

## A APPENDIX

### A.1 DIRICHLET-BASED UNCERTAINTY MODELS

In this section, we provide details on the losses used by each DBU model. PostNet uses a Bayesian loss which can be expressed as follows:

$$L_{\text{PostNet}} = \frac{1}{N} \sum_i \mathbb{E}_{q(p^{(i)})} [\text{CE}(p^{(i)}, y^{(i)})] - H(q^{(i)}) \quad (3)$$

where CE denotes the cross-entropy. Both the expectation term (i.e.  $\mathbb{E}_{q(p^{(i)})} [\text{CE}(p^{(i)}, y^{(i)})]$ ) and the entropy term (i.e.  $H(q^{(i)})$ ) can be computed in closed-form (Charpentier et al., 2020). PriorNet uses a loss composed of two KL divergence terms for ID and OOD data:

$$L_{\text{PriorNet}} = \frac{1}{N} \left[ \sum_{\mathbf{x}^{(i)} \in \text{ID data}} [\text{KL}[\text{Dir}(\alpha^{\text{ID}}) || q^{(i)}]] + \sum_{\mathbf{x}^{(i)} \in \text{OOD data}} [\text{KL}[\text{Dir}(\alpha^{\text{OOD}}) || q^{(i)}]] \right]. \quad (4)$$

Both KL divergences terms can be computed in closed-form (Malinin & Gales, 2019). The precision  $\alpha^{\text{ID}}$  and  $\alpha^{\text{OOD}}$  are hyper-parameters. The precision  $\alpha^{\text{ID}}$  is usually set to  $1e^1$  for the correct class and 1 otherwise. The precision  $\alpha^{\text{OOD}}$  is usually set to 1. DDNet uses the Dirichlet likelihood of soft labels produce by an ensemble of  $M$  neural networks:

$$L_{\text{DDNet}} = -\frac{1}{N} \sum_i \sum_{m=1}^M [\ln q^{(i)}(\pi^{im})] \quad (5)$$

where  $\pi^{im}$  denotes the soft-label of  $m$ th neural network. The Dirichlet likelihood can be computed in closed-form (Malinin et al., 2019). EvNet uses the expected mean square error between the one-hot encoded label and the predicted categorical distribution:

$$L_{\text{EvNet}} = \frac{1}{N} \sum_i \mathbb{E}_{\mathbf{p}^{(i)} \sim \text{Dir}(\alpha^{(i)})} ||\mathbf{y}^{*(i)} - \mathbf{p}^{(i)}||^2 \quad (6)$$

where  $\mathbf{y}^{*(i)}$  denotes the one-hot encoded label. The expected MSE loss can also be computed in closed form (Sensoy et al., 2018). For more details please have a look at the original paper on PriorNet (Malinin & Gales, 2018a), PostNet (Charpentier et al., 2020), DDNet (Malinin & Gales, 2019) and EvNet (Sensoy et al., 2018).

### A.2 CLOSED-FORM COMPUTATION OF UNCERTAINTY MEASURES & UNCERTAINTY ATTACKS

Dirichlet-based uncertainty models allow to compute several uncertainty measures in closed form (see (Malinin & Gales, 2018a) for a derivation). As proposed by Malinin & Gales (2018a), we use precision  $m_{\alpha_0}$ , differential entropy  $m_{\text{diffE}}$  and mutual information  $m_{\text{MI}}$  to estimate uncertainty on predictions.

The differential entropy  $m_{\text{diffE}}$  of a DBU model reaches its maximum value for equally probable categorical distributions and thus, a on flat Dirichlet distribution. It is a measure for distributional uncertainty and expected to be low on ID data, but high on OOD data.

$$m_{\text{diffE}} = \sum_c^K \ln \Gamma(\alpha_c) - \ln \Gamma(\alpha_0) - \sum_c^K (\alpha_c - 1) \cdot (\Psi(\alpha_c) - \Psi(\alpha_0)) \quad (7)$$

where  $\alpha$  are the parameters of the Dirichlet-distribution,  $\Gamma$  is the Gamma function and  $\Psi$  is the Digamma function.

The mutual information  $m_{\text{MI}}$  is the difference between the total uncertainty (entropy of the expected distribution) and the expected uncertainty on the data (expected entropy of the distribution). This uncertainty is expected to be low on ID data and high on OOD data.

$$m_{\text{MI}} = -\sum_{c=1}^K \frac{\alpha_c}{\alpha_0} \left( \ln \frac{\alpha_c}{\alpha_0} - \Psi(\alpha_c + 1) + \Psi(\alpha_0 + 1) \right) \quad (8)$$

Furthermore, we use the precision  $\alpha_0$  to measure uncertainty, which is expected to be high on ID data and low on OOD data.

$$m_{\alpha_0} = \alpha_0 = \sum_{c=1}^K \alpha_c \quad (9)$$

As these uncertainty measures are computed in closed form and it is possible to obtain their gradients, we use them (i.e.  $m_{\text{diffE}}$ ,  $m_{\text{MI}}$ ,  $m_{\alpha_0}$ ) as target function of our uncertainty attacks. Changing the attacked target function allows us to use a wide range of gradient-based attacks such as FGSM attacks, PGD attacks, but also more sophisticated attacks such as Carlini-Wagner attacks.

### A.3 DETAILS OF THE EXPERIMENTAL SETUP

**Models.** We trained all models with a similar based architecture. We used namely 3 linear layers for vector data sets, 3 convolutional layers with size of  $5 + 3$  linear layers for MNIST and the VGG16 Simonyan & Zisserman (2015) architecture with batch normalization for CIFAR10. All the implementation are performed using Pytorch (Paszke et al., 2019). We optimized all models using Adam optimizer. We performed early stopping by checking for loss improvement every 2 epochs and a patience of 10. The models were trained on GPUs (1 TB SSD).

We performed a grid-search for hyper-parameters for all models. The learning rate grid search was done in  $[1e^{-5}, 1e^{-3}]$ . For PostNet, we used Radial Flows with a depth of 6 and a latent space equal to 6. Further, we performed a grid search for the regularizing factor in  $[1e^{-7}, 1e^{-4}]$ . For PriorNet, we performed a grid search for the OOD loss weight in  $[1, 10]$ . For DDNet, we distilled the knowledge of 5 neural networks after a grid search in  $[2, 5, 10, 20]$  neural networks. Note that it already implied a significant overhead at training compare to other models.

**Metrics.** For all experiments, we focused on using AUC-PR scores since it is well suited to imbalance tasks (Saito & Rehmsmeier, 2015) while bringing theoretically similar information than AUC-ROC scores (Davis & Goadrich, 2006). We scaled all scores from  $[0, 1]$  to  $[0, 100]$ . All results are average over 5 training runs using the best hyper-parameters found after the grid search.

**Data sets.** For vector data sets, we use 5 different random splits to train all models. We split the data in training, validation and test sets (60%, 20%, 20%).

We use the segment vector data set Dua & Graff (2017), where the goal is to classify areas of images into 7 classes (window, foliage, grass, brickface, path, cement, sky). We remove class window from ID training data to provide OOD training data to PriorNet. Further, We remove the class 'sky' from training and instead use it as the OOD data set for OOD detection experiments. Each input is composed of 18 attributes describing the image area. The data set contains 2,310 samples in total.

We further use the Sensorless Drive vector data set Dua & Graff (2017), where the goal is to classify extracted motor current measurements into 11 different classes. We remove class 9 from ID training data to provide OOD training data to PriorNet. We remove classes 10 and 11 from training and use them as the OOD dataset for OOD detection experiments. Each input is composed of 49 attributes describing motor behaviour. The data set contains 58,509 samples in total.

Additionally, we use the MNIST image data set LeCun & Cortes (2010) where the goal is to classify pictures of hand-drawn digits into 10 classes (from digit 0 to digit 9). Each input is composed of a  $1 \times 28 \times 28$  tensor. The data set contains 70,000 samples. For OOD detection experiments, we use FashionMNIST Xiao et al. (2017) and KMNIST Clanuwat et al. (2018) containing images of Japanese characters and images of clothes, respectively. FashionMNIST was used as training OOD for PriorNet while KMNIST is used as OOD at test time.

Finally, we use the CIFAR10 image data set Krizhevsky et al. (2009) where the goal is to classify a picture of objects into 10 classes (airplane, automobile, bird, cat, deer, dog, frog, horse, ship, truck). Each input is a  $3 \times 32 \times 32$  tensor. The data set contains 60,000 samples. For OOD detection experiments, we use street view house numbers (SVHN) Netzer et al. (2011) and CIFAR100 (Krizhevsky et al., 2009) containing images of numbers and objects respectively. CIFAR100 was used as training OOD for PriorNet while SVHN is used as OOD at test time.

**Perturbations.** For all label and uncertainty attacks, we used Fast Gradient Sign Methods and Project Gradient Descent. We tried 6 different radii  $[0.0, 0.1, 0.2, 0.5, 1.0, 2.0, 4.0]$ . These radii operate on

the input space after data normalization. We bound perturbations by  $L_\infty$ -norm or by  $L_2$ -norm, with

$$L_\infty(x) = \max_{i=1,\dots,D} |x_i| \quad \text{and} \quad L_2(x) = \left( \sum_{i=1}^D x_i^2 \right)^{0.5}. \quad (10)$$

For  $L_\infty$ -norm it is obvious how to relate perturbation size  $\varepsilon$  with perturbed input images, because all inputs are standardized such that the values of their features are between 0 and 1. A perturbation of size  $\varepsilon = 0$  corresponds to the original input, while a perturbation of size  $\varepsilon = 1$  corresponds to the whole input space and allows to change all features to any value.

For  $L_2$ -norm the relation between perturbation size  $\varepsilon$  and perturbed input images is less obvious. To justify our choice for  $\varepsilon$  w.r.t. this norm, we relate perturbations size  $\varepsilon_2$  corresponding to  $L_2$ -norm with perturbations size  $\varepsilon_\infty$  corresponding to  $L_\infty$ -norm. First, we compute  $\varepsilon_2$ , such that the  $L_2$ -norm is the smallest super-set of the  $L_\infty$ -norm. Let us consider a perturbation of  $\varepsilon_\infty$ . The largest  $L_2$ -norm would be obtained if each feature is perturbed by  $\varepsilon_\infty$ . Thus, perturbation  $\varepsilon_2$ , such that  $L_2$  encloses  $L_\infty$  is  $\varepsilon_2 = \left( \sum_{i=1}^D \varepsilon_\infty^2 \right)^{0.5} = \sqrt{D} \varepsilon_\infty$ . For the MNIST-data set, with  $D = 28 \times 28$  input features  $L_2$ -norm with  $\varepsilon_2 = 28$  encloses  $L_\infty$ -norm with  $\varepsilon_\infty = 1$ .

Alternatively,  $\varepsilon_2$  can be computed such that the volume spanned by  $L_2$ -norm is equivalent to the one spanned by  $L_\infty$ -norm. Using that the volume spanned by  $L_\infty$ -norm is  $\varepsilon_\infty^D$  and the volume spanned by  $L_2$ -norm is  $\frac{\pi^{0.5D} \varepsilon_2^D}{\Gamma(0.5D+1)}$  (where  $\Gamma$  is the Gamma-function), we obtain volume equivalence if  $\varepsilon_2 = \Gamma(0.5D+1)^{\frac{1}{D}} \sqrt{\pi} \varepsilon_\infty$ . For the MNIST-data set, with  $D = 28 \times 28$  input features  $L_2$ -norm with  $\varepsilon_2 \approx 21.39$  is volume equivalent to  $L_\infty$ -norm with  $\varepsilon_\infty = 1$ .

## A.4 ADDITIONAL EXPERIMENTS

Table 7 and 8 illustrate that no DBU model maintains high accuracy under gradient-based label attacks. Accuracy under PGD attacks decreases more than under FGSM attacks, since PGD is stronger. Interestingly Noise attacks achieve also good performances with increasing Noise standard deviation. Note that the attack is not constraint to be with a given radius for Noise attacks.

Table 7: Accuracy under PGD label attacks.

Att. Rad.	0.0	0.1	0.2	0.5	1.0	2.0	4.0		0.0	0.1	0.2	0.5	1.0	2.0	4.0
MNIST								CIFAR10							
PostNet	<b>99.4</b>	<b>99.2</b>	<b>98.8</b>	96.8	89.6	53.8	13.0		89.5	73.5	51.7	13.2	2.2	0.8	0.3
PriorNet	99.3	99.1	<b>98.8</b>	97.4	<b>93.9</b>	<b>75.3</b>	4.8		88.2	<b>77.8</b>	<b>68.4</b>	<b>54.0</b>	<b>37.9</b>	<b>17.5</b>	<b>5.1</b>
DDNet	<b>99.4</b>	99.1	<b>98.8</b>	<b>97.5</b>	91.6	48.8	0.2		86.1	73.9	59.1	20.5	1.5	0.0	0.0
EvNet	99.2	98.9	98.4	96.8	92.4	73.1	<b>40.9</b>		<b>89.8</b>	71.7	48.8	11.5	2.7	1.5	0.4
Sensorless								Segment							
PostNet	98.3	13.1	6.4	4.0	<b>7.0</b>	<b>9.8</b>	<b>11.3</b>		98.9	82.8	<b>50.1</b>	<b>19.2</b>	<b>8.8</b>	<b>5.1</b>	<b>8.6</b>
PriorNet	<b>99.3</b>	16.5	5.6	1.2	0.4	0.2	1.6		<b>99.5</b>	90.7	47.6	7.8	0.2	0.0	0.4
DDNet	<b>99.3</b>	12.4	2.4	0.6	0.3	0.1	0.1		99.2	<b>90.8</b>	45.7	6.9	0.0	0.0	0.0
EvNet	99.0	<b>35.3</b>	<b>22.3</b>	<b>11.2</b>	<b>7.0</b>	5.2	4.0		99.3	91.8	54.0	10.3	0.8	0.5	0.6

Table 8: Accuracy under FGSM label attacks.

Att. Rad.	0.0	0.1	0.2	0.5	1.0	2.0	4.0		0.0	0.1	0.2	0.5	1.0	2.0	4.0
MNIST								CIFAR10							
PostNet	<b>99.4</b>	<b>99.2</b>	<b>98.9</b>	97.7	95.2	<b>90.1</b>	<b>79.2</b>		89.5	72.3	54.9	31.2	21.0	16.8	15.6
PriorNet	99.3	99.1	<b>98.9</b>	97.7	<b>95.8</b>	93.2	76.7		88.2	<b>77.3</b>	<b>70.1</b>	<b>59.4</b>	<b>52.3</b>	<b>48.5</b>	<b>46.8</b>
DDNet	<b>99.4</b>	<b>99.2</b>	<b>98.9</b>	<b>97.8</b>	94.7	79.2	25.2		86.1	73.0	60.2	32.5	14.6	7.1	6.0
EvNet	99.2	98.9	98.6	97.6	<b>95.8</b>	<b>90.1</b>	74.4		<b>89.8</b>	71.4	54.5	29.6	18.1	14.4	13.4
Sensorless								Segment							
PostNet	98.3	19.6	10.9	10.9	11.9	12.4	12.5		98.9	79.6	<b>57.3</b>	<b>31.5</b>	<b>18.4</b>	<b>20.6</b>	<b>19.9</b>
PriorNet	<b>99.3</b>	24.7	11.8	8.6	8.5	8.1	8.3		<b>99.5</b>	85.5	40.5	8.9	0.4	0.3	0.2
DDNet	<b>99.3</b>	18.0	8.2	6.5	5.4	6.7	7.8		99.2	86.4	36.2	11.9	0.9	0.0	0.0
EvNet	99.0	<b>42.0</b>	<b>28.0</b>	<b>17.5</b>	<b>13.7</b>	<b>13.6</b>	<b>14.9</b>		99.3	<b>90.6</b>	55.2	14.2	2.4	0.5	0.1

Table 9: Accuracy under Noise label attacks.

Noise Std	0.0	0.1	0.2	0.5	1.0	2.0	4.0		0.0	0.1	0.2	0.5	1.0	2.0	4.0
MNIST								CIFAR10							
PostNet	<b>99.4</b>	<b>98.6</b>	91.8	<b>14.9</b>	<b>1.3</b>	<b>0.1</b>	0.0		<b>91.7</b>	21.5	10.1	0.1	1.2	0.0	1.9
PriorNet	99.3	98.5	<b>95.7</b>	14.4	0.0	0.0	0.0		87.7	<b>28.1</b>	<b>11.2</b>	9.7	5.0	<b>8.5</b>	<b>9.0</b>
DDNet	<b>99.4</b>	<b>98.6</b>	92.4	13.3	0.7	0.0	0.0		81.7	23.0	<b>11.2</b>	<b>11.2</b>	<b>11.0</b>	7.8	6.7
EvNet	99.3	96.9	81.6	11.7	0.5	0.0	0.0		89.5	20.7	11.1	5.2	0.5	2.3	3.9
Sensorless								Segment							
PostNet	98.1	0.1	<b>3.7</b>	<b>11.7</b>	<b>11.7</b>	<b>11.7</b>	<b>11.7</b>		98.5	39.4	3.9	<b>1.8</b>	<b>12.1</b>	<b>20.3</b>	<b>22.1</b>
PriorNet	<b>99.3</b>	0.2	0.0	0.0	0.0	0.3	2.4		<b>99.4</b>	47.9	8.8	0.0	0.0	0.0	0.0
DDNet	99.0	<b>0.4</b>	0.1	0.0	0.0	0.0	0.0		99.1	50.0	<b>10.3</b>	0.0	0.0	0.3	0.0
EvNet	98.6	0.2	0.0	0.1	1.4	4.6	8.8		99.1	<b>50.3</b>	<b>10.3</b>	1.2	0.3	0.0	1.5

## A.4.1 UNCERTAINTY ESTIMATION UNDER LABEL ATTACKS

## Is high certainty a reliable indicator of correct predictions?

On non-perturbed data uncertainty estimates are an indicator of correctly classified samples, but if the input data is perturbed none of the DBU models maintains its high performance. Thus, uncertainty estimates are not a robust indicator of correctly labeled inputs.

Table 10: Certainty based on differential entropy under PGD label attacks (AUC-PR).

Att. Rad.	MNIST								Segment						
	0.0	0.1	0.2	0.5	1.0	2.0	4.0		0.0	0.1	0.2	0.5	1.0	2.0	4.0
PostNet	99.9	99.9	99.8	98.7	89.5	43.5	9.0		99.9	77.6	31.6	<b>11.1</b>	<b>5.3</b>	<b>4.4</b>	8.7
PriorNet	99.9	99.8	99.6	97.7	90.5	<b>69.1</b>	6.4		<b>100.0</b>	<b>96.8</b>	44.5	4.5	0.4	0.0	<b>15.2</b>
DDNet	<b>100.0</b>	<b>100.0</b>	<b>99.9</b>	<b>99.7</b>	<b>97.6</b>	50.2	0.1		<b>100.0</b>	<b>96.8</b>	<b>54.0</b>	4.3	0.0	0.0	0.0
EvNet	99.6	99.3	98.7	96.1	88.8	63.1	<b>31.7</b>		<b>100.0</b>	95.9	44.3	5.9	0.8	0.6	0.7

Table 2, 10, 11, and 12 illustrate that neither differential entropy nor precision, nor mutual information are a reliable indicator of correct predictions under PGD attacks. DBU-models achieve significantly

Table 11: Certainty based on precision  $\alpha_0$  under PGD label attacks (AUC-PR).

Att. Rad.	0.0	0.1	0.2	0.5	1.0	2.0	4.0		0.0	0.1	0.2	0.5	1.0	2.0	4.0
MNIST								CIFAR10							
PostNet	<b>100.0</b>	99.9	99.7	98.2	87.9	39.1	6.9		<b>98.7</b>	88.6	56.2	7.8	1.2	0.4	0.3
PriorNet	99.9	99.8	99.6	97.7	90.4	<b>69.1</b>	6.6		92.9	77.7	60.5	<b>37.6</b>	<b>24.9</b>	<b>11.3</b>	<b>3.0</b>
DDNet	<b>100.0</b>	<b>100.0</b>	<b>100.0</b>	<b>99.8</b>	<b>98.2</b>	51.1	0.1		97.6	<b>91.8</b>	<b>78.3</b>	18.1	0.8	0.0	0.0
EvNet	99.6	99.2	98.6	95.7	88.6	63.6	<b>32.6</b>		97.9	85.9	57.2	10.2	4.0	2.4	0.3
Sensorless								Segment							
PostNet	99.6	7.0	3.3	3.1	<b>6.9</b>	<b>9.8</b>	<b>11.3</b>		99.9	74.2	31.6	<b>11.1</b>	<b>5.0</b>	<b>4.2</b>	<b>8.6</b>
PriorNet	99.8	10.5	3.2	0.6	0.2	0.2	1.8		<b>100.0</b>	96.9	<b>45.2</b>	4.4	0.4	0.0	1.2
DDNet	99.8	8.7	1.3	0.3	0.2	0.1	0.2		<b>100.0</b>	<b>97.1</b>	45.0	4.1	0.0	0.0	0.0
EvNet	<b>99.9</b>	<b>23.2</b>	<b>13.2</b>	<b>6.0</b>	3.7	2.7	2.1		<b>100.0</b>	95.7	44.5	5.9	0.8	0.6	0.7

Table 12: Certainty based on mutual information under PGD label attacks (AUC-PR).

Att. Rad.	0.0	0.1	0.2	0.5	1.0	2.0	4.0		0.0	0.1	0.2	0.5	1.0	2.0	4.0
MNIST								CIFAR10							
PostNet	99.7	99.7	99.6	99.2	92.4	40.0	6.9		<b>97.3</b>	84.5	56.2	12.2	2.4	0.7	0.3
PriorNet	99.9	99.8	99.6	97.7	90.3	<b>68.9</b>	6.4		82.7	65.6	51.4	<b>35.5</b>	<b>24.4</b>	<b>11.0</b>	<b>2.9</b>
DDNet	<b>100.0</b>	<b>99.9</b>	<b>99.9</b>	<b>99.7</b>	<b>97.4</b>	50.2	0.1		96.9	<b>90.8</b>	<b>77.2</b>	18.8	0.8	0.0	0.0
EvNet	97.8	97.0	95.7	92.6	86.1	62.3	<b>28.9</b>		91.3	72.4	47.9	11.4	1.6	0.9	1.6
Sensorless								Segment							
PostNet	99.3	7.0	3.3	3.3	<b>7.0</b>	<b>9.8</b>	11.3		99.9	73.2	31.5	<b>11.1</b>	<b>5.0</b>	<b>4.3</b>	<b>8.7</b>
PriorNet	<b>99.8</b>	10.5	3.2	0.6	0.2	0.1	<b>11.8</b>		<b>100.0</b>	<b>96.6</b>	<b>45.2</b>	4.5	0.4	0.0	1.1
DDNet	99.6	8.6	1.3	0.3	0.2	0.1	0.1		<b>100.0</b>	96.5	42.4	4.1	0.0	0.0	0.0
EvNet	99.1	<b>22.0</b>	<b>12.6</b>	<b>5.9</b>	3.7	2.7	2.2		<b>100.0</b>	90.5	41.0	5.9	0.8	0.6	0.7

Table 13: Certainty based on differential entropy under FGSM label attacks (AUC-PR).

Att. Rad.	0.0	0.1	0.2	0.5	1.0	2.0	4.0		0.0	0.1	0.2	0.5	1.0	2.0	4.0
MNIST								CIFAR10							
PostNet	99.9	99.9	99.8	99.4	97.8	<b>92.1</b>	<b>83.2</b>		<b>98.5</b>	88.7	68.9	31.0	18.6	15.5	16.7
PriorNet	99.9	99.9	99.7	98.3	94.1	88.5	78.6		90.1	73.6	61.6	<b>46.1</b>	<b>38.5</b>	<b>35.6</b>	<b>37.3</b>
DDNet	<b>100.0</b>	<b>100.0</b>	<b>99.9</b>	<b>99.8</b>	<b>98.7</b>	86.4	23.0		97.3	<b>90.6</b>	<b>78.7</b>	39.4	13.7	6.0	5.1
EvNet	99.6	99.4	99.1	97.8	95.8	90.4	76.8		98.0	86.2	67.4	32.7	19.9	18.2	19.7
Sensorless								Segment							
PostNet	99.7	11.7	7.3	9.3	11.8	12.5	12.5		99.9	73.6	40.6	<b>23.7</b>	<b>17.2</b>	<b>19.8</b>	<b>20.2</b>
PriorNet	99.8	21.4	10.4	8.5	9.0	9.2	10.3		<b>100.0</b>	93.7	37.7	5.8	1.1	0.9	0.8
DDNet	99.7	18.5	5.4	4.3	4.2	5.7	7.9		<b>100.0</b>	<b>94.1</b>	42.9	7.2	1.0	0.0	0.0
EvNet	<b>99.9</b>	<b>44.8</b>	<b>29.2</b>	<b>18.2</b>	<b>15.1</b>	<b>14.9</b>	<b>15.5</b>		<b>100.0</b>	93.7	<b>48.7</b>	8.7	2.4	1.6	0.5

Table 14: Certainty based on differential entropy under Noise label attacks (AUC-PR).

Noise Std	0.0	0.1	0.2	0.5	1.0	2.0	4.0		0.0	0.1	0.2	0.5	1.0	2.0	4.0
MNIST								CIFAR10							
PostNet	99.9	99.8	99.6	<b>74.2</b>	<b>7.4</b>	<b>0.2</b>	0.0		<b>98.7</b>	<b>76.3</b>	24.3	0.4	4.9	0.0	1.7
PriorNet	99.9	99.9	<b>99.8</b>	73.4	0.0	0.0	0.0		85.0	27.8	15.9	<b>20.4</b>	7.0	<b>7.7</b>	<b>8.3</b>
DDNet	<b>100.0</b>	<b>99.9</b>	99.4	51.1	0.6	0.1	0.0		96.1	61.0	<b>39.8</b>	14.2	<b>11.3</b>	6.9	6.9
EvNet	99.5	98.4	88.5	20.2	0.9	0.0	0.0		97.5	66.1	21.4	7.7	2.3	3.0	3.8
Sensorless								Segment							
PostNet	99.7	0.3	<b>3.2</b>	<b>13.3</b>	<b>12.0</b>	<b>11.7</b>	<b>11.7</b>		99.9	53.9	4.8	1.8	<b>11.2</b>	<b>21.7</b>	<b>21.6</b>
PriorNet	<b>100.0</b>	0.3	0.0	0.0	0.0	7.8	11.5		<b>100.0</b>	<b>84.5</b>	15.6	0.0	0.0	0.0	0.0
DDNet	99.7	<b>0.9</b>	0.6	0.0	0.0	0.0	0.0		<b>100.0</b>	82.7	<b>23.9</b>	0.0	0.0	0.6	0.0
EvNet	99.8	0.3	0.0	0.1	1.7	5.5	10.0		<b>100.0</b>	78.3	19.0	<b>3.5</b>	0.5	0.0	1.7

better results when they are attacked by FGSM-attacks (Table 13), but as FGSM attacks provide much weaker adversarial examples than PGD attacks, this cannot be seen as real advantage.

### Can we use uncertainty estimates to detect attacks against the classification decision?

PGD attacks do not explicitly consider uncertainty during the computation of adversarial examples, but they seem to provide perturbed inputs with similar uncertainty as the original input.

Table 15: Attack-Detection based on differential entropy under PGD label attacks (AUC-PR).

Att. Rad.	MNIST						Segment					
	0.1	0.2	0.5	1.0	2.0	4.0	0.1	0.2	0.5	1.0	2.0	4.0
PostNet	57.7	66.3	83.4	90.5	79.0	50.1	<b>95.6</b>	73.5	<b>47.0</b>	<b>42.3</b>	<b>53.4</b>	<b>82.7</b>
PriorNet	<b>67.7</b>	<b>83.2</b>	<b>97.1</b>	<b>96.7</b>	92.1	82.9	86.7	83.3	38.0	31.3	30.8	31.5
DDNet	53.4	57.1	68.5	83.9	<b>96.0</b>	<b>86.3</b>	76.1	<b>83.5</b>	45.4	32.4	30.8	30.8
EvNet	54.8	59.0	68.5	75.9	72.6	59.8	94.9	80.9	41.5	32.5	31.1	31.1

Table 16: Attack-Detection based on precision  $\alpha_0$  under PGD label attacks (AUC-PR).

Att. Rad.	0.1	0.2	0.5	1.0	2.0	4.0	0.1	0.2	0.5	1.0	2.0	4.0
	MNIST						CIFAR10					
PostNet	63.3	75.7	92.6	95.1	75.3	39.5	<b>63.4</b>	<b>66.9</b>	42.1	32.9	31.6	31.2
PriorNet	<b>67.6</b>	<b>83.2</b>	<b>97.1</b>	<b>96.9</b>	<b>92.7</b>	<b>84.7</b>	53.3	56.0	55.6	<b>49.2</b>	42.2	35.4
DDNet	52.7	55.7	64.7	78.4	91.9	80.9	55.8	60.5	<b>57.3</b>	38.7	32.3	31.4
EvNet	49.1	48.0	45.1	42.7	41.8	39.2	48.4	46.9	46.3	46.3	<b>44.5</b>	<b>42.5</b>
	Sensorless						Segment					
	0.1	0.2	0.5	1.0	2.0	4.0	0.1	0.2	0.5	1.0	2.0	4.0
PostNet	39.8	35.8	35.4	<b>52.0</b>	<b>88.2</b>	<b>99.0</b>	<b>94.6</b>	70.3	<b>46.3</b>	<b>42.6</b>	<b>54.9</b>	<b>84.0</b>
PriorNet	40.9	35.1	32.0	31.1	30.7	30.7	82.7	82.6	39.4	31.6	30.8	30.8
DDNet	<b>47.7</b>	<b>40.3</b>	35.3	32.8	31.3	30.8	80.0	<b>86.0</b>	43.3	33.6	31.0	30.8
EvNet	45.4	39.7	<b>36.1</b>	34.8	34.7	36.0	90.9	72.4	40.4	32.4	31.1	31.1

Table 17: Attack-Detection based on mutual information under PGD label attacks (AUC-PR).

Att. Rad.	0.1	0.2	0.5	1.0	2.0	4.0	0.1	0.2	0.5	1.0	2.0	4.0
	MNIST						CIFAR10					
PostNet	42.2	37.5	36.7	54.5	70.5	70.3	52.2	52.1	50.0	<b>65.9</b>	<b>76.3</b>	<b>80.7</b>
PriorNet	<b>67.7</b>	<b>83.3</b>	<b>97.1</b>	<b>96.9</b>	92.6	<b>84.5</b>	54.0	56.9	56.3	49.7	42.4	35.5
DDNet	53.1	56.3	66.5	81.0	<b>94.0</b>	82.9	<b>56.0</b>	<b>60.8</b>	<b>57.4</b>	38.2	32.1	31.3
EvNet	49.1	48.0	45.2	42.9	41.9	39.3	48.7	47.3	46.3	46.0	44.1	42.2
	Sensorless						Segment					
	0.1	0.2	0.5	1.0	2.0	4.0	0.1	0.2	0.5	1.0	2.0	4.0
PostNet	<b>75.3</b>	<b>76.6</b>	<b>66.5</b>	<b>57.7</b>	<b>85.6</b>	<b>98.7</b>	<b>94.8</b>	73.5	<b>55.9</b>	<b>47.9</b>	<b>58.0</b>	<b>84.0</b>
PriorNet	40.7	35.0	32.0	31.0	30.7	30.7	83.5	82.7	39.2	31.6	30.8	30.8
DDNet	48.0	40.0	35.2	32.6	31.2	30.8	82.4	<b>88.1</b>	43.4	33.4	30.9	30.8
EvNet	45.5	39.7	36.1	34.8	34.7	36.0	91.7	72.9	40.5	32.4	31.1	31.1

FGSM and Noise attacks are easier to detect, but also weaker than PGD attacks. This suggests that DBU models are capable of detecting weak attacks by using uncertainty estimation.

Table 18: Attack-Detection based on differential entropy under FGSM label attacks (AUC-PR).

Att. Rad.	0.1	0.2	0.5	1.0	2.0	4.0	0.1	0.2	0.5	1.0	2.0	4.0
	MNIST						CIFAR10					
PostNet	55.9	61.8	74.8	84.0	88.9	89.9	<b>62.1</b>	<b>67.2</b>	65.7	63.1	65.4	73.8
PriorNet	<b>67.4</b>	<b>82.4</b>	<b>96.9</b>	<b>98.3</b>	<b>98.9</b>	<b>99.6</b>	58.4	63.1	68.5	<b>70.1</b>	68.5	62.5
DDNet	53.6	57.3	68.3	82.6	95.6	98.7	57.2	62.9	<b>69.1</b>	68.7	<b>69.7</b>	<b>76.5</b>
EvNet	54.1	57.4	63.8	67.6	68.6	69.9	57.8	61.7	63.3	62.9	65.7	72.5
	Sensorless						Segment					
	0.1	0.2	0.5	1.0	2.0	4.0	0.1	0.2	0.5	1.0	2.0	4.0
PostNet	<b>98.4</b>	<b>99.8</b>	<b>99.9</b>	<b>99.9</b>	<b>99.9</b>	<b>99.9</b>	<b>96.9</b>	<b>93.9</b>	<b>99.5</b>	<b>99.9</b>	<b>100.0</b>	<b>100.0</b>
PriorNet	48.7	38.6	32.7	32.9	38.6	44.3	89.0	80.8	46.7	37.2	33.7	32.4
DDNet	61.5	47.8	37.1	33.1	32.4	33.2	79.6	86.2	60.2	47.5	36.6	31.6
EvNet	67.3	65.5	72.3	73.4	75.3	79.1	95.7	87.2	59.3	51.7	51.1	53.5

Table 19: Attack-Detection based on differential entropy under Noise label attacks (AUC-PR).

Noise Std.	0.1	0.2	0.5	1.0	2.0	4.0		0.1	0.2	0.5	1.0	2.0	4.0
MNIST								CIFAR10					
PostNet	51.3	65.3	93.8	95.1	95.2	95.2		<b>80.8</b>	<b>84.5</b>	<b>97.6</b>	<b>99.5</b>	99.3	98.2
PriorNet	32.5	36.8	88.9	99.6	99.7	92.7		34.7	32.3	34.3	60.3	95.5	<b>100.0</b>
DDNet	<b>60.7</b>	<b>87.6</b>	<b>99.8</b>	<b>100.0</b>	<b>99.9</b>	<b>99.8</b>		59.1	62.6	81.5	98.6	<b>99.8</b>	98.7
EvNet	51.2	55.7	66.9	70.3	68.0	67.1		75.7	78.6	88.2	97.8	96.4	95.6
Sensorless								Segment					
PostNet	<b>99.8</b>	<b>100.0</b>	<b>100.0</b>	<b>100.0</b>	<b>100.0</b>	<b>100.0</b>		<b>95.6</b>	<b>99.4</b>	<b>100.0</b>	<b>100.0</b>	<b>100.0</b>	<b>100.0</b>
PriorNet	42.0	33.8	31.5	34.7	43.7	47.0		56.7	56.7	39.8	33.7	31.9	33.7
DDNet	53.4	43.5	34.3	31.6	32.5	36.1		57.0	58.9	43.1	33.7	31.5	31.3
EvNet	67.1	78.8	88.3	95.4	96.9	97.8		60.8	63.5	61.2	64.8	73.7	85.2

## A.4.2 ATTACKING UNCERTAINTY ESTIMATION

## Are uncertainty estimates a robust feature for OOD detection?

Using uncertainty estimation to distinguish between ID and OOD data is not robust as shown in the following tables.

Table 20: OOD detection based on differential entropy under PGD uncertainty attacks against differential entropy on ID data and OOD data (AUC-PR).

Att. Rad.	ID-Attack (non-attacked OOD)							OOD-Attack (non-attacked ID)						
	0.0	0.1	0.2	0.5	1.0	2.0	4.0	0.0	0.1	0.2	0.5	1.0	2.0	4.0
MNIST – KMNIST														
PostNet	94.5	94.1	93.9	91.1	77.1	44.0	31.9	94.5	93.1	91.4	82.1	62.2	50.7	48.8
PriorNet	<b>99.6</b>	<b>99.4</b>	<b>99.1</b>	<b>97.8</b>	<b>93.8</b>	<b>77.6</b>	<b>32.0</b>	<b>99.6</b>	<b>99.4</b>	<b>99.1</b>	98.0	94.6	85.5	<b>73.9</b>
DDNet	99.3	99.1	98.9	<b>97.8</b>	93.5	63.3	30.7	99.3	99.1	99.0	<b>98.3</b>	<b>96.7</b>	<b>91.3</b>	73.8
EvNet	69.0	67.1	65.6	61.8	57.4	50.9	43.6	69.0	55.8	48.0	39.4	36.2	34.9	34.4
Seg. – Seg. class sky														
PostNet	<b>99.0</b>	<b>80.7</b>	<b>53.5</b>	<b>38.0</b>	<b>34.0</b>	<b>41.6</b>	<b>49.5</b>	<b>99.0</b>	<b>88.4</b>	69.2	45.1	<b>36.4</b>	<b>42.6</b>	<b>75.4</b>
PriorNet	34.8	31.4	30.9	30.8	30.8	30.8	30.8	34.8	31.8	31.0	30.8	30.8	30.8	32.1
DDNet	31.5	30.9	30.8	30.8	30.8	30.8	30.8	31.5	31.0	30.8	30.8	30.8	30.8	30.8
EvNet	92.5	67.2	43.2	31.6	30.9	30.9	31.2	92.5	86.1	<b>82.7</b>	<b>48.9</b>	32.7	30.9	30.9

Table 21: OOD detection under PGD uncertainty attacks against differential entropy on ID data and OOD data (AUC-ROC).

Att. Rad.	ID-Attack (non-attacked OOD)							OOD-Attack (non-attacked ID)						
	0.0	0.1	0.2	0.5	1.0	2.0	4.0	0.0	0.1	0.2	0.5	1.0	2.0	4.0
MNIST – KMNIST														
PostNet	91.6	91.3	91.9	91.5	80.2	38.8	9.2	91.6	90.4	89.0	81.6	62.6	45.0	43.1
PriorNet	<b>99.8</b>	<b>99.7</b>	<b>99.5</b>	<b>99.0</b>	<b>97.1</b>	<b>81.1</b>	8.7	<b>99.8</b>	<b>99.7</b>	<b>99.6</b>	<b>99.1</b>	<b>97.7</b>	<b>93.0</b>	<b>84.9</b>
DDNet	99.2	98.9	98.6	97.3	92.1	58.2	1.2	99.2	99.0	98.8	97.9	95.8	89.1	69.3
EvNet	81.2	79.6	78.2	74.6	69.5	58.7	<b>43.0</b>	81.2	67.2	54.8	35.4	25.5	20.7	18.5
CIFAR10 – SVHN														
PostNet	87.0	71.9	56.3	<b>30.2</b>	<b>20.2</b>	<b>15.0</b>	<b>9.7</b>	87.0	71.0	54.3	33.5	30.3	26.2	19.4
PriorNet	62.4	48.2	35.9	13.8	3.6	0.9	0.3	62.4	48.0	35.6	14.8	6.6	3.4	1.6
DDNet	<b>87.0</b>	<b>76.0</b>	<b>63.6</b>	29.3	6.1	1.1	0.4	<b>87.0</b>	<b>78.1</b>	<b>66.1</b>	26.2	5.1	0.7	0.1
EvNet	<b>88.0</b>	69.1	51.7	24.6	15.5	9.5	4.2	<b>88.0</b>	72.0	60.7	<b>47.9</b>	<b>42.1</b>	<b>33.3</b>	<b>24.0</b>
Sens. – Sens. class 10, 11														
PostNet	<b>85.3</b>	<b>49.1</b>	<b>38.1</b>	<b>7.8</b>	<b>8.2</b>	8.2	8.2	<b>85.3</b>	<b>57.2</b>	<b>54.0</b>	<b>27.3</b>	<b>31.5</b>	<b>86.7</b>	<b>99.5</b>
PriorNet	28.1	0.8	0.3	0.4	1.6	<b>8.4</b>	<b>26.8</b>	28.1	2.5	0.7	0.2	2.3	18.9	41.0
DDNet	21.0	3.0	0.9	0.4	0.6	2.1	7.3	21.0	4.4	2.1	1.9	2.2	2.2	4.1
EvNet	74.2	21.4	12.2	4.3	1.4	0.6	0.3	74.2	45.3	38.5	19.6	9.6	12.1	26.0
Seg. – Seg. class sky														
PostNet	<b>99.2</b>	<b>84.7</b>	<b>55.5</b>	<b>23.0</b>	<b>9.7</b>	<b>4.4</b>	<b>4.7</b>	<b>99.2</b>	<b>92.1</b>	<b>77.1</b>	41.5	<b>24.9</b>	<b>41.0</b>	<b>80.8</b>
PriorNet	17.1	4.4	1.3	0.0	0.0	0.0	0.1	17.1	5.9	1.5	0.1	0.0	0.1	5.8
DDNet	4.1	1.1	0.0	0.0	0.0	0.0	0.0	4.1	1.8	0.4	0.0	0.0	0.0	0.0
EvNet	91.2	54.5	23.3	3.9	0.9	0.4	0.2	91.2	82.9	76.4	<b>42.2</b>	9.7	0.8	0.6



Table 22: OOD detection (AU-PR) under PGD uncertainty attacks against precision  $\alpha_0$  on ID data and OOD data.

Att. Rad.	ID-Attack (non-attacked OOD)							OOD-Attack (non-attacked ID)						
	0.0	0.1	0.2	0.5	1.0	2.0	4.0	0.0	0.1	0.2	0.5	1.0	2.0	4.0
<b>MNIST – KMNIST</b>														
PostNet	98.4	97.4	96.0	88.8	70.9	39.3	31.3	98.4	97.2	95.2	82.8	52.6	34.3	32.1
PriorNet	<b>99.6</b>	<b>99.5</b>	<b>99.2</b>	<b>98.0</b>	<b>94.1</b>	<b>76.0</b>	31.1	<b>99.6</b>	<b>99.5</b>	<b>99.2</b>	<b>98.2</b>	<b>95.3</b>	<b>87.5</b>	<b>75.6</b>
DDNet	97.2	96.7	96.1	93.8	86.4	53.2	31.0	97.2	96.7	96.2	94.5	91.1	82.9	64.6
EvNet	39.8	39.2	38.8	37.9	37.1	36.3	<b>35.4</b>	39.8	34.5	32.5	31.2	31.0	30.9	31.0
<b>CIFAR10 – SVHN</b>														
PostNet	<b>82.4</b>	63.8	46.1	22.3	17.4	16.7	16.4	<b>82.4</b>	61.8	41.5	21.8	<b>19.8</b>	<b>17.5</b>	<b>15.8</b>
PriorNet	37.9	25.0	19.2	15.8	15.4	15.4	15.4	37.9	25.9	19.4	15.6	15.4	15.4	15.4
DDNet	81.1	<b>70.1</b>	<b>58.4</b>	<b>30.0</b>	16.7	15.5	15.4	81.1	<b>71.2</b>	<b>59.9</b>	<b>27.8</b>	16.5	15.5	15.4
EvNet	34.7	27.4	25.4	22.0	<b>19.7</b>	<b>18.1</b>	<b>17.1</b>	34.7	19.4	18.1	17.1	16.8	16.2	15.7
<b>Sens. – Sens. class 10, 11</b>														
PostNet	<b>77.4</b>	<b>39.6</b>	<b>35.9</b>	<b>31.7</b>	<b>44.4</b>	<b>44.4</b>	<b>44.4</b>	<b>77.4</b>	40.3	<b>38.6</b>	29.5	<b>34.0</b>	<b>79.4</b>	<b>97.4</b>
PriorNet	35.9	27.0	26.8	26.8	27.5	36.2		35.9	27.7	27.0	26.7	26.6	26.5	26.5
DDNet	55.6	34.4	31.7	30.4	29.5	30.2	33.4	55.6	<b>40.9</b>	34.1	28.0	26.9	26.6	26.5
EvNet	66.3	33.3	29.7	27.0	27.1	29.2	33.9	66.3	39.3	37.1	<b>31.3</b>	28.3	28.4	29.7
<b>Seg. – Seg. class sky</b>														
PostNet	<b>98.4</b>	74.8	51.0	<b>37.2</b>	<b>32.8</b>	<b>43.5</b>	<b>49.9</b>	<b>98.4</b>	84.7	66.1	42.4	34.8	<b>40.9</b>	<b>71.2</b>
PriorNet	32.1	30.9	30.8	30.8	30.8	30.8	30.8	32.1	31.0	30.8	30.8	30.8	30.8	30.8
DDNet	31.0	30.8	30.8	30.8	30.8	30.8	30.8	31.0	30.8	30.8	30.8	30.8	30.8	30.8
EvNet	98.3	<b>83.0</b>	<b>60.5</b>	34.0	31.0	30.8	30.8	98.3	<b>94.4</b>	<b>88.8</b>	<b>65.6</b>	<b>37.0</b>	31.4	30.9

Table 23: OOD detection (AUC-ROC) under PGD uncertainty attacks against precision  $\alpha_0$  on ID data and OOD data.

Att. Rad.	ID-Attack (non-attacked OOD)							OOD-Attack (non-attacked ID)						
	0.0	0.1	0.2	0.5	1.0	2.0	4.0	0.0	0.1	0.2	0.5	1.0	2.0	4.0
<b>MNIST – KMNIST</b>														
PostNet	98.4	97.6	96.4	90.9	74.0	28.9	6.3	98.4	97.6	96.3	89.0	61.3	19.6	9.7
PriorNet	<b>99.8</b>	<b>99.7</b>	<b>99.6</b>	<b>99.1</b>	<b>97.2</b>	<b>79.4</b>	4.4	<b>99.8</b>	<b>99.7</b>	<b>99.6</b>	<b>99.2</b>	<b>98.0</b>	<b>93.9</b>	<b>85.8</b>
DDNet	96.5	95.9	95.1	92.0	82.6	44.3	3.5	96.5	95.9	95.2	92.9	88.6	78.7	59.4
EvNet	35.9	34.1	32.8	30.1	27.4	24.6	<b>21.4</b>	35.9	18.7	10.4	3.7	2.0	1.7	2.0
<b>CIFAR10 – SVHN</b>														
PostNet	<b>87.4</b>	71.2	54.8	29.2	19.0	14.0	9.4	<b>87.4</b>	71.4	54.1	30.1	<b>25.8</b>	<b>17.5</b>	<b>5.8</b>
PriorNet	45.6	31.1	20.4	6.3	1.4	0.3	0.1	45.6	32.2	21.7	5.4	1.0	0.3	0.1
DDNet	84.9	<b>73.8</b>	<b>61.8</b>	30.2	9.3	3.0	0.8	84.9	<b>76.6</b>	<b>66.2</b>	<b>34.6</b>	10.4	2.3	0.3
EvNet	61.2	49.4	45.2	<b>37.6</b>	<b>30.5</b>	<b>23.4</b>	<b>17.0</b>	61.2	29.4	23.0	16.8	14.2	10.2	5.5
<b>Sens. – Sens. class 10, 11</b>														
PostNet	<b>87.2</b>	<b>48.8</b>	<b>37.3</b>	4.1	0.7	0.7	0.7	<b>87.2</b>	<b>50.0</b>	<b>45.4</b>	16.5	<b>27.6</b>	<b>81.9</b>	<b>98.0</b>
PriorNet	37.3	3.5	2.4	2.2	2.9	6.3	<b>19.2</b>	37.3	8.0	3.6	1.4	0.6	0.1	0.0
DDNet	55.2	23.7	17.7	<b>14.1</b>	<b>12.5</b>	<b>12.7</b>	15.7	55.2	37.1	27.7	9.4	2.5	0.6	0.1
EvNet	75.5	30.8	18.2	5.8	1.6	0.6	0.2	75.5	47.8	41.9	<b>24.1</b>	10.2	10.2	15.6
<b>Seg. – Seg. class sky</b>														
PostNet	<b>98.6</b>	77.7	<b>50.8</b>	<b>20.3</b>	<b>8.2</b>	<b>1.3</b>	<b>0.5</b>	<b>98.6</b>	88.9	73.4	36.2	19.4	<b>36.7</b>	<b>75.2</b>
PriorNet	8.5	1.3	0.2	0.0	0.0	0.0	0.1	8.5	2.0	0.4	0.0	0.0	0.0	0.0
DDNet	2.2	0.3	0.0	0.0	0.0	0.0	0.0	2.2	0.5	0.1	0.0	0.0	0.0	0.0
EvNet	97.7	<b>78.4</b>	47.7	9.9	1.2	0.2	0.1	97.7	<b>93.5</b>	<b>86.9</b>	<b>62.2</b>	<b>21.5</b>	3.7	1.0

Table 24: OOD detection (AU-PR) under PGD uncertainty attacks against distributional uncertainty on ID data and OOD data.

Att. Rad.	ID-Attack (non-attacked OOD)							OOD-Attack (non-attacked ID)						
	0.0	0.1	0.2	0.5	1.0	2.0	4.0	0.0	0.1	0.2	0.5	1.0	2.0	4.0
<b>MNIST – KMIST</b>														
PostNet	80.5	76.2	73.4	69.1	66.6	65.4	<b>60.2</b>	80.5	72.1	63.9	43.9	33.0	30.9	30.8
PriorNet	<b>99.6</b>	<b>99.4</b>	<b>99.2</b>	<b>98.0</b>	<b>94.1</b>	<b>76.3</b>	31.2	<b>99.6</b>	<b>99.4</b>	<b>99.2</b>	<b>98.2</b>	<b>95.2</b>	<b>87.2</b>	<b>75.2</b>
DDNet	98.4	98.1	97.7	95.8	89.5	56.2	30.9	98.4	98.1	97.8	96.5	93.8	86.3	67.7
EvNet	40.1	39.5	39.1	38.2	37.3	36.5	35.6	40.1	34.6	32.6	31.3	31.0	31.0	31.1
<b>CIFAR10 – SVHN</b>														
PostNet	64.2	44.7	37.5	<b>31.1</b>	<b>28.5</b>	<b>25.0</b>	<b>19.3</b>	64.2	31.0	19.5	16.3	16.4	<b>16.5</b>	<b>16.3</b>
PriorNet	40.8	27.4	20.4	15.9	15.4	15.4	15.4	40.8	28.3	21.1	15.9	15.4	15.4	15.4
DDNet	<b>82.0</b>	<b>71.0</b>	<b>59.1</b>	29.9	16.6	15.5	15.4	<b>82.0</b>	<b>72.2</b>	<b>60.3</b>	<b>26.3</b>	16.2	15.4	15.4
EvNet	36.4	28.7	26.5	22.8	20.2	18.4	17.2	36.4	19.8	18.3	17.2	<b>16.9</b>	16.2	15.7
<b>Sens. – Sens. class 10, 11</b>														
PostNet	<b>79.1</b>	<b>40.3</b>	<b>35.9</b>	<b>33.0</b>	<b>45.5</b>	<b>45.5</b>	45.5	<b>79.1</b>	<b>47.3</b>	<b>43.7</b>	<b>36.5</b>	<b>37.9</b>	<b>74.6</b>	<b>96.5</b>
PriorNet	35.5	26.8	26.7	26.9	29.6	43.7	<b>68.7</b>	35.5	27.5	26.9	26.7	26.6	26.5	26.5
DDNet	52.9	31.7	29.8	29.1	28.4	30.1	37.6	52.9	38.4	31.5	27.5	26.8	26.6	26.5
EvNet	66.3	33.3	29.6	27.0	27.2	29.3	35.2	66.3	39.3	37.1	31.3	28.3	28.4	29.7
<b>Seg. – Seg. class sky</b>														
PostNet	98.0	76.3	53.1	<b>37.4</b>	<b>32.9</b>	<b>44.6</b>	<b>50.2</b>	98.0	83.5	64.8	41.8	35.4	<b>43.1</b>	<b>71.3</b>
PriorNet	32.3	30.9	30.8	30.8	30.8	32.5	45.0	32.3	31.0	30.8	30.8	30.8	30.8	30.8
DDNet	30.9	30.8	30.8	30.8	30.8	30.8	30.8	30.9	30.8	30.8	30.8	30.8	30.8	30.8
EvNet	<b>98.1</b>	<b>82.1</b>	<b>59.1</b>	33.8	31.0	30.8	30.8	<b>98.1</b>	<b>93.8</b>	<b>88.2</b>	<b>64.5</b>	<b>36.4</b>	31.3	31.0

Table 25: OOD detection (AUC-ROC) under PGD uncertainty attacks against distributional uncertainty on ID data and OOD data.

Att. Rad.	ID-Attack (non-attacked OOD)							OOD-Attack (non-attacked ID)						
	0.0	0.1	0.2	0.5	1.0	2.0	4.0	0.0	0.1	0.2	0.5	1.0	2.0	4.0
<b>MNIST – KMIST</b>														
PostNet	90.1	88.0	86.2	82.2	79.0	77.1	<b>66.1</b>	90.1	84.5	77.2	46.4	12.9	2.7	2.4
PriorNet	<b>99.8</b>	<b>99.7</b>	<b>99.6</b>	<b>99.1</b>	<b>97.2</b>	<b>79.7</b>	4.7	<b>99.8</b>	<b>99.7</b>	<b>99.6</b>	<b>99.2</b>	<b>97.9</b>	<b>93.7</b>	<b>85.6</b>
DDNet	98.1	97.7	97.2	94.8	87.0	48.7	3.0	98.1	97.8	97.3	95.8	92.3	83.3	63.3
EvNet	36.8	35.0	33.7	30.9	28.2	25.3	22.1	36.8	19.3	10.7	3.9	2.1	1.8	2.2
<b>CIFAR10 – SVHN</b>														
PostNet	82.9	67.7	59.2	<b>51.3</b>	<b>47.7</b>	<b>40.1</b>	<b>24.2</b>	82.9	51.9	26.2	8.9	9.5	<b>11.1</b>	<b>9.9</b>
PriorNet	48.0	33.6	22.5	7.1	1.6	0.3	0.1	48.0	34.8	24.0	6.7	1.6	0.6	0.2
DDNet	<b>85.9</b>	<b>74.9</b>	<b>62.7</b>	30.1	8.3	2.3	0.6	<b>85.9</b>	<b>77.6</b>	<b>66.9</b>	<b>32.1</b>	8.0	1.5	0.2
EvNet	63.3	51.4	47.1	39.3	32.1	24.9	17.9	63.3	31.1	24.4	17.7	<b>15.0</b>	10.7	5.7
<b>Sens. – Sens. class 10, 11</b>														
PostNet	<b>87.1</b>	<b>50.9</b>	<b>37.8</b>	5.5	4.5	4.5	4.5	<b>87.1</b>	<b>55.3</b>	<b>51.1</b>	<b>34.4</b>	<b>38.9</b>	<b>79.7</b>	<b>97.9</b>
PriorNet	36.5	2.9	1.8	1.8	5.2	<b>21.5</b>	<b>52.8</b>	36.5	7.3	3.0	1.3	0.5	0.1	0.0
DDNet	52.3	18.7	13.1	<b>10.3</b>	<b>9.3</b>	10.8	18.4	52.3	33.1	22.0	6.7	2.2	0.6	0.1
EvNet	75.5	30.7	18.1	5.8	1.6	0.6	0.8	75.5	47.7	41.8	23.8	10.3	10.2	15.8
<b>Seg. – Seg. class sky</b>														
PostNet	<b>98.6</b>	<b>78.3</b>	<b>51.9</b>	<b>20.5</b>	<b>8.3</b>	<b>2.1</b>	1.7	<b>98.6</b>	88.8	73.1	35.9	<b>21.4</b>	<b>39.9</b>	<b>75.9</b>
PriorNet	9.4	1.6	0.3	0.0	0.0	1.8	<b>15.4</b>	9.4	2.4	0.4	0.0	0.0	0.0	0.0
DDNet	1.3	0.2	0.0	0.0	0.0	0.0	0.0	1.3	0.2	0.0	0.0	0.0	0.0	0.0
EvNet	97.4	77.1	45.9	9.4	1.3	0.2	0.1	97.4	<b>92.9</b>	<b>86.1</b>	<b>60.9</b>	20.4	3.0	1.2

Table 26: OOD detection (AU-PR) under FGSM uncertainty attacks against differential entropy on ID data and OOD data.

Att. Rad.	ID-Attack (non-attacked OOD)							OOD-Attack (non-attacked ID)						
	0.0	0.1	0.2	0.5	1.0	2.0	4.0	0.0	0.1	0.2	0.5	1.0	2.0	4.0
<b>MNIST – KMNIST</b>														
PostNet	94.5	94.2	94.1	93.5	89.9	81.2	<b>71.6</b>	94.5	93.3	92.0	87.6	81.1	75.7	75.7
PriorNet	<b>99.6</b>	<b>99.4</b>	<b>99.2</b>	<b>98.1</b>	<b>95.6</b>	<b>90.0</b>	65.3	<b>99.6</b>	<b>99.4</b>	<b>99.2</b>	<b>98.6</b>	97.5	<b>95.9</b>	<b>94.4</b>
DDNet	99.3	99.1	98.9	98.0	95.4	80.9	48.2	99.3	99.2	99.0	98.5	<b>97.6</b>	95.5	92.0
EvNet	69.0	67.4	66.2	64.0	61.9	59.8	56.70	9.0	60.1	56.5	53.4	52.7	52.9	53.5
<b>CIFAR10 – SVHN</b>														
PostNet	81.8	66.2	61.6	<b>64.2</b>	<b>65.7</b>	61.3	48.4	81.8	63.1	51.9	43.4	46.6	<b>61.7</b>	<b>77.0</b>
PriorNet	54.4	40.6	33.8	27.0	25.5	27.2	35.5	54.4	42.3	36.8	30.6	28.3	29.5	32.1
DDNet	<b>82.8</b>	<b>71.9</b>	<b>64.6</b>	53.8	50.2	47.8	41.0	<b>82.8</b>	<b>71.5</b>	<b>60.5</b>	39.1	31.4	41.2	66.6
EvNet	80.3	67.8	64.0	61.9	61.6	57.4	<b>49.6</b>	80.3	59.2	51.5	<b>46.7</b>	<b>49.0</b>	56.3	64.6
<b>Sens. – Sens. class 10, 11</b>														
PostNet	<b>74.5</b>	40.6	37.2	31.4	38.1	44.9	45.9	<b>74.5</b>	<b>99.6</b>	<b>99.8</b>	<b>99.9</b>	<b>99.9</b>	<b>99.9</b>	<b>99.9</b>
PriorNet	32.3	35.7	<b>57.6</b>	<b>83.1</b>	<b>88.8</b>	79.7	70.0	32.3	28.3	28.1	27.6	28.0	32.7	38.5
DDNet	31.7	31.3	44.4	70.3	87.9	<b>92.5</b>	<b>91.9</b>	31.7	28.8	29.3	29.1	27.7	27.9	28.01
EvNet	66.5	<b>45.7</b>	46.8	42.3	42.0	41.4	41.8	66.5	54.7	66.5	76.2	71.1	75.3	75.8
<b>Seg. – Seg. class sky</b>														
PostNet	<b>99.0</b>	<b>80.8</b>	<b>66.4</b>	43.6	37.0	35.5	43.0	<b>99.0</b>	<b>94.8</b>	<b>92.0</b>	<b>98.5</b>	<b>99.7</b>	<b>100.0</b>	<b>100.0</b>
PriorNet	34.8	31.2	31.4	46.3	<b>74.0</b>	<b>88.8</b>	<b>94.5</b>	34.8	31.6	31.0	31.2	30.9	30.8	30.8
DDNet	31.5	30.8	30.8	30.9	37.9	56.2	84.3	31.5	30.9	30.8	30.8	30.8	30.8	30.8
EvNet	92.5	64.9	54.6	<b>66.6</b>	69.5	69.6	64.6	92.5	85.9	83.0	66.3	66.1	61.1	56.8

Table 27: OOD detection (AU-PR) under Noise uncertainty attacks against differential entropy on ID data and OOD data.

Noise Std	ID-Attack (non-attacked OOD)							OOD-Attack (non-attacked ID)						
	0.0	0.1	0.2	0.5	1.0	2.0	4.0	0.0	0.1	0.2	0.5	1.0	2.0	4.0
<b>MNIST – KMNIST</b>														
PostNet	93.0	94.2	82.3	34.4	31.6	31.0	30.9	92.2	91.8	91.5	92.3	92.7	93.2	93.5
PriorNet	<b>99.7</b>	<b>99.6</b>	<b>96.7</b>	<b>40.0</b>	<b>40.6</b>	<b>45.7</b>	<b>55.6</b>	<b>99.5</b>	97.3	96.5	99.4	<b>100.0</b>	99.5	72.4
DDNet	99.1	97.5	81.2	31.3	31.0	30.9	31.2	99.0	<b>98.8</b>	<b>99.2</b>	<b>99.8</b>	99.9	<b>99.8</b>	<b>99.1</b>
EvNet	65.5	60.5	51.4	35.3	34.5	35.5	35.0	62.5	47.2	40.9	35.1	34.6	33.5	34.9
<b>CIFAR10 – SVHN</b>														
PostNet	88.5	41.4	39.8	31.0	30.7	31.6	33.9	88.5	<b>86.6</b>	<b>81.9</b>	<b>93.0</b>	<b>98.5</b>	98.6	97.3
PriorNet	73.3	88.3	<b>95.3</b>	<b>92.4</b>	<b>70.4</b>	30.9	30.8	73.3	31.6	30.9	31.7	51.8	94.3	<b>100.0</b>
DDNet	87.3	69.3	78.4	55.2	31.6	30.7	31.4	87.3	55.8	57.9	73.9	97.3	<b>99.5</b>	97.2
EvNet	<b>92.4</b>	<b>56.8</b>	53.8	33.4	30.9	<b>32.9</b>	<b>36.6</b>	<b>92.4</b>	73.7	73.5	77.7	93.7	92.5	92.1
<b>Sens. – Sens. class 10, 11</b>														
PostNet	<b>85.3</b>	<b>30.8</b>	<b>39.4</b>	50.0	50.0	50.0	50.0	<b>85.3</b>	<b>98.9</b>	<b>100.0</b>	<b>100.0</b>	<b>100.0</b>	<b>100.0</b>	<b>100.0</b>
PriorNet	32.3	<b>30.8</b>	34.9	<b>83.7</b>	<b>77.7</b>	49.8	<b>80.3</b>	32.3	30.7	30.7	32.5	40.1	49.9	47.6
DDNet	31.1	30.7	30.7	32.4	58.8	<b>88.1</b>	74.3	31.1	30.7	30.7	30.7	30.8	31.6	39.1
EvNet	80.3	<b>30.8</b>	31.2	37.9	46.3	50.0	50.0	80.3	34.6	38.4	53.9	69.3	78.8	81.5
<b>Seg. – Seg. class sky</b>														
PostNet	<b>99.9</b>	<b>41.8</b>	30.8	<b>34.5</b>	<b>49.1</b>	50.0	50.0	<b>99.9</b>	<b>97.4</b>	<b>96.6</b>	<b>99.5</b>	<b>100.0</b>	<b>100.0</b>	<b>100.0</b>
PriorNet	31.0	30.8	30.8	30.8	32.7	<b>69.0</b>	78.3	31.0	30.8	30.8	30.8	30.9	31.1	32.4
DDNet	30.8	30.8	30.8	30.8	30.8	58.2	<b>91.3</b>	30.8	30.8	30.8	30.8	30.8	30.8	31.9
EvNet	99.1	38.1	<b>32.2</b>	30.8	30.8	32.2	37.5	99.1	95.6	87.6	58.0	44.9	46.6	53.8

## A.5 ROBUST TRAINING FOR DBU MODELS &amp; ID/OOD VERIFICATION

Table 5 and 29 on adversarial training illustrate that there is a jump between ID-verification and OOD-verification, where robustness on ID data drops while robustness on OOD data increases. These jumps are observed for each model and each training (normal, noise-based, adversarial with label attacks, adversarial with uncertainty attacks). Thus, either ID-verification or OOD-verification perform well, depending on the chosen threshold.

In contrast to that, adversarial training improves robustness w.r.t. the predicted class label for most pair model/data set (Fig. 7, 32).

Table 28: Randomized smoothing verification of CIFAR10 (ID data) and SVHN (OOD data) harmonic mean.

	0.1	0.2	0.5
<b>adv. train. loss: None</b>			
PriorNet	26.7	3.7	0.0
PostNet	35.9	34.1	0.0
DDNet	45.2	18.1	<b>46.6</b>
EvNet	<b>47.6</b>	<b>45.4</b>	22.6
<b>adv. train. loss: crossentropy</b>			
PriorNet	0.2	0.0	<b>41.4</b>
PostNet	34.4	<b>47.9</b>	0.0
DDNet	<b>49.2</b>	44.3	0.0
EvNet	41.1	22.4	4.7
<b>adv. train. loss: diffe</b>			
PriorNet	2.2	0.0	0.0
PostNet	41.9	11.4	0.0
DDNet	46.2	8.4	0.0
EvNet	<b>47.3</b>	<b>34.6</b>	<b>2.0</b>

Table 29: Randomized smoothing verification of MNIST (ID data) and KMNIST (OOD data): percentage of samples that is certifiably correct (cc) and mean certified radius (R).

$\sigma$	ID-Verification						OOD-Verification					
	0.1		0.2		0.5		0.1		0.2		0.5	
	cc	R	cc	R	cc	R	cc	R	cc	R	cc	R
<b>adv. train. loss: None</b>												
PriorNet	<b>97.0</b>	<b>0.36</b>	<b>88.2</b>	<b>0.52</b>	3.0	0.20	<b>98.7</b>	<b>0.37</b>	<b>99.5</b>	<b>0.74</b>	100.0	1.88
PostNet	93.2	0.32	68.4	0.31	0.8	0.11	98.4	0.36	99.5	0.68	100.0	1.55
DDNet	90.6	0.35	52.3	0.46	0.0	0.00	97.8	0.37	<b>99.5</b>	<b>0.74</b>	<b>100.0</b>	<b>1.90</b>
EvNet	95.0	0.31	83.0	0.30	<b>17.3</b>	<b>0.21</b>	77.3	0.17	82.7	0.24	88.6	0.39
<b>adv. train. loss: crossentropy</b>												
PriorNet	<b>97.0</b>	<b>0.36</b>	<b>94.3</b>	<b>0.58</b>	1.0	0.15	<b>99.8</b>	<b>0.38</b>	99.5	0.74	100.0	1.89
PostNet	94.4	0.31	57.7	0.32	3.2	0.13	97.2	0.33	95.6	0.51	99.6	1.02
DDNet	82.6	0.34	55.5	0.46	0.0	0.00	99.6	0.38	<b>100.0</b>	<b>0.75</b>	<b>100.0</b>	<b>1.90</b>
EvNet	96.8	0.34	70.1	0.27	<b>18.8</b>	<b>0.25</b>	58.7	0.11	85.2	0.24	89.1	0.26
<b>adv. train. loss: diffe</b>												
PriorNet	<b>98.0</b>	<b>0.37</b>	<b>83.4</b>	<b>0.49</b>	0.7	0.10	<b>99.7</b>	<b>0.38</b>	<b>100.0</b>	<b>0.76</b>	<b>100.0</b>	<b>1.90</b>
PostNet	93.5	0.33	47.1	0.23	0.6	0.15	95.8	0.34	98.8	0.63	100.0	1.38
DDNet	93.6	0.36	52.7	0.43	0.0	0.00	97.7	0.37	99.7	0.75	<b>100.0</b>	<b>1.90</b>
EvNet	95.4	0.33	81.6	0.34	<b>23.1</b>	<b>0.63</b>	81.7	0.20	82.8	0.28	99.1	1.70

Table 30: Randomized smoothing verification of MNIST (ID data) and KMNIST (OOD data): percentage of samples that is certifiably wrong (cw) and mean certified radius (R).

	0.1		0.2		0.5	
	cw	R	cw	R	cw	R
<b>adv. train. loss: None</b>						
PriorNet	<b>2.8</b>	<b>0.16</b>	<b>10.7</b>	<b>0.21</b>	96.0	0.97
PostNet	6.4	0.17	28.8	0.22	99.0	1.15
DDNet	9.1	0.24	46.3	0.42	100.0	1.81
EvNet	4.5	0.10	15.1	0.13	<b>78.8</b>	<b>0.31</b>
<b>adv. train. loss: crossentropy</b>						
PriorNet	<b>2.9</b>	<b>0.20</b>	<b>4.9</b>	<b>0.24</b>	98.4	1.05
PostNet	5.3	0.17	38.8	0.23	95.2	0.93
DDNet	16.4	0.25	43.5	0.41	100.0	1.74
EvNet	3.0	0.08	26.3	0.13	<b>76.3</b>	<b>0.27</b>
<b>adv. train. loss: diffe</b>						
PriorNet	<b>2.0</b>	<b>0.19</b>	<b>15.7</b>	<b>0.25</b>	98.8	1.10
PostNet	6.3	0.17	49.8	0.25	99.1	1.10
DDNet	6.2	0.22	46.2	0.42	100.0	1.81
EvNet	4.2	0.14	17.0	0.16	<b>73.9</b>	<b>0.94</b>

Table 31: Randomized smoothing verification of MNIST (ID data) and KMNIST (OOD data) harmonic mean.

	0.1	0.2	0.5
<b>adv. train. loss: None</b>			
PriorNet	5.5	19.1	5.9
PostNet	12.0	40.5	1.5
DDNet	<b>16.5</b>	<b>49.2</b>	0.0
EvNet	8.7	25.6	<b>28.4</b>
<b>adv. train. loss: crossentropy</b>			
PriorNet	5.6	9.3	2.0
PostNet	10.0	46.4	6.2
DDNet	<b>27.4</b>	<b>48.8</b>	0.0
EvNet	5.8	38.2	<b>30.2</b>
<b>adv. train. loss: diffe</b>			
PriorNet	3.9	26.4	1.4
PostNet	<b>11.8</b>	48.4	1.2
DDNet	11.6	<b>49.2</b>	0.0
EvNet	8.0	28.1	<b>35.2</b>

Table 32: Adversarial training with CE: Accuracy under PGD label attacks (AUC-PR).

Att. Rad.	0.1	0.2	0.5	1.0	2.0	4.0		0.1	0.2	0.5	1.0	2.0	4.0
<b>MNIST</b>								<b>CIFAR10</b>					
PostNet	<b>99.1</b>	98.7	96.7	89.3	62.4	14.8		72.1	50.6	12.8	2.4	0.2	0.1
PriorNet	<b>99.1</b>	98.8	<b>97.6</b>	<b>94.8</b>	<b>91.1</b>	<b>79.5</b>		69.6	<b>63.9</b>	<b>55.8</b>	<b>46.4</b>	<b>30.9</b>	<b>11.1</b>
DDNet	<b>99.1</b>	<b>98.9</b>	97.4	90.6	47.0	0.2		74.7	<b>63.9</b>	23.6	2.0	0.0	0.0
EvNet	80.3	98.3	96.7	90.5	60.0	52.7		<b>78.3</b>	51.0	14.6	2.6	2.9	0.7
<b>Sensorless</b>								<b>Segment</b>					
PostNet	15.5	6.4	4.7	<b>6.7</b>	<b>11.1</b>	<b>11.7</b>		84.5	52.4	<b>21.1</b>	<b>7.6</b>	<b>5.0</b>	<b>6.3</b>
PriorNet	31.3	15.8	0.2	0.0	0.3	5.3		<b>94.0</b>	<b>65.2</b>	19.1	0.6	0.0	0.0
DDNet	12.4	4.2	0.2	0.3	0.2	0.1		91.4	46.2	7.4	0.2	0.0	0.0
EvNet	<b>33.6</b>	<b>19.4</b>	<b>8.3</b>	5.4	2.6	1.7		93.0	55.2	15.5	2.0	1.4	1.4

Table 33: Adversarial training with CE: Accuracy under FGSM label attacks (AUC-PR).

Att. Rad.	0.1	0.2	0.5	1.0	2.0	4.0		0.1	0.2	0.5	1.0	2.0	4.0
<b>MNIST</b>								<b>CIFAR10</b>					
PostNet	<b>99.1</b>	98.8	97.6	95.3	90.1	81.3		71.0	54.4	30.4	19.9	17.0	17.7
PriorNet	<b>99.1</b>	<b>98.9</b>	<b>97.9</b>	<b>96.2</b>	<b>94.7</b>	<b>90.0</b>		69.7	<b>65.3</b>	<b>58.8</b>	<b>52.1</b>	<b>38.9</b>	<b>22.2</b>
DDNet	<b>99.1</b>	<b>98.9</b>	97.8	94.6	79.1	27.9		<b>73.7</b>	64.7	34.8	15.5	8.0	4.8
EvNet	80.3	98.5	97.6	94.0	72.7	81.3		48.0	56.9	27.0	17.1	17.8	15.7
<b>Sensorless</b>								<b>Segment</b>					
PostNet	20.6	10.6	11.0	11.8	12.5	12.5		82.9	60.1	<b>27.5</b>	<b>22.7</b>	<b>19.1</b>	<b>24.3</b>
PriorNet	35.0	20.8	0.5	0.1	0.9	7.5		91.2	55.2	19.0	0.7	0.0	0.0
DDNet	16.4	9.7	7.0	4.6	6.4	7.5		86.8	36.4	10.5	0.8	0.0	0.0
EvNet	<b>41.1</b>	<b>27.4</b>	<b>20.1</b>	<b>15.2</b>	<b>14.9</b>	<b>12.6</b>		<b>93.8</b>	<b>64.2</b>	25.0	1.5	0.0	0.4

Table 34: Adversarial training with CE: Accuracy under Noise label attacks (AUC-PR).

Att. Rad.	0.1	0.2	0.5	1.0	2.0	4.0		0.1	0.2	0.5	1.0	2.0	4.0
MNIST							CIFAR10						
PostNet	<b>97.8</b>	84.3	8.4	<b>3.2</b>	0.0	0.0		19.3	7.3	3.4	0.0	<b>10.2</b>	<b>10.1</b>
PriorNet	97.7	<b>94.6</b>	<b>32.7</b>	0.0	0.0	0.0		<b>39.1</b>	<b>16.8</b>	<b>3.8</b>	<b>10.5</b>	3.7	0.4
DDNet	97.6	89.1	16.5	2.1	0.0	<b>0.3</b>		25.8	15.1	2.2	0.2	9.8	8.9
EvNet	94.1	74.3	5.0	0.0	0.0	0.0		27.1	7.1	0.1	2.8	8.6	9.7
Sensorless							Segment						
PostNet	1.0	<b>4.7</b>	<b>11.4</b>	<b>11.7</b>	<b>11.7</b>	<b>11.7</b>		27.6	1.5	<b>3.6</b>	<b>15.2</b>	<b>20.9</b>	<b>21.2</b>
PriorNet	<b>4.7</b>	0.0	0.0	0.1	0.0	2.2		<b>56.4</b>	<b>16.1</b>	2.1	0.9	0.0	0.0
DDNet	0.3	0.0	0.0	0.0	0.0	0.1		49.4	6.1	0.0	0.0	0.0	0.0
EvNet	0.9	0.0	0.9	0.2	3.5	3.1		51.2	10.6	0.3	0.0	0.0	0.6

Table 35: Randomized smoothing verification of CIFAR10: percentage of samples that is certifiably correct (cc) w.r.t. the predicted class label and mean certified radius (R) w.r.t. class labels.

	0.1		0.2		0.5	
	cc	R	cc	R	cc	R
adv. train. loss: None						
PriorNet	<b>42.8</b>	<b>0.25</b>	<b>21.2</b>	<b>0.42</b>	<b>11.8</b>	<b>1.30</b>
PostNet	35.0	0.22	12.3	0.51	9.4	0.12
DDNet	31.7	0.26	12.2	0.69	10.8	1.91
EvNet	34.3	0.22	15.4	0.42	11.0	0.63
adv. train. loss: crossentropy						
PriorNet	<b>56.2</b>	<b>0.25</b>	<b>25.4</b>	<b>0.48</b>	<b>13.0</b>	<b>0.35</b>
PostNet	34.7	0.22	15.6	0.45	11.0	0.32
DDNet	41.7	0.24	19.6	0.44	9.1	1.30
EvNet	34.3	0.16	11.1	0.55	10.8	0.74
adv. train. loss: diffe						
PriorNet	48.1	0.23	<b>28.0</b>	<b>0.40</b>	8.4	0.22
PostNet	45.5	0.21	18.0	0.36	5.4	0.18
DDNet	<b>49.2</b>	<b>0.25</b>	26.3	0.34	9.6	0.27
EvNet	21.9	0.30	15.2	0.24	<b>10.8</b>	<b>1.06</b>

Table 36: Randomized smoothing verification of MNIST: percentage of samples that is certifiably correct (cc) w.r.t. the predicted class label and mean certified radius (R) w.r.t. class labels.

	0.1		0.2		0.5	
	cc	R	cc	R	cc	R
adv. train. loss: None						
PriorNet	99.2	0.38	<b>98.8</b>	<b>0.71</b>	<b>61.4</b>	<b>0.45</b>
PostNet	99.2	0.38	98.1	0.66	51.2	0.51
DDNet	<b>99.3</b>	<b>0.38</b>	98.0	0.68	47.3	0.52
EvNet	98.9	0.37	96.2	0.56	57.1	0.42
adv. train. loss: crossentropy						
PriorNet	99.1	0.38	<b>99.0</b>	<b>0.72</b>	50.4	0.53
PostNet	<b>99.4</b>	<b>0.38</b>	97.4	0.62	28.8	0.51
DDNet	99.3	0.38	98.6	0.69	<b>75.4</b>	<b>0.64</b>
EvNet	99.1	0.37	92.1	0.43	35.0	0.40
adv. train. loss: diffe						
PriorNet	99.5	0.38	<b>98.3</b>	<b>0.71</b>	64.0	0.48
PostNet	99.1	0.38	96.8	0.62	48.1	0.44
DDNet	<b>99.6</b>	<b>0.38</b>	98.1	0.69	32.4	0.64
EvNet	99.1	0.37	96.7	0.59	<b>89.5</b>	<b>0.93</b>

Table 37: Adversarial training with CE: Certainty based on differential entropy under PGD label attacks (AUC-PR).

Att. Rad.	0.1	0.2	0.5	1.0	2.0	4.0		0.1	0.2	0.5	1.0	2.0	4.0
MNIST							CIFAR10						
PostNet	99.9	99.8	98.5	88.7	47.6	9.0		88.1	54.1	7.5	1.3	0.1	0.1
PriorNet	99.8	99.4	97.7	92.4	<b>79.7</b>	<b>67.5</b>		54.0	45.5	<b>37.9</b>	<b>29.7</b>	<b>18.1</b>	<b>6.1</b>
DDNet	<b>100.0</b>	<b>99.9</b>	<b>99.7</b>	<b>96.9</b>	46.1	0.1		<b>92.1</b>	<b>83.1</b>	24.8	1.1	0.0	0.0
EvNet	81.2	98.4	95.5	90.4	53.2	38.3		62.9	59.2	13.1	1.5	1.6	0.4
Sensorless							Segment						
PostNet	8.8	4.2	<b>4.6</b>	<b>6.7</b>	<b>11.1</b>	<b>11.7</b>		76.1	35.8	12.6	<b>4.9</b>	<b>4.9</b>	<b>6.3</b>
PriorNet	<b>22.6</b>	<b>11.7</b>	0.2	0.0	0.3	3.6		<b>98.1</b>	<b>66.2</b>	<b>12.8</b>	0.6	0.0	0.0
DDNet	10.9	3.0	0.1	0.2	0.1	0.1		95.9	52.6	4.5	0.5	0.0	0.0
EvNet	21.4	11.0	4.4	3.1	1.7	1.4		94.8	42.7	8.9	1.2	1.3	1.3

Table 38: Adversarial training with CE: Certainty based on differential entropy under FGSM label attacks (AUC-PR).

Att. Rad.	0.1	0.2	0.5	1.0	2.0	4.0		0.1	0.2	0.5	1.0	2.0	4.0
MNIST							CIFAR10						
PostNet	99.9	99.8	99.5	97.4	<b>90.7</b>	<b>83.5</b>		88.0	66.6	28.8	18.0	15.5	<b>20.4</b>
PriorNet	<b>99.8</b>	<b>99.5</b>	<b>98.2</b>	94.6	89.3	<b>83.5</b>		54.7	48.1	42.3	<b>36.8</b>	<b>27.3</b>	19.8
DDNet	100.0	99.9	99.7	<b>98.7</b>	87.3	27.4		<b>91.3</b>	<b>83.4</b>	<b>45.4</b>	15.0	6.5	3.6
EvNet	81.2	98.7	97.4	95.9	73.4	80.7		62.4	68.7	28.6	15.1	20.5	23.1
Sensorless							Segment						
PostNet	12.5	7.5	9.9	11.6	12.5	12.5		75.1	46.2	<b>20.5</b>	<b>19.5</b>	<b>20.2</b>	<b>26.3</b>
PriorNet	30.0	22.2	0.7	0.5	3.2	8.7		<b>96.1</b>	<b>55.6</b>	13.5	2.1	0.0	0.0
DDNet	17.5	6.7	6.2	3.1	4.8	5.5		92.0	42.0	6.4	2.3	0.0	0.0
EvNet	<b>41.6</b>	<b>25.1</b>	<b>18.9</b>	<b>13.2</b>	<b>14.6</b>	<b>13.9</b>		93.1	55.1	15.5	1.6	0.0	1.3

Table 39: Adversarial training with CE: Certainty based on differential entropy under Noise label attacks (AUC-PR).

Att. Rad.	0.1	0.2	0.5	1.0	2.0	4.0		0.1	0.2	0.5	1.0	2.0	4.0
MNIST							CIFAR10						
PostNet	<b>100.0</b>	98.8	54.7	12.1	0.0	0.0		52.1	22.7	6.6	0.0	9.8	9.0
PriorNet	<b>100.0</b>	<b>99.9</b>	<b>88.5</b>	0.0	0.0	0.0		34.9	11.2	<b>7.8</b>	<b>8.7</b>	<b>10.0</b>	0.4
DDNet	99.8	98.5	77.2	<b>15.1</b>	0.0	<b>0.4</b>		<b>81.6</b>	<b>45.3</b>	4.2	0.2	9.6	8.6
EvNet	98.4	86.9	13.3	0.0	0.0	0.0		54.5	17.6	0.1	3.7	8.3	<b>10.5</b>
Sensorless							Segment						
PostNet	0.6	<b>5.1</b>	<b>12.2</b>	<b>11.7</b>	<b>11.7</b>	<b>11.7</b>		36.7	2.0	3.6	<b>17.2</b>	<b>20.8</b>	<b>21.3</b>
PriorNet	<b>8.5</b>	0.0	0.0	0.2	0.0	2.0		<b>90.5</b>	<b>32.8</b>	<b>7.1</b>	1.2	0.0	0.0
DDNet	1.5	0.0	0.0	0.0	0.0	0.0		79.6	21.8	0.0	0.0	0.0	0.0
EvNet	1.5	0.0	1.0	0.2	4.9	4.8		75.7	22.0	3.2	0.0	0.0	0.7

Table 40: Adversarial training with CE: Attack-Detection based on differential entropy under PGD label attacks (AUC-PR).

Att. Rad.	0.1	0.2	0.5	1.0	2.0	4.0		0.1	0.2	0.5	1.0	2.0	4.0
MNIST							CIFAR10						
PostNet	57.8	67.0	84.1	91.0	76.8	47.9		<b>62.5</b>	<b>66.7</b>	41.6	35.0	37.5	36.6
PriorNet	<b>71.7</b>	<b>83.8</b>	<b>96.5</b>	<b>96.0</b>	90.0	79.3		54.4	55.2	54.8	<b>51.1</b>	<b>45.9</b>	40.6
DDNet	54.4	57.4	69.9	86.4	<b>96.2</b>	<b>86.3</b>		56.7	62.4	<b>60.8</b>	39.3	32.9	31.8
EvNet	52.9	59.7	67.7	71.9	66.5	58.5		52.4	59.0	48.9	41.7	40.5	<b>40.7</b>
Sensorless							Segment						
PostNet	43.7	41.1	<b>38.4</b>	<b>53.0</b>	<b>83.5</b>	<b>98.7</b>		94.2	73.5	47.7	42.7	<b>56.8</b>	<b>70.7</b>
PriorNet	<b>60.9</b>	<b>47.5</b>	35.8	31.1	30.8	34.5		86.2	<b>90.1</b>	<b>59.5</b>	<b>47.6</b>	34.0	30.8
DDNet	53.1	43.3	34.7	33.0	31.1	32.6		76.6	83.0	45.7	32.7	30.8	30.8
EvNet	48.3	42.1	37.7	36.6	39.2	48.5		<b>95.9</b>	79.6	43.3	33.4	31.3	31.2

Table 41: Adversarial training with CE: Attack-Detection based on differential entropy under FGSM label attacks (AUC-PR).

Att. Rad.	0.1	0.2	0.5	1.0	2.0	4.0		0.1	0.2	0.5	1.0	2.0	4.0
MNIST							CIFAR10						
PostNet	56.1	62.8	75.5	86.3	90.4	92.6		<b>63.0</b>	<b>68.9</b>	<b>66.1</b>	63.2	66.2	<b>77.5</b>
PriorNet	<b>68.6</b>	<b>80.6</b>	<b>96.8</b>	<b>98.0</b>	<b>98.4</b>	98.2		55.9	59.0	63.0	65.4	63.4	58.4
DDNet	54.4	57.4	69.5	84.0	95.5	<b>99.0</b>		57.6	63.7	70.3	<b>69.0</b>	<b>73.4</b>	76.4
EvNet	52.6	57.9	62.9	66.0	64.0	70.0		52.4	60.0	59.0	61.2	62.9	72.3
Sensorless							Segment						
PostNet	<b>98.3</b>	<b>99.8</b>	<b>99.9</b>	<b>100.0</b>	<b>99.9</b>	<b>100.0</b>		<b>98.0</b>	<b>99.8</b>	<b>99.9</b>	<b>100.0</b>	<b>99.9</b>	<b>100.0</b>
PriorNet	78.6	68.1	37.6	32.0	30.7	49.1		68.8	60.0	42.9	31.4	30.7	32.3
DDNet	60.9	55.5	41.0	34.6	31.7	32.7		61.3	51.7	39.0	33.8	31.5	32.5
EvNet	70.0	70.4	67.5	63.0	77.2	76.6		69.5	70.0	66.9	62.4	77.2	76.4

Table 42: Adversarial training with CE: Attack-Detection based on differential entropy under Noise label attacks (AUC-PR).

Att. Rad.	0.1	0.2	0.5	1.0	2.0	4.0		0.1	0.2	0.5	1.0	2.0	4.0
MNIST							CIFAR10						
PostNet	<b>59.0</b>	66.4	97.2	95.7	95.9	99.6		<b>80.5</b>	<b>89.2</b>	95.2	<b>99.5</b>	85.7	99.7
PriorNet	31.8	33.8	61.3	99.5	<b>100.0</b>	95.8		52.2	50.2	31.2	54.4	<b>99.8</b>	<b>100.0</b>
DDNet	51.8	<b>86.4</b>	<b>99.8</b>	<b>100.0</b>	<b>100.0</b>	99.6		80.3	88.4	<b>99.7</b>	98.8	99.4	68.4
EvNet	51.7	58.6	85.3	84.9	66.3	<b>100.0</b>		46.9	68.1	93.4	94.9	71.4	77.8
Sensorless							Segment						
PostNet	<b>99.9</b>	<b>100.0</b>	<b>100.0</b>	<b>100.0</b>	<b>99.9</b>	<b>100.0</b>		<b>93.2</b>	<b>99.3</b>	<b>99.9</b>	<b>100.0</b>	<b>100.0</b>	<b>100.0</b>
PriorNet	56.7	42.9	32.0	30.8	31.0	30.7		68.8	60.0	42.9	31.4	30.7	32.3
DDNet	51.2	43.9	33.8	32.4	31.7	35.0		61.3	51.7	39.0	33.8	31.5	32.5
EvNet	69.2	58.6	71.7	52.3	70.9	77.6		69.5	70.0	66.9	62.4	77.2	76.4

Table 43: Adversarial training with CE: OOD detection based on differential entropy under PGD uncertainty attacks against differential entropy on ID data and OOD data (AUC-PR).

Att. Rad.	ID-Attack (non-attacked OOD)						OOD-Attack (non-attacked ID)					
	0.1	0.2	0.5	1.0	2.0	4.0	0.1	0.2	0.5	1.0	2.0	4.0
<b>MNIST – KMNIST</b>												
PostNet	95.7	93.1	88.3	78.1	46.9	32.1	94.8	90.3	78.6	58.7	46.4	41.2
PriorNet	<b>99.6</b>	<b>99.3</b>	<b>98.1</b>	<b>95.4</b>	<b>86.7</b>	<b>62.6</b>	<b>99.7</b>	<b>99.3</b>	98.3	90.7	77.7	37.3
DDNet	99.0	98.9	97.8	91.7	58.6	30.7	99.1	99.0	<b>98.4</b>	<b>96.2</b>	<b>90.8</b>	<b>75.7</b>
EvNet	71.3	66.9	60.6	64.4	50.4	42.7	66.3	47.8	37.4	46.7	37.3	33.3
<b>CIFAR10 – SVHN</b>												
PostNet	65.1	45.6	21.0	<b>17.7</b>	<b>16.4</b>	15.5	63.8	41.1	19.6	19.4	17.0	16.1
PriorNet	17.0	16.6	16.0	15.9	16.0	<b>16.1</b>	17.1	16.4	15.8	15.8	15.6	15.7
DDNet	<b>70.8</b>	<b>63.5</b>	<b>34.0</b>	16.8	15.5	15.4	<b>72.7</b>	<b>64.8</b>	28.3	17.9	15.4	15.4
EvNet	53.9	43.7	24.2	16.6	16.1	15.5	55.8	34.7	<b>29.6</b>	<b>21.5</b>	<b>22.0</b>	<b>22.5</b>
<b>Sens. – Sens. class 10, 11</b>												
PostNet	<b>40.5</b>	<b>37.3</b>	<b>43.8</b>	<b>46.7</b>	<b>47.3</b>	<b>45.8</b>	<b>42.6</b>	<b>41.7</b>	<b>31.7</b>	<b>38.5</b>	<b>81.9</b>	<b>99.3</b>
PriorNet	26.6	26.6	26.5	26.5	30.8	40.0	27.9	27.7	26.5	26.5	26.5	26.5
DDNet	26.6	26.6	26.5	26.5	26.6	28.2	26.6	26.6	26.8	26.7	26.6	26.7
EvNet	31.8	29.7	27.2	28.0	32.8	37.8	36.5	38.1	27.8	27.4	30.0	38.3
<b>Seg. – Seg. class sky</b>												
PostNet	61.2	<b>50.8</b>	<b>53.3</b>	<b>32.7</b>	<b>45.3</b>	<b>49.2</b>	<b>79.9</b>	<b>61.6</b>	<b>62.7</b>	<b>32.6</b>	<b>46.0</b>	<b>66.7</b>
PriorNet	31.1	30.8	30.8	30.8	30.8	30.8	31.4	30.8	30.8	30.8	30.8	30.8
DDNet	30.8	30.8	30.8	30.8	30.8	30.8	30.8	30.8	30.8	30.8	30.8	30.8
EvNet	<b>67.0</b>	34.9	30.9	30.8	30.8	31.6	75.5	52.1	31.2	31.2	30.8	30.8

Table 44: Adversarial training with CE: OOD detection based on differential entropy under FGSM uncertainty attacks against differential entropy on ID data and OOD data (AUC-PR).

Att. Rad.	ID-Attack (non-attacked OOD)						OOD-Attack (non-attacked ID)					
	0.1	0.2	0.5	1.0	2.0	4.0	0.1	0.2	0.5	1.0	2.0	4.0
<b>MNIST – KMNIST</b>												
PostNet	95.8	93.7	91.7	90.2	80.3	73.5	95.0	91.6	84.9	80.8	73.9	81.4
PriorNet	<b>99.5</b>	<b>99.2</b>	<b>98.3</b>	<b>95.7</b>	<b>87.2</b>	<b>77.0</b>	<b>99.5</b>	<b>99.2</b>	<b>98.8</b>	96.8	92.4	76.2
DDNet	99.0	99.0	98.1	94.6	80.3	47.3	99.1	99.1	98.6	<b>97.3</b>	<b>95.3</b>	<b>92.4</b>
EvNet	71.4	67.7	62.6	68.6	56.9	53.0	68.3	57.1	50.7	62.6	50.7	46.4
<b>CIFAR10 – SVHN</b>												
PostNet	67.4	61.4	<b>62.9</b>	<b>70.7</b>	<b>65.0</b>	44.4	65.6	54.6	39.0	<b>45.1</b>	<b>62.7</b>	<b>77.6</b>
PriorNet	17.0	16.9	17.1	18.1	24.0	37.5	17.2	17.3	17.4	18.9	22.2	29.8
DDNet	<b>71.0</b>	<b>66.9</b>	56.1	55.7	48.7	44.7	<b>72.2</b>	<b>66.4</b>	<b>48.0</b>	42.2	48.7	69.1
EvNet	56.5	62.2	51.6	53.9	64.4	<b>46.6</b>	55.9	42.1	35.0	37.4	56.2	68.7
<b>Sens. – Sens. class 10, 11</b>												
PostNet	41.2	37.9	34.3	41.4	45.7	45.7	<b>99.2</b>	<b>99.8</b>	<b>99.8</b>	<b>100.0</b>	<b>99.9</b>	<b>100.0</b>
PriorNet	27.2	28.0	29.3	37.5	96.5	77.7	28.8	29.5	26.6	26.5	26.5	26.5
DDNet	30.7	32.8	<b>65.6</b>	<b>72.7</b>	<b>92.9</b>	<b>94.4</b>	27.5	29.5	28.7	26.6	26.5	27.2
EvNet	<b>44.3</b>	<b>47.2</b>	47.7	46.3	38.8	40.0	51.6	69.3	50.4	48.5	65.4	72.3
<b>Seg. – Seg. class sky</b>												
PostNet	61.9	<b>54.3</b>	<b>57.4</b>	34.0	37.5	43.3	<b>92.9</b>	<b>92.2</b>	<b>91.1</b>	<b>99.9</b>	<b>99.9</b>	<b>100.0</b>
PriorNet	31.0	30.8	30.8	30.8	30.8	36.9	31.2	30.8	30.8	30.8	30.8	30.8
DDNet	30.8	30.8	33.0	37.8	59.4	<b>92.1</b>	30.8	30.8	30.8	30.8	30.8	30.8
EvNet	<b>66.7</b>	42.3	45.6	<b>49.0</b>	<b>61.5</b>	50.1	74.6	57.3	51.2	45.8	60.0	63.2



Table 45: Adversarial training with CE: OOD detection based on differential entropy under Noise uncertainty attacks against differential entropy on ID data and OOD data (AUC-PR).

Att. Rad.	ID-Attack (non-attacked OOD)						OOD-Attack (non-attacked ID)					
	0.1	0.2	0.5	1.0	2.0	4.0	0.1	0.2	0.5	1.0	2.0	4.0
<b>MNIST – KMNIST</b>												
PostNet	91.7	72.7	32.4	<b>36.2</b>	33.2	30.7	93.0	85.5	89.9	90.5	73.1	99.1
PriorNet	<b>99.7</b>	<b>98.4</b>	<b>71.5</b>	34.3	<b>37.3</b>	<b>38.7</b>	<b>99.2</b>	96.8	96.8	99.3	<b>100.0</b>	77.4
DDNet	96.6	79.5	31.9	31.2	30.9	35.0	98.6	<b>99.5</b>	<b>99.9</b>	<b>100.0</b>	99.8	98.5
EvNet	87.4	49.9	32.8	32.0	33.2	36.3	87.4	46.8	48.6	45.1	33.3	<b>100.0</b>
<b>CIFAR10 – SVHN</b>												
PostNet	43.9	31.9	30.7	30.7	<b>56.2</b>	31.5	85.1	<b>84.8</b>	85.2	97.0	82.2	99.5
PriorNet	31.4	32.8	<b>85.8</b>	<b>37.1</b>	30.7	30.8	35.2	40.1	30.8	42.1	<b>99.0</b>	<b>100.0</b>
DDNet	50.8	<b>42.8</b>	30.7	32.8	30.7	<b>94.9</b>	<b>82.3</b>	80.3	<b>99.2</b>	<b>97.8</b>	98.9	63.9
EvNet	<b>56.2</b>	34.4	32.1	32.4	37.9	46.7	59.8	63.2	82.5	92.0	41.6	51.2
<b>Sens. – Sens. class 10, 11</b>												
PostNet	<b>30.8</b>	<b>47.8</b>	<b>50.0</b>	50.0	50.0	50.0	<b>98.7</b>	<b>99.8</b>	<b>100.0</b>	<b>100.0</b>	<b>99.8</b>	<b>100.0</b>
PriorNet	30.7	30.7	30.9	<b>87.7</b>	<b>100.0</b>	<b>100.0</b>	31.0	30.7	32.4	31.4	34.2	30.7
DDNet	30.7	30.7	47.4	75.0	92.9	79.7	30.7	30.7	30.7	41.3	34.3	39.6
EvNet	<b>30.8</b>	30.8	40.8	36.3	50.7	34.2	34.5	31.0	34.4	38.8	47.4	51.1
<b>Seg. – Seg. class sky</b>												
PostNet	34.2	<b>31.0</b>	<b>42.6</b>	<b>49.9</b>	50.0	50.0	<b>97.7</b>	<b>93.7</b>	<b>99.8</b>	<b>100.0</b>	<b>100.0</b>	<b>100.0</b>
PriorNet	30.8	30.8	30.8	30.8	30.9	<b>100.0</b>	30.9	30.8	30.8	30.9	31.2	32.9
DDNet	30.8	30.8	30.8	31.3	<b>77.8</b>	93.3	30.8	30.8	30.8	30.8	30.8	32.3
EvNet	<b>63.3</b>	30.8	30.8	30.8	32.7	36.0	98.8	40.2	40.4	34.8	50.0	32.5

Table 46: Adversarial training with Diff. Ent.: Accuracy based on differential entropy under PGD label attacks (AUC-PR).

Att. Rad.	0.1	0.2	0.5	1.0	2.0	4.0	0.1	0.2	0.5	1.0	2.0	4.0
	MNIST						CIFAR10					
PostNet	<b>99.1</b>	98.8	96.8	88.8	65.9	15.6	72.9	50.7	12.9	3.1	0.4	0.4
PriorNet	99.1	98.8	<b>97.7</b>	<b>94.7</b>	<b>88.8</b>	<b>73.4</b>	66.5	62.8	<b>52.9</b>	<b>35.8</b>	<b>23.0</b>	<b>9.6</b>
DDNet	99.1	<b>98.9</b>	97.4	91.9	48.7	0.3	<b>78.9</b>	<b>63.1</b>	22.0	1.9	0.0	0.0
EvNet	98.3	98.1	95.2	91.0	72.7	40.1	65.6	48.9	14.8	8.4	3.8	1.8
<b>Sensorless</b>							<b>Segment</b>					
PostNet	16.1	7.4	5.8	<b>7.5</b>	<b>9.4</b>	<b>12.5</b>	84.7	47.1	<b>22.3</b>	<b>6.4</b>	<b>10.8</b>	<b>3.8</b>
PriorNet	33.3	15.6	3.7	0.0	0.0	0.0	<b>93.9</b>	<b>65.9</b>	18.1	2.9	0.0	0.0
DDNet	12.9	3.0	0.5	0.3	0.2	0.2	90.6	47.5	8.4	0.1	0.0	0.0
EvNet	<b>36.1</b>	<b>22.1</b>	<b>10.8</b>	3.8	1.7	3.1	92.0	56.2	11.9	2.1	0.4	2.8

Table 47: Adversarial training with Diff. Ent.: Accuracy based on differential entropy under FGSM label attacks (AUC-PR).

Att. Rad.	0.1	0.2	0.5	1.0	2.0	4.0	0.1	0.2	0.5	1.0	2.0	4.0
	MNIST						CIFAR10					
PostNet	<b>99.1</b>	98.9	97.6	95.0	91.2	80.9	23.5	11.2	0.2	0.1	9.0	1.6
PriorNet	<b>99.1</b>	98.8	<b>97.9</b>	<b>95.9</b>	<b>93.5</b>	<b>87.4</b>	31.6	<b>14.4</b>	6.2	7.8	0.2	1.4
DDNet	<b>99.1</b>	<b>99.0</b>	<b>97.9</b>	94.9	78.4	23.3	<b>36.8</b>	13.9	<b>9.2</b>	<b>10.3</b>	<b>10.0</b>	<b>10.0</b>
EvNet	98.3	98.3	95.7	95.4	88.9	63.7	24.4	4.8	0.5	9.0	11.2	3.9
<b>Sensorless</b>							<b>Segment</b>					
PostNet	21.4	10.6	10.3	12.3	12.4	12.5	83.9	53.7	<b>27.8</b>	<b>15.4</b>	<b>20.1</b>	<b>19.1</b>
PriorNet	<b>40.2</b>	21.7	8.1	0.0	0.0	2.3	91.3	57.2	18.1	3.2	0.0	0.0
DDNet	17.6	5.0	4.9	8.0	7.0	5.7	86.1	39.2	12.1	0.4	0.0	2.9
EvNet	43.1	<b>29.3</b>	<b>21.3</b>	<b>14.4</b>	<b>13.4</b>	<b>13.5</b>	<b>91.9</b>	<b>65.3</b>	17.9	2.9	0.1	0.8

Table 48: Adversarial training with Diff. Ent.: Accuracy based on differential entropy under Noise label attacks (AUC-PR).

Att. Rad.	0.1	0.2	0.5	1.0	2.0	4.0	0.1	0.2	0.5	1.0	2.0	4.0
	MNIST						CIFAR10					
PostNet	97.0	88.5	10.0	2.2	0.0	<b>0.4</b>	71.8	54.2	30.7	21.6	16.8	14.6
PriorNet	<b>98.1</b>	88.5	<b>24.5</b>	<b>4.1</b>	0.0	0.0	67.2	<b>65.2</b>	<b>58.9</b>	<b>48.5</b>	<b>40.5</b>	<b>31.4</b>
DDNet	97.8	<b>92.5</b>	6.4	2.2	<b>0.1</b>	0.1	<b>78.1</b>	63.8	34.2	16.0	8.0	6.0
EvNet	96.2	87.0	2.3	0.1	0.0	0.0	65.2	54.8	29.2	18.7	17.4	16.2
<b>Sensorless</b>							<b>Segment</b>					
PostNet	1.0	<b>4.7</b>	<b>11.4</b>	<b>11.7</b>	<b>11.7</b>	<b>11.7</b>	27.6	1.5	<b>3.6</b>	<b>15.2</b>	<b>20.9</b>	<b>21.2</b>
PriorNet	<b>4.7</b>	0.0	0.0	0.1	0.0	2.2	<b>56.4</b>	<b>16.1</b>	2.1	0.9	0.0	0.0
DDNet	0.3	0.0	0.0	0.0	0.0	0.1	49.4	6.1	0.0	0.0	0.0	0.0
EvNet	0.9	0.0	0.9	0.2	3.5	3.1	51.2	10.6	0.3	0.0	0.0	0.6

Table 49: Adversarial training with Diff. Ent.: Certainty based on differential entropy under PGD label attacks (AUC-PR).

Att. Rad.	0.1	0.2	0.5	1.0	2.0	4.0		0.1	0.2	0.5	1.0	2.0	4.0
MNIST							CIFAR10						
PostNet	99.9	99.8	98.5	86.8	53.0	10.2		88.5	56.6	7.5	1.7	0.3	0.2
PriorNet	99.7	99.6	98.0	91.4	<b>76.2</b>	<b>54.8</b>		51.0	44.7	<b>36.3</b>	<b>23.7</b>	<b>13.8</b>	<b>5.5</b>
DDNet	<b>100.0</b>	<b>99.9</b>	<b>99.7</b>	<b>97.6</b>	47.9	0.1		<b>94.7</b>	<b>82.4</b>	21.2	1.1	0.0	0.0
EvNet	99.2	98.9	96.8	86.5	60.8	33.73		80.6	50.4	14.1	9.1	9.7	2.2
Sensorless							Segment						
PostNet	10.6	5.3	5.8	<b>7.5</b>	<b>9.5</b>	<b>12.5</b>		76.1	30.1	<b>13.4</b>	<b>4.9</b>	<b>13.2</b>	<b>3.8</b>
PriorNet	<b>22.6</b>	10.3	3.8	0.0	0.0	0.0		<b>97.9</b>	<b>63.8</b>	11.4	1.7	0.0	0.0
DDNet	13.2	2.2	0.4	0.3	0.1	0.2		95.8	51.1	5.0	0.3	0.0	0.0
EvNet	<b>22.6</b>	<b>12.8</b>	<b>5.9</b>	2.0	1.1	2.9		94.5	44.3	7.0	1.2	0.4	2.0

Table 50: Adversarial training with Diff. Ent.: Certainty based on differential entropy under FGSM label attacks (AUC-PR).

Att. Rad.	0.1	0.2	0.5	1.0	2.0	4.0		0.1	0.2	0.5	1.0	2.0	4.0
MNIST							CIFAR10						
PostNet	99.9	<b>99.9</b>	99.4	97.1	<b>92.3</b>	<b>82.2</b>		88.3	67.8	30.9	19.2	15.9	16.7
PriorNet	99.7	99.6	98.4	93.8	87.0	79.5		52.9	49.1	<b>44.1</b>	<b>36.3</b>	<b>28.6</b>	<b>20.9</b>
DDNet	<b>100.0</b>	<b>99.9</b>	<b>99.8</b>	<b>98.8</b>	85.8	21.0		<b>94.1</b>	<b>82.8</b>	43.1	15.4	6.8	5.1
EvNet	99.2	99.1	97.8	95.3	87.8	67.0		80.2	61.1	32.5	21.9	21.6	22.2
Sensorless							Segment						
PostNet	13.8	7.2	9.7	12.3	12.4	12.5		74.4	37.9	20.2	<b>16.7</b>	<b>19.5</b>	<b>19.9</b>
PriorNet	32.3	19.1	9.1	0.0	0.0	5.2		<b>96.2</b>	56.6	11.4	10.7	0.0	0.0
DDNet	19.7	7.0	3.7	7.6	7.6	7.1		92.2	40.5	7.0	0.4	0.0	4.9
EvNet	<b>37.8</b>	<b>30.5</b>	<b>26.0</b>	<b>14.7</b>	<b>13.9</b>	<b>13.6</b>		91.9	<b>57.7</b>	<b>11.8</b>	2.3	0.7	1.1

Table 51: Adversarial training with Diff. Ent.: Certainty based on differential entropy under Noise label attacks (AUC-PR).

Att. Rad.	0.1	0.2	0.5	1.0	2.0	4.0		0.1	0.2	0.5	1.0	2.0	4.0
MNIST							CIFAR10						
PostNet	99.7	99.2	64.5	8.7	0.0	<b>0.5</b>		77.4	14.6	0.3	0.1	9.1	3.7
PriorNet	<b>99.9</b>	98.8	<b>84.5</b>	<b>17.5</b>	0.0	0.0		44.7	11.3	13.6	9.0	0.4	2.6
DDNet	99.7	<b>99.4</b>	39.5	5.9	<b>0.6</b>	0.3		<b>86.1</b>	<b>46.5</b>	<b>15.5</b>	<b>17.4</b>	<b>10.2</b>	<b>10.1</b>
EvNet	99.2	97.0	19.4	0.1	0.0	0.0		67.1	12.7	4.5	12.8	13.2	3.3
Sensorless							Segment						
PostNet	0.6	<b>5.1</b>	<b>12.2</b>	<b>11.7</b>	<b>11.7</b>	<b>11.7</b>		36.7	2.0	3.6	<b>17.2</b>	<b>20.8</b>	<b>21.3</b>
PriorNet	<b>8.5</b>	0.0	0.0	0.2	0.0	2.0		<b>90.5</b>	<b>32.8</b>	<b>7.1</b>	1.2	0.0	0.0
DDNet	1.5	0.0	0.0	0.0	0.0	0.0		79.6	21.8	0.0	0.0	0.0	0.0
EvNet	1.5	0.0	1.0	0.2	4.9	4.8		75.7	22.0	3.2	0.0	0.0	0.7

Table 52: Adversarial training with Diff. Ent.: Attack-Detection based on differential entropy under PGD label attacks (AUC-PR).

Att. Rad.	0.1	0.2	0.5	1.0	2.0	4.0		0.1	0.2	0.5	1.0	2.0	4.0
MNIST							CIFAR10						
PostNet	57.9	65.4	87.1	93.6	79.3	47.7		<b>63.1</b>	<b>67.1</b>	41.7	34.4	35.0	36.5
PriorNet	<b>66.9</b>	<b>76.0</b>	<b>95.1</b>	<b>96.4</b>	88.7	74.8		55.7	55.8	53.1	<b>48.9</b>	<b>43.3</b>	37.8
DDNet	53.7	58.5	69.3	85.5	<b>96.1</b>	<b>87.7</b>		56.7	62.5	<b>60.4</b>	41.2	32.6	31.8
EvNet	54.3	58.9	63.2	72.3	69.4	59.1		55.9	60.2	49.7	44.6	41.4	<b>39.4</b>
Sensorless							Segment						
PostNet	49.8	41.5	36.3	<b>51.0</b>	<b>85.9</b>	<b>99.0</b>		95.0	77.6	48.9	<b>42.9</b>	<b>45.2</b>	<b>68.5</b>
PriorNet	50.4	39.4	31.6	30.8	30.7	30.7		86.1	<b>89.7</b>	<b>50.8</b>	37.5	32.9	30.8
DDNet	<b>52.2</b>	41.5	35.5	32.3	31.5	35.9		77.9	87.3	43.4	32.4	30.9	30.8
EvNet	48.0	<b>44.3</b>	<b>38.8</b>	35.2	39.1	48.9		<b>95.4</b>	77.9	42.6	33.7	31.3	31.6

Table 53: Adversarial training with Diff. Ent.: Attack-Detection based on differential entropy under FGSM label attacks (AUC-PR).

Att. Rad.	0.1	0.2	0.5	1.0	2.0	4.0		0.1	0.2	0.5	1.0	2.0	4.0
MNIST							CIFAR10						
PostNet	56.2	61.6	78.0	89.5	92.1	90.7		<b>63.8</b>	<b>68.8</b>	65.8	62.9	67.0	73.6
PriorNet	<b>67.3</b>	<b>76.9</b>	<b>95.3</b>	<b>98.0</b>	<b>98.0</b>	98.0		58.0	62.5	<b>67.3</b>	<b>67.3</b>	<b>65.9</b>	<b>62.2</b>
DDNet	53.8	58.6	68.7	83.9	95.8	<b>98.9</b>		57.7	63.8	71.0	72.8	74.8	79.3
EvNet	53.7	57.2	59.8	65.2	71.5	72.4		56.2	62.4	58.9	59.0	63.2	70.5
Sensorless							Segment						
PostNet	<b>98.3</b>	<b>99.8</b>	<b>99.9</b>	<b>99.9</b>	<b>99.9</b>	<b>100.0</b>		<b>96.0</b>	<b>93.3</b>	<b>93.1</b>	<b>97.4</b>	<b>99.8</b>	<b>99.6</b>
PriorNet	67.0	56.0	35.8	30.8	30.7	30.7		88.6	89.3	55.1	45.2	37.6	30.8
DDNet	60.1	50.0	40.3	32.5	31.5	33.8		81.6	89.2	56.9	49.1	35.8	31.8
EvNet	68.7	69.8	69.9	68.6	68.8	73.2		95.9	86.6	66.2	63.2	74.2	77.2

Table 54: Adversarial training with Diff. Ent.: Attack-Detection based on differential entropy under Noise label attacks (AUC-PR).

Att. Rad.	0.1	0.2	0.5	1.0	2.0	4.0		0.1	0.2	0.5	1.0	2.0	4.0
MNIST								CIFAR10					
PostNet	53.7	67.7	94.5	94.1	98.1	<b>100.0</b>		<b>85.8</b>	79.1	93.0	<b>98.5</b>	<b>96.3</b>	96.6
PriorNet	33.1	42.4	80.8	93.8	78.8	56.0		52.3	38.2	51.8	38.1	87.4	<b>99.9</b>
DDNet	<b>56.0</b>	<b>86.7</b>	<b>99.9</b>	<b>100.0</b>	<b>100.0</b>	99.4		79.5	<b>90.9</b>	<b>99.6</b>	98.2	96.2	85.0
EvNet	50.8	71.7	83.1	88.2	68.4	88.0		67.9	87.0	93.3	89.6	93.7	97.3
Sensorless								Segment					
PostNet	<b>99.9</b>	<b>100.0</b>	<b>100.0</b>	<b>100.0</b>	<b>99.9</b>	<b>100.0</b>		<b>93.2</b>	<b>99.3</b>	<b>99.9</b>	<b>100.0</b>	<b>100.0</b>	<b>100.0</b>
PriorNet	56.7	42.9	32.0	30.8	31.0	30.7		55.9	63.1	41.7	33.3	30.8	30.8
DDNet	51.2	43.9	33.8	32.4	31.7	35.0		55.9	58.9	42.1	33.4	31.1	32.7
EvNet	69.2	58.6	71.7	52.3	70.9	77.6		63.0	62.7	63.9	54.4	74.6	61.2

Table 55: Adversarial training with Diff. Ent.: OOD detection based on differential entropy under PGD uncertainty attacks against differential entropy on ID data and OOD data (AUC-PR).

Att. Rad.	ID-Attack (non-attacked OOD)							OOD-Attack (non-attacked ID)					
	0.1	0.2	0.5	1.0	2.0	4.0		0.1	0.2	0.5	1.0	2.0	4.0
MNIST – KMNIST													
PostNet	92.6	94.0	89.5	69.1	49.6	15.6		91.9	91.4	79.8	54.2	42.4	46.0
PriorNet	<b>99.6</b>	<b>99.3</b>	<b>97.6</b>	<b>93.7</b>	<b>81.4</b>	<b>50.8</b>		<b>99.6</b>	<b>99.2</b>	97.4	92.1	66.3	37.6
DDNet	99.1	98.9	<b>97.6</b>	93.6	60.7	30.7		99.1	99.0	<b>98.3</b>	<b>97.0</b>	<b>91.4</b>	<b>77.5</b>
EvNet	73.4	66.7	72.9	57.7	49.3	45.0		63.7	51.3	58.8	35.1	33.4	36.4
CIFAR10 – SVHN													
PostNet	68.6	46.1	21.7	17.5	16.2	15.6		63.3	37.4	19.0	17.4	16.7	16.8
PriorNet	17.3	15.9	17.4	15.5	15.4	15.4		17.1	15.7	16.6	15.4	15.4	15.4
DDNet	<b>77.5</b>	<b>66.0</b>	<b>34.5</b>	16.4	15.4	15.4		<b>79.6</b>	<b>67.4</b>	<b>33.2</b>	16.9	15.4	15.4
EvNet	57.8	35.2	22.0	<b>21.5</b>	<b>17.5</b>	<b>16.0</b>		52.7	30.7	30.3	<b>31.1</b>	<b>20.8</b>	<b>18.0</b>
Sens. – Sens. class 10, 11													
PostNet	<b>39.6</b>	<b>34.8</b>	<b>41.8</b>	<b>46.0</b>	<b>44.9</b>	<b>47.6</b>		<b>41.1</b>	<b>40.6</b>	<b>30.8</b>	<b>35.6</b>	<b>83.0</b>	<b>99.5</b>
PriorNet	26.6	26.5	26.5	26.5	26.5	26.6		28.8	27.0	26.6	26.9	26.5	26.5
DDNet	26.8	26.5	26.5	26.6	26.6	28.0		26.8	26.6	26.7	26.7	26.6	26.6
EvNet	31.0	29.4	27.2	29.1	32.4	36.5		39.1	35.1	28.9	28.7	30.0	38.3
Seg. – Seg. class sky													
PostNet	<b>91.7</b>	<b>45.3</b>	<b>44.6</b>	<b>38.8</b>	<b>46.0</b>	<b>49.4</b>		<b>98.7</b>	67.3	<b>44.1</b>	<b>47.7</b>	<b>37.5</b>	<b>59.4</b>
PriorNet	31.2	30.8	30.8	30.8	30.8	30.8		31.7	30.8	30.8	30.8	30.8	30.8
DDNet	31.0	30.8	30.8	30.8	30.8	30.8		31.2	30.8	30.8	30.8	30.8	30.8
EvNet	58.0	39.4	31.0	30.8	30.8	31.5		84.1	<b>71.7</b>	36.1	31.1	30.8	30.8

Table 56: Adversarial training with Diff. Ent.: OOD detection based on differential entropy under FGSM uncertainty attacks against differential entropy on ID data and OOD data (AUC-PR).

Att. Rad.	ID-Attack (non-attacked OOD)							OOD-Attack (non-attacked ID)					
	0.1	0.2	0.5	1.0	2.0	4.0		0.1	0.2	0.5	1.0	2.0	4.0
MNIST – KMNIST													
PostNet	92.8	94.2	93.1	84.7	81.2	<b>72.9</b>		92.1	92.2	88.7	79.0	78.2	75.9
PriorNet	<b>99.6</b>	<b>99.3</b>	97.7	94.5	<b>86.6</b>	71.0		<b>99.6</b>	<b>99.3</b>	98.1	96.5	90.2	71.5
DDNet	99.1	98.9	<b>97.9</b>	<b>95.5</b>	79.2	46.3		99.1	99.0	<b>98.5</b>	<b>97.8</b>	<b>95.4</b>	<b>92.1</b>
EvNet	73.7	67.6	74.4	62.2	56.5	59.5		67.7	57.6	67.5	49.8	49.0	55.0
CIFAR10 – SVHN													
PostNet	71.3	60.8	<b>67.1</b>	<b>68.7</b>	<b>61.8</b>	<b>54.4</b>		65.7	49.9	41.7	37.8	<b>61.0</b>	<b>78.4</b>
PriorNet	17.5	16.3	20.8	17.4	19.8	19.6		17.5	16.4	21.7	17.2	21.0	20.1
DDNet	<b>77.7</b>	<b>69.9</b>	61.4	53.8	45.7	37.2		<b>79.3</b>	<b>69.3</b>	<b>54.2</b>	47.5	52.4	73.9
EvNet	62.9	51.2	48.9	58.7	58.1	46.3		53.9	35.0	37.9	<b>51.4</b>	56.4	59.6
Sens. – Sens. class 10, 11													
PostNet	40.1	35.5	33.8	43.0	43.6	47.5		<b>99.7</b>	<b>99.9</b>	<b>99.9</b>	<b>99.9</b>	<b>99.9</b>	<b>100.0</b>
PriorNet	28.4	29.5	39.4	53.3	93.0	<b>98.3</b>		30.2	28.3	27.2	27.8	26.9	26.5
DDNet	32.2	34.4	<b>53.5</b>	<b>84.2</b>	<b>93.8</b>	92.2		26.8	26.6	26.7	26.7	26.6	26.6
EvNet	<b>43.1</b>	<b>41.0</b>	37.0	43.5	42.9	41.9		56.8	56.4	57.8	60.0	66.0	67.8
Seg. – Seg. class sky													
PostNet	<b>91.4</b>	<b>54.2</b>	<b>53.1</b>	45.4	38.1	40.6		<b>98.8</b>	<b>75.9</b>	<b>87.7</b>	<b>99.1</b>	<b>99.9</b>	<b>100.0</b>
PriorNet	31.1	30.8	30.8	30.8	30.8	39.2		31.5	30.8	30.8	30.8	30.8	30.8
DDNet	30.9	30.8	30.8	34.9	<b>58.2</b>	<b>85.4</b>		31.1	30.8	30.8	30.8	30.8	30.8
EvNet	55.3	46.6	51.8	<b>51.0</b>	48.9	41.4		86.6	71.6	49.5	57.4	76.6	58.6

Table 57: Adversarial training with Diff. Ent.: OOD detection based on differential entropy under Noise uncertainty attacks against differential entropy on ID data and OOD data (AUC-PR).

Att. Rad.	ID-Attack (non-attacked OOD)						OOD-Attack (non-attacked ID)					
	0.1	0.2	0.5	1.0	2.0	4.0	0.1	0.2	0.5	1.0	2.0	4.0
<b>MNIST – KMNIST</b>												
PostNet	90.8	81.0	39.1	38.3	30.8	30.7	88.1	88.9	81.0	70.8	95.4	<b>100.0</b>
PriorNet	99.8	<b>93.6</b>	<b>41.0</b>	<b>42.1</b>	<b>40.0</b>	<b>51.7</b>	<b>100.0</b>	99.3	98.1	94.0	44.8	32.1
DDNet	<b>97.2</b>	87.1	31.0	31.1	31.6	40.6	98.8	<b>99.5</b>	<b>99.9</b>	<b>99.9</b>	<b>100.0</b>	98.1
EvNet	68.6	52.2	31.9	33.1	32.9	31.3	53.2	62.2	39.7	45.9	33.7	50.3
<b>CIFAR10 – SVHN</b>												
PostNet	46.7	<b>67.3</b>	<b>33.9</b>	30.7	<b>39.4</b>	44.0	88.0	78.2	84.8	92.1	93.6	95.9
PriorNet	41.9	37.0	<b>33.9</b>	<b>42.5</b>	32.6	30.7	35.4	32.4	38.2	33.3	76.1	<b>99.5</b>
DDNet	<b>58.7</b>	37.3	30.7	33.5	33.6	<b>61.0</b>	<b>90.6</b>	<b>88.3</b>	<b>99.0</b>	<b>97.5</b>	<b>93.9</b>	80.9
EvNet	48.9	34.8	30.9	31.8	37.7	34.4	69.6	77.7	85.3	87.4	93.1	95.7
<b>Sens. – Sens. class 10, 11</b>												
PostNet	<b>30.8</b>	<b>47.8</b>	<b>50.0</b>	50.0	50.0	50.0	<b>98.7</b>	<b>99.8</b>	<b>100.0</b>	<b>100.0</b>	<b>99.8</b>	<b>100.0</b>
PriorNet	30.7	30.7	30.9	<b>87.7</b>	<b>100.0</b>	<b>100.0</b>	31.0	30.7	32.4	31.4	34.2	30.7
DDNet	30.7	30.7	47.4	75.0	92.9	79.7	30.7	30.7	30.7	41.3	34.3	39.6
EvNet	<b>30.8</b>	30.8	40.8	36.3	50.7	34.2	34.5	31.0	34.4	38.8	47.4	51.1
<b>Seg. – Seg. class sky</b>												
PostNet	34.2	<b>31.0</b>	<b>42.6</b>	<b>49.9</b>	50.0	50.0	97.7	<b>93.7</b>	<b>99.8</b>	<b>100.0</b>	<b>100.0</b>	<b>100.0</b>
PriorNet	30.8	30.8	30.8	30.8	30.9	<b>100.0</b>	30.9	30.8	30.8	30.9	31.2	32.9
DDNet	30.8	30.8	30.8	31.3	<b>77.8</b>	93.3	30.8	30.8	30.8	30.8	30.8	32.3
EvNet	<b>63.3</b>	30.8	30.8	30.8	32.7	36.0	<b>98.8</b>	40.2	40.4	34.8	50.0	32.5

#### A.6 VISUALIZATION OF DIFFERENTIAL ENTROPY DISTRIBUTIONS ON ID DATA AND OOD DATA

The following Figures visualize the differential entropy distribution for ID data and OOD data for all models with standard training. We used label attacks and uncertainty attacks for CIFAR10 and MNIST. Thus, they show how well the DBU models separate on clean and perturbed ID data and OOD data.

Figures 4 and 5 visualizes the differential entropy distribution of ID data and OOD data under label attacks. On CIFAR10, PriorNet and DDNet can barely distinguish between clean ID and OOD data. We observe a better ID/OOD distinction for PostNet and EvNet for clean data. However, we do not observe for any model an increase of the uncertainty estimates on label attacked data. Even worse, PostNet, PriorNet and DDNet seem to assign higher confidence on class label attacks. On MNIST, models show a slightly better behavior. They are capable to assign a higher uncertainty to label attacks up to some attack radius.

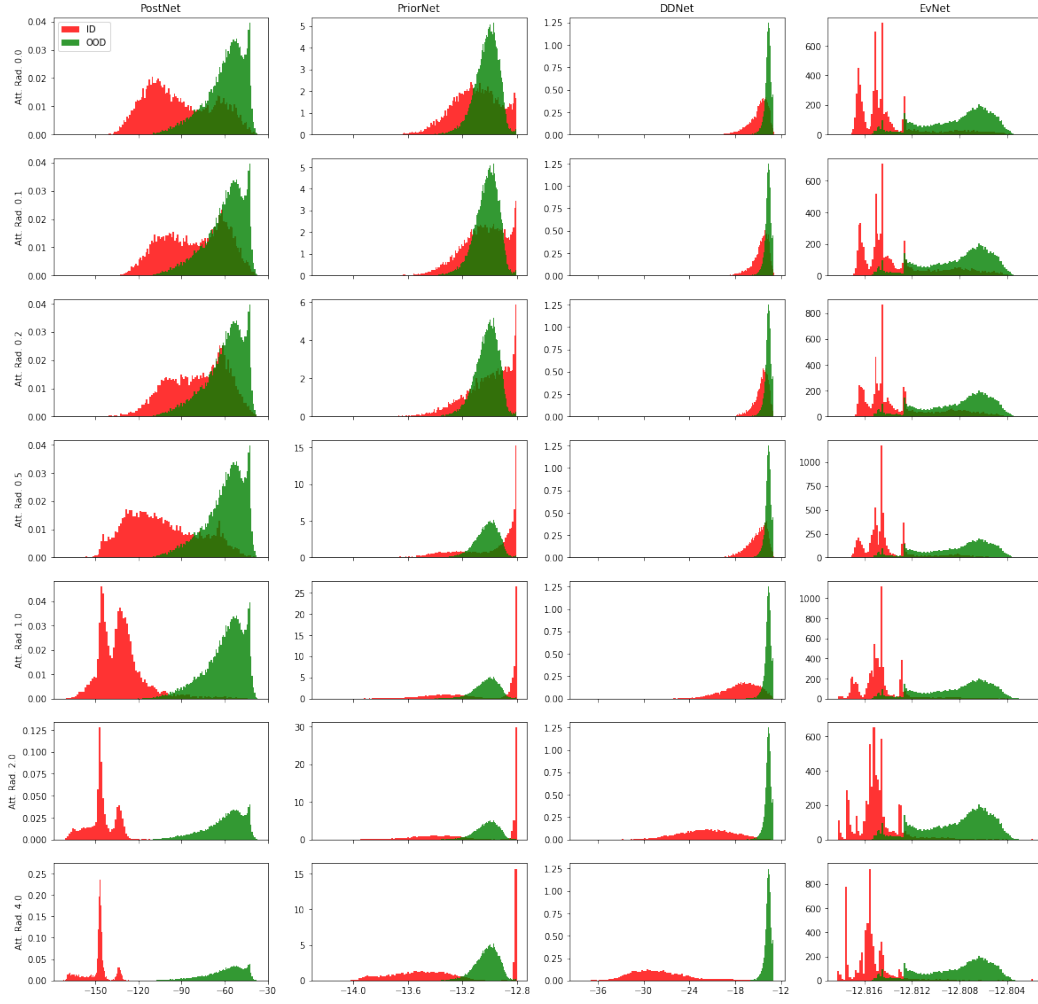


Figure 4: Visualization of the differential entropy distribution of ID data (CIFAR10) and OOD data (SVHN) under label attack. The first row corresponds to no attack. The other rows correspond to increasingly stronger attack strength.

Figures 6, 7, 8 and 9 visualize the differential entropy distribution of ID data and OOD data under uncertainty attacks. For both CIFAR10 and MNIST data sets, we observed that uncertainty estimations of all models can be manipulated. That is, OOD uncertainty attacks can shift the OOD uncertainty distribution to more certain predictions, and ID uncertainty attacks can shift the ID uncertainty distribution to less certain predictions.

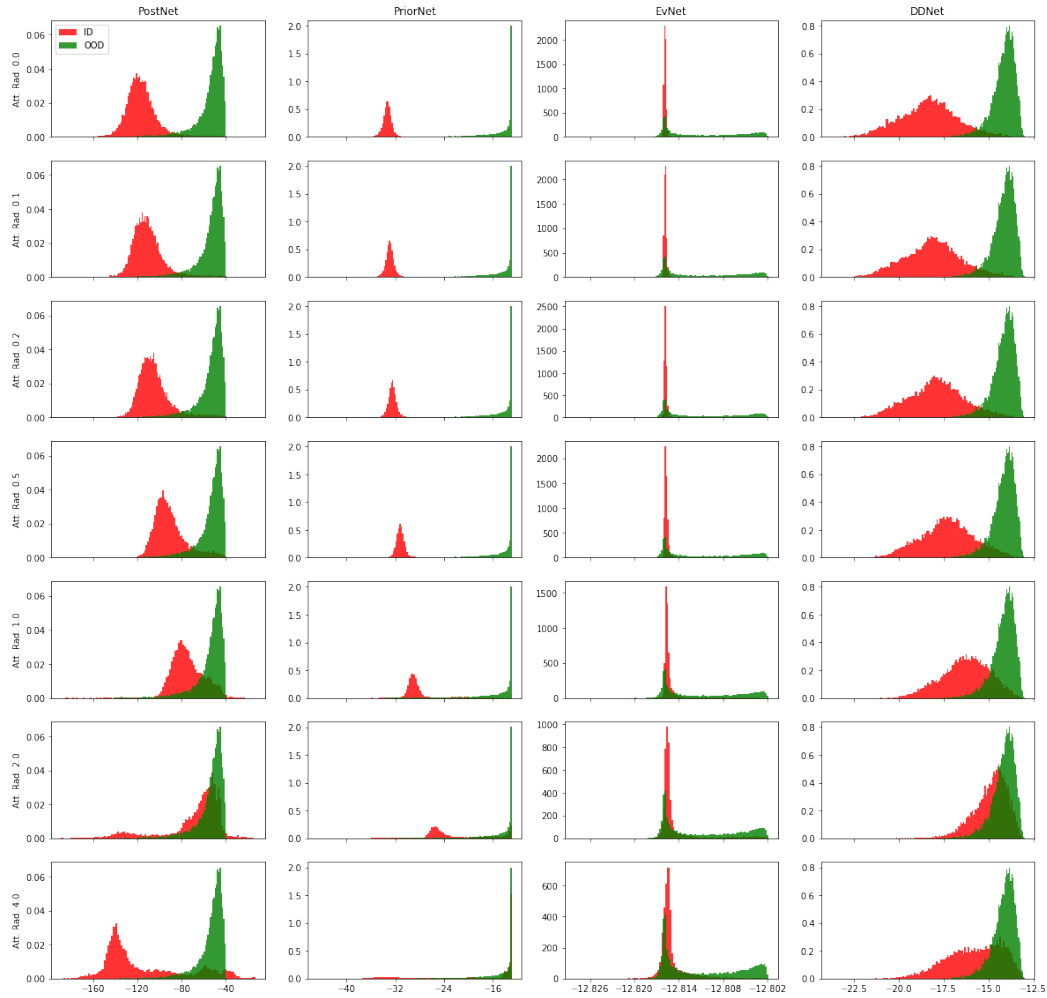


Figure 5: Visualization of the differential entropy distribution of ID data (MNIST) and OOD data (KMNIST) under label attack. The first row corresponds to no attack. The other rows correspond to increasingly stronger attack strength.

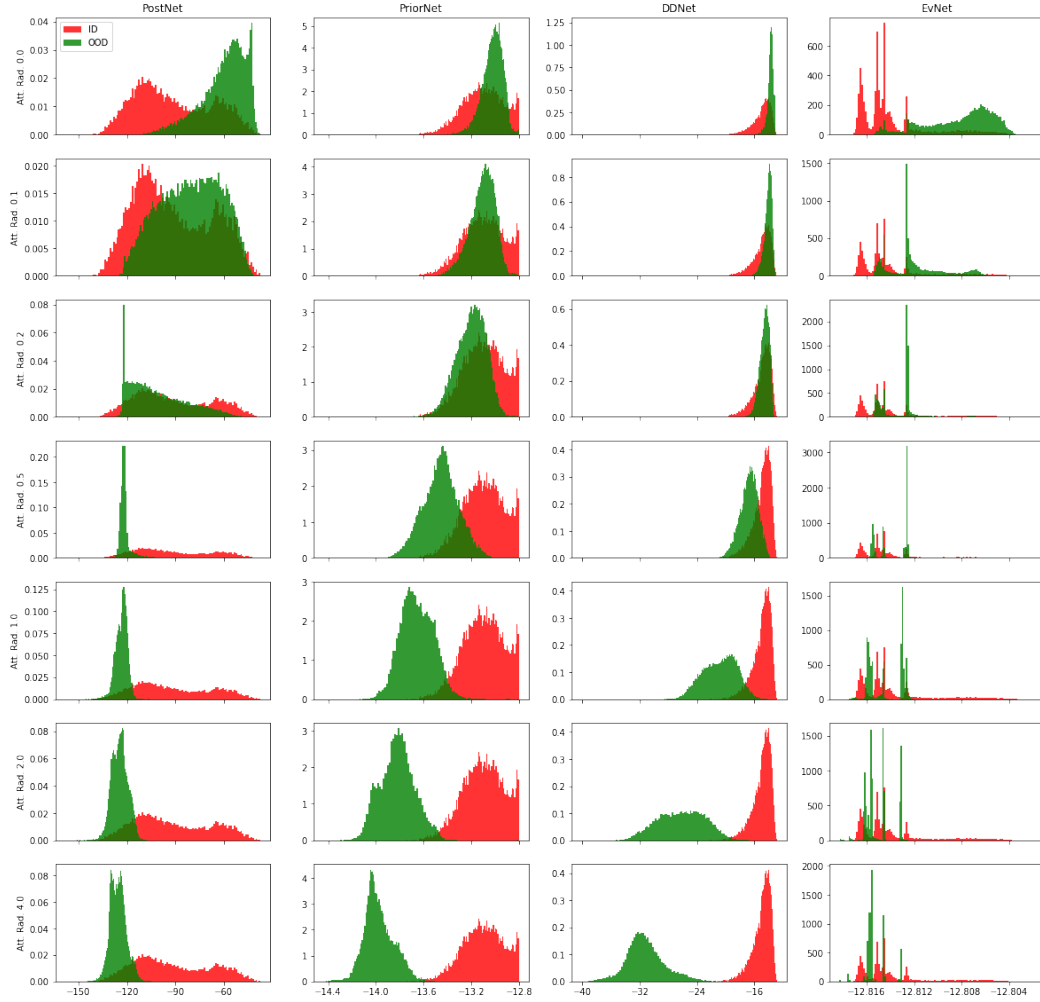


Figure 6: Visualization of the differential entropy distribution of ID data (CIFAR10) and OOD data (SVHN) under OOD uncertainty attack. The first row corresponds to no attack. The other rows correspond to increasingly stronger attack strength.

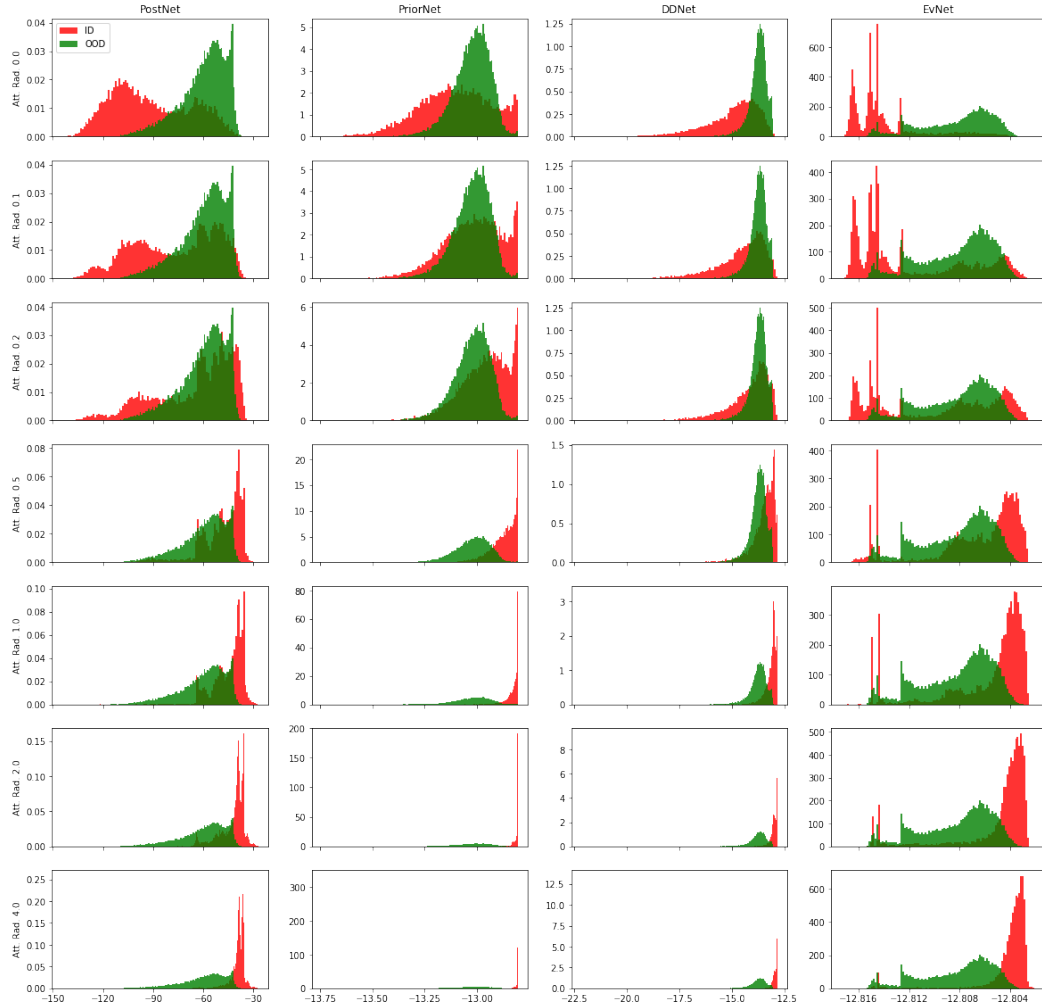


Figure 7: Visualization of the differential entropy distribution of ID data (CIFAR10) and OOD data (SVHN) under ID uncertainty attack. The first row corresponds to no attack. The other rows correspond to increasingly stronger attack strength.



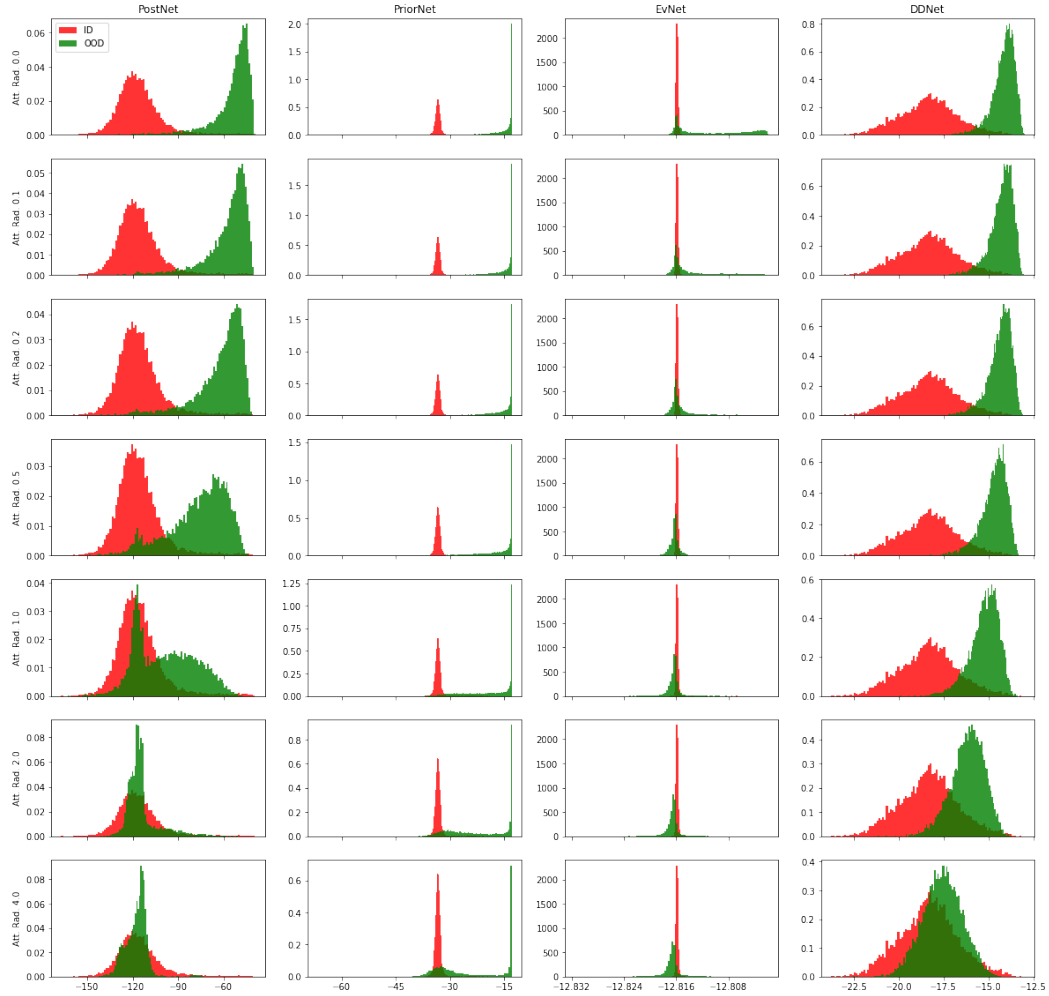


Figure 8: Visualization of the differential entropy distribution of ID data (MNIST) and OOD data (KMNIST) under OOD uncertainty attack. The first row corresponds to no attack. The other rows correspond to increasingly stronger attack strength.

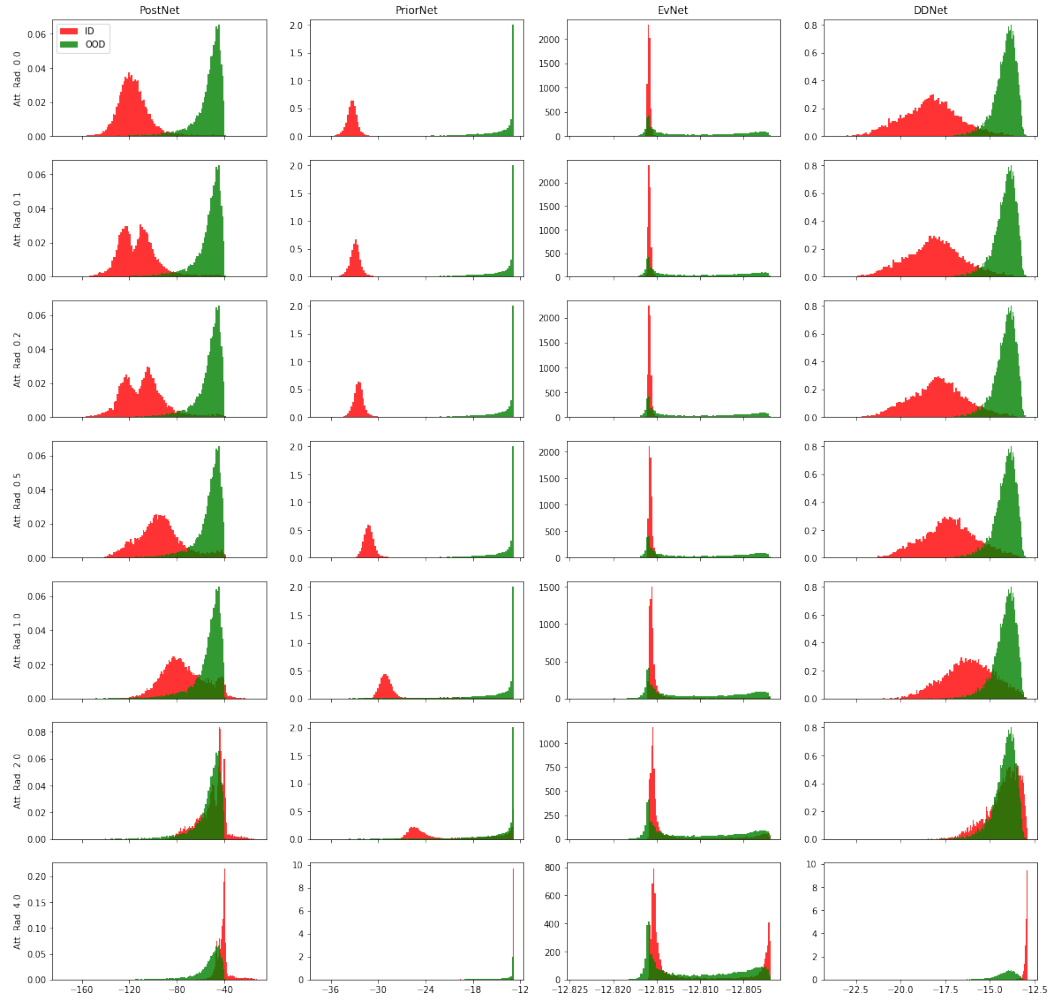


Figure 9: Visualization of the differential entropy distribution of ID data (MNIST) and OOD data (KMNIST) under ID uncertainty attack. The first row corresponds to no attack. The other rows correspond to increasingly stronger attack strength.