

---

# Natural Color Fool: Towards Boosting Black-box Unrestricted Attacks (Supplementary Material)

---

Anonymous Author(s)

Affiliation

Address

email

---

## CONTENTS

|                       |   |          |
|-----------------------|---|----------|
| <b>A</b>              | <b>Color Distribution Library</b>                     | <b>2</b> |
| <b>B</b>              | <b>Algorithm</b>                                      | <b>2</b> |
| <b>C</b>              | <b>Initialization Reset</b>                           | <b>3</b> |
| <b>D</b>              | <b>Transferability on ViTs</b>                        | <b>3</b> |
| <sup>1</sup> <b>E</b> | <b>Neighborhood Search Impact on Human Perception</b> | <b>4</b> |
| <b>F</b>              | <b>Visualization of Attention</b>                     | <b>4</b> |
| <b>G</b>              | <b>Visualization of Adversarial Examples</b>          | <b>5</b> |
| <b>H</b>              | <b>The Influence of Segmentation Models</b>           | <b>5</b> |
| <b>I</b>              | <b>The Effect of Ensemble Attack</b>                  | <b>6</b> |
| <b>J</b>              | <b>Attack by Selecting Random Colors</b>              | <b>7</b> |

---

## 2 A Color Distribution Library

3 In this section, we explain how to build our color distribution library. First, all semantic classes are  
 4 extracted based on the images and annotations of the ADE20K training set [1] (ADE20K contains 150  
 5 semantic classes). The semantic classes are clustered into a set of 20 different dominant distribution  
 6 sets according to their color styles using hierarchical clustering, which is the “distribution of color  
 7 distribution” (DoD) [2]. Then, since the color distributions of semantic classes in each cluster are  
 8 similar, we select one object in each cluster to represent that style. Finally, we extract the color  
 9 distributions to form our color distribution library. The pipeline is shown in Figure 1.

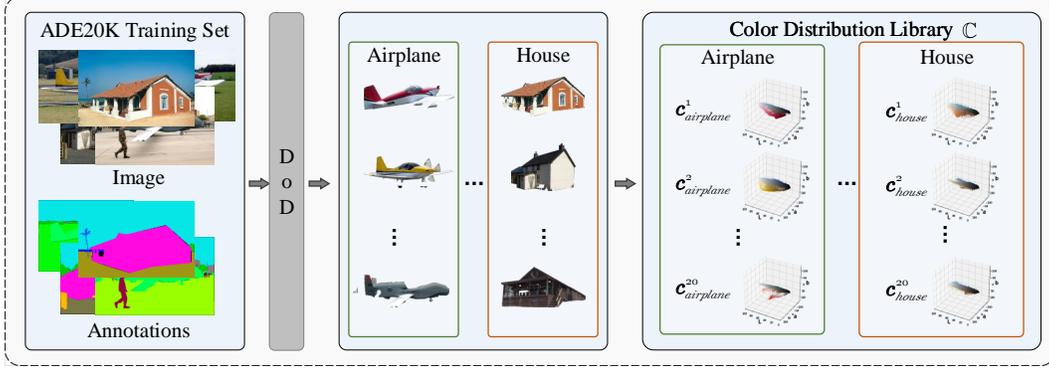


Figure 1: **The pipeline of color distribution library.** We select 20 objects with different color styles for each semantic class and extract their color distributions to build a color distribution library.

## 10 B Algorithm

11 In this section, we provide the algorithm of NCF (see Algorithm 1).

---

### 13 Algorithm 1 Natural Color Fool

---

#### 14 Input:

15  $x$ : clean image,  $y$ : label,  $\tilde{y}$ : the semantic class,  $\mathcal{S}$ : semantic segmentation model,  
 16  $\mathcal{F}_\phi(\cdot)$ : the substitute model,  $\mathbb{C}$ : color distribution library,  $K$ : the restart number,  
 17  $N$ : the iteration of NS,  $u$ : the momentum,  $\alpha$ : the step size,  
 18  $\Sigma$ : covariance matrix,  $\mu$ : mean

#### 19 Output: $x'$ : adversarial image

20 1: Obtain all semantic classes  $\tilde{Y}$  of image  $x$  by  $\mathcal{S}$ .  
 21 2: Calculate each semantic class’s pixel ratio  $w$  in the  $x$ .  
 22 3:  $\bar{x} \leftarrow \text{rgb2lab}(x)$   
 23 4: **for**  $i \leftarrow 1$  to  $K$  **do** ▷ Initialization Reset.  
 24 5:   **for**  $j \leftarrow 1$  to  $\eta$  **do**  
 25 6:     Randomly sample the color distribution  $c_{\tilde{y}}$  of each semantic class  $\tilde{y}$  from the library  $\mathbb{C}_{\tilde{y}}$   
 26 7:      $H_j \leftarrow \sum_{\tilde{y} \in \tilde{Y}} w_{\tilde{y}} * c_{\tilde{y}}$   
 27 8:      $\bar{x}_{H_j} \leftarrow$  Convert  $H_j$  to an image without spatial information.  
 28 9:      $\Sigma_{\bar{x}_{H_j}}, \mu_{\bar{x}_{H_j}} \leftarrow$  Calculate covariance matrix and channel means for  $\bar{x}_{H_j}$ .  
 29 10:      $T \leftarrow \text{MKSolution}(\Sigma_{\bar{x}}, \Sigma_{\bar{x}_{H_j}})$  (see Eq.(6))  
 30 11:      $\bar{x}'_{H_j} \leftarrow T(\bar{x} - \mu_{\bar{x}}) + \mu_{\bar{x}_{H_j}}$   
 31 12:   **end for**  
 32 13:  $\mathbb{H} \leftarrow \{H_j\}_{j=1}^\eta$   
 33 14:  $H^* \leftarrow \arg \max_{H \in \mathbb{H}} \mathcal{L}_{adv}(\mathcal{F}_\phi(x'_H), y)$   
 34 15:  $\Sigma_{\bar{x}_{H^*}}, \mu_{\bar{x}_{H^*}} \leftarrow$  Calculate covariance matrix and channel means for  $\bar{x}_{H^*}$ .  
 35 16:  $T^* \leftarrow \text{MKSoulution}(\Sigma_{\bar{x}}, \Sigma_{\bar{x}_{H^*}})$  (see Eq.(6))

```

36 17:  $\mu^* \leftarrow \mu_{\bar{x}_{H^*}}$ 
37 18: Initialize  $T'_0 \leftarrow T^*$ 
38 19: for  $n \leftarrow 1$  to  $N$  do ▷ Neighborhood Search.
39 20:    $\bar{x}'_i \leftarrow T'_{n-1}(\bar{x} - \mu_{\bar{x}}) + \mu^*$ 
40 21:    $g \leftarrow \nabla_{T'} \mathcal{L}_{adv}(\mathcal{F}_\phi(\text{lab2rgb}(\bar{x}'_i)), y)$ 
41 22:    $g_n \leftarrow u \cdot g_{n-1} + \frac{g}{\|g\|_1}$ 
42 23:    $T'_n \leftarrow T'_{n-1} + \alpha \cdot \text{sign}(g_n)$ 
43 24: end for
44 25:  $x'_i \leftarrow \text{lab2rgb}(\bar{x}'_i)$ 
45 26: end for
46 27:  $i' \leftarrow \arg \max_{i \in \{1, 2, \dots, K\}} \mathcal{L}_{adv}(\mathcal{F}_\phi(x'_i), y)$ 
47 28:  $x' \leftarrow x'_{i'}$ 
48 29: return  $x'$ 

```

---

## 49 C Initialization Reset

50 As mentioned in Section 3.3, initialization reset (IR) avoids the optimization of convert matrix to  
51 fall into local optimum. Therefore, we study the appropriate value for the number of resets in this  
52 section. Our adversarial examples are crafted via Inc-v3 by NCF with different resets, i.e., from 1 to  
53 20. We report the attack success rate on the victim model in Figure 2. As demonstrated in the figure,  
54 in the beginning, there is an approximate positive relationship between the attack success rate and  
55 the number of resets. Once  $K = 10$  or larger, the attack success rate maintains stable. We choose  
56  $K = 10$  to reduce the computational overhead, where the attack success rate is 84.0%, 57.7%,  
57 57.6%, 56.8%, 40.1%, 47.6% for Inc-v3 [3], Res-18 [4], VGG-19 [5], Mobile-v2 [6], Dense-121 [7]  
58 and Res-50 [4], respectively.

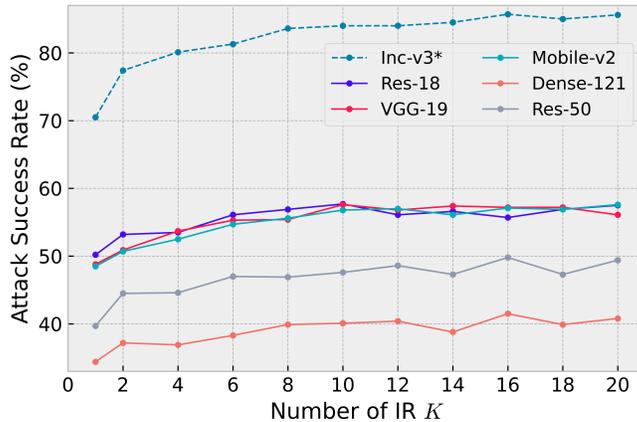


Figure 2: Attack success rates (%) w.r.t. the number of IR  $K$ . The substitute model is Inc-v3.

## 59 D Transferability on ViTs

60 Recent work [8] points out that the vision transformer models are more concerned with low-frequency  
61 information, thus we further evaluate attack success rate of black-box when Transformer models are  
62 used as substitute models. Specifically, we generate adversarial examples using ViT-S/16 (ViT-S) [9]  
63 and XCiT-N12/16 (XCiT-N12) [10] as substitute models and evaluate the adversarial examples on  
64 multiple Transformer models (including ViT-S, XCiT-N12, DeiT-S [9], ViT-B [11], Swin-T [12], PiT-  
65 Ti [13]) and multiple CNNs (including Res-18 [4], VGG-19 [5], Mobile-v2 [6]). Table 1 summarizes  
66 the results on different black-box models. We can observe that our NCF usually achieves the highest  
67 transferability on black-box. In particular, when transferring from ViTs to CNNs, NCF achieves an  
68 average success rate of 59.8%, but other attacks only achieve 48.3% (SAE), 29.9% (ReColorAdv),  
69 36.8% (cAdv), 35.1% (ColorFool), 27.3% (ACE).

70 We observe that the transferability of adversarial examples between ViTs is weaker than different  
71 CNNs (see Section 4.2.1). This is consistent with the observation in [14]. It may be since that  
72 ViTs have multiple self-attention blocks that can generate class tokens independently. However,  
73 our loss only utilizes the last logit (which is equivalent to utilizing only the last class token). The  
74 discriminative information in the previous tokens is not directly utilized, which results the poor  
75 transferability of the adversarial examples between different ViTs.

Table 1: **Transferability on ViTs.** We report attack success rates (%) of each method and the leftmost model column denotes the substitute model (“\*” means white-box attack results).

| Model    | Attacks    | Transformers  |               |             |             |             |             | CNNs        |             |             |
|----------|------------|---------------|---------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
|          |            | ViT-S         | XCiT-N12      | DeiT-S      | ViT-B       | Swin-T      | PiT-Ti      | Res-18      | VGG-19      | Mobile-v2   |
|          | Clean      | 13.3          | 13.7          | 5.8         | 10.7        | 5.0         | 11.6        | 16.1        | 11.4        | 12.8        |
| ViT-S    | SAE        | 98.2*         | 38.7          | 24.2        | 37.5        | 19.6        | 34.0        | 49.9        | 47.6        | 46.5        |
|          | ReColorAdv | 96.2*         | 27.4          | 20.4        | 30.6        | 13.8        | 24.5        | 27.6        | 23.7        | 26.3        |
|          | cAdv       | <b>100.0*</b> | 36.5          | <b>32.8</b> | 42.4        | 19.3        | 35.4        | 35.5        | 27.1        | 33.4        |
|          | ColorFool  | 99.2*         | 23.6          | 11.0        | 24.3        | 7.6         | 21.5        | 35.0        | 26.7        | 29.3        |
|          | ACE        | 99.7*         | 21.1          | 10.5        | 20.0        | 7.1         | 19.6        | 29.8        | 23.8        | 28.4        |
|          | NCF (Ours) | 93.9*         | <b>42.4</b>   | 25.9        | <b>62.7</b> | <b>20.4</b> | <b>39.9</b> | <b>57.5</b> | <b>55.4</b> | <b>54.4</b> |
| XCiT-N12 | SAE        | 49.3          | 86.5*         | 25.9        | 36.9        | 18.9        | 39.9        | 51.2        | 47.5        | 47.0        |
|          | ReColorAdv | 21.2          | 95.9*         | 20.2        | 16.4        | 14.7        | 37.4        | 36.3        | 29.5        | 35.9        |
|          | cAdv       | 39.4          | <b>100.0*</b> | <b>38.5</b> | 33.0        | 26.5        | 53.9        | 46.9        | 36.1        | 41.7        |
|          | ColorFool  | 47.6          | 85.9          | 12.6        | 34.4        | 9.7         | 26.5        | 43.9        | 36.5        | 39.4        |
|          | ACE        | 19.5          | 98.7          | 9.8         | 15.0        | 6.3         | 21.0        | 29.5        | 24.9        | 27.3        |
|          | NCF (Ours) | <b>54.1</b>   | 89.1*         | 36.1        | <b>38.5</b> | <b>27.3</b> | <b>55.8</b> | <b>65.1</b> | <b>63.7</b> | <b>62.6</b> |

## 76 E Neighborhood Search Impact on Human Perception

77 In this section, we study the impact of neighborhood search (NS) on human perception. Figure 3  
78 shows the adversarial examples with and without NS, and we observe that our NS strategy has almost  
79 no effect on human perception.



Figure 3: **Visualizations of NCF-IR and NCF-IR-NS.** We use Inc-v3 as an alternative model to generate adversarial examples. We observe from the results that there is no perceptual difference after using NS.

## 80 F Visualization of Attention

81 In this section, we show the Gradient-weighted Class Activation Mapping (Grad-CAM) [15] of  
82 different adversarial examples. Figure 4 illustrates that all of SAE, ReColorAdv, cAdv, ColorFool,  
83 and ACE have difficulty shifting the attention of the model, while our NCF can dramatically it, i.e.,  
84 no longer on the target object.

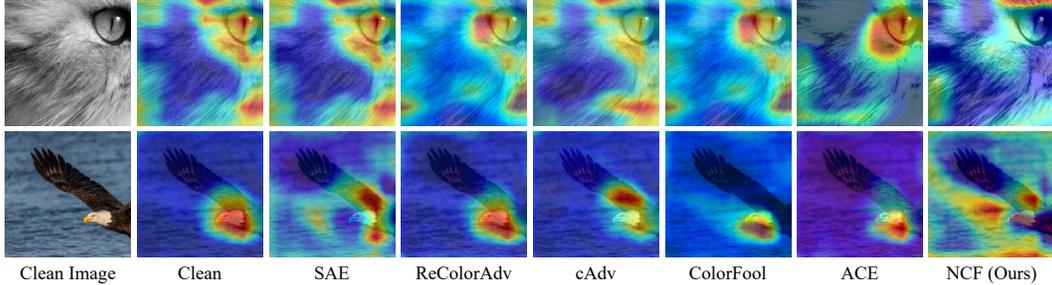


Figure 4: **Visualization of attention.** We use Inc-v3 as the alternative model to generate adversarial examples, and then visualize the attention on VGG-19.

## 85 G Visualization of Adversarial Examples

86 In this section, we show some adversarial examples. In Figure 5, we observe that SAE is prone to  
 87 global color distortion due to the unbounded range of variation. ReColorAdv limits the perturbation  
 88 to a small range, thus its adversarial examples are less perceptually different from the clean images.  
 89 However, some areas of color gradients (e.g., the “clouds” in the second row of Figure 5) may produce  
 90 some perceptual anomalies. cAdv relies on the colorization model, and some color anomalies are  
 91 likely to occur locally in the image. ColorFool does not set the perturbation range in non-sensitive  
 92 areas, resulting in non-sensitive areas prone to color distortion. ACE is similar to ReColorAdv in that  
 93 sometimes details are lost. (e.g., the “plants” in the third row of Figure 5). NCF transforms the color  
 94 distribution so that the images’ details are preserved.

## 95 H The Influence of Segmentation Models

96 In this section, we compare the performance of our NCF under different semantic segmentation  
 97 models pre-trained on ADE20K (including Swin-T [12], OCRNet [16] and Deeplabv3+ [17]). As  
 98 indicated in Table 2, segmentation models have impact on the attack success rates of resulting  
 99 adversarial examples. Among these models, Swin-T is usually the best choice for our NCF. Therefore,  
 100 in our paper, we choose it to segment inputs. Note that even if the segmentation model affects  
 101 our method, the lowest black-box attack success rate of NCF is still much higher than the existing  
 102 methods.

Table 2: **The influence of different segmentation models on attack success rates.** (“\*” denotes the white-box attack)

| Segm       | Res-18*      | VGG-19      | Mobile-v2   | Inc-v3      | Dense-121   | Res-50      | ViT-S       |
|------------|--------------|-------------|-------------|-------------|-------------|-------------|-------------|
| Swin-T     | <b>92.9*</b> | <b>72.1</b> | <b>72.7</b> | <b>48.3</b> | <b>55.3</b> | <b>66.7</b> | 53.0        |
| OCRNet     | 89.9*        | 69.1        | 67.1        | 44.2        | 50.6        | 61.1        | <b>56.5</b> |
| Deeplabv3+ | 91.0*        | 68.0        | 68.6        | 45.3        | 49.2        | 62.0        | 54.0        |

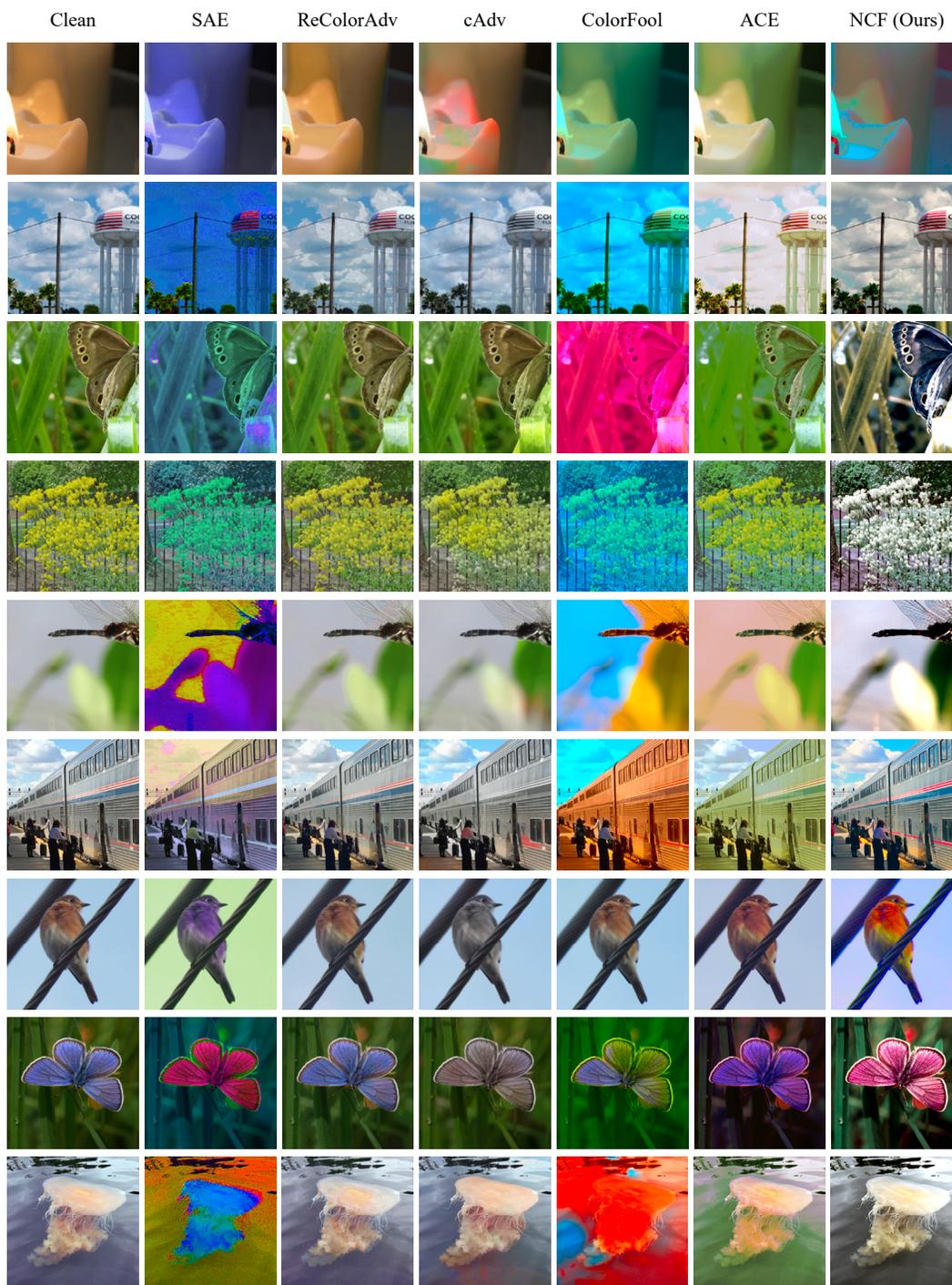


Figure 5: Visualization of adversarial examples. All adversarial examples are generated on Inc-v3.

103 **I The Effect of Ensemble Attack**

104 As for ensemble model attack (fusing the logits of multiple models like [18]), here we report the  
 105 result of NCF. As indicated in Table 3, the attack success rate of NCF can be further improved when  
 106 crafting via an ensemble of models.

Table 3: Comparison of ensemble attack and single model attack. We report attack success rates (%) of NCF and the leftmost model column denotes the substitute model, where Ensemble means an ensemble of Res-18, VGG-19 and Mobile-v2.

| Model     | Dense-121   | Res-50      | ViT-S       | XCiT-N12    | DeiT-S      |
|-----------|-------------|-------------|-------------|-------------|-------------|
| Clean     | 7.9         | 7.5         | 13.3        | 13.7        | 5.8         |
| Res-18    | 55.3        | 66.7        | 53.0        | 55.3        | 32.8        |
| VGG-19    | 53.6        | 64.3        | 56.5        | 53.5        | 30.7        |
| Mobile-V2 | 54.4        | 66.2        | 55.4        | 56.4        | 32.6        |
| Ensemble  | <b>63.5</b> | <b>71.6</b> | <b>59.7</b> | <b>61.7</b> | <b>37.0</b> |

## 107 J Attack by Selecting Random Colors

108 In this section, we discuss the difference between NCF-IR-NS and random color attack. Formally,  
 109 NCF-IR-NS does not mean selecting random colors to attack. Specifically, it first generates a set of  
 110 adversarial examples with different color distributions and then selects the best example from them  
 111 based on the loss of the white-box model to attack. Therefore, NCF-IR-NS is close to a white-box  
 112 attack.

113 To support our claim, we evaluate the performance of random color attack (NCF-IR-NS-), i.e.,  
 114 randomly select colors for each semantic class and use the resulting adversarial examples to attack.  
 115 As demonstrated in Table 4, the performance of NCF-IR-NS- is much lower than NCF-IR-NS. For  
 116 example, NCF-IR-NS- only achieves a 27.3% (degraded from 51.6%) success rate on Inc-v3. Thus,  
 117 directly using random colors to generate adversarial examples is ineffective.

Table 4: The attack success rate of using white-box information and not using it. NCF-IR-NS using Inc-v3 as the substitute model. (“\*” denotes the white-box attack)

| Methods     | Inc-v3*      | Res-18      | VGG-19      | Mobile-v2   | Dense-121   | Res-50      | ViT-S       | XCiT-N12    | DeiT-S      |
|-------------|--------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| Clean       | 19.2*        | 16.1        | 11.4        | 12.8        | 7.9         | 7.5         | 13.3        | 13.7        | 5.8         |
| NCF-IR-NS   | <b>51.6*</b> | <b>43.8</b> | <b>42.2</b> | <b>42.4</b> | <b>28.0</b> | <b>33.0</b> | <b>38.3</b> | <b>32.0</b> | <b>14.8</b> |
| NCF-IR-NS-* | 27.3         | 34.8        | 30.9        | 31.1        | 20.5        | 24.2        | 32.7        | 25.0        | 11.6        |

## References

- 118
- 119 [1] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba.  
120 Scene parsing through ADE20K dataset. In *CVPR*, 2017. 2
- 121 [2] Mahmoud Afifi, Brian L. Price, Scott Cohen, and Michael S. Brown. Image recoloring based  
122 on object color distributions. In *Eurographics*, 2019. 2
- 123 [3] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna.  
124 Rethinking the inception architecture for computer vision. In *CVPR*, 2016. 3
- 125 [4] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
126 recognition. In *CVPR*, 2016. 3
- 127 [5] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale  
128 image recognition. In *ICLR*, 2015. 3
- 129 [6] Mark Sandler, Andrew G. Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen.  
130 Mobilenetv2: Inverted residuals and linear bottlenecks. In *CVPR*, 2018. 3
- 131 [7] Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. Densely connected  
132 convolutional networks. In *CVPR*, 2017. 3
- 133 [8] Namuk Park and Songkuk Kim. How do vision transformers work? In *ICLR*, 2022. 3
- 134 [9] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and  
135 Hervé Jégou. Training data-efficient image transformers & distillation through attention. In  
136 *ICML*, 2021. 3
- 137 [10] Alaaeldin Ali, Hugo Touvron, Mathilde Caron, Piotr Bojanowski, Matthijs Douze, Armand  
138 Joulin, Ivan Laptev, Natalia Neverova, Gabriel Synnaeve, Jakob Verbeek, and Hervé Jégou.  
139 Xcit: Cross-covariance image transformers. In *NeurIPS*, 2021. 3
- 140 [11] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,  
141 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,  
142 Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image  
143 recognition at scale. In *ICLR*, 2021. 3
- 144 [12] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining  
145 Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, pages  
146 9992–10002. IEEE, 2021. 3, 5
- 147 [13] Byeongho Heo, Sangdoon Yun, Dongyoon Han, Sanghyuk Chun, Junsuk Choe, and Seong Joon  
148 Oh. Rethinking spatial dimensions of vision transformers. In *ICCV*, 2021. 3
- 149 [14] Muzammal Naseer, Kanchana Ranasinghe, Salman Khan, Fahad Khan, and Fatih Porikli. On  
150 improving adversarial transferability of vision transformers. In *ICLR*, 2022. 4
- 151 [15] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi  
152 Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based  
153 localization. In *ICCV*, 2017. 4
- 154 [16] Yuhui Yuan, Xilin Chen, and Jingdong Wang. Object-contextual representations for semantic  
155 segmentation. 2020. 5
- 156 [17] Hang Zhang, Chongruo Wu, Zhongyue Zhang, Yi Zhu, Zhi Zhang, Haibin Lin, Yue Sun,  
157 Tong He, Jonas Muller, R. Manmatha, Mu Li, and Alexander Smola. Resnest: Split-attention  
158 networks. *arXiv preprint arXiv:2004.08955*, 2020. 5
- 159 [18] Yinpeng Dong, Fangzhou Liao, Tianyu Pang, Hang Su, Jun Zhu, Xiaolin Hu, and Jianguo Li.  
160 Boosting adversarial attacks with momentum. In *CVPR*, 2018. 6