

Table 1. Comparison of the baseline method (Hendrycks &amp; Gimpel, 2017) and confidence-based thresholding. All results are averaged over 5 runs. All values are shown in percentages. ↓ indicates that lower values are better, while ↑ indicates that higher scores are better.

Model In-distribution Dataset	Out-of-distribution Dataset	Classification Error ↓	FPR (95% TPR) ↓	Detection Error ↓	AUROC ↑	AUPR In ↑	AUPR Out ↑
<b>Baseline (Hendrycks &amp; Gimpel, 2017)/Confidence Thresholding</b>							
<b>DenseNet-BC</b> SVHN	TinyImageNet (resize)	<b>2.89/2.77</b>	7.2/ <b>1.5</b>	5.3/ <b>2.8</b>	98.4/ <b>99.5</b>	99.4/ <b>99.8</b>	95.6/ <b>98.7</b>
	LSUN (resize)		6.0/ <b>1.0</b>	4.9/ <b>2.3</b>	98.6/ <b>99.7</b>	99.5/ <b>99.9</b>	96.0/ <b>99.0</b>
	iSUN		6.0/ <b>0.9</b>	4.9/ <b>2.3</b>	98.6/ <b>99.7</b>	99.5/ <b>99.9</b>	95.7/ <b>98.8</b>
	All Images		12.2/ <b>4.2</b>	7.2/ <b>4.5</b>	97.3/ <b>98.9</b>	95.1/ <b>97.4</b>	98.4/ <b>99.4</b>
<b>WRN-16-8</b> SVHN	TinyImageNet (resize)	<b>2.77/2.66</b>	10.6/ <b>1.5</b>	6.1/ <b>2.7</b>	97.8/ <b>99.6</b>	99.2/ <b>99.8</b>	93.6/ <b>99.2</b>
	LSUN (resize)		9.5/ <b>0.6</b>	5.8/ <b>1.8</b>	98.0/ <b>99.8</b>	99.3/ <b>99.9</b>	94.0/ <b>99.5</b>
	iSUN		9.6/ <b>0.8</b>	5.9/ <b>2.1</b>	98.0/ <b>99.8</b>	99.3/ <b>99.9</b>	93.4/ <b>99.4</b>
	All Images		15.7/ <b>5.3</b>	7.9/ <b>5.0</b>	96.7/ <b>98.7</b>	94.1/ <b>96.8</b>	97.9/ <b>99.4</b>
<b>VGG13</b> SVHN	TinyImageNet (resize)	<b>3.05/2.98</b>	11.4/ <b>1.8</b>	6.2/ <b>3.1</b>	97.8/ <b>99.6</b>	99.2/ <b>99.8</b>	93.7/ <b>99.1</b>
	LSUN (resize)		9.4/ <b>0.8</b>	5.7/ <b>2.0</b>	98.1/ <b>99.8</b>	99.3/ <b>99.9</b>	94.3/ <b>99.6</b>
	iSUN		10.0/ <b>1.0</b>	6.0/ <b>2.2</b>	98.0/ <b>99.8</b>	99.3/ <b>99.9</b>	93.7/ <b>99.5</b>
	All Images		14.2/ <b>4.3</b>	7.1/ <b>4.6</b>	97.3/ <b>99.2</b>	95.9/ <b>98.5</b>	98.2/ <b>99.6</b>
<b>DenseNet-BC</b> CIFAR-10	TinyImageNet (resize)	<b>4.17/4.39</b>	44.9/ <b>33.8</b>	12.8/ <b>12.3</b>	93.2/ <b>94.2</b>	94.6/ <b>95.0</b>	91.2/ <b>93.0</b>
	LSUN (resize)		38.6/ <b>30.7</b>	10.8/ <b>10.3</b>	94.6/ <b>95.4</b>	95.9/ <b>96.4</b>	92.8/ <b>93.9</b>
	iSUN		41.4/ <b>31.6</b>	11.6/ <b>11.0</b>	94.1/ <b>95.0</b>	95.8/ <b>96.3</b>	91.3/ <b>93.0</b>
	All Images		40.9/ <b>28.9</b>	11.6/ <b>10.9</b>	94.1/ <b>95.3</b>	87.6/ <b>88.1</b>	98.3/ <b>98.7</b>
<b>WRN-28-10</b> CIFAR-10	TinyImageNet (resize)	<b>3.25/3.46</b>	41.0/ <b>26.6</b>	14.3/ <b>11.6</b>	91.0/ <b>94.5</b>	88.9/ <b>94.1</b>	90.5/ <b>94.0</b>
	LSUN (resize)		34.7/ <b>24.0</b>	11.7/ <b>9.1</b>	93.7/ <b>96.0</b>	93.4/ <b>96.6</b>	92.7/ <b>94.5</b>
	iSUN		36.7/ <b>24.9</b>	12.6/ <b>9.8</b>	92.8/ <b>95.7</b>	92.6/ <b>96.5</b>	91.1/ <b>94.0</b>
	All Images		36.1/ <b>23.3</b>	12.4/ <b>9.7</b>	92.9/ <b>95.7</b>	73.3/ <b>86.7</b>	98.1/ <b>98.8</b>
<b>VGG13</b> CIFAR-10	TinyImageNet (resize)	<b>5.28/5.44</b>	43.8/ <b>18.4</b>	12.0/ <b>9.4</b>	93.5/ <b>97.0</b>	94.6/ <b>97.3</b>	91.7/ <b>96.9</b>
	LSUN (resize)		41.9/ <b>16.4</b>	11.5/ <b>8.3</b>	94.0/ <b>97.5</b>	95.1/ <b>97.8</b>	92.2/ <b>97.2</b>
	iSUN		41.2/ <b>16.3</b>	11.4/ <b>8.5</b>	94.0/ <b>97.5</b>	95.5/ <b>98.0</b>	91.5/ <b>96.9</b>
	All Images		41.6/ <b>19.2</b>	11.7/ <b>9.1</b>	93.9/ <b>97.1</b>	85.5/ <b>92.0</b>	98.2/ <b>99.3</b>

we equally weight  $P_{in}$  and  $P_{out}$  as if they have the same probability of appearing in the test set.

**AUROC:** Measures the Area Under the Receiver Operating Characteristic curve. The Receiver Operating Characteristic (ROC) curve plots the relationship between TPR and FPR. The area under the ROC curve can be interpreted as the probability that a positive example (in-distribution) will have a higher detection score than a negative example (out-of-distribution).

**AUPR:** Measures the Area Under the Precision-Recall (PR) curve. The PR curve is made by plotting precision =  $TP/(TP + FP)$  versus recall =  $TP/(TP + FN)$ . In our tests, AUPR-In indicates that in-distribution examples are used as the positive class, while AUPR-Out indicates that out-of-distribution examples are used as the positive class.

#### 4.2.4. MODEL TRAINING

We evaluate our confidence estimation technique by applying it to several different neural network architectures: DenseNet (Huang et al., 2017), WideResNet, (Zagoruyko & Komodakis, 2016), and VGGNet (Simonyan & Zisserman, 2015). Following Liang et al. (2018), we train a DenseNet with depth  $L = 100$  and growth rate  $k = 12$  (hereafter referred to as DenseNet-BC), as well as a WideResNet. We use a depth of 16 and a widening factor of 8 for the SVHN

dataset, and a depth of 28 and widening factor of 10 for CIFAR-10 (referred to as WRN-16-8 and WRN-28-10 respectively). We also train a VGG13 model to evaluate performance on network architectures without skip connections.

All models are trained using stochastic gradient descent, with Nesterov momentum of 0.9. We also apply standard data augmentation (cropping and flipping) and Cutout. Following Huang et al. (2017), DenseNet-BC is trained for 300 epochs with batches of 64 images, and a weight decay of  $1e-4$ . The learning rate is initialized to 0.1 at the beginning of training and is reduced by a factor of  $10\times$  after the 150th and 225th epochs. We use the settings of Zagoruyko & Komodakis (2016) for both WideResNet and VGG13, training for 200 epochs with batches of size 128, and a weight decay of  $5e-4$ . The learning rate is initialized to 0.1 at the beginning of training, and reduced by a factor of  $5\times$  after the 60<sup>th</sup>, 120<sup>th</sup>, and 160<sup>th</sup> epochs.

We train two sets of models, each with 5 runs per model. The first set of models are trained without any confidence branch, and are used to evaluate the baseline method (Hendrycks & Gimpel, 2017) and ODIN (Liang et al., 2018). The second set of models is trained with a confidence estimation branch, and are used to evaluate confidence thresholding with and without input preprocessing. A budget value of  $\beta = 0.3$  is