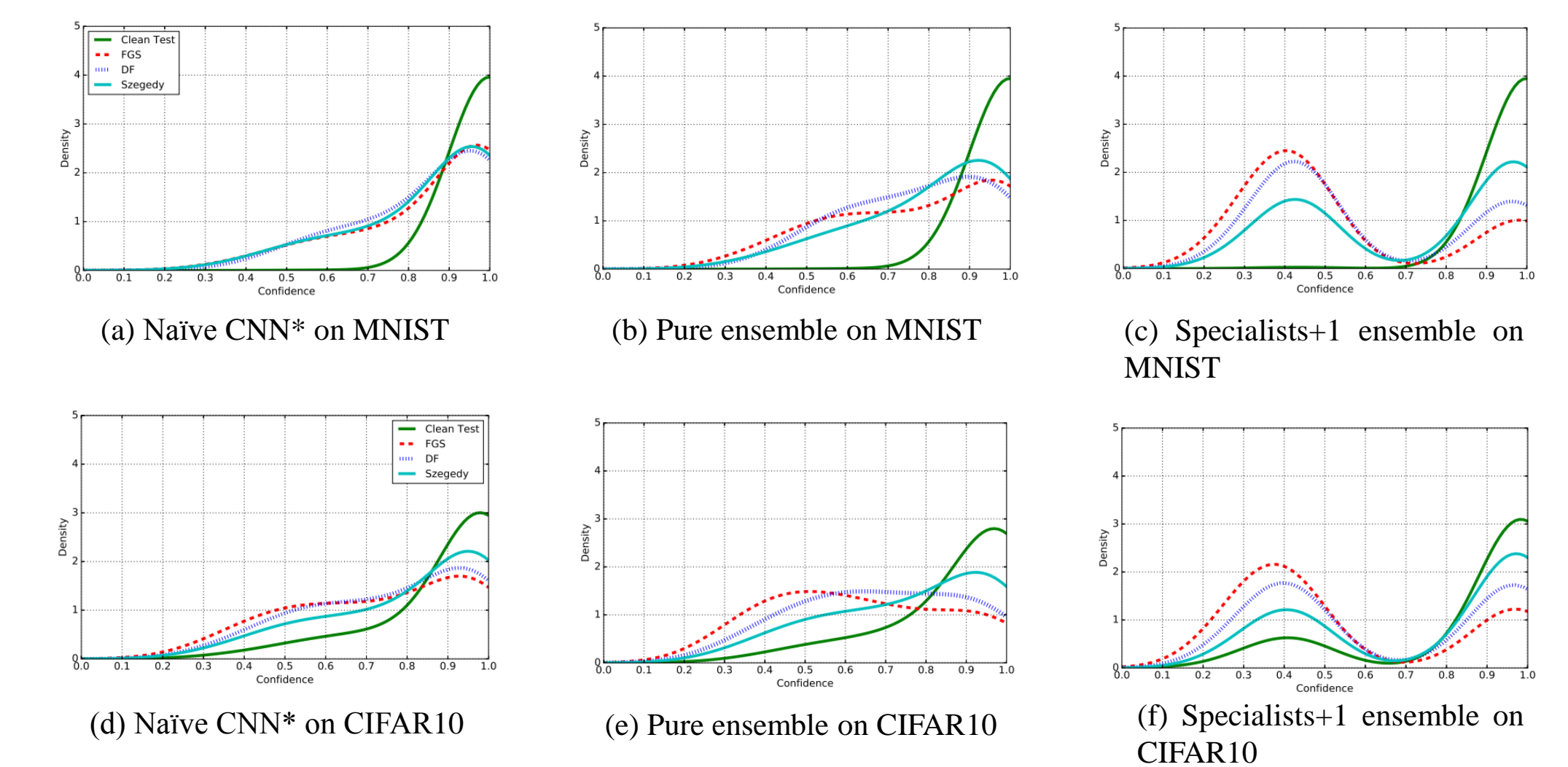


Figure 4: Confidence densities on MNIST (the first row) and CIFAR-10 (the second row)

- Error rate**

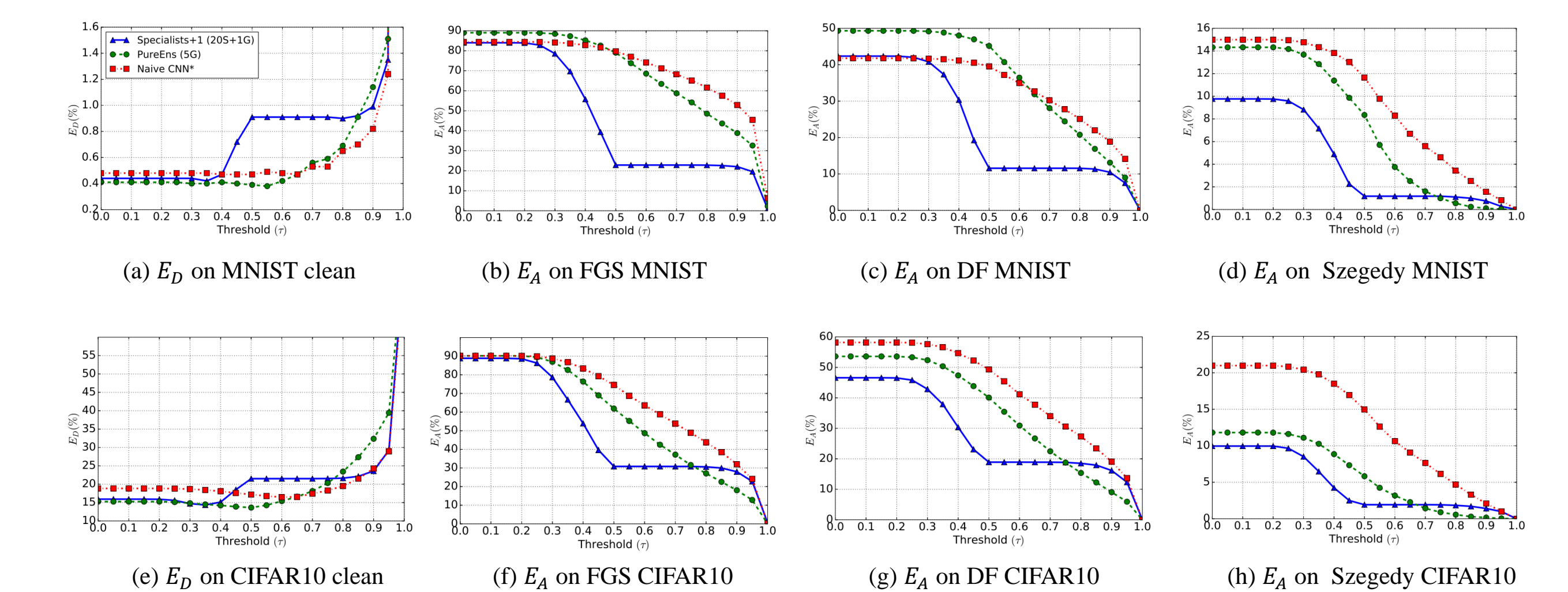


Figure 5: Error rates E_D on clean test samples and error rates E_A on their corresponding adversaries as a function of threshold (τ) for MNIST and CIFAR10 datasets.

9 CONCLUSION

- Without training from adversaries and by leveraging diversity in specialists ensemble, clean samples are discriminated from adversaries.
- Increasing the robustness of CNNs by refusing the suspicious samples.