# Supplementary Material for: SymDNN: Simple & Effective Adversarial Robustness for Embedded Systems

Swarnava DeyPallab DasguptaPartha P ChakrabartiIndian Institute of Technology Kharagpur, Kharagpur, India - 721302

swarnava.dey@kgpian.iitkgp.ac.in, {pallab, ppchak}@cse.iitkgp.ac.in

The supplementary material is organized as follows:

- Sec. 1 presents the detailed results of MNIST and ImageNet dataset, as referred from the main paper. Further details and visualization of CIFAR-10 experiments on adversarial robustness are also presented here.
- Sec. 2 presents the, computation load for clustering and some ablation studies on the most important hyperparameters of SymDNN.
- Sec. 3 presents the details about all the building blocks used by us in this work with their licensing details.

### **1. Extended Results and Visualizations**

This section presents more elaborate results supporting SymDNN's inference mechanism over multiple datasets.

For the experiments in the current section, a model white-box / defense black-box setting is used, where SymDNN works best (see Sec. 5.2 of main paper). The gradient obfuscation [3] problem is mostly manifested in a complete white-box setting. We use 100 iterations for finding adversarial examples in the gradient based attacks, and 10 iteration for the Expectation over Transformation (EoT) [4]. These are standard iteration values used in the literature, and default in TorchAttacks [17]. Using very high values for iterations may imply that an adversary has access to an enormous amount of computing power and significant time at hand for crafting adversarial examples [22].

It is important to note that with increase in the strength  $\epsilon$  of an attack, the attacked images begin to show visually identifiable distortions. Inherent robustness to adversarial attacks is essential within a bound, where the changes are visually imperceptible. The visualizations in the current section highlight this aspect. **This part is referred in Sec. 5.2 of the main paper.** 

#### 1.1. Adversarial Robustness of MNIST

Tab. 1 provides extended results on the attacks on the MNIST dataset and the corresponding LeNet-5 model [9].

For most of the attacks, when the attack is *effective*, namely the attacked image is misclassified but no visually identifiable patterns are present, SymDNN successfully defends those. This happens when the maximum changes per pixel is limited to  $\epsilon = 8/255$ . For higher attack strengths SymDNN's inference accuracy drops, albeit remaining much higher than the standard non-symbolic inference accuracy. Fig. 1 and Fig. 2 provide visual illustrations.

Table 1. MNIST symbolic inference accuracy after attack. Results are reported on randomly selected 2000 images from MNIST test set. 's' indicates SymDNN inference that uses a 64 symbol codebook.

Attacks	Strength	Acc.	Acc. <sup>s</sup>
Auto Attack [8]	$\epsilon = 8/255$	95.60	98.05
	$\epsilon = 16/255$	85.50	93.65
	$\epsilon = 32/255$	22.35	68.60
EOTPGD [33]	$\epsilon = 8/255$	95.65	98.05
	$\epsilon = 16/255$	85.90	93.65
	$\epsilon = 32/255$	29.80	69.75
DIFGSM [30]	$\epsilon = 8/255$	96.45	98.05
	$\epsilon = 16/255$	86.95	94.00
	$\epsilon = 32/255$	31.00	69.20
RFGSM [29]	$\epsilon = 8/255$	95.60	98.05
	$\epsilon = 16/255$	85.95	93.65
	$\epsilon = 32/255$	29.09	68.70
MIFGSM [12]	$\epsilon = 8/255$	95.60	98.05
	$\epsilon = 16/255$	85.95	93.60
	$\epsilon = 32/255$	29.60	69.10
PGD [19]	$ \begin{aligned} \boldsymbol{\epsilon} &= 8/255 \\ \boldsymbol{\epsilon} &= 16/255 \\ \boldsymbol{\epsilon} &= 32/255 \end{aligned} $	95.75 86.15 29.60	98.05 93.65 69.10
OnePixel [26]	pixels = 1	98.60	98.60
	pixels = 5	94.15	94.60
	pixels = 10	94.55	94.65

# 1.2. Adversarial Robustness of CIFAR-10

We have reported the efficacy of SymDNN on the CIFAR-10 dataset and the corresponding model [32] in the main paper. Some of the attacks (AutoAttack [8],



(a) Autoattack with  $\epsilon = 16/255$ : there is neither identifiable visual distortion, nor missclassification on attack



(b) Autoattack with  $\epsilon = 32/255$ : pattern for 2 faded out at places and is identifiable. SymDNN provides resistance and survives the attack.



(c) DIFGSM with  $\epsilon = 32/255$ : pattern for 7 faded out at places and is identifiable. SymDNN provides resistance and survives the attack.

Figure 1. A visual comparison. For each set of 3 images Column 1 (Orig:) shows the clean original image, Column 2 (Sym:) shows the image reconstructed by SymDNN, Column 3 (orig perturb:) shows the original image after attack and Column 4 (Sym perturb:) shows the symbolically reconstructed image after attack.

For  $\epsilon = 16/255$ , the adversarial attack is effective on the original image, but fails to cause misclassification in SymDNN. At  $\epsilon > 16/255$ , we observe that SymDNN misclassified, but the attacks are not adversarially *effective*, that is, the attacked image has visually identifiable changes. High  $\epsilon$  is used here to highlight the patch replacement by representative symbols in SymDNN.

DIFGSM [30], RFGSM [29], APGDT [8], and variants of APGD [8]), using a strength greater than  $\epsilon$ =2/255 ( $L_{\infty}$ ) results in visually identifiable changes on CIFAR-10 images. Fig. 3 illustrates one such example where AutoAttack [8] is applied at different strengths. It may be noted that both standard and SymDNN robust accuracy drops heavily in these cases. This is consistent with the observation in [1], that k-centroid clustering is inherently robust to small set of changes, if such a clustering is well separated. We find most of our centroids are well separated in terms of the dissimilarity measure used in [16]. As it is difficult to visualize a 2048 × 2048 matrix here, we present a heatmap of cluster distances in Fig. 4 for MNIST similarity index, with 32 centroid patches.

We report the results on CIFAR-10 using L2 distance measure in Tab. 2.

To compare the SymDNN inference and the standard non-symbolic inference in terms of numbers, we consider the following outcomes:

Table 2. CIFAR-10 symbolic inference accuracy after attack using  $L_2$  distance. Results are reported on randomly selected 2000 images from CIFAR-10 test set. <sup>(s)</sup> indicates SymDNN inference that uses a 2048 symbol codebook. MSR<sup>s</sup> is the Multi-Sym Randomized inference mechanism defined in the Sec. 4.2.2 of the main paper.

Attacks	Strength	Acc.	Acc. <sup>s</sup>
Auto Attack [8]	$\epsilon = 8/255$	57	81
Square [2]	$\epsilon = 8/255$	79	86
FAB [7]	$\epsilon = 8/255$	58	85
APGDT [8]	$\epsilon = 8/255$	57	82
APGD(CE) [8]	$\epsilon = 8/255$	57	81
APGD(DLR) [8]	$\epsilon = 8/255$	58	82

1. *Clean non-symbolic*. The original image is correctly classified.



Figure 2. Efficacy of SymDNN on MNIST digits under PGD attack: for each set of 4 images the first column (Orig:) is the clean original image, second column (Sym:) is the symbolically abstracted clean image, third column (orig perturb:) is the image after attack and the fourth column (Sym perturb:) is the symbolically abstracted image after attack.

At  $\epsilon = 8/255$ , the attack is very *effective*, i.e., the attacked image is misclassified but there are no visual changes in the image (orig perturb). SymDNN not only results in a correct classification in most cases (except for 8), but keeps a clear mark in the symbolically abstracted perturbed image (Sym perturb), depicting that some undesirable changes have happened to the image. The reader is requested to zoom in for finding the pattern in the images labeled *Sym ( perturb)*. One such image is presented in the last row.



(c) At  $\epsilon = 2/255$  changes are not visually identifiable

Figure 3. Visually identifiable changes in some CIFAR-10 images under Autoattack: For each set of 4 images the first column (Orig:) shows the clean original image, the second column (Sym:) shows the symbolically reconstructed clean image, the third column (orig perturb:) shows the image after attack, and the fourth column (Sym perturb:) shows the symbolically reconstructed image after attack. Autoattack results in visually identifiable patterns on the images for changes per pixel greater than 2/255

- 2. *Clean symbolic*. The SymDNN reconstructed image is correctly classified.
- 3. *Robust non-symbolic (post attack).* The original image is correctly classified even after attack.
- 4. *Robust symbolic (post attack).* The SymDNN reconstructed image is correctly classified even after attack.

Based on these outcomes, we partition the results into the following cases:

- Case-1 This refers to the instances where all the above correctly infer the ground truth label.
- Case-2 This refers to the rare cases where only the

SymDNN reconstructed image is misclassified after attack.

- Case-3 This refers to the instances where the clean images were correctly classified, the SymDNN reconstructed images were correctly classified after attack, but the original images were misclassified after the attack. These instances are most relevant for demonstrating the benefit of SymDNN.
- Case-4 These are instances where both the original and SymDNN reconstructed images are misclassified after the attack. As discussed earlier, many of these cases show significant visual distortions after attack and thereby rendering the attacks ineffective from an

o-0.00 0.43 0.38 1.12	0.73 0.16 0.12	.95 0.21 0.29 0.	09 1.28 0.46 (	0.49 0.33 0.07 0.2	24 0.96 1.70 1	1.00 0.47 1.47 0.96	0.45 0.37 1.35 1.69 0	.47 0.46 0.72 2.31 0.7	3.0
	0.66.0.21.0.73.0	93 1 06 1 13 0	15 0 34 0 19	1 06 0 98 0 75 0 6	52 1 13 0 68 0	90 0 77 0 99 0 70	1 10 1 33 0 60 0 96 0	65 1 11 1 11 1 09 0 1	
		56 0.05 0.07 0.				70 0 20 1 62 1 69	0.12, 0.02, 1, 60, 2, 22, 1	11 0.06 0.24 2.20 1.7	
112 040 217 0.00	1.26 0.47 1.12 0		00 0.22 0.07				1.62.2.20.0.75.0.24.6		
m - 1.12 0.40 2.17 0.00	1.36 0.47 1.13 0	.41 2.02 1.84 0.	80 0.22 0.87 .	1.22 1.97 1.59 1.5	59 0.98 0.11 1	1.15 1.61 0.66 0.21	1.63 2.38 0.75 0.34 0	1.47 1.93 1.54 0.35 0.1	1
<del>√</del> -0.73 0.66 0.56 1.36	0.00 0.92 0.86 1	.12 0.72 0.73 0.	60 0.92 0.20 (	0.74 0.39 0.70 0.3	35 0.65 1.44 0	0.16 0.09 0.66 1.27	0.44 0.77 0.34 1.25 1	.27 0.32 0.27 1.74 0.9	9
ю-0.16 0.21 0.89 0.47	0.92 0.00 0.25 1	.30 0.69 0.68 0.	13 0.74 0.50 (	0.59 0.81 0.41 0.5	59 0.86 0.94 1	1.02 0.82 1.11 0.45	0.76 0.95 1.09 1.05 0	0.18 0.89 0.95 1.44 0.2	8 2.5
0.12 0.73 0.33 1.13 و.	0.86 0.25 0.00 1	.88 0.23 0.15 0.	33 1.47 0.74 (	0.17 0.38 0.12 0.4	40 0.61 1.69 0	0.90 0.55 1.25 0.74	0.26 0.36 1.48 1.51 0	0.26 0.40 0.53 2.28 0.9	3
⊾ - 1.95 0.93 2.56 0.41	1.12 1.30 1.88 0	.00 2.67 2.40 1.	52 0.23 1.07	1.55 2.33 2.34 2.0	04 0.82 0.14 0	0.72 1.62 0.20 0.54	1.80 2.96 0.30 0.11 1	.22 2.05 1.36 0.11 0.6	2
∞-0.21 1.06 0.05 2.02	0.72 0.69 0.23 2	.67 0.00 0.06 0.	46 2.11 0.73 (	0.44 0.07 0.04 0.1	16 1.05 2.64 1	1.02 0.32 1.82 1.62	0.23 0.02 1.75 2.42 0	.96 0.18 0.54 3.32 1.5	4
σn - 0.29 1.13 0.07 1.84	0.73 0.68 0.15 2	.40 0.06 0.00 0.	57 2.05 0.81 0	0.20 0.13 0.11 0.2	24 0.73 2.40 0	0.86 0.34 1.51 1.34	0.11 0.07 1.66 2.11 0	.76 0.17 0.39 3.05 1.5	1
음-0.09 0.15 0.64 0.80	0.60 0.13 0.33 1	.52 0.46 0.57 0.	00 0.84 0.21 0	0.71 0.48 0.26 0.2	25 1.06 1.27 0	0.92 0.49 1.28 0.89	0.68 0.66 0.97 1.44 0	.55 0.66 0.86 1.82 0.4	- 2.0
긐-1.28 0.34 <mark>2.15</mark> 0.22	0.92 0.74 1.47 0	.23 2.11 2.05 0.	84 0.00 0.61	1.52 1.88 1.72 1.5	51 1.05 0.15 0	0.82 1.32 0.51 0.53	1.67 2.45 0.29 0.32 0	.97 1.77 1.37 0.26 0.1	8
္ဌ-0.46 0.19 0.76 0.87	0.20 0.50 0.74 1	.07 0.73 0.81 0.	21 0.61 0.00 0	0.88 0.51 0.59 0.2	27 0.95 1.09 0	0.50 0.29 0.87 1.07	0.71 0.88 0.40 1.23 0	.98 0.62 0.67 1.50 0.4	6
<u>m</u> -0.49 1.06 0.39 1.22	0.74 0.59 0.17 1	.55 0.44 0.20 0.	71 1.52 0.88 (	0.00 0.47 0.41 0.5	57 0.23 1.59 0	0.55 0.51 0.81 0.67	0.12 0.51 1.25 1.22 0	.38 0.35 0.26 2.07 1.1	8
<u></u> - 0.33 0.98 0.04 1.97	0.39 0.81 0.38 2	.33 0.07 0.13 0.	48 1.88 0.51 0	0.47 0.00 0.14 0.0	07 0.92 2.45 0	0.70 0.11 1.51 1.65	0.18 0.07 1.33 2.22 1	13 0.08 0.36 3.05 1.4	8
ហ្ន-0.07 0.75 0.14 1.59	0.70 0.41 0.12 2	.34 0.04 0.11 0.	26 1.72 0.59 0	0.41 0.14 0.00 0.1	16 0.96 2.20 0	0.98 0.35 1.64 1.28	0.28 0.13 1.56 2.06 0	.68 0.26 0.57 2.85 1.1	6
ទុ-0.24 0.62 0.19 1.59	0.35 0.59 0.40 2	.04 0.16 0.24 0.	25 1.51 0.27 0	0.57 0.07 0.16 0.0	00 1.03 2.06 0	0.74 0.15 1.44 1.50	0.34 0.20 1.12 2.03 1	.02 0.24 0.52 2.67 1.0	- 1.5
ːː-0.96 1.13 0.87 0.98	0.65 0.86 0.61 0	.82 1.05 0.73 1.	06 1.05 0.95 0	0.23 0.92 0.96 1.0	03 0.00 1.05 0	0.24 0.66 0.25 0.44	0.35 1.15 0.74 0.57 0	.55 0.61 0.22 1.27 1.0	9
<u>∞</u> - 1.70 0.68 2.68 0.11	1.44 0.94 1.69 0	.14 2.64 2.40 1.	27 0.15 1.09	1.59 2.45 2.20 2.0	06 1.05 0.00 1	1.10 1.87 0.48 0.35	1.98 3.00 0.57 0.16 0	.89 2.28 1.69 0.07 0.3	1
ត្ន-1.00 0.90 0.79 1.15	0.16 1.02 0.90 0	.72 1.02 0.86 0.	92 0.82 0.50 0	0.55 0.70 0.98 0.7	74 0.24 1.10 0	0.00 0.31 0.24 0.86	0.41 1.09 0.28 0.73 1	.05 0.45 0.15 1.28 1.0	4
<sub>ର</sub> - 0.47 0.77 0.20 1.61	0.09 0.82 0.55 1	.62 0.32 0.34 0.4	49 1.32 0.29 0	0.51 0.11 0.35 0.1	15 0.66 1.87 0	0.31 0.00 0.97 1.37	0.20 0.34 0.74 1.63 1		9
<mark>5 -</mark> 1.47 0.99 1.62 0.66	0.66 1.11 1.25 0	.20 1.82 1.51 1.	28 0.51 0.87 (	0.81 1.51 1.64 1.4	44 0.25 0.48 0	0.24 0.97 0.00 0.45	0.95 1.99 0.27 0.19 0	.92 1.18 0.59 0.53 0.8	- 1.0
ಜ್ಞ-0.96 0.70 1.68 0.21	1.27 0.45 0.74 0	.54 1.62 1.34 0.	89 0.53 1.07 (	0.67 1.65 1.28 1.5	50 0.44 0.35 0	0.86 1.37 0.45 0.00	1.08 1.92 0.90 0.24 0	.19 1.44 1.01 0.58 0.4	2
ຕູ-0.45 1.10 0.12 1.63	0.44 0.76 0.26 1	.80 0.23 0.11 0.	68 1.67 0.71 (	0.12 0.18 0.28 0.3	34 0.35 1.98 0	0.41 0.20 0.95 1.08	0.00 0.24 1.16 1.55 0	.75 0.07 0.09 2.46 1.4	0
± - 0.37 1.33 0.03 2.38	0.77 0.95 0.36 2	.96 0.02 0.07 0.	66 2.45 0.88 (	0.51 0.07 0.13 0.2	20 1.15 3.00 1	1.09 0.34 1.99 1.92	0.24 0.00 1.94 2.72 1	.21 0.17 0.55 3.71 1.8	8
∾ - 1.35 0.60 1.60 0.75	0.34 1.09 1.48 0	130 1.75 1.66 <b>0</b> .	97 0.29 0.40	1.25 1.33 1.56 1.1	12 0.74 0.57 0	0.28 0.74 0.27 0.90	1.16 1.94 0.00 0.53 1	.34 1.17 0.78 0.66 0.6	3
Q 1 69 0 96 2 33 0 34	1 25 1 05 1 51 0	11 2 42 2 11 1	44 0 32 1 23	1 22 2 22 2 06 2 0	03 0 57 0 16 0	73 1 63 0 19 0 24	1 55 2 72 0 53 0 00 0	81 1 87 1 20 0 16 0 6	1
A 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	1 27 0 18 0 26 1	22 0.96 0.76 0	55 0.97 0.98 (		0.55 0.89 1		0.75 1.21 1.34 0.81 0		- 0.5
× 0.47 0.03 1.11 0.47	0.32 0.89 0.40	05 0 18 0 17 0	66 1 77 0 62 0	0.35 0.08 0.26 0.2	24 0.61 2.22 0	45 0 09 1 19 1 44	0.07 0.17 1.17 1.07 1	07 0.00 0.12 2.76 1.5	4
N -0.40 1.11 0.00 1.95	0.32 0.05 0.40 2	26 0 54 0 20 0	96 1 37 0 C7	0.35 0.06 0.28 0.2	E2 0.22 1 60 0		0.00 0.55 0.78 1.20	07 0.00 0.13 2.70 1.5	7
$\sim -0.72$ 1.11 0.34 1.54	0.27 0.95 0.53 1		00 1.57 0.07 (	0.20 0.30 0.57 0.5	52 0.22 <b>1.</b> 59 0	0.19 0.19 0.59 1.01	0.09 0.55 0.78 1.20 0	24 2 76 2 62 0 6 6 7	
<u>≈</u> -2.31 1.09 <b>3.29</b> 0.35	1.74 1.44 2.28 0	.11 3.32 3.05 1.	oz 0.26 1.50 i	2.07 3.05 2.85 2.6	57 1.27 0.07 1		2.46 3.71 0.66 0.16 1		4
m - 0.73 0.10 1.72 0.11	0.99 0.28 0.93 0		40 0.18 0.46 : 0 11 12	1.18 1.48 1.16 1.0	6 17 18	1.04 1.19 0.83 0.42	1.40 1.88 0.63 0.61 0	1.53 1.54 1.37 0.64 0.0 27 28 29 30 31	- 0.0

Figure 4. A heatmap of cluster centroid separation for MNIST clustering model with 32 centroids .

adversarial point of view.

- Case-5 These are instances where the SymDNN reconstructed images are misclassified without any attack.
- Case-6 These are instances where the original clean images are misclassified.
- The case distributions and accuracy of SymDNN and non-

symbolic inference against some adversarial attacks are shown in Fig. 5. We present these results in two different sets:

- 1. Set 1 (Fig. 5b) illustrates the attacks against which SymDNN performs very well, and
- 2. Set 2 (Fig. 5d)the attack for which the robust accuracy



Figure 5. The distribution of inference outcomes for the 6 different cases defined in Sec. 1.2. We segregate the adversarial attacks into two different sets, based on their effect on robust accuracy. Figures (a) and (b) illustrate the accuracy and distribution of cases respectively, for the attacks against which SymDNN results in high robust accuracy gain. Figures (c) and (d) illustrate the accuracies and distribution of cases respectively, for the attacks against which SymDNN results in relatively lower robust accuracy gain.

#### gain is limited.

The attacks in *Set 2*, mostly generate adversarial examples with improved transferability. The perturbations in such examples are relatively higher as these attacks try to maximize the confidence instead of minimizing the distance [10]. We have already discussed this issue in Sec. 5.2 & Sec. 5.4 of the main paper, and we believe that a more strongly separable clustering may help SymDNN to thwart adversarial examples with larger perturbations.

#### **1.3. ImageNet Adversarial Robustness.**

We perform ImageNet adversarial tests with Fool-Box [23]. Tab. 3 presents the results of these preliminary experiments.

In this section we report some interesting observations regarding SymDNN inference on Imagenet images. Fig. 6 illustrates examples where we find that SymDNN preserves inference against adversarial attacks, however the symbol-

Table 3. ImageNet robust symbolic Foolbox.

Attack	Acc.	$\mathrm{Acc.}^{s}$
FGSM [14,27]	0	62.5
Adam-PGD [23]	0	68.5
CW [5]	0	68.5
PGD [19]	0	68.75
DeepFool $L_2$ [20]	0	31.25

ically reconstructed images contain spurious patterns that are not present in the original images. Fig. 7 illustrates the different cases of symbolic and non-symbolic inference on ImageNet.

#### **1.4. Illustrative Example on Compaction.**

We provide an illustration below for data reduction results, as presented in the main paper. The first part of this



(a) The symbolic image showing a pattern not present in the original image: leafhopper under C & W attack



(b) The symbolic image of loafer showing a pattern inside the shadow region, absent in the original image: PGD



(d) The symbolic image showing marks in the upper left and right, absent in the original image: PGD

Figure 6. Some symbolically abstracted images have visual distortion, although the symbolic inference is correct. This shows that SymDNN inference may not be used for visual reconstruction, however it is inference preserving from a CNN perspective. For each set of 4 images the first column (Orig:) is the clean original image, second column (Sym:) is the symbolically abstracted clean image, third column (orig perturb:) is the image after attack and the fourth column (Sym perturb:) is the symbolically abstracted image after attack.



(b) Both symbolic and non-symbolic inferences are affected by attack



'pickup







(d) Only non-symbolic inference is affected by attack

Figure 7. Visualization of different inference outcomes of ImageNet images under different adversarial attacks: for each set of 4 images the first column (Orig:) is the clean original image, second column (Sym:) is the symbolically abstracted clean image, third column (orig perturb:) is the image after attack and the fourth column (Sym perturb:) is the symbolically abstracted image after attack.

Sym (perturb): 'jeep



Sym (perturb):



reduction,  $\Delta^a$ , comes from the symbolic abstraction. We can calculate this guaranteed compaction by finding out the reduction for representing each patch. For a 2 × 2 patch we need 2 × 2 × 8 = 32 bits, assuming 8 bit per pixel representation. For representing a 2 × 2 patch we use one symbol. Considering a 512 symbol codebook, this requires 9 bits to encode one symbol. This results in the  $\Delta^a$  compaction of 71.8%.

To employ entropy encoding over this compaction, we assume a distribution similar to the 10000 image testset of CIFAR-10. We apply Huffman encoding and obtain a distribution illustrated in Fig. 8. We find that using a 512 symbol codebook results in extra 22.67% compaction. The combined compaction of 78.3% is reported as the  $\Delta^e$  in the main paper. This compaction can be useful when the image has to be communicated, say, from an edge device to an edge server in a computation offloading scenario [11].



Figure 8. Huffman encoding of symbols: Variable number of bits assigned to encode the symbols based on their frequency of occurrence in the CIFAR-10 testset.

# 2. Patch Clustering

All the experiments in this section are conducted on a 32 core Intel(R) Xeon(R) Silver 4108 CPU 1.80GHz workstation with 128 Gibibytes of RAM, and NVIDIA P1000 GPU (4 Gigabytes memory). We use the GPU for unsupervised clustering of ImageNet patches.

#### 2.1. Clustering Computation Load

This part supplements Sec. 4.1.1 and Sec. 5.4 of the main paper.

Fig. 9 illustrates the increasing time required for clustering patches from the CIFAR-10 dataset, as the number of clusters are increased. The patch extraction overhead becomes comparatively less with increasing number of centroids. It may be noted that here we extract patches from each channel of the image, collate the patches from all the images, and pass these patches to FAISS *k-means* for building the index at one go.



Figure 9. Compute load of CIFAR-10 unsupervised clustering of patches.

Collating all the patches and passing those to the FAISS library has a peak memory load of 22GB for CIFAR-10 dataset. For ImageNet, index building from all patches is not feasible in our setup as the number of patches generated is very high. In this case, we extract patches from each channel of the image, collate the patches from several images, and pass these patches to FAISS *k-means* for building the index incrementally. Building an index on ImageNet training set (images of resolution  $224 \times 224$ ), with a patch size of  $2 \times 2$  and 2048 patch centroids generates a patch dataset of  $112 \times 112 \times 3 \times 150000 = 5.6448$  billion patches, each having a dimensionality of 4. Here we have M = 2048, N = 5.6448 billion and  $P^2 = 4$ . This takes 15 days in total, of which 4 days are needed for patch extraction.

With a patch size of  $4 \times 4$ , and 512 centroids a total time of 98 hours is needed of which 48 hours are needed for patch extraction.

We plan to explore the possibilities of vectorizing the patch extraction code or creating and storing a patch dataset in advance in the future, to reduce this overhead.

#### 2.2. Optimal Number of Centroids

To find the optimal number of cluster centroids for CIFAR-10, we use the *knee detection* method of [25]. Fig. 10 illustrates the change in the within cluster sum

square (WCSS) error for using an increasing number of centroids. The knee point detection algorithm of Fig. 10 aptly finds out the maximum curvature for 4 centroids. However, we find that using 4 centroids the SymDNN inference accuracy is 17%. This shows that the traditional method for finding optimal centroid may not be suitable in the context of symbolic inference.



Figure 10. Knee point detection for finding the optimal number of clusters for symbolic inference on CIFAR-10 dataset. Y-axis values of within cluster sum square (WCSS) error values are obtained from the FAISS kmeans inertia metric.

# 2.3. Optimal Patch Sizes

This part is referred from Sec. 5.4 of the main paper. Fig. 11 illustrates that  $2 \times 2$  patches result in best accuracy on MNIST dataset. The accuracy reduces with the increase in patch size, and drops by almost 10% for an  $8 \times 8$  patch size. We also observe that  $8 \times 8$  patch size gives better accuracy with a larger number of symbols.

For ImageNet, using a  $4 \times 4$  patch size and 512 symbols, the SymDNN Top-1 accuracy of ResNet-152 comes down to 46% from 70% with a  $2 \times 2$  patch size and 2048 symbols.

For MNIST, the robust accuracy remains highest on using a lower number of symbols. This is illustrated in Fig. 12.

# 2.4. Conclusion about the Ablation studies

We conclude that  $2 \times 2$  patches results in inference preservation for all the datasets we tested. There is some correlation between the combination of a particular patch size and number of symbols, with the inference accuracy, which could not be established with the limited set of experiments we performed.

We believe that Neural Architecture Search (NAS) [13, 18, 34], which has helped in designing compact DNNs can be applied to SymDNN as well. Instead of *brute force* trials, evolutionary or reinforcement learning based optimization in NAS can be applied in this context. The purpose of such



Figure 11. Inference accuracies corresponding to different patch sizes.



Figure 12. A comparison of robust accuracy for different numbers of symbols used for SymDNN inference, under different attack strengths (maximum change per pixel).

a search will be to find the optimal value of SymDNN hyperparameters, with respect to the non-smooth gradient of the SymDNN inference accuracy.

# 3. Details of Assets Used

# This part provides additional details over Sec. 5.1 of the main paper.

Architectures. For ImageNet we use publicly available pre-trained models: WideResNet [31], ResNet-152 [15] and MnasNet [28] from PyTorch model repository<sup>1</sup>. We did not find any separate licensing for these models. We point the reader to the *Facebook Open Source – Terms of Use*: https://pytorch.org/assets/tos-oss-privacy-policy/fb-tos-privacy-policy.pdf. All these models are trained on the ImageNet dataset.

<sup>&</sup>lt;sup>1</sup>https://pytorch.org/vision/stable/models.html

The licence for ImageNet can be found at: https://www.image-net.org/.

For CIFAR-10 we train the models on our own infrastructure following the ResNet-20 architecture [32]. We could not find explicit licensing terms for this repository https://github.com/zhuchen03/gradinit or the paper [32].

We use the public LeNet-5 [9] architecture for MNIST experiments.

Adversarial Attacks. We choose Torchattacks [17] and the corresponding repository *Adversarial-Attacks-PyTorch*: https://github.com/Harry24k/ adversarial-attacks-pytorch for attack implementations. This repository has a MIT license: https: //github.com/Harry24k/adversarialattacks-pytorch/blob/master/LICENSE.

We select *Adversarial-Attacks-PyTorch* for primarily two reasons: (A) It implements most of the successful and the most recent attacks. (B) Unlike many other frameworks, it does not expose an integrated API for adversarial image generation and testing of accuracy. This is important for an easy integration of SymDNN. The symbolic inference needs to modify the image after the adversarial attack, before it is fed to the DNN for inference. SymDNN can also be integrated to frameworks that do not return the adversarial examples directly, however that would require a modification in such adversarial attack frameworks.

For demonstrating efficacy of SymDNN on already robust models, we use RobustBench [6]. These robust models that we use for evaluation purpose are available from: https://robustbench.github.io/ released with a MIT license: https://github.com/ RobustBench / robustbench / blob / master / LICENSE.

We compare and integrate our work two state-of-the-art input defenses:

- NRP [21], available from: https://github. com/Muzammal-Naseer/NRP with a MIT license: https://github.com/Muzammal-Naseer/NRP/blob/master/LICENSE.
- DefenseGAN [24] PyTorch implementation, available from: https://github.com/sky4689524/ DefenseGAN-Pytorch with no specified license.

# References

 Margareta Ackerman, Shai Ben-David, David Loker, and Sivan Sabato. Clustering oligarchies. In Carlos M. Carvalho and Pradeep Ravikumar, editors, *Proceedings of the Sixteenth International Conference on Artificial Intelligence* and Statistics, volume 31 of *Proceedings of Machine Learn*- *ing Research*, pages 66–74, Scottsdale, Arizona, USA, 29 Apr–01 May 2013. PMLR. 2

- [2] Maksym Andriushchenko, Francesco Croce, Nicolas Flammarion, and Matthias Hein. Square attack: a query-efficient black-box adversarial attack via random search. September 2020. 2
- [3] Anish Athalye, Nicholas Carlini, and David A. Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In *ICML*, pages 274–283, 2018. 1
- [4] Anish Athalye, Logan Engstrom, Andrew Ilyas, and Kevin Kwok. Synthesizing robust adversarial examples. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 284–293. PMLR, 10–15 Jul 2018. 1
- [5] Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In 2017 IEEE Symposium on Security and Privacy (SP), pages 39–57, 2017. 6
- [6] Francesco Croce, Maksym Andriushchenko, Vikash Sehwag, Edoardo Debenedetti, Nicolas Flammarion, Mung Chiang, Prateek Mittal, and Matthias Hein. Robustbench: a standardized adversarial robustness benchmark, 2021. 11
- [7] Francesco Croce and Matthias Hein. Minimally distorted adversarial examples with a fast adaptive boundary attack. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 2196–2205. PMLR, 13–18 Jul 2020. 2
- [8] Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference* on Machine Learning, volume 119 of Proceedings of Machine Learning Research, pages 2206–2216. PMLR, 13–18 Jul 2020. 1, 2
- [9] Y. Le Cun, B. Boser, J. S. Denker, R. E. Howard, W. Habbard, L. D. Jackel, and D. Henderson. *Handwritten Digit Recognition with a Back-Propagation Network*, page 396–404. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1990. 1, 11
- [10] Ambra Demontis, Marco Melis, Maura Pintor, Matthew Jagielski, Battista Biggio, Alina Oprea, Cristina Nita-Rotaru, and Fabio Roli. Why do adversarial attacks transfer? explaining transferability of evasion and poisoning attacks. In 28th USENIX Security Symposium (USENIX Security 19), pages 321–338, Santa Clara, CA, Aug. 2019. USENIX Association. 6
- [11] Swarnava Dey, Jayeeta Mondal, and A. Mukherjee. Offloaded execution of deep learning inference at edge: Challenges and insights. 2019 IEEE International Conference on Pervasive Computing and Communications Workshops (Per-Com Workshops), pages 855–861, 2019. 9
- [12] Yinpeng Dong, Fangzhou Liao, Tianyu Pang, Hang Su, Jun Zhu, Xiaolin Hu, and Jianguo Li. Boosting adversarial attacks with momentum. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 1

- [13] Thomas Elsken, Jan Hendrik Metzen, and Frank Hutter. Neural architecture search: A survey. *Journal of Machine Learning Research*, 20(55):1–21, 2019. 10
- [14] Ian Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *International Conference on Learning Representations*, 2015. 6
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 770–778, 2016. 10
- [16] Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billionscale similarity search with gpus. *IEEE Transactions on Big Data*, 7(3):535–547, 2021. 2
- [17] Hoki Kim. Torchattacks: A pytorch repository for adversarial attacks, 2021. 1, 11
- [18] Chenxi Liu, Barret Zoph, Maxim Neumann, Jonathon Shlens, Wei Hua, Li-Jia Li, Li Fei-Fei, Alan Yuille, Jonathan Huang, and Kevin Murphy. Progressive neural architecture search. In *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018. 10
- [19] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings. OpenReview.net, 2018. 1, 6
- [20] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. Deepfool: a simple and accurate method to fool deep neural networks, 2016. 6
- [21] Muzammal Naseer, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, and Fatih Porikli. A self-supervised approach for adversarial robustness. In *IEEE/CVF Conference* on Computer Vision and Pattern Recognition (CVPR), June 2020. 11
- [22] Edward Raff, Jared Sylvester, Steven Forsyth, and Mark McLean. Barrage of random transforms for adversarially robust defense. In 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 6521–6530, 2019. 1
- [23] Jonas Rauber, Wieland Brendel, and Matthias Bethge. Foolbox: A python toolbox to benchmark the robustness of machine learning models. In *Reliable Machine Learning in the Wild Workshop, 34th International Conference on Machine Learning*, 2017. 6
- [24] Pouya Samangouei, Maya Kabkab, and Rama Chellappa. Defense-GAN: Protecting classifiers against adversarial attacks using generative models. In *International Conference* on Learning Representations, 2018. 11
- [25] Ville Satopaa, Jeannie Albrecht, David Irwin, and Barath Raghavan. Finding a "kneedle" in a haystack: Detecting knee points in system behavior. In 2011 31st International Conference on Distributed Computing Systems Workshops, pages 166–171, 2011. 9
- [26] Jiawei Su, Danilo Vasconcellos Vargas, and Kouichi Sakurai. One pixel attack for fooling deep neural networks. *IEEE Transactions on Evolutionary Computation*, 23(5):828–841, Oct 2019. 1

- [27] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In Yoshua Bengio and Yann LeCun, editors, 2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings, 2014. 6
- [28] M. Tan, B. Chen, R. Pang, V. Vasudevan, M. Sandler, A. Howard, and Q. V. Le. Mnasnet: Platform-aware neural architecture search for mobile. In 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 2815–2823, Los Alamitos, CA, USA, jun 2019. IEEE Computer Society. 10
- [29] Florian Tramèr, Alexey Kurakin, Nicolas Papernot, Ian J. Goodfellow, Dan Boneh, and Patrick D. McDaniel. Ensemble adversarial training: Attacks and defenses. In 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings. OpenReview.net, 2018. 1, 2
- [30] Cihang Xie, Zhishuai Zhang, Yuyin Zhou, Song Bai, Jianyu Wang, Zhou Ren, and Alan Yuille. Improving transferability of adversarial examples with input diversity, 2019. 1, 2
- [31] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. In Edwin R. Hancock Richard C. Wilson and William A. P. Smith, editors, *Proceedings of the British Machine Vision Conference (BMVC)*, pages 87.1–87.12. BMVA Press, September 2016. 10
- [32] Chen Zhu, Renkun Ni, Zheng Xu, Kezhi Kong, W. Ronny Huang, and Tom Goldstein. Gradinit: Learning to initialize neural networks for stable and efficient training. *CoRR*, abs/2102.08098, 2021. 1, 11
- [33] Roland S. Zimmermann. Comment on "adv-bnn: Improved adversarial defense through robust bayesian neural network", 2019. 1
- [34] Barret Zoph and Quoc V. Le. Neural architecture search with reinforcement learning. *CoRR*, abs/1611.01578, 2016. 10