

## A Combining the Additive and Functional Threat Models

Here we provide a proof of Theorem [1](#).

**Threat model** Let  $\mathbf{x}$  be a grayscale image with  $n \geq 2$  pixels, i.e.  $\mathbf{x} \in [0, 1]^n = \mathcal{X}^n$ . Let  $t_{\text{add}}$  be an additive threat model where the  $\ell_\infty$  distance between input and adversarial example is bounded by  $\epsilon_1$ , i.e.  $\|(\delta_1, \dots, \delta_n)\|_\infty \leq \epsilon_1$ . Let  $t_{\text{func}}$  be a functional threat model where  $f(x) = cx$  for some  $c \in [1 - \epsilon_2, 1 + \epsilon_2]$  and let  $\epsilon_2 > \epsilon_1 > 0$ . The additive threat model allows individually changing each pixel’s value by up to  $\epsilon_1$ ; the functional threat model allows darkening or lightening the entire image by up to a proportion of  $\epsilon_2$ . Both of these are arguably imperceptible perturbations for small enough  $\epsilon_1$  and  $\epsilon_2$ . We also consider  $t_{\text{combined}} = t_{\text{add}} \circ t_{\text{func}}$ :

$$t_{\text{combined}}(\mathcal{S}) \triangleq \left\{ (cx_1 + \delta_1, \dots, cx_n + \delta_n) \left| \begin{array}{l} (x_1, \dots, x_n) \in \mathcal{S} \\ |\delta_i| \leq \epsilon_1 \\ c \in [1 - \epsilon_2, 1 + \epsilon_2] \end{array} \right. \right\} \quad (3)$$

This combined threat model allows darkening or lightening the image, followed by changing each pixel value individually by a small amount.

**Theorem 1** (restated). *Let  $\mathcal{S} \in \mathcal{P}(\mathcal{X}^n)$  be a set of inputs such that  $\mathcal{S}$  contains an image that is not too dark; that is,  $\exists \mathbf{x} \in \mathcal{S}$  for which  $\exists x_i$  s.t.  $x_i > \epsilon_1/\epsilon_2$ . Then*

$$t_{\text{combined}}(\mathcal{S}) \supsetneq t_{\text{add}}(\mathcal{S}) \cup t_{\text{func}}(\mathcal{S}) \quad \text{or equivalently} \quad \exists \tilde{\mathbf{x}} \text{ s.t. } \begin{array}{l} \tilde{\mathbf{x}} \in t_{\text{combined}}(\mathcal{S}) \\ \tilde{\mathbf{x}} \notin t_{\text{add}}(\mathcal{S}) \cup t_{\text{func}}(\mathcal{S}) \end{array}$$

*Proof.* The above two statements are equivalent, so we focus on the formulation on the right. We calculate  $\tilde{\mathbf{x}}$  and show that it satisfies the given criteria. Let  $\mathbf{x} \in \mathcal{S}$  such that  $\exists x_i$  s.t.  $x_i > \epsilon_1/\epsilon_2$ . Without loss of generality, assume that in particular  $x_2 > \epsilon_1/\epsilon_2$ . Then let

$$\tilde{\mathbf{x}} = ((1 - \epsilon_2)x_1 + \epsilon_1, (1 - \epsilon_2)x_2, \dots, (1 - \epsilon_2)x_n)$$

First, we show that  $\tilde{\mathbf{x}} \in t_{\text{combined}}(\mathcal{S})$ . Using the definition of  $t_{\text{combined}}$  in [\(3\)](#), we set  $c = 1 - \epsilon_2$ ,  $\delta_1 = \epsilon_1$ , and  $\delta_2 = \dots = \delta_n = 0$ , which generates  $\tilde{\mathbf{x}}$ . These values clearly satisfy the constraints in [\(3\)](#).

Second, we prove that  $\tilde{\mathbf{x}} \notin t_{\text{add}}(\mathcal{S})$  by contradiction. Say that  $\tilde{\mathbf{x}} \in t_{\text{add}}(\mathcal{S})$ . Then  $\exists \delta_1, \delta_2, \dots, \delta_n$  such that  $\tilde{x}_i = x_i + \delta_i$  and  $\|\delta_i\| \leq \epsilon_1$ . Consider  $\delta_2$ , which must satisfy  $\tilde{x}_2 = (1 - \epsilon_2)x_2 = x_2 + \delta_2$ , or alternatively  $\delta_2 = x_2 - (1 - \epsilon_2)x_2 = \epsilon_2 x_2$ . However,  $x_2 > \epsilon_1/\epsilon_2$  implies that  $\delta_2 > \epsilon_1$ , which is a contradiction since the constraints on  $t_{\text{add}}$  specify that  $|\delta_2| \leq \epsilon_1$ . Thus,  $\tilde{\mathbf{x}} \notin t_{\text{add}}(\mathcal{S})$ .

Third, we prove that  $\tilde{\mathbf{x}} \notin t_{\text{func}}(\mathcal{S})$ , again by contradiction. Say that  $\tilde{\mathbf{x}} \in t_{\text{func}}(\mathcal{S})$ . Then  $\exists c \in [1 - \epsilon_2, 1 + \epsilon_2]$  such that  $\tilde{x}_i = cx_i$  for all  $i$ . Considering  $i = 1, 2$ , we have the following system of equations:

$$\begin{aligned} \tilde{x}_1 &= cx_1 = (1 - \epsilon_2)x_1 + \epsilon_1 \\ \tilde{x}_2 &= cx_2 = (1 - \epsilon_2)x_2 \end{aligned}$$

From the second equation, we have  $c = 1 - \epsilon_2$ . However, using this in the first equation gives  $(1 - \epsilon_2)x_1 = (1 - \epsilon_2)x_1 + \epsilon_1$ , which implies  $0 = \epsilon_1$ . This is a contradiction since  $\epsilon_1 > 0$ , showing that  $\tilde{\mathbf{x}} \notin t_{\text{func}}(\mathcal{S})$ .  $\blacksquare$

## B Experimental Setup

We implement ReColorAdv using the `mister_ed` library [\[9\]](#) and PyTorch [\[16\]](#). Adversarial examples are generated by 100 iterations of PGD using the Adam optimizer [\[10\]](#) with learning rate 0.001. After all iterations have completed, we choose the result of the iteration with the lowest loss as the adversarial example.

When combining attacks, we apply multiple attacks sequentially to the input example and optimize over the parameters of all attacks simultaneously, similarly to Jordan et al. [9].

In all adversarial training experiments on CIFAR-10, we begin with a trained ResNet32 [8] and then train it further on batches which are half original training data and half adversarial examples. We adversarially train with a batch size of 500 for 50 epochs. We preprocess images after adversarial perturbation, but before classification, by standardizing them based on the mean and standard deviation of each channel for all images in the dataset. The CIFAR-10 dataset can be obtained from <https://www.cs.toronto.edu/~kriz/cifar.html>.

In CIELUV color space (see section 4.1), we define

$$(c_1, c_2, c_3) = \left( \frac{L}{100}, \frac{U + 100}{200}, \frac{V + 100}{200} \right) \quad (4)$$

so that  $(c_1, c_2, c_3) \in [0, 1]^3$ .

For the experiments described in section 5.3, we use LPIPS v0.1 with AlexNet.

## B.1 Regularization Parameters

The objective function and constraints described in section 4 include a number of constants that can be used to regularize the outputs of the ReColorAdv attack. Changing these constants alters the strength of the attack and the perceptual similarity of a generated adversarial example to the input.

First,  $\epsilon_1$ ,  $\epsilon_2$ , and  $\epsilon_3$  control the maximum amount by which a color in  $\mathbf{x}$  can be changed to produce  $\tilde{\mathbf{x}}$ . For RGB color space, we set  $\epsilon_1 = \epsilon_2 = \epsilon_3 = 0.1$ ; that is, each channel of a color can change by up to  $\sim 25/255$ . This is greater than the usual  $\epsilon = 8/255$  allowed for adversarial examples, but we find that the uniform perturbation used by the functional threat model allows each pixel to change by a greater amount while remaining almost indistinguishable. For the CIELUV color space, we let  $\epsilon_1 = \epsilon_2 = \epsilon_3 = 0.06$ . This corresponds to a maximum change of 6 in  $L$  and a maximum change of 3 in  $U$  and  $V$ , since we find that changes in luma are usually less noticeable than changes in chroma. The  $\epsilon_i$  values for RGB and CIELUV color spaces result in similar total amounts of perturbation, but the CIELUV color space allows the perturbation to be greater in areas where it is less noticeable.

Second, we can control the resolution of the grid  $\mathcal{G}$  over which the perturbation function  $f(\cdot)$  is parameterized. Let  $R_1 \times R_2 \times R_3$  be the resolution of  $\mathcal{G}$ . Lowering the resolution in a particular dimension acts as a regularizer because it allows less variation in how colors are transformed along that dimension. For RGB color space, we use  $R_1 = R_2 = R_3 = 25$ . However, for CIELUV color space, we use  $R_1 = 16$  and  $R_2 = R_3 = 32$ . With a high  $R_1$  value, we find that the attack sometimes recolors different values of a particular hue very differently. For instance, the attack might make the light parts of a white car green and the dark parts purple. Lowering  $R_1$  forces the attack to alter these colors more similarly.

Finally,  $\lambda$  controls the importance of the smoothness optimization term  $\mathcal{L}_{\text{smooth}}$ . We always set  $\lambda = 0.05$ .

## C Learning Rate Experiments

We consistently use Adam with a learning rate of 0.001 throughout the main paper to craft adversarial examples. However, we also experimented with a learning rate of 0.01. The results of these experiments are shown below, similar to table 1. All numbers reported are accuracy over the CIFAR-10 test set. Each column corresponds to an attack and each row corresponds to a model trained against a particular attack. C(-RGB) is ReColorAdv using CIELUV (RGB) color space, D is delta attack, and S is StAdv attack. TRADES is the method of Zhang et al. [26]. For classifiers marked (B&W), the images are converted to black-and-white before classification. The learning rate used in an attack is marked above that attack or to the right when the attack is used in adversarial training. There are a couple interesting conclusions that can be drawn from this experiment:

- The higher learning rate (0.01) is stronger against TRADES and undefended networks. A ReColorAdv + StAdv + delta (C+S+D) attack with learning rate 0.01 against a TRADES-trained classifier reduces its accuracy to just 6.0%, compared with 10.1% at learning rate 0.001.

- Adversarial training against an attack at the higher learning rate (0.01) increases robustness against that attack but lowers it against other attacks. For instance, consider the network defended against C+S+D with learning rate 0.01. This network achieves 15.0% accuracy against attacks of the same type, but the accuracy decreases to 7.1% against some other attacks. In contrast, adversarial training against attacks at the lower learning rate (0.001) leads to more robustness across different attacks.

Defense	LR	Attack (learning rate = 0.01)								
		None	C-RGB	C	D	S	C+S	C+D	S+D	C+S+D
<b>Undefended</b>		92.3	5.1	3.8	<b>0.0</b>	1.5	1.5	<b>0.0</b>	<b>0.0</b>	<b>0.0</b>
<b>C</b>	0.01	87.8	37.4	45.5	4.7	3.2	2.9	1.2	<b>0.2</b>	0.4
<b>D</b>	0.01	88.8	40.4	22.7	32.7	4.2	4.3	15.0	5.0	<b>4.0</b>
<b>S</b>	0.01	89.3	11.5	9.8	0.3	29.0	9.5	0.4	0.4	<b>0.3</b>
<b>C+S</b>	0.01	90.5	27.0	24.3	2.8	31.4	23.0	<b>2.1</b>	2.8	<b>2.1</b>
<b>C+D</b>	0.01	88.3	46.3	32.7	34.4	6.0	5.4	22.6	<b>4.7</b>	4.8
<b>S+D</b>	0.01	88.0	25.3	17.4	28.4	9.3	<b>8.1</b>	22.6	17.0	13.9
<b>C+S+D</b>	0.01	89.0	32.3	23.8	29.8	13.4	<b>11.4</b>	26.1	17.4	15.0
<b>C</b>	0.001	89.2	37.2	46.6	5.1	3.4	3.0	1.1	<b>0.3</b>	<b>0.3</b>
<b>D</b>	0.001	84.7	72.9	57.4	30.8	12.2	11.2	12.6	2.4	<b>1.8</b>
<b>S</b>	0.001	82.7	14.9	11.9	0.5	22.2	6.7	<b>0.1</b>	0.2	<b>0.1</b>
<b>C+S</b>	0.001	82.3	37.5	40.4	5.9	18.5	13.1	1.9	0.9	<b>0.8</b>
<b>C+D</b>	0.001	84.3	70.8	60.0	33.8	9.4	8.7	18.1	<b>1.8</b>	1.9
<b>S+D</b>	0.001	82.0	65.2	49.9	35.0	18.5	14.0	16.5	5.5	<b>4.5</b>
<b>C+S+D</b>	0.001	82.3	65.8	53.0	34.8	16.8	14.7	18.3	5.1	<b>5.0</b>
<b>TRADES</b>		84.2	79.7	69.2	53.5	21.0	17.8	33.8	6.6	<b>6.0</b>
<b>Undefended (B&amp;W)</b>		87.9	4.7	4.8	<b>0.0</b>	1.6	1.5	<b>0.0</b>	0.1	<b>0.0</b>
<b>C (B&amp;W)</b>	0.01	84.7	40.4	41.7	4.5	2.4	2.4	1.0	<b>0.2</b>	0.3
<b>C (B&amp;W)</b>	0.001	85.6	37.8	40.7	4.0	2.5	2.5	0.7	<b>0.3</b>	<b>0.3</b>

Defense	LR	Attack (learning rate = 0.001)								
		None	C-RGB	C	D	S	C+S	C+D	S+D	C+S+D
<b>Undefended</b>		92.3	8.3	5.3	<b>0.0</b>	2.2	1.8	<b>0.0</b>	<b>0.0</b>	<b>0.0</b>
<b>C</b>	0.01	87.8	46.2	48.4	5.9	4.5	4.4	1.6	<b>0.3</b>	0.7
<b>D</b>	0.01	88.8	43.7	25.4	26.4	4.1	<b>3.8</b>	15.7	8.3	7.9
<b>S</b>	0.01	89.3	18.7	13.9	0.4	13.8	8.4	0.8	<b>0.6</b>	0.9
<b>C+S</b>	0.01	90.5	39.1	32.3	4.3	22.7	17.5	<b>2.8</b>	3.5	3.3
<b>C+D</b>	0.01	88.3	49.1	35.6	29.2	<b>5.3</b>	5.4	20.3	8.7	8.3
<b>S+D</b>	0.01	88.0	25.8	17.7	10.7	4.5	<b>4.1</b>	9.3	6.1	5.9
<b>C+S+D</b>	0.01	89.0	33.9	24.9	15.7	7.5	<b>7.1</b>	13.0	8.4	8.5
<b>C</b>	0.001	89.2	47.4	50.3	5.9	4.6	4.6	1.7	<b>0.5</b>	0.9
<b>D</b>	0.001	84.7	77.3	61.9	32.8	18.6	17.2	17.3	4.3	<b>4.2</b>
<b>S</b>	0.001	82.7	20.3	15.7	0.8	29.9	10.7	<b>0.2</b>	<b>0.2</b>	<b>0.2</b>
<b>C+S</b>	0.001	82.3	47.2	44.6	7.5	26.2	20.2	3.5	2.2	<b>2.0</b>
<b>C+D</b>	0.001	84.3	74.7	63.5	35.2	14.0	13.4	22.2	4.5	<b>4.2</b>
<b>S+D</b>	0.001	82.0	69.3	53.9	36.5	26.4	21.1	21.7	9.6	<b>8.0</b>
<b>C+S+D</b>	0.001	82.3	70.1	56.4	35.5	25.5	21.4	23.4	10.0	<b>8.5</b>
<b>TRADES</b>		84.2	81.6	72.8	53.7	31.2	27.5	39.3	11.1	<b>10.1</b>
<b>Undefended (B&amp;W)</b>		87.9	7.3	6.1	<b>0.0</b>	1.6	1.6	<b>0.0</b>	<b>0.0</b>	<b>0.0</b>
<b>C (B&amp;W)</b>	0.01	84.7	49.3	44.9	5.4	4.1	3.8	1.6	0.5	0.7
<b>C (B&amp;W)</b>	0.001	85.6	46.7	43.8	5.0	3.5	3.8	1.3	<b>0.4</b>	0.7

## D Non-Additive Threat Models

Here, we discuss some other non-additive adversarial threat models that have been explored in the literature and how our work differs from them.

**Spatial Threat Models** Some recent work has focused on *spatial threat models*, which allow for slight perturbations of the locations of features in an input rather than perturbations of the features themselves. Xiao et al. [23] propose StAdv, which optimizes the parameters of a smooth flow field that moves each pixel of an input image by a small, bounded distance to generate an example that fools the classifier. Wong et al. [22] bound the Wasserstein distance between the original input and the adversarial example. Engstrom et al. [3] apply a small rotation and translation to an input image to generate a misclassification.

**Other Threat Models** A few papers have focused on threat models that are neither additive or spatial. Zeng et al. [25] perturb the properties of a 3D renderer to render an image of an object which is unrecognizable to a classifier or other machine learning algorithm. Hosseini and Poovendran [6] propose "Semantic Adversarial Examples," which allow modifications of the input image's hue and saturation. Hosseini et al. [7] also explore inverting images to cause misclassification. These latter two papers can be considered as special examples of functional threat models. In the first, each pixel's hue and saturation is shifted by the same amount; that is, each pixel is transformed by the function  $f(h, s, v) = (h + \delta_h, s + \delta_s, v)$ . In the second, each pixel is inverted, i.e. each pixel channel is transformed by the function  $f(x_i) = 1 - x_i$ . However, the authors do not propose a general framework for these types of attacks, as we do. Furthermore, the adversarial examples generated by these attacks are often not realistic and not imperceptible. For example, their crafted adversarial examples include green skies, purple fields of grass, and inverted street signs—unlike our proposed ReColorAdv attack, which results in imperceptible changes.

## E Additional Images

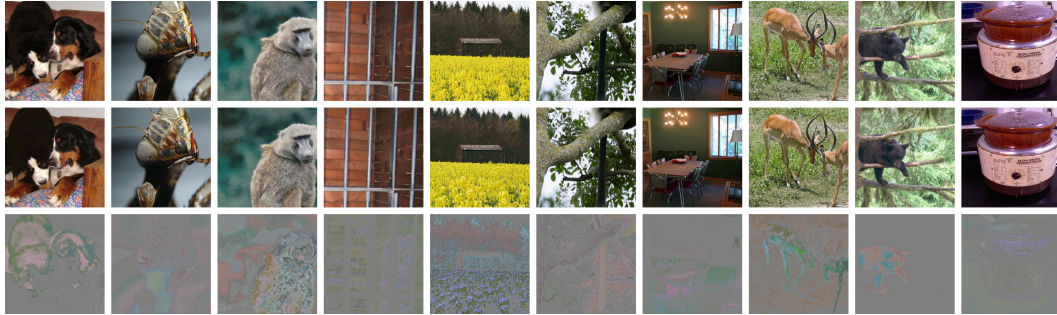


Figure 7: More adversarial examples like those in figure 2, generated by ReColorAdv against an Inception-v4 classifier on ImageNet. Top row: original images; middle row: adversarial examples; bottom row: magnified difference.

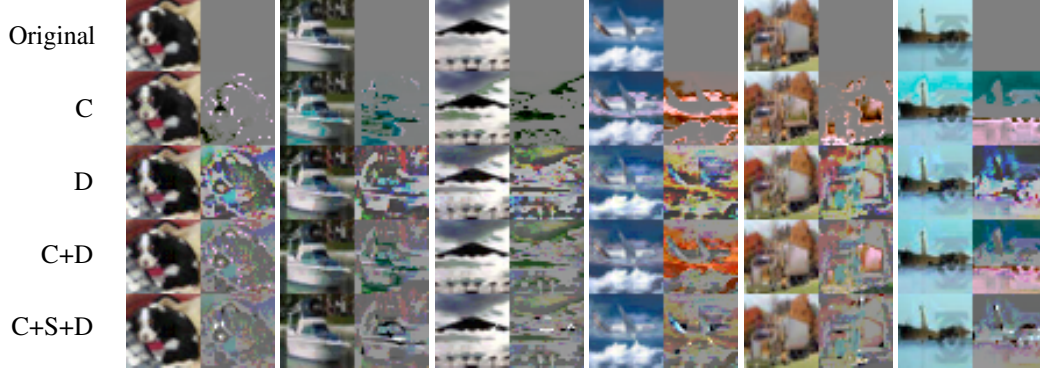


Figure 8: More adversarial examples like those in figure 5 generated with combinations of attacks against a CIFAR-10 WideResNet trained using TRADES. C is ReColorAdv, D is delta attack, and S is StAdv attack [23]. The difference from the original is shown to the right of each example. Combinations of attacks tend to produce less perceptible changes than the attacks do separately.

## F Lipschitz Regularization

In addition to the regularizations defined in section 3.1, we can also enforce that the perturbation function  $f(\cdot)$  in a functional threat model is Lipschitz for some suitably small  $\kappa$ :

$$\mathcal{F}_{\text{lips}} \triangleq \{f : \mathcal{X} \rightarrow \mathcal{X} \mid \forall x_1, x_2 \in \mathcal{X} \|f(x_1) - f(x_2)\| \leq \kappa \|x_1 - x_2\|\} \quad (5)$$

$\mathcal{F}_{\text{lips}}$  requires some smoothness in the perturbation function  $f(\cdot)$ , ensuring that similar features in the input are mapped to similar features in the adversarial example. However, one disadvantage of  $\mathcal{F}_{\text{lips}}$  is that it includes constant functions  $f(x) = c$ , i.e. functions which map every feature to a single value, removing salient features from the input. Thus, we ultimately use  $\mathcal{F}_{\text{smooth}}$  instead.