

# NAT: Learning to Attack Neurons for Enhanced Adversarial Transferability

Krishna Kanth Nakka  
VITA Lab, EPFL, Switzerland  
krishkanth.92@gmail.com

Alexandre Alahi  
VITA Lab, EPFL, Switzerland  
alexandre.alahi@epfl.ch

## 1. Implementation

We provide below additional details for the reproducibility of our experiments.

**Libraries.** We conducted our experiments on NVIDIA GeForce RTX 3090 using PyTorch 2.1.1 [28], CUDA 11.8, Timm 1.0.9 [44], and Torchvision 0.16.1 [25].

**Generator Modifications.** Due to the non-deterministic behavior<sup>1</sup> of the `ReflectionPad2d` operation, which was used inside the generator architecture in prior works [27, 33, 50], we removed it in the `block1` layer and set padding to 1 in the subsequent `conv2d` operation within the same layer. We also removed `ReflectionPad2d` in the `blockf` layer and set the padding to 5 in the `conv2d` operation. Furthermore, the `residualblock` implementation contains two instances of `ReflectionPad2d`, both of which were removed, with padding set to 1 in the subsequent `conv2d` operations. These changes allowed us to eliminate `ReflectionPad2d` from the generator while compensating for changes in feature size with additional padding in the subsequent convolution operations.

**Top- $k$  Neurons After Initial Lightweight Training.** In Section 3.1, we outlined that all 512 generators, corresponding to the 512 neurons in layer 18 of VGG16, undergo an initial lightweight training phase using 3.125% of the training set. After this phase, we selected the top- $k = 40$  generators based on their performance, either on the source VGG16 [35] or by using a held-out model such as DenseNet121 [13] or ResNet [11] in our experiments. We provide the top-ranked neurons from this step across different held-out models in Table 2. Additionally, for full transparency, we present the performance of all 512 neurons after the lightweight training phase on two heldout models in Figures 1, and 2 for DenseNet121 [13] and VGG16 [35] heldout models.

**Target Models.** We used 41 target models, publicly avail-

able in Torchvision [25] and Timm libraries [44]. We provide in Table 1 the exact version details used for each model and also report the clean accuracy on the evaluation set of 5K images [33], sampled from the ImageNet validation set [31]. We used 9 models for cross-domain provided by the authors of BIA [50]

## 2. Additional Quantitative Results

**Performance of each neuron-specific generator.** In Tables 3 and 4, we provide the performance of all fully trained top-40 generators selected using DenseNet121 [13] as the held-out model. We observed that our neuron-specific generators outperform the baselines in 37 out of 40 cases in the single-query setting. Moreover, the best-performing generator  $G_{391}$ , attacking neuron 391, outperforms the highest-ranked generator  $G_{250}$  obtained at the end of lightweight training by more than 4%. We believe that with better neuron ranking schemes, the fooling rate in the single-query setting can be further enhanced.

**Choice of Heldout model for top- $k$  neuron selection.** In Figure 3, we show the performance of top-40 generators trained with different models for neuron selection versus the number of queries. We observe that our top- $k$  neuron position is robust to the choice of heldout model selection. And in all cases, top-1 generator consisted outperformed baselines in the single-query setting.

## 3. Additional Qualitative Results

We present the visualizations of adversarial images generated by different neuron-specific generators in Figures 4, 5, 6, and 7. Next to each adversarial image, we show the synthesized images that is optimized to have large activation magnitudes for the attacked neuron using activation maximization algorithm [34]. A clear visual correlation can be observed between the synthesized images and the adversarial patterns in the generated images, indicating that the perturbations are effectively disrupting neurons associated with specific concepts.

<sup>1</sup><https://github.com/pytorch/pytorch/issues/98925>

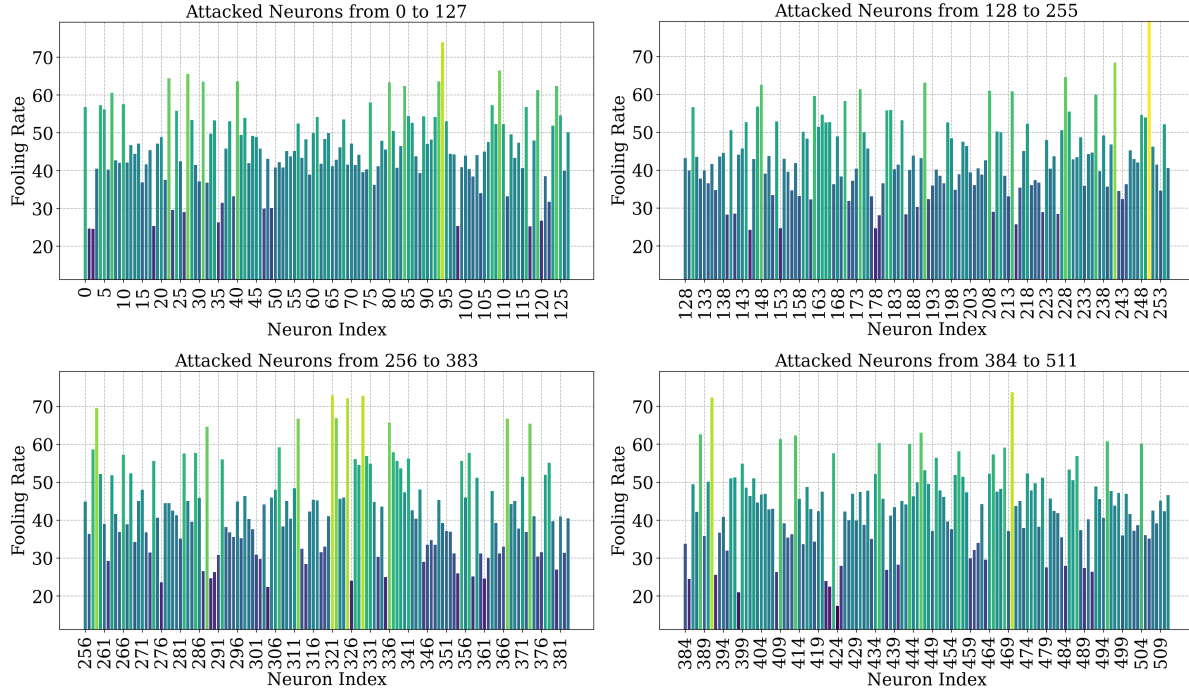


Figure 1. Evaluation of 512 Neuron specific generators on DenseNet121 [13] obtained after end of lightweight training on VGG16 [35]. We observe that some filters are more transferable than others and we choose top- $k$  generators that has highest transferability.

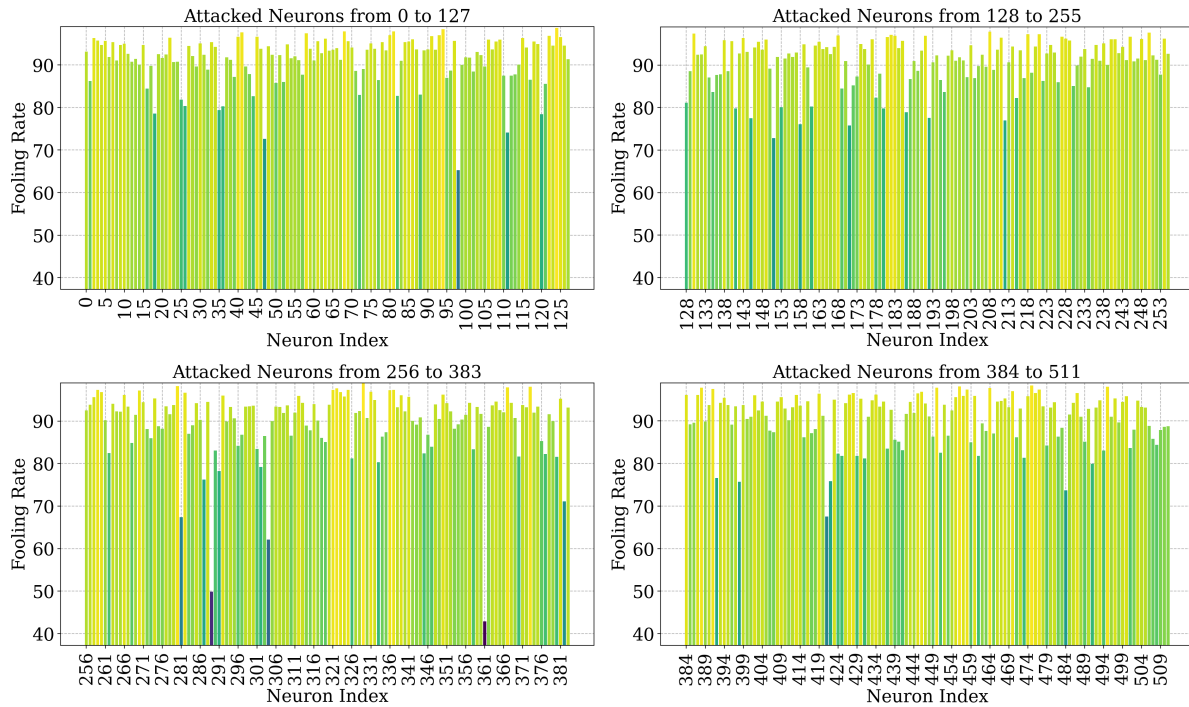


Figure 2. Evaluation of 512 Neuron specific generators on VGG16 [35] obtained after end of lightweight training on the same VGG16 [35]. We observe that all generators have high fooling rates suggesting that overfitting to the source VGG16 [35] model.

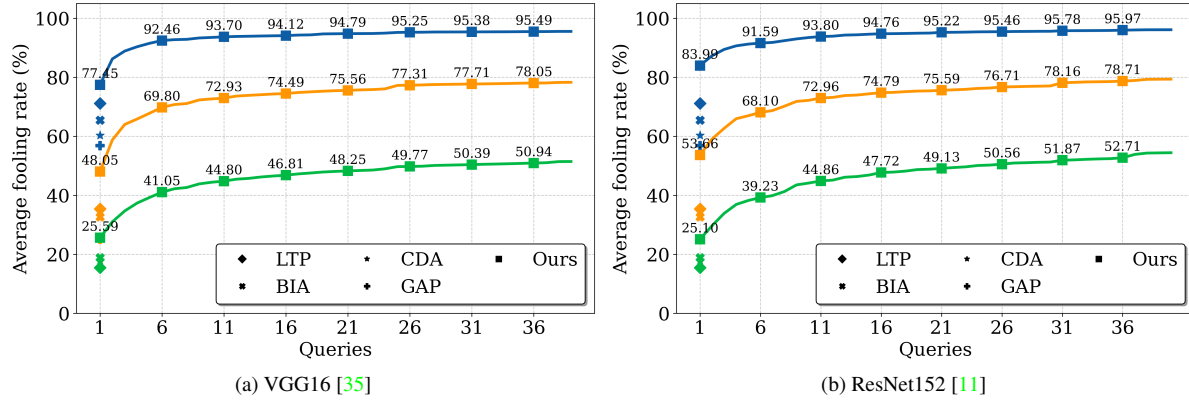


Figure 3. **Impact of heldout model for top- $k$  neuron selection after lightweight training** We use the same source model VGG16 [35] as the heldout model for top- $k$  neuron selection after lightweight training in the left and similarly using ResNet152 [11] on the right. We observe that for both the cases, the fooling rate with  $k$ -queries differ marginally and perform along the same lines as reported in the main paper with DenseNet121 [13] model.

Model	Version	Accuracy	Source
VGG [35]	vgg16	71.07%	TorchVision [25]
Resnet [11]	resnet152	78.05%	TorchVision [25]
Densenet [13]	densenet121	74.72%	TorchVision [25]
Squeezenet1 [14]	squeezenet1_1	58.58%	TorchVision [25]
Shufflenet [23]	shufflenetv2	69.52%	TorchVision [25]
Mobilenet [12]	mobilenet_v3	74.84%	TorchVision [25]
Mnasnet [36]	mnasnet1_0	73.84%	TorchVision [25]
Wide Resnet [49]	wide_resnet50_2	78.74%	TorchVision [25]
Convnext [22]	convnext_base	84.05%	TorchVision [25]
Efficientnet [16]	efficientnet_v2_m	80.00%	TorchVision [25]
Regnet [30]	regnet_x_32gf	83.62%	TorchVision [25]
Alexnet [17]	alexnet	56.88%	TorchVision [25]
ViT [1]	vit_b_16	81.06%	TorchVision [25]
ViT [1]	vit_l_16	80.01%	TorchVision [25]
<hr/>			
Swin [21]	swin_base_patch4_window7_224	85.62%	timm [44]
BeiT [2]	beit_base_patch16_224	85.68%	timm [44]
DeiT [38]	deit3_base_patch16_224	84.01%	timm [44]
Mixer [37]	mixer_b16_224	77.02%	timm [44]
ConvMixer [40]	convmixer_768_32	80.78%	timm [44]
Efficientvit [20]	efficientvit_m3	73.54%	timm [44]
Xception [6]	xception	75.54%	timm [44]
Crossvit [4]	crossvit_base_240	82.56%	timm [44]
Csp [42]	csresnet50	78.66%	timm [44]
Davit [7]	davit	84.86%	timm [44]
Edgenext [24]	edgenext	83.07%	timm [44]
Gcvit [10]	gcvit_base	85.32%	timm [44]
Ghostnet [9]	ghostnetv2_100	75.92%	timm [44]
Visformer [5]	visformer_small	82.22%	timm [44]
Focalnet [46]	focalnet_base_lrf	84.54%	timm [44]
Hiera [32]	hiera_base_224	85.02%	timm [44]
Hrnet [43]	hrnet_w32	79.03%	timm [44]
Maxvit [41]	maxvit_base_tf_224	85.52%	timm [44]
Convformer [48]	convformer_s36	84.56%	timm [44]
Mobilevit [26]	mobilevitv2_150	79.08%	timm [44]
Pnasnet [19]	pnasnet5large	78.52%	timm [44]
Rdnet [15]	rdnet_base	85.14%	timm [44]
Levit [8]	levit_conv_256	82.18%	timm [44]
Mvit [18]	mvitv2_base	85.08%	timm [44]
Nf [3]	nf_resnet50	79.86%	timm [44]
Coat [45]	coat_lite_small	82.76%	timm [44]
Cait [39]	cait_s24_224	84.03%	timm [44]
Dla [47]	dla102	78.54%	timm [44]

Table 1. Details of all 41 target models used in our evaluation, along with their clean accuracy on 5K evaluation ImageNet dataset provided by LTP [33].

Heldout model	Top-40 Neuron locations
DenseNet121	250, 94, 470, 321, 329, 391, 325, 259, 241, 322, 312, 367, 109, 336, 27, 373, 288, 228, 22, 93, 40, 31, 80, 191, 446, 388, 148, 124, 84, 413, 409, 174, 119, 208, 495, 214, 7, 435, 504, 443
ResNet152	250, 94, 22, 259, 93, 321, 367, 325, 391, 27, 470, 329, 228, 119, 504, 322, 288, 214, 307, 292, 130, 40, 109, 413, 312, 435, 208, 336, 5, 495, 373, 236, 446, 84, 80, 124, 75, 266, 241, 218
VGG16	329, 124, 94, 475, 280, 456, 373, 495, 367, 68, 81, 450, 208, 388, 464, 322, 41, 250, 391, 130, 58, 458, 321, 477, 259, 337, 325, 336, 221, 218, 270, 182, 470, 93, 168, 80, 331, 191, 446, 183

Table 2. **Attacked Neuron Positions.** We provide the list of the top-40 neurons selected based on their transferability performance on the held-out model. Note that all 512 generators were initially trained on just 3.12% of the dataset, with each generator specifically targeting an individual neuron in layer 18 of the source model VGG16.

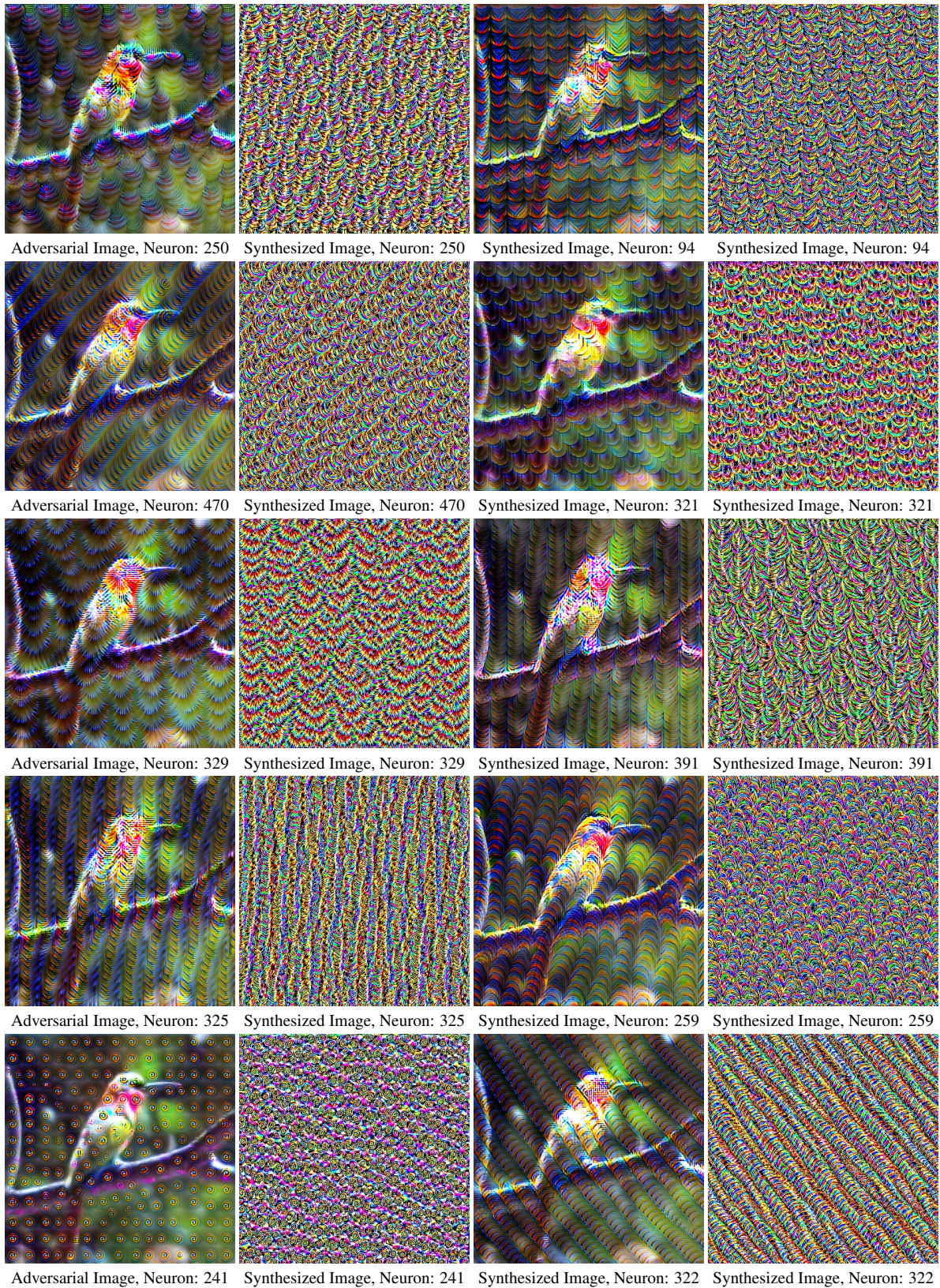


Figure 4. Generated unbounded adversarial images along with synthesized neuron visualizations for the top 10 attacked neurons. The positions of the neurons are listed below each image. Note that, the top-k neurons were selected based on the transferability to the DenseNet121 [13] as the held-out model after initial lightweight training.

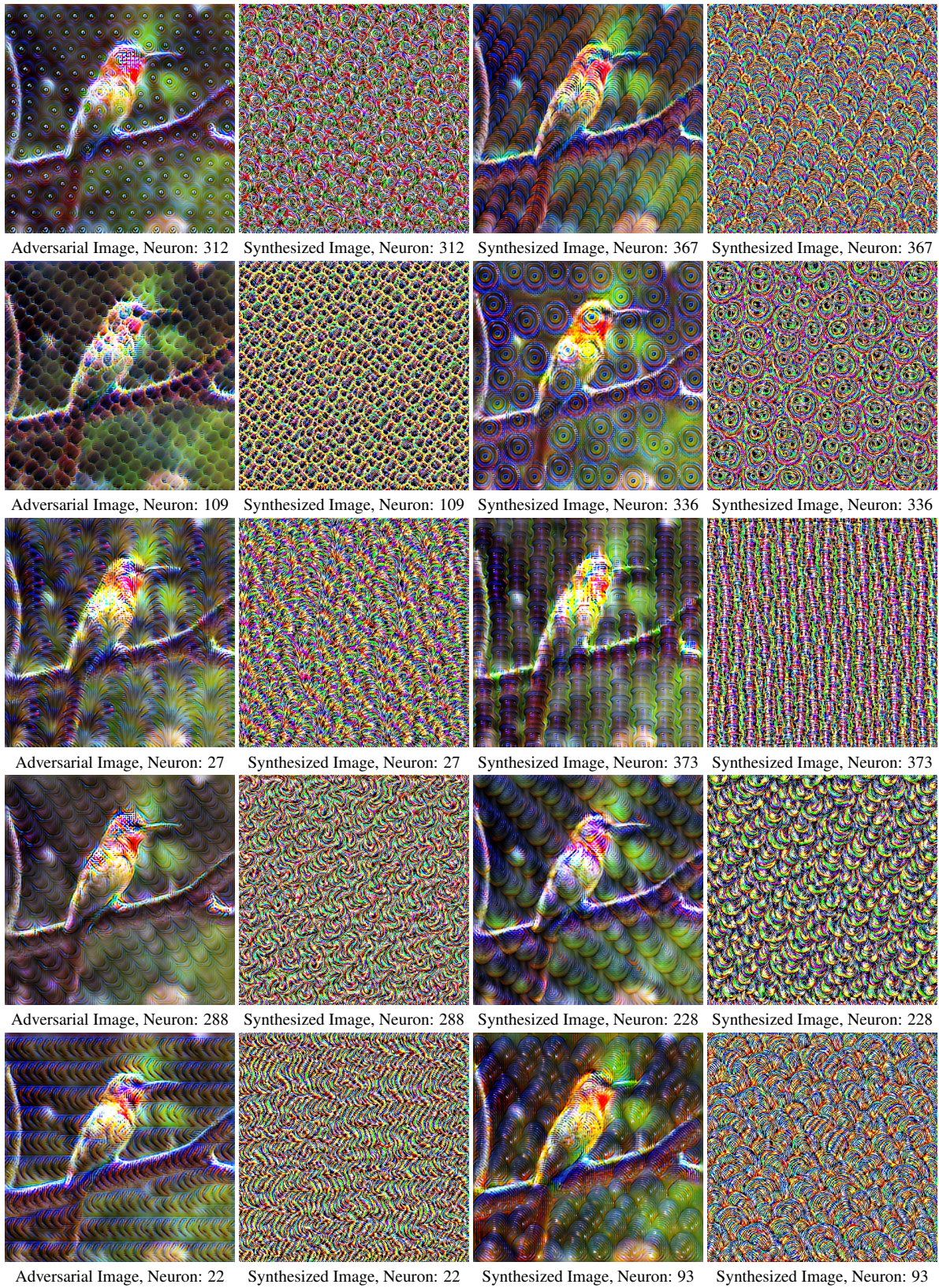


Figure 5. Generated unbounded adversarial images along with synthesized neuron visualizations for the top 10 to 20 attacked neurons. The positions of the neurons are listed below each image. Note that, the top-k neurons were selected based on the transferability to the DenseNet121 [13] as the held-out model after initial lightweight training

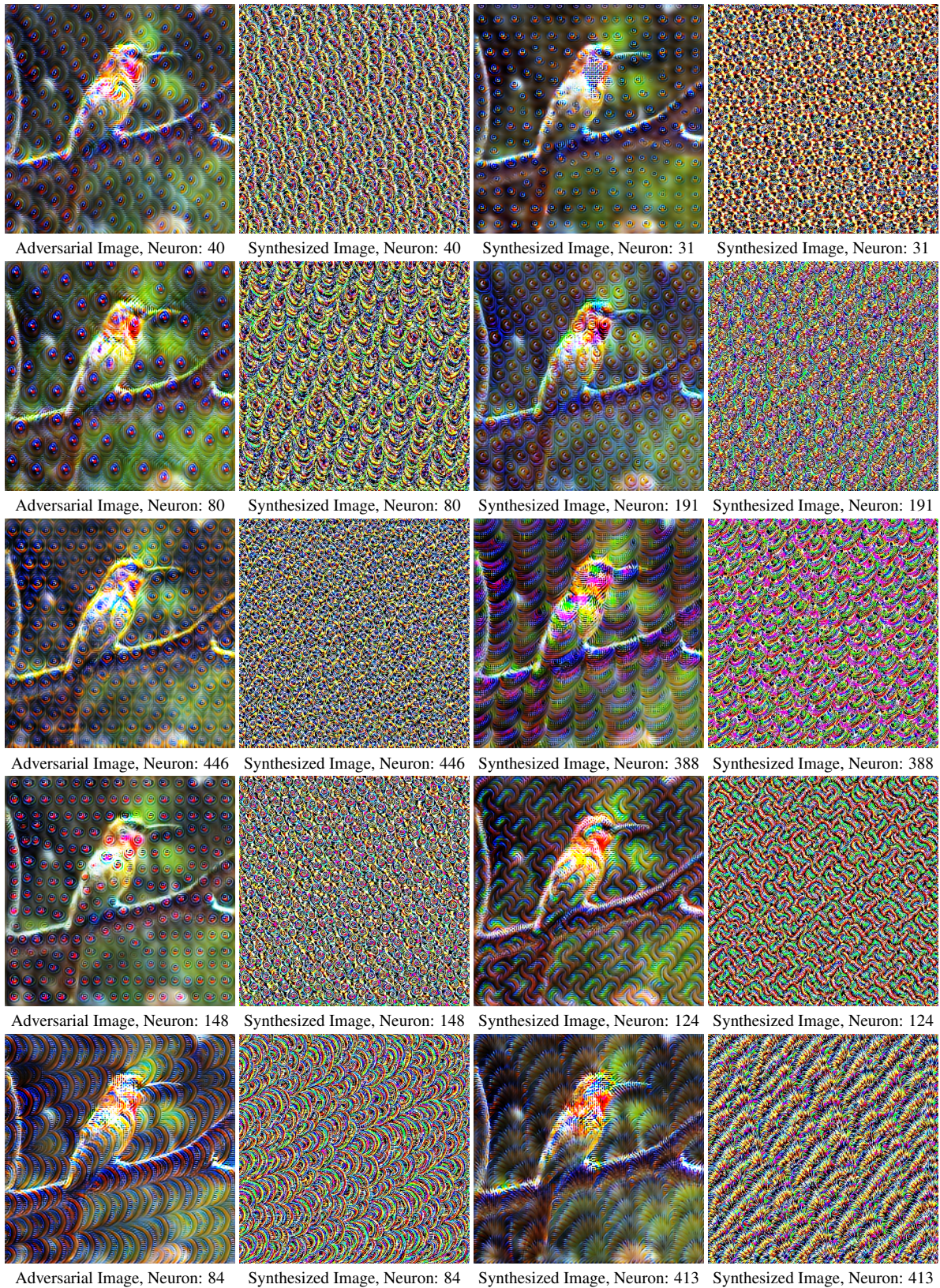


Figure 6. Generated unbounded adversarial images along with synthesized neuron visualizations for the top 20 to 30 attacked neurons. The positions of the neurons are listed below each image. Note that, the top-k neurons were selected based on the transferability to the DenseNet121 [13] as the held-out model after initial lightweight training

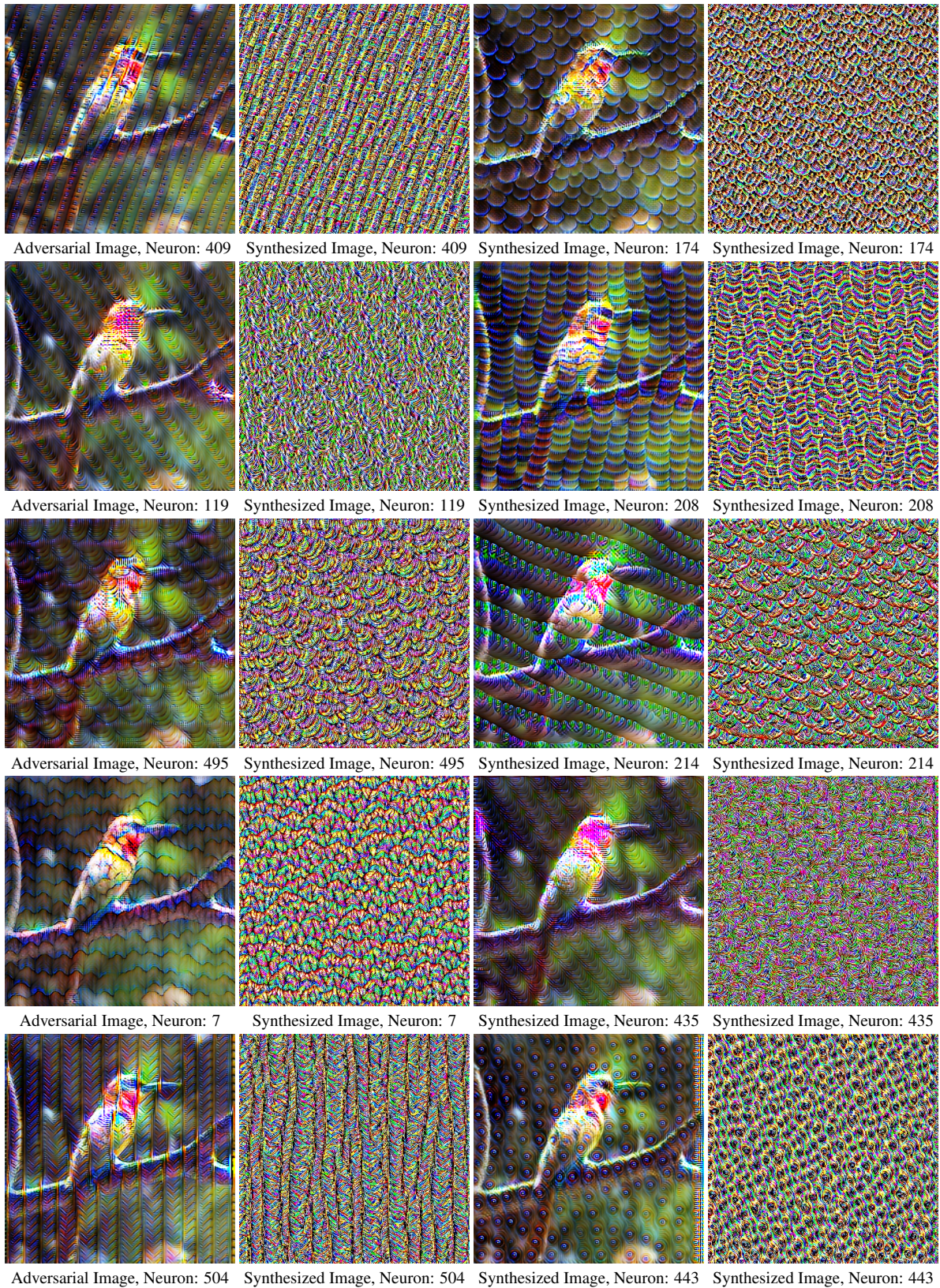


Figure 7. Generated unbounded adversarial images along with synthesized neuron visualizations for the top 30 to 40 neurons. . The positions of the neurons are listed below each image. Note that, the top-k neurons were selected based on the transferability to the DenseNet121 [13] as the held-out model after initial lightweight training



Model	250	94	470	321	329	391	325	259	241	322	312	367	109	336	27	373	288	228	22	93
ResNet152 [11]	89.10	<u>90.66</u>	84.34	88.26	80.44	<b>94.28</b>	79.14	90.06	66.70	80.22	73.14	85.76	77.84	75.66	79.68	65.40	77.28	79.16	81.54	78.00
DenseNet121 [13]	91.02	89.56	91.44	90.70	89.44	<b>94.68</b>	82.22	<u>91.94</u>	85.70	84.52	88.28	86.84	84.74	89.40	81.32	75.64	85.98	83.94	80.62	77.86
SqueezeNet1 [14]	89.46	83.04	88.24	85.22	89.56	86.94	82.46	81.16	87.04	72.98	<u>90.44</u>	78.38	79.50	<b>95.20</b>	88.36	74.78	83.28	87.34	87.06	84.90
ShuffleNet [23]	89.02	82.86	<b>91.64</b>	80.36	84.56	<u>91.10</u>	83.48	90.50	75.28	80.48	86.70	84.34	70.12	84.88	87.02	65.64	84.58	84.62	82.04	87.52
MNASNet [36]	93.32	86.32	<b>95.58</b>	87.28	87.72	93.12	85.76	88.54	85.28	90.60	93.16	88.58	84.48	92.30	84.86	79.54	89.20	<u>94.16</u>	88.90	89.22
PNasNet [19]	63.64	36.88	58.72	52.46	56.18	57.08	43.82	52.04	30.94	<u>66.60</u>	43.68	45.38	52.10	65.52	55.24	36.56	62.66	<b>69.38</b>	57.30	54.72
MobileNet [12]	84.26	66.86	<b>86.66</b>	71.22	80.98	<u>86.22</u>	75.94	76.26	57.00	71.88	71.58	72.46	54.34	80.42	82.08	55.66	78.64	79.08	79.72	79.48
WideResNet [49]	92.20	90.86	91.82	93.14	87.60	<b>93.32</b>	85.16	92.12	81.22	90.94	83.68	<u>93.26</u>	80.00	85.94	84.56	76.22	82.68	85.48	84.60	85.48
RegNet [30]	91.20	94.30	93.78	<b>96.78</b>	82.66	<u>96.10</u>	83.62	86.86	72.66	91.88	81.74	88.72	72.50	82.98	80.08	76.28	83.04	82.72	79.32	78.28
AlexNet [17]	49.82	47.28	46.68	43.86	48.40	41.98	42.94	45.28	38.04	48.92	41.12	42.20	49.22	<b>56.68</b>	50.06	49.24	46.12	<u>55.66</u>	49.94	46.90
NF ResNet [3]	76.86	75.34	77.98	77.94	65.46	<b>82.92</b>	56.56	<u>80.38</u>	45.62	76.88	59.26	76.34	63.22	68.26	63.34	54.32	69.86	63.04	70.76	64.86
DLA	97.86	96.74	96.72	97.10	95.16	<b>98.64</b>	94.70	<u>98.00</u>	92.26	96.46	94.98	97.96	89.98	91.98	95.16	94.28	93.28	93.88	93.66	94.64
GhostNetV2 [9]	82.68	75.54	<b>89.84</b>	84.04	76.50	<u>86.80</u>	75.06	83.50	58.90	83.82	75.10	83.08	65.74	79.78	79.52	65.00	80.42	82.24	80.54	77.70
CSPResNet [42]	<u>83.38</u>	78.54	83.10	78.60	67.74	<b>88.98</b>	64.08	82.80	64.46	75.36	72.68	78.08	61.24	77.24	70.92	52.32	73.28	71.20	74.10	72.32
HRNet [43]	97.10	95.20	94.14	<b>97.46</b>	90.32	<u>97.34</u>	88.60	89.08	90.66	93.78	94.30	92.36	89.48	95.48	85.06	92.10	90.82	92.52	91.22	87.44
Xception [6]	<u>73.82</u>	62.84	<u>66.22</u>	<b>74.48</b>	66.10	71.66	51.88	70.46	50.14	71.38	57.42	57.56	70.32	68.82	60.60	49.60	64.40	70.86	64.68	56.52
ViT_B_16 [1]	17.96	14.60	14.30	14.88	<u>18.30</u>	16.32	14.22	15.10	11.48	16.98	13.40	13.92	<b>18.38</b>	15.48	17.32	15.30	16.40	18.06	17.20	15.56
ViT_L_16 [1]	18.08	14.90	17.18	16.10	<b>19.14</b>	17.88	14.22	17.82	12.40	17.58	14.76	15.30	17.78	17.96	18.14	15.36	18.24	<u>18.74</u>	17.28	17.06
Swin [21]	28.78	24.52	<u>30.94</u>	28.86	27.26	<b>33.76</b>	21.34	28.20	17.62	21.58	22.16	25.46	18.32	21.66	25.46	19.32	25.52	25.02	24.54	24.54
BeiT [2]	<b>30.48</b>	16.22	25.60	22.04	23.46	27.94	18.36	24.70	12.00	24.02	16.46	17.50	21.40	19.64	<u>29.20</u>	18.30	21.44	28.32	25.38	21.94
DeiT [38]	20.54	17.42	19.02	19.64	20.32	<u>22.14</u>	17.48	19.84	11.90	20.38	14.86	17.08	18.44	17.84	<b>24.08</b>	16.92	17.50	18.68	18.60	15.76
Cait [39]	27.62	20.42	24.52	22.40	<b>31.72</b>	<u>29.08</u>	18.22	24.24	13.48	22.52	18.46	20.54	27.36	20.46	25.22	18.24	23.92	22.66	22.92	22.14
Davit [7]	31.62	31.70	37.70	36.50	<b>41.20</b>	<u>40.84</u>	24.46	34.48	23.46	31.14	27.16	37.10	19.66	19.94	24.90	21.22	32.60	27.98	27.72	29.50
ConvNext [22]	70.90	72.64	<u>75.02</u>	71.64	63.38	<b>83.56</b>	73.76	75.00	47.54	74.94	51.76	66.60	73.76	57.08	63.24	69.12	49.94	60.04	55.50	54.38
Mixer [37]	32.38	28.32	35.92	31.82	<b>42.82</b>	35.18	28.66	30.82	24.72	33.50	29.48	29.34	29.74	30.98	35.88	25.40	<u>36.88</u>	31.44	32.92	30.48
ConvMixer [40]	<u>76.12</u>	52.04	75.24	54.44	72.02	<b>84.60</b>	66.78	67.82	41.76	71.50	55.02	73.14	61.04	55.74	67.46	51.90	58.08	70.04	60.02	66.34
CrossViT [4]	19.18	16.50	17.92	17.44	17.70	<b>21.40</b>	15.36	17.36	12.28	16.16	16.00	14.50	18.26	17.02	17.78	15.94	<u>19.82</u>	17.24	17.68	15.62
Edgenext [24]	73.26	64.30	69.48	74.38	64.24	<b>80.80</b>	58.90	70.72	44.98	<u>74.54</u>	49.44	68.72	58.90	52.40	65.96	54.38	60.10	62.02	55.52	59.90
GCViT [10]	44.54	38.94	44.48	40.98	35.30	<b>49.38</b>	30.38	39.60	24.32	34.20	31.36	<u>45.92</u>	22.00	24.22	27.88	22.92	37.62	33.92	30.86	36.12
Visformer [5]	71.10	62.96	68.78	59.32	<u>79.26</u>	<b>80.60</b>	55.74	67.10	39.36	66.38	47.78	69.04	40.48	40.36	58.78	41.48	53.32	57.72	59.92	56.92
FocalNet [46]	45.66	36.48	<u>47.98</u>	44.82	37.76	<b>54.32</b>	35.60	38.32	31.02	40.14	35.12	39.82	27.88	30.84	37.34	25.24	42.38	38.36	33.40	36.38
Hiera [32]	43.16	33.72	<u>43.56</u>	36.40	38.92	<b>44.78</b>	31.18	41.14	19.82	34.50	24.22	36.46	21.30	24.92	37.04	24.58	31.56	38.72	32.80	34.06
MaxViT [41]	<u>22.94</u>	17.98	21.00	20.96	17.40	<b>25.60</b>	17.54	19.22	15.38	17.42	19.30	19.86	16.46	15.56	16.34	15.30	20.04	18.32	17.18	18.70
Conformer [48]	47.72	38.44	45.98	42.24	<b>53.34</b>	<u>51.72</u>	33.54	37.30	27.68	43.56	32.64	42.04	29.26	29.14	39.54	28.68	38.18	46.40	42.96	42.60
MobileViT [26]	<u>92.06</u>	91.96	88.10	91.12	68.42	<b>94.42</b>	81.82	90.68	72.34	87.40	79.12	88.36	61.04	78.04	76.06	85.48	80.62	84.84	80.42	79.52
RDNet [15]	<u>50.72</u>	42.44	44.28	47.88	36.58	<b>58.78</b>	33.00	41.88	25.68	43.76	31.04	43.74	21.82	24.98	28.28	28.96	36.62	37.52	34.72	37.34
LeViT [8]	67.52	58.20	66.52	62.00	63.98	<b>78.00</b>	61.14	<u>75.12</u>	37.40	62.40	46.94	66.24	43.46	46.66	62.48	39.78	55.86	58.26	60.10	57.04
MViT [18]	35.38	30.04	37.50	34.04	36.40	<b>43.42</b>	24.24	36.56	19.64	30.90	27.14	<u>38.20</u>	20.28	23.00	25.56	21.78	32.66	28.44	29.96	31.32
Coat [45]	52.14	<u>55.50</u>	51.10	52.22	36.20	<b>70.70</b>	42.32	45.70	32.08	39.02	34.52	46.00	30.94	30.28	34.56	30.14	45.92	38.98	34.86	39.74
EfficientViT [20]	60.04	44.78	66.24	49.00	62.44	<u>68.62</u>	51.62	65.74	41.48	52.02	54.46	59.30	41.32	61.32	<b>70.62</b>	43.32	58.10	57.48	62.30	57.30
EfficientNet [16]	60.64	48.00	<b>64.40</b>	49.48	55.08	<u>63.26</u>	43.48	55.00	36.96	51.92	44.70	47.54	35.40	54.68	47.72	36.64	50.08	54.62	54.70	48.40
Average	60.62	54.29	60.24	57.06	56.38	64.35	50.95	58.01	43.38	56.22	49.87	55.93	47.30	51.73	53.77	45.08	53.88	55.36	53.56	52.79

Table 3. **Fooling rate of top-20 generators on the 41 target models.** The attack neuron position of each generator is listed along with the columns. Note that the ranking here is obtained by using DenseNet121 [13] as the held-out model after lightweight training. The best performing generator for each target model is highlighted in bold, and the second-best is underlined for each target model. For comparison, baselines LTP [33], BIA [50], CDA [27], GAP [29] achieves average fooling rate of 44.6%, 42.0%, 40.7% and 34.7% over 41 architectures. All the top-20 generators obtained with our method outperform the baselines in all cases.



## References

- [1] Dosovitskiy Alexey. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. [4](#), [9](#), [10](#)
- [2] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254*, 2021. [4](#), [9](#), [10](#)
- [3] Andrew Brock, Soham De, and Samuel L Smith. Characterizing signal propagation to close the performance gap in unnormalized resnets. *arXiv preprint arXiv:2101.08692*, 2021. [4](#), [9](#), [10](#)
- [4] MDN Chandrasiri and Priyanga Dilini Talagala. Cross-vit: Cross-attention vision transformer for image duplicate detection. In *2023 8th International Conference on Information Technology Research (ICITR)*, pages 1–6. IEEE, 2023. [4](#), [9](#), [10](#)
- [5] Zhengsu Chen, Lingxi Xie, Jianwei Niu, Xuefeng Liu, Longhui Wei, and Qi Tian. Visformer: The vision-friendly transformer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 589–598, 2021. [4](#), [9](#), [10](#)
- [6] François Chollet. Xception: Deep learning with depthwise separable convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1251–1258, 2017. [4](#), [9](#), [10](#)
- [7] Mingyu Ding, Bin Xiao, Noel Codella, Ping Luo, Jingdong Wang, and Lu Yuan. Davit: Dual attention vision transformers. In *European conference on computer vision*, pages 74–92. Springer, 2022. [4](#), [9](#), [10](#)
- [8] Benjamin Graham, Alaeldin El-Nouby, Hugo Touvron, Pierre Stock, Armand Joulin, Hervé Jégou, and Matthijs Douze. Levit: a vision transformer in convnet’s clothing for faster inference. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 12259–12269, 2021. [4](#), [9](#), [10](#)
- [9] Kai Han, Yunhe Wang, Qi Tian, Jianyuan Guo, Chunjing Xu, and Chang Xu. Ghostnet: More features from cheap operations. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1580–1589, 2020. [4](#), [9](#), [10](#)
- [10] Ali Hatamizadeh, Hongxu Yin, Greg Heinrich, Jan Kautz, and Pavlo Molchanov. Global context vision transformers. In *International Conference on Machine Learning*, pages 12633–12646. PMLR, 2023. [4](#), [9](#), [10](#)
- [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. [1](#), [3](#), [4](#), [9](#), [10](#)
- [12] Andrew G Howard. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017. [4](#), [9](#), [10](#)
- [13] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [8](#), [9](#), [10](#)
- [14] Forrest N Iandola. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and 0.5 mb model size. *arXiv preprint arXiv:1602.07360*, 2016. [4](#), [9](#), [10](#)
- [15] Donghyun Kim, Byeongho Heo, and Dongyoon Han. Densenets reloaded: Paradigm shift beyond resnets and vits. *arXiv preprint arXiv:2403.19588*, 2024. [4](#), [9](#), [10](#)
- [16] Brett Koonce and Brett Koonce. Efficientnet. *Convolutional neural networks with swift for Tensorflow: image recognition and dataset categorization*, pages 109–123, 2021. [4](#), [9](#), [10](#)
- [17] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012. [4](#), [9](#), [10](#)
- [18] Yanghao Li, Chao-Yuan Wu, Haoqi Fan, Kartikeya Mangalam, Bo Xiong, Jitendra Malik, and Christoph Feichtenhofer. Mvitv2: Improved multiscale vision transformers for classification and detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4804–4814, 2022. [4](#), [9](#), [10](#)
- [19] 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)*, pages 19–34, 2018. [4](#), [9](#), [10](#)
- [20] Xinyu Liu, Houwen Peng, Ningxin Zheng, Yuqing Yang, Han Hu, and Yixuan Yuan. Efficientvit: Memory efficient vision transformer with cascaded group attention. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14420–14430, 2023. [4](#), [9](#), [10](#)
- [21] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021. [4](#), [9](#), [10](#)
- [22] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11976–11986, 2022. [4](#), [9](#), [10](#)
- [23] Ningning Ma, Xiangyu Zhang, Hai-Tao Zheng, and Jian Sun. Shufflenet v2: Practical guidelines for efficient cnn architecture design. In *Proceedings of the European conference on computer vision (ECCV)*, pages 116–131, 2018. [4](#), [9](#), [10](#)
- [24] Muhammad Maaz, Abdelrahman Shaker, Hisham Cholakkal, Salman Khan, Syed Waqas Zamir, Rao Muhammad Anwer, and Fahad Shahbaz Khan. Edgenext: efficiently amalgamated cnn-transformer architecture for mobile vision applications. In *European conference on computer vision*, pages 3–20. Springer, 2022. [4](#), [9](#), [10](#)
- [25] TorchVision maintainers and contributors. Torchvision: Pytorch’s computer vision library. <https://github.com/pytorch/vision>, 2016. [1](#), [4](#)
- [26] Sachin Mehta and Mohammad Rastegari. Mobilevit: lightweight, general-purpose, and mobile-friendly vision transformer. *arXiv preprint arXiv:2110.02178*, 2021. [4](#), [9](#), [10](#)
- [27] Muhammad Muzammal Naseer, Salman H Khan, Muhammad Haris Khan, Fahad Shahbaz Khan, and Fatih Porikli.

- Cross-domain transferability of adversarial perturbations. *Advances in Neural Information Processing Systems*, 32, 2019. 1, 9, 10
- [28] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019. 1
- [29] Omid Poursaeed, Isay Katsman, Bicheng Gao, and Serge Belongie. Generative adversarial perturbations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4422–4431, 2018. 9, 10
- [30] Ilija Radosavovic, Raj Prateek Kosaraju, Ross Girshick, Kaiming He, and Piotr Dollár. Designing network design spaces. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10428–10436, 2020. 4, 9, 10
- [31] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115:211–252, 2015. 1
- [32] Chaitanya Ryali, Yuan-Ting Hu, Daniel Bolya, Chen Wei, Haoqi Fan, Po-Yao Huang, Vaibhav Aggarwal, Arkabandhu Chowdhury, Omid Poursaeed, Judy Hoffman, et al. Hiera: A hierarchical vision transformer without the bells-and-whistles. In *International Conference on Machine Learning*, pages 29441–29454. PMLR, 2023. 4, 9, 10
- [33] Mathieu Salzmann et al. Learning transferable adversarial perturbations. *Advances in Neural Information Processing Systems*, 34:13950–13962, 2021. 1, 4, 9, 10
- [34] Karen Simonyan. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*, 2013. 1
- [35] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 1, 2, 3, 4
- [36] Mingxing Tan, Bo Chen, Ruoming Pang, Vijay Vasudevan, Mark Sandler, Andrew Howard, and Quoc V Le. Mnasnet: Platform-aware neural architecture search for mobile. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2820–2828, 2019. 4, 9, 10
- [37] Ilya O Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Andreas Steiner, Daniel Keysers, Jakob Uszkoreit, et al. Mlp-mixer: An all-mlp architecture for vision. *Advances in neural information processing systems*, 34:24261–24272, 2021. 4, 9, 10
- [38] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International conference on machine learning*, pages 10347–10357. PMLR, 2021. 4, 9, 10
- [39] Hugo Touvron, Matthieu Cord, Alexandre Sablayrolles, Gabriel Synnaeve, and Hervé Jégou. Going deeper with image transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 32–42, 2021. 4, 9, 10
- [40] Asher Trockman and J Zico Kolter. Patches are all you need? *arXiv preprint arXiv:2201.09792*, 2022. 4, 9, 10
- [41] Zhengzhong Tu, Hossein Talebi, Han Zhang, Feng Yang, Peyman Milanfar, Alan Bovik, and Yinxiao Li. Maxvit: Multi-axis vision transformer. In *European conference on computer vision*, pages 459–479. Springer, 2022. 4, 9, 10
- [42] Chien-Yao Wang, Hong-Yuan Mark Liao, Yueh-Hua Wu, Ping-Yang Chen, Jun-Wei Hsieh, and I-Hau Yeh. Cspnet: A new backbone that can enhance learning capability of cnn. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 390–391, 2020. 4, 9, 10
- [43] Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu, Yadong Mu, Mingkui Tan, Xinggang Wang, et al. Deep high-resolution representation learning for visual recognition. *IEEE transactions on pattern analysis and machine intelligence*, 43(10):3349–3364, 2020. 4, 9, 10
- [44] Ross Wightman. Pytorch image models. <https://github.com/rwightman/pytorch-image-models>, 2019. 1, 4
- [45] Weijian Xu, Yifan Xu, Tyler Chang, and Zhuowen Tu. Co-scale conv-attentional image transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9981–9990, 2021. 4, 9, 10
- [46] Jianwei Yang, Chunyuan Li, Xiyang Dai, and Jianfeng Gao. Focal modulation networks. *Advances in Neural Information Processing Systems*, 35:4203–4217, 2022. 4, 9, 10
- [47] Fisher Yu, Dequan Wang, Evan Shelhamer, and Trevor Darrell. Deep layer aggregation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2403–2412, 2018. 4
- [48] Weihao Yu, Chenyang Si, Pan Zhou, Mi Luo, Yichen Zhou, Jiashi Feng, Shuicheng Yan, and Xinchao Wang. Metaformer baselines for vision. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023. 4, 9, 10
- [49] Sergey Zagoruyko. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016. 4, 9, 10
- [50] Qilong Zhang, Xiaodan Li, Yuefeng Chen, Jingkuan Song, Lianli Gao, Yuan He, and Hui Xue. Beyond imagenet attack: Towards crafting adversarial examples for black-box domains. *arXiv preprint arXiv:2201.11528*, 2022. 1, 9, 10