ImageNet-E: Benchmarking Neural Network Robustness via Attribute Editing

Xiaodan Li¹, Yuefeng Chen^{1*}, Yao Zhu², Shuhui Wang^{3*}, Rong Zhang¹, Hui Xue¹ ¹Alibaba Group ²Zhejiang University ³Inst. of Comput. Tech., CAS, China

{fiona.lxd, yuefeng.chenyf, stone.zhangr, hui.xueh}@alibaba-inc.com ee_zhuy@zju.edu.cn, wangshuhui@ict.ac.cn



ImageNet-C

Figure 1. Examples of the proposed ImageNet-E dataset. In contrast to adversarial examples or datasets like ImageNet-C [21] who add perturbation or corruptions to original images, we edit the object attributes with controls of backgrounds, sizes, positions and directions.

Abstract

Recent studies have shown that higher accuracy on ImageNet usually leads to better robustness against different corruptions. Therefore, in this paper, instead of following the traditional research paradigm that investigates new out-of-distribution corruptions or perturbations deep models may encounter, we conduct model debugging in indistribution data to explore which object attributes a model may be sensitive to. To achieve this goal, we create a toolkit for object editing with controls of backgrounds, sizes, positions, and directions, and create a rigorous benchmark named ImageNet-E(diting) for evaluating the image classifier robustness in terms of object attributes. With our ImageNet-E, we evaluate the performance of current deep learning models, including both convolutional neural networks and vision transformers. We find that most models are quite sensitive to attribute changes. A small change in the background can lead to an average of 9.23% drop

on top-1 accuracy. We also evaluate some robust models including both adversarially trained models and other robust trained models and find that some models show worse robustness against attribute changes than vanilla models. Based on these findings, we discover ways to enhance attribute robustness with preprocessing, architecture designs, and training strategies. We hope this work can provide some insights to the community and open up a new avenue for research in robust computer vision. The code and dataset are available at https://github.com/ alibaba/easyrobust.

1. Introduction

Deep learning has triggered the rise of artificial intelligence and has become the workhorse of machine intelligence. Deep models have been widely applied in various fields such as autonomous driving [27], medical science [32], and finance [37]. With the spread of these techniques, the robustness and safety issues begin to be essential, especially after the finding that deep models can be easily fooled by negligible noises [15]. As a result, more researchers contribute to building datasets for benchmark-

^{*} Corresponding author.

This research is supported in part by the National Key Research and Development Progrem of China under Grant No.2020AAA0140000.

ing model robustness to spot vulnerabilities in advance.

Most of the existing work builds datasets for evaluating the model robustness and generalization ability on outof-distribution data [6, 21, 29] using adversarial examples and common corruptions. For example, the ImageNet-C(orruption) dataset conducts visual corruptions such as Gaussian noise to input images to simulate the possible processors in real scenarios [21]. ImageNet-R(enditions) contains various renditions (e.g., paintings, embroidery) of ImageNet object classes [20]. As both studies have found that higher accuracy on ImageNet usually leads to better robustness against different domains [21,50]. However, most previous studies try to achieve this in a top-down way, such as architecture design, exploring a better training strategy, etc. We advocate that it is also essential to manage it in a bottom-up way, that is, conducting model debugging with the in-distribution dataset to provide clues for model repairing and accuracy improvement. For example, it is interesting to explore whether a bird with a water background can be recognized correctly even if most birds appear with trees or grasses in the training data. Though this topic has been investigated in studies such as causal and effect analysis [8], the experiments and analysis are undertaken on domain generalization datasets. How a deep model generalizes to different backgrounds is still unknown due to the vacancy of a qualified benchmark. Therefore, in this paper, we provide a detached object editing tool to conduct the model debugging from the perspective of object attribute and construct a dataset named ImageNet-E(diting).

The ImageNet-E dataset is a compact but challenging test set for object recognition that contains controllable object attributes including backgrounds, sizes, positions and directions, as shown in Fig. 1. In contrast to ObjectNet [5] whose images are collected by their workers via posing objects according to specific instructions and differ from the target data distribution. This makes it hard to tell whether the degradation comes from the changes of attribute or distribution. Our ImageNet-E is automatically generated with our object attribute editing tool based on the original ImageNet. Specifically, to change the object background, we provide an object background editing method that can make the background simpler or more complex based on diffusion models [24, 46]. In this way, one can easily evaluate how much the background complexity can influence the model performance. To control the object size, position, and direction to simulate pictures taken from different distances and angles, an object editing method is also provided. With the editing toolkit, we apply it to the large-scale ImageNet dataset [41] to construct our ImageNet-E(diting) dataset. It can serve as a general dataset for benchmarking robustness evaluation on different object attributes.

With the ImageNet-E dataset, we evaluate the performance of current deep learning models, including both convolutional neural networks (CNNs), vision transformers as well as the large-scale pretrained CLIP [39]. We find that deep models are quite sensitive to object attributes. For example, when editing the background towards high complexity (see Fig. 1, the 3rd row in the background part), the drop in top-1 accuracy reaches 9.23% on average. We also find that though some robust models share similar top-1 accuracy on ImageNet, the robustness against different attributes may differ a lot. Meanwhile, some models, being robust under certain settings, even show worse results than the vanilla ones on our dataset. This suggests that improving robustness is still a challenging problem and the object attributes should be taken into account. Afterward, we discover ways to enhance robustness against object attribute changes. The main contributions are summarized as follows:

- We provide an object editing toolkit that can change the object attributes for manipulated image generation.
- We provide a new dataset called ImageNet-E that can be used for benchmarking robustness to different object attributes. It opens up new avenues for research in robust computer vision against object attributes.
- We conduct extensive experiments on ImageNet-E and find that models that have good robustness on adversarial examples and common corruptions may show poor performance on our dataset.

2. Related Work

The literature related to attribute robustness benchmarks can be broadly grouped into the following themes: robustness benchmarks and attribute editing datasets. Existing robustness benchmarks such as ImageNet-C(orruption) [21], ImageNet-R(endition) [20], ImageNet-Stylized [13] and ImageNet-3DCC [29] mainly focus on the exploration of the corrupted or out-of-distribution data that models may encounter in reality. For instance, the ImageNet-R dataset contains various renditions (e.g., paintings, embroidery) of ImageNet object classes. ImageNet-C analyzes image models in terms of various simulated image corruptions (e.g., noise, blur, weather, JPEG compression, etc.). Attribute editing dataset creation is a new topic and few studies have explored it before. Among them, ObjectNet [5] and ImageNet-9 (a.k.a. background challenge) [50] can be the representative. Specifically, ObjectNet collects a large realworld test set for object recognition with controls where object backgrounds, rotations, and imaging viewpoints are random. The images in ObjectNet are collected by their workers who image objects in their homes. It consists of 313 classes which are mainly household objects. ImageNet-9 mainly creates a suit of datasets that help disentangle the impact of foreground and background signals on classification. To achieve this goal, it uses coarse-grained classes



Figure 2. Attribute editing with DDPMs. Give an input image and its corresponding object mask, the object is firstly removed with inpainting operation to get the pure background image. Then, we leverage the diffusion process to edit the background image x_0 and object image coherently. \odot denotes the element-wise blending of these two images using the object mask. For background editing, the background complexity objective function is added during the diffusion process (Alg. 1, line 5). For other object attributes editing, the object image needs to be transformed first (Alg. 2, line 1).

with corresponding rectangular bounding boxes to remove the foreground and then paste the cut area with other backgrounds. It can be observed that there lacks a dataset that can smoothly edit the object attribute.

3. Preliminaries

Since the editing tool is developed based on diffusion models, let us first briefly review the theory of denoising diffusion probabilistic models (DDPM) [24,46] and analyze how it can be used to generate images.

According to the definition of the Markov Chain, one can always reach a desired stationary distribution from a given distribution along with the Markov Chain [14]. To get a generative model that can generate images from random Gaussian noises, one only needs to construct a Markov Chain whose stationary distribution is Gaussian distribution. This is the core idea of DDPM. In DDPM, given a data distribution $\mathbf{x}_0 \sim q(\mathbf{x}_0)$, a forward noising process produces a series of latents $\mathbf{x}_1, ..., \mathbf{x}_T$ of the same dimensionality as the data \mathbf{x}_0 by adding Gaussian noise with variance $\beta_t \in (0, 1)$ at time t:

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}), s.t. \ 0 < \beta_t < 1,$$
(1)

where β_t is the diffusion rate. Then the distribution $q(\mathbf{x}_t|\mathbf{x}_0)$ at any time t is:

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t}, (1 - \bar{\alpha}_t)\mathbf{I}), \ \mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon$$
(2)

where $\bar{\alpha}_t = \prod_{s=1}^t (1 - \beta_t)$, $\epsilon \sim \mathcal{N}(0, \mathbf{I})$. It can be proved that $\lim_{t\to\infty} q(\mathbf{x}_t) = \mathcal{N}(0, \mathbf{I})$. In other words, we can map the original data distribution into a Gaussian distribution with enough iterations. Such a stochastic forward process is named as diffusion process since what the process $q(\mathbf{x}_t|\mathbf{x}_{t-1})$ does is adding noise to \mathbf{x}_{t-1} .

To draw a fresh sample from the distribution $q(\mathbf{x}_0)$, the Markov process is reversed. That is, beginning from a

Gaussian noise sample $\mathbf{x}_T \sim \mathcal{N}(0, \mathbf{I})$, a reverse sequence is constructed by sampling the posteriors $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$. To approximate the unknown function $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$, in DDPMs, a deep model p_{θ} is trained to predict the mean and the covariance of \mathbf{x}_{t-1} given \mathbf{x}_t instead. Then the \mathbf{x}_{t-1} can be sampled from the normal distribution defined as:

$$p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mu_{\theta}(\mathbf{x}_t, t), \Sigma_{\theta}(\mathbf{x}_t, t)).$$
(3)

In stead of inferring $\mu_{\theta}(\mathbf{x}_t, t)$ directly, [24] propose to predict the noise $\epsilon_{\theta}(\mathbf{x}_t, t)$ which was added to \mathbf{x}_0 to get \mathbf{x}_t with Eq. (2). Then $\mu_{\theta}(\mathbf{x}_t, t)$ is:

$$\mu_{\theta}(\mathbf{x}_t, t) = \frac{1}{\sqrt{\bar{\alpha}_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right).$$
(4)

[24] keep the value of $\Sigma_{\theta}(\mathbf{x}_t, t)$ to be constant. As a result, given a sample \mathbf{x}_t at time t, with a trained model that can predict the noise $\epsilon_{\theta}(\mathbf{x}_t, t)$, we can get $\mu_{\theta}(\mathbf{x}_t, t)$ according to Eq. (4) to reach the \mathbf{x}_{t-1} with Equation (3) and eventually we can get to \mathbf{x}_0 .

Previous studies have shown that diffusion models can achieve superior image generation quality compared to the current state-of-the-art generative models [1]. Besides, there have been plenty of works on utilizing the DDPMs to generate samples with desired properties, such as semantic image translation [36], high fidelity data generation from low-density regions [44], *etc.* In this paper, we also choose the DDPM adopted in [1] as our generator.

4. Attribute Editing with Diffusion Models and ImageNet-E

Most previous robustness-related work has focused on the important challenges of robustness on adversarial examples [6], common corruptions [21]. They have found that higher clean accuracy usually leads to better robustness. Therefore, instead of exploring a new corruption that models may encounter in reality, we pay attention to the model debugging in terms of object attributes, hoping to provide new insights to clean accuracy improvement. In the following, we describe our object attribute editing tool and the generated ImageNet-E dataset in detail. The whole pipeline can be found in Fig. 2.

4.1. Object Attribute Editing with Diffusion Models

Background editing. Most existing corruptions conduct manipulations on the whole image, as shown in Fig. 1. Compared to adding global corruptions that may hinder the visual quality, a more likely-to-happen way in reality is to manipulate the backgrounds to fool the model. Besides, it is shown that there exists a spurious correlation between labels and image backgrounds [12]. From this point, a background corruption benchmark is needed to evaluate the model's robustness. However, the existing background challenge dataset achieves background editing with copy-paste operation, resulting an obvious artifacts in generated images [50]. This may leave some doubts about whether the evaluation is precise since the dataset's distribution may have changed. To alleviate this concern, we adopt DDPM approach to incorporate background editing by adding a guiding loss that can lead to backgrounds with desired properties to make the generated images stay in/close to the original distribution. Specifically, we choose to manipulate the background in terms of texture complexity due to the hypothesis that an object should be observed more easily from simple backgrounds than from complicated ones. In general, the texture complexity can be evaluated with the gray-level cooccurrence matrix (GLCM) [16], which calculates the graylevel histogram to show the texture characteristic. However, the calculation of GLCM is non-differentiable, thus it cannot serve as the conditional guidance of image generation. We hypothesize that a complex image should contain more frequency components in its spectrum and higher amplitude indicates greater complexity. Thus, we define the objective of complexity as:

$$\mathcal{L}_{c} = \sum \left| \mathcal{A}(\mathcal{F}(\mathbf{x})) \right|, \qquad (5)$$

where \mathcal{F} is the Fourier transform [45], \mathcal{A} extracts the amplitude of the input spectrum. **x** is the evaluated image. Since minimizing this loss helps us generate an image with desired properties and should be conducted on the **x**₀, we need a way of estimating a clean image **x**₀ from each noisy latent representation **x**_t during the denoising diffusion process. Recall that the process estimates at each step the noise $\epsilon_{\theta}(\mathbf{x}_t, t)$ added to **x**₀ to obtain **x**_t. Thus, $\hat{\mathbf{x}}_0$ can be estimated via Equation (6) [1]. The whole optimization procedure is shown in Algorithm 1.

$$\hat{\mathbf{x}}_0 = \frac{\mathbf{x}_t}{\sqrt{\bar{\alpha}_t}} - \frac{\sqrt{1 - \bar{\alpha}_t}\epsilon_\theta(\mathbf{x}_t, t)}{\sqrt{\bar{\alpha}_t}}.$$
 (6)

As shown in Fig. 3(a), with the proposed method, when we guide the generation procedure with the proposed objective towards the complex direction, it will return images with visually complex backgrounds. We also provide the GLCM dissimilarity and contrast of each image to make a quantitative analysis of the generated images. A higher dissimilarity/contrast score indicates a more complex image background [16]. It can be observed that the complexity is consistent with that calculated with GLCM, indicating the effectiveness of the proposed method.

Controlling object size, position and direction. In general, the human vision system is robust to position, direction and small size changes. Whether the deep models are also robust to these object attribute changes is still unknown to researchers. Therefore, we conduct the image editing with controls of object sizes, positions and directions to find the answer. For a valid evaluation on different attributes, all other variables should remain unchanged, especially the background. Therefore, we first disentangle the object and background with the in-painting strategy provided by [54]. Specifically, we mask the object area in input image x. Then we conduct in-painting to remove the object and get the pure background image \mathbf{x}^b , as shown in Fig. 3(b) column 3. To realize the aforementioned object attribute controlling, we adopt the orthogonal transformation. Denote P as the pixel locations of object in image \mathbf{x} where $P \in \mathbb{R}^{3 \times N_o}$. N_o is the number of pixels belong to object and $p_i = [x_i, y_i, 1]^T$ is the position of object's *i*-th pixel. $h' \in [0, H - h], w' \in [0, W - w]$ where [x, y, w, h]stand for the enclosing rectangle of the object with mask M. Then the newly edited $\mathbf{x}[T_{\text{attribute}} \cdot P] = \mathbf{x}[P]$ and $M[T_{\text{attribute}} \cdot P] = M[P]$, where

$$T_{\text{size}} = \begin{bmatrix} s \ 0 \ \Delta x \\ 0 \ s \ \Delta y \\ 0 \ 0 \ 1 \end{bmatrix}, T_{\text{position}} = \begin{bmatrix} 1 \ 0 \ w' \\ 0 \ 1 \ h' \\ 0 \ 0 \ 1 \end{bmatrix}, T_{\text{direction}} = \begin{bmatrix} \cos \theta \ \sin \theta \ 0 \\ -\sin \theta \cos \theta \ 0 \\ 0 \ 0 \ 1 \end{bmatrix}.$$
(7)

where s is the resize scale. θ is the rotation angle. $\Delta x = (1-s) \cdot (x+w/2), \Delta y = (1-s) \cdot (y+h/2).$

With the background image \mathbf{x}^b and edited object \mathbf{x}^o , a naive way is to place the object in the original image to the corresponding area of background image \mathbf{x}^b as $M \odot \mathbf{x}^o + (1 - M) \odot \mathbf{x}^b$. However, the result generated in this manner may look disharmonic, lacking a delicate adjustment to blending them together. Besides, as shown in Fig. 3(b) column 3, the object-removing operation may leave some artifacts behind, failing to produce a coherent and seamless result. To deal with this problem, we leverage DDPM models to blend them at different noise levels along the diffusion process. Denote the image with desired object attribute as \mathbf{x}^o . Starting from the pure background image \mathbf{x}^b at time t_0 , at each stage, we perform a guided diffusion step with a latent \mathbf{x}_t to obtain the \mathbf{x}_{t-1} and at the same time, obtain a noised version of object image



Figure 3. (a) Images generated with the proposed background complexity editing method. (b) Edited images with size changing. The Fréchet inception distance (FID) for pasting is 50.64 while it is 32.59 for ours, indicating the effectiveness of the leveraging of DDPMs.

Algorithm 1: Background editing	Algorithm 2: Object size controlling
input : source image x , mask M , diffusion model $(\mu_{\theta}(\mathbf{x}_t), \Sigma_{\theta}(\mathbf{x}_t)), \bar{\alpha}_t, \lambda$, iteration steps t_0 output: edited image \mathbf{x}_0	input : source image \mathbf{x} , mask M , diffusion model $(\mu_{\theta}(\mathbf{x}_t), \Sigma_{\theta}(\mathbf{x}_t)), \bar{\alpha}_t$, iteration steps t_0 , ratio s output: edited image \mathbf{x}_0
$ \begin{array}{c c c} \mathbf{i} & \mathbf{x}_{t_0} \sim \mathcal{N}(\sqrt{\bar{\alpha}_{t_0}}\mathbf{x}, (1 - \bar{\alpha}_{t_0})\mathbf{I}); \\ 2 & \mathbf{for} \ t \leftarrow t_0 \ \mathbf{to} \ 1 \ \mathbf{do} \\ 3 & \mathbf{\hat{x}}_0 \leftarrow \frac{\mathbf{x}_t}{\sqrt{\bar{\alpha}_t}} - \frac{\sqrt{1 - \bar{\alpha}_t}\epsilon_{\theta}(\mathbf{x}_t, t)}{\sqrt{\bar{\alpha}_t}}; \\ 4 & \nabla_{bg} \leftarrow \nabla_{\hat{\mathbf{x}}_0}\mathcal{L}_c(\hat{\mathbf{x}}_0); \\ 5 & \mathbf{x}_{t-1}^b \sim \mathcal{N}(\mu_{\theta}(\mathbf{x}_t) + \lambda \Sigma_{\theta}(\mathbf{x}_t) \nabla_{bg}, \Sigma_{\theta}(\mathbf{x}_t)); \\ 6 & \mathbf{x}^o \sim \mathcal{N}(\sqrt{\bar{\alpha}_t}\mathbf{x}, (1 - \bar{\alpha}_t)\mathbf{I}); \\ 7 & \mathbf{x}_{t-1} \leftarrow M \odot \mathbf{x}^o + (1 - M) \odot \mathbf{x}_{t-1}^b; \\ 8 \ \mathbf{end} \end{array} $	$ 1 \ \mathbf{x}^{b} \leftarrow ObjectRemoving(\mathbf{x}, M); 2 \ \mathbf{x}, M \leftarrow Rescale \ (\mathbf{x}, M, s); 3 \ \mathbf{x}_{t_{0}} \sim \mathcal{N}(\sqrt{\overline{\alpha}_{t_{0}}}\mathbf{x}^{b}, (1 - \overline{\alpha}_{t_{0}})\mathbf{I}); 4 \ \mathbf{for} \ t \leftarrow t_{0} \ \mathbf{to} \ 1 \ \mathbf{do} \\ 5 \ \ \left \begin{array}{c} \mathbf{x}_{t-1}^{b} \sim \mathcal{N}(\mu_{\theta}(\mathbf{x}_{t}), \Sigma_{\theta}(\mathbf{x}_{t})); \\ \mathbf{x}^{o} \sim \mathcal{N}(\sqrt{\overline{\alpha}_{t}}\mathbf{x}, (1 - \overline{\alpha}_{t})\mathbf{I}); \\ 7 \ \ \right \ \mathbf{x}_{t-1} \leftarrow M \odot \mathbf{x}^{o} + (1 - M) \odot \mathbf{x}_{t-1}^{b}; \\ 8 \ \mathbf{end} \end{array} $
8 end	8 end

 \mathbf{x}_{t-1}^{o} . Then the two latents are blended with the mask M as $\mathbf{x}_{t-1} = M \odot \mathbf{x}_{t-1}^{o} + (1 - M) \odot \mathbf{x}_{t-1}$. The DDPM denoising procedure may change the background. Thus a proper initial timing is required to maintain a high resemblance to the original background. We set the iteration steps t_0 as 50 and 25 in Algorithm 1 and 2 respectively.

4.2. ImageNet-E dataset

With the tool above, we conduct object attribute editing including background, size, direction and position changes based on the large-scale ImageNet dataset [41] and ImageNet-S [11], which provides the mask annotation. To guarantee the dataset quality, we choose the animal classes from ImageNet classes such as dogs, fishes and birds, since they appear more in nature without messy backgrounds. Classes such as stove and mortarboard are removed. Finally, our dataset consists of 47872 images with 373 classes based on the initial 4352 images, each of which is applied 11 transforms. Detailed information can be found in Appendix A. For background editing, we choose three levels of the complexity, including $\lambda = -20, \lambda = 20$ and $\lambda =$ 20-adv with adversarial guidance (see Sec.B for details) instead of complexity. Larger λ indicates stronger guidance towards high complexity. For the object size, we design

four levels of sizes in terms of the object pixel rates (= $sum(M > 0.5)/sum(M \ge 0)$): [Full, 0.1, 0.08, 0.05] where 'Full' indicates making the object as large as possible while maintaining its whole body inside the image. Smaller rates indicate smaller objects. For object position, we find that some objects hold a high object pixel rate in the whole image, resulting in a small H - h. Take the first picture in Fig. 3 for example, the dog is big and it will make little visual differences after position changing. Thus, we adopt the data whose pixel rate is 0.05 as the initial images for the position-changing operation.

In contrast to benchmarks like ImageNet-C [21] giving images from different domains so that the model robustness in these situations may be assessed, our effort aims to give an editable image tool that can conduct model debugging with in-distribution (ID) data, in order to identify specific shortcomings of different models and provide some insights for clean accuracy improving. Thus, the data distribution should not differ much from the original ImageNet. We choose the out-of-distribution (OOD) detection methods Energy [33] and GradNorm [26] to evaluate whether our editing tool will move the edited image out of its original distribution. These OOD detection methods aim to distinguish the OOD examples from the ID exam-



Figure 4. Distributions of ID score of different datasets in terms of the quantities in Energy (the first row) and GradNorm (the second row) for in-distribution (ImageNet) and other datasets. Higher overlap indicates greater proximity to ImageNet.

ples. The results are shown in Fig. 4. x-axis is the ID score in terms of the quantities in Energy and GradNorm and y-axis is the frequency of each ID score. A high ID score indicates the detection method takes the input sample as the ID data. Compared to other datasets, our method barely changes the data distribution under both Energy (the 1st row) and GradNorm (the 2nd row) evaluation methods. Besides, the Fréchet inception distance (FID) [23] for our ImageNet-E is 15.57 under the random background setting, while it is 34.99 for ImageNet-9 (background challenge). These all imply that our editing tool can ensure the proximity to the original ImageNet, thus can give a controlled evaluation on object attribute changes. To find out whether the DDPM will induce some degradation to our evaluation, we have conducted experiment in Tab. 1 with the setting $\lambda = 0$ during background editing. This operation will first add noises to the original and then denoise them. It can be found in "Inver" column that the degradation is negligible compared to degradation induced by attribute changes.

5. Experiments

We conduct evaluation experiments on various architectures including both CNNs (ResNet (RN) [19], DenseNet [25], EfficientNet (EF) [47], ResNest [53], ConvNeXt [35]) and transformer-based models (Vision-Transformer (ViT) [9], Swin-Transformer (Swin) [34]). Other state-of-the-art models that trained with extra data such as CLIP [39], EfficientNet-L2-Noisy-Student [51] are also evaluated in the Appendix. Apart from different sizes of these models, we have also evaluated their adversarially trained versions for comprehensive studies. We report the drop of top-1 accuracy as metric based on the idea that the attribute changes should induce little influence to a robust trained model. More experimental details and results of top-1 accuracy can be found in the Appendix.

5.1. Robustness evaluation

Normally trained models. To find out whether the widely used models in computer vision have gained robustness against changes on different object attributes, we conduct extensive experiments on different models. As shown in Tab. 1, when only the background is edited towards high complexity, the average drop in top-1 accuracy is 9.23% ($\lambda = 20$). This indicates that most models are sensitive to object background changes. Other attribute changes such as size and position can also lead to model performance degradation. For example, when changing the object pixel rate to 0.05, as shown in Fig. 1 row 4 in the 'size' column, while we can still recognize the image correctly, the performance drop is 18.34% on average. We also find that the robustness under different object attributes is improved along with improvements in terms of clean accuracy (Original) on different models. Accordingly, a switch from an RN50 (92.69% top-1 accuracy) to a Swin-S (96.21%) leads to the drop in accuracy decrease from 15.72% to 10.20% on average. By this measure, models have become more and more capable of generalizing to different backgrounds, which implies that they indeed learn some robust features. This shows that object attribute robustness can be a good way to measure future progress in representation learning. We also observe that larger networks possess better robustness on the attribute editing. For example, swapping a Swin-S (96.21% top-1 accuracy) with the larger Swin-B (95.96% top-1 accuracy) leads to the decrease of the dropped accuracy from 10.20% to 8.99% when $\lambda = 20$. In a similar fashion, a ConvNeXt-T (9.32% drop) is less robust than the giant ConvNeXt-B (7.26%). Consequently, models with even more depth, width, and feature aggregation may attain further attribute robustness. Previous studies [30] have shown that zero-shot CLIP exhibits better outof-distribution robustness than the finetuned CLIP, which is opposite to our ImageNet-E as shown in Tab. 1. This may serve as the evidence that our ImageNet-E has a good proximity to ImageNet. We also find that compared with fullysupervised trained model under the same backbone (ViT-B), the CLIP fails to show a better attribute robustness. We think this may be caused by that the CLIP has spared some capacity for OOD robustness.



Figure 5. Comparisons between vanilla models and adversarially trained models across different architectures in terms of size changes (left). Evaluation of adversarial models (RN50) trained with different perturbation budgets is provided in the right figure.

Adversarially trained models. Adversarial training [42] is one of the state-of-the-art methods for improving the adversarial robustness of deep models and has been widely studied [2]. To find out whether they can boost the attribute robustness, we conduct extensive experiments in terms of different architectures and perturbation budgets (constraints of l_2 norm bound). As shown in Fig. 5, the adversarially trained ones are not robust against attribute changes including both backgrounds and size-changing situations. The dropped accuracies are much greater compared to normally trained models. As the perturbation budget grows, the situation gets worse. This indicates that adversarial training can do harm to robustness against attributes.

5.2. Robustness enhancements

Based on the above evaluations, we step further to discover ways to enhance the attribute robustness in terms of preprocessing, network design and training strategies. More details including training setting and numerical experimental results can be found in Appendix C.5.

Preprocessing. Given that an object can be inconspicuous due to its small size or subtle position, viewing an object at several different locations may lead to a more stable prediction. Having this intuition in mind, we perform the classical Ten-Crop strategy to find out if this operation can help to get a robustness boost. The Ten-Crop operation is executed by cropping all four corners and the center of the input image. We average the predictions of these crops together with their horizontal mirrors as the final result. We find this operation can contribute a 0.69% and 1.24% performance boost on top-1 accuracy in both background and size changes scenarios on average respectively.

Network designs. Intuitively, a robust model should tend to focus more on the object of interest instead of the background. Therefore, recent models begin to enhance the model by employing attention modules. Of these, the ResNest [53] can be a representative. The ResNest is a modularized architecture, which applies channel-wise attention on different network branches to leverage their success in capturing cross-feature interactions and learning diverse representations. As it has achieved a great boost in the ImageNet dataset, it also shows superiority on ImageNet-E compared to ResNet. For example, a switch from RN50 decreases the average dropped accuracy from 15.72% to 12.57%. This indicates that the channel-wise attention module can be a good choice to improve the attribute robustness. Another representative model can be the vision transformer, which consists of multiple self-attention modules. To study whether incorporating transformer's self-attention-like architecture into the model design can help attribute robustness generalization, we establish a hybrid architecture by directly feeding the output of res_3 block in RN50 into ViT-S as the input feature like [3]. The dropped accuracy decreases by 1.04% compared to the original RN50, indicating the effectiveness of the self-attention-like architectures.

Training strategy. a) Robust trained. There have been plenty of studies focusing on the robust training strategy to improve model robustness. To find out whether these works can boost the robustness on our dataset, we further evaluate these state-of-the-art models including SIN [13], Debiased-CNN [31], Augmix [22], ANT [40], DeepAugment [20] and model trained with lots of standard augmentations (RN50-T) [48]. As shown in Tab. 2, apart from the RN50-T, while the Augmix model shows the best performance against the background change scenario, the Debiased model holds the best in the object size change scenario. What we find unexpectedly is the SIN performance. The SIN method features the novel data augmentation scheme where ImageNet images are stylized with style transfer as the training data to force the model to rely less on textural cues for classification. Though the robustness boost is achieved on ImageNet-C (mCE 69.32%) compared to its vanilla model (mCE 76.7%), it fails to improve the robustness in both object background and size-changing scenarios. The drops of top-1 accuracy for vanilla RN50 and RN50-SIN are 21.26% and 24.23% respectively, when the object size rate is 0.05, though they share similar accuracy on original ImageNet. This indicates that existing benchmarks cannot reflect the real robustness in object attribute changing. Therefore, a dataset like ImageNet-E is necessary for comprehensive evaluations on deep models. b) Masked image modeling. Considering that masked image modeling has demonstrated impressive results in self-supervised representation learning by recovering corrupted image patches [4], it may be robust to the attribute changes. Therefore, we choose the Masked AutoEncoder (MAE) [17] as the training strategy since its objective is recovering images with only 25% patches. Specifically, we adopt the MAE training strategy with ViT-B backbone and then finetune it with ImageNet

Table 1. Evaluations with different state-of-the-art models in terms of Top-1 accuracy and the corresponding drop of accuracy under background changes, size changes, random position (rp) and random direction (rd).

Madala	Original		Ba	ckground c	changes			Size c	hanges		Position	Direction	4.00
widdels	Original	Inver	$\lambda = -20$	$\lambda = 20$	$\lambda = 20\text{-adv}$	Random	Full	0.1	0.08	0.05	rp	rd	Avg.
RN50	92.69%	1.97%	7.30%	13.35%	29.92%	13.34%	2.71%	7.25%	10.51%	21.26%	26.46%	25.12%	15.72%
DenseNet121	92.10%	1.49%	6.29%	9.00%	29.20%	12.43%	3.50%	7.00%	10.68%	21.55%	26.53%	23.64%	14.98%
EF-B0	92.85%	1.07%	7.10%	10.71%	34.88%	15.64%	3.03%	8.00%	11.57%	23.28%	27.91%	19.11%	16.12%
ResNest50	95.38%	1.44%	6.33%	8.98%	26.62%	11.28%	2.53%	5.27%	8.01%	18.03%	21.37%	17.32%	12.57%
ViT-S	94.14%	0.82%	6.42%	8.98%	31.12%	13.06%	0.80%	5.37%	8.59%	17.37%	22.86%	17.13%	13.17%
Swin-S	96.21%	1.13%	5.18%	7.33%	23.50%	9.31%	1.27%	4.21%	6.29%	14.16%	17.35%	13.42%	10.20%
ConvNeXt-T	96.07%	1.43%	4.69%	6.26%	19.83%	7.93%	1.75%	3.28%	5.18%	12.76%	15.71%	15.78%	9.32%
	94.00%	2.11%	7.05%	11.62%	29.47%	13.57%	2.57%	6.81%	10.12%	20.65%	25.85%	24.42%	15.21%
DenseNet169	92.37%	1.12%	5.81%	8.43%	27.51%	11.61%	2.25%	6.90%	10.41%	20.59%	24.93%	20.68%	13.91%
EF-B3	94.97%	1.87%	7.77%	8.40%	29.90%	12.92%	1.36%	6.80%	10.16%	21.36%	24.98%	17.24%	14.09%
ResNest101	95.54%	1.10%	5.58%	6.65%	23.03%	10.40%	1.35%	3.97%	6.53%	15.44%	19.11%	14.31%	10.64%
ViT-B	95.38%	0.83%	5.32%	8.43%	26.60%	10.98%	0.62%	4.00%	6.30%	14.51%	18.82%	14.95%	11.05%
Swin-B	95.96%	0.79%	4.46%	6.23%	21.44%	8.25%	0.99%	3.16%	5.04%	12.34%	15.38%	12.60%	8.99%
ConvNeXt-B	96.42%	0.69%	3.75%	4.86%	16.49%	6.04%	0.99%	2.25%	3.36%	9.47%	12.40%	13.01%	7.26%
CLIP-zeroshot	80.01%	4.88%	11.56%	15.28%	36.14%	20.09%	3.33%	12.67%	15.77%	25.31%	28.87%	21.57%	19.06%
CLIP-finetuned	93.68%	2.17%	9.82%	11.83%	38.33%	18.19%	9.06%	9.25%	12.67%	23.32%	28.56%	22.00%	18.30%

Table 2. Evaluations with different robust models in terms of Top-1 accuracy and the corresponding dropped accuracy.

Architacturas	Ori	Background changes					Size changes				Position	Direction	Ava
Architectures	UII	Inver	$\lambda = -20$	$\lambda = 20$	$\lambda = 20$ -adv	Random	Full	0.1	0.08	0.05	rp	rd	Avg.
RN50	92.69%	1.97%	7.30%	13.35%	29.92%	13.34%	2.71%	7.25%	10.51%	21.26%	26.46%	25.12%	15.72%
RN50-Adversarial	81.96%	0.66%	4.75%	13.62%	37.87%	15.25%	4.87%	9.62%	13.94%	25.51%	32.51%	31.96%	18.99%
RN50-SIN	91.57%	2.23%	7.61%	12.19%	33.16%	13.58%	1.68%	8.30%	12.60%	24.23%	29.16%	27.24%	16.98%
RN50-Debiased	93.34%	1.43%	6.09%	11.45%	27.99%	12.12%	1.98%	5.53%	8.76%	19.27%	24.01%	24.97%	14.22%
RN50-Augmix	93.50%	0.98%	6.26%	8.38%	30.49%	12.94%	1.61%	6.40%	9.97%	21.42%	27.14%	22.42%	14.70%
RN50-ANT	91.87%	1.68%	6.62%	11.94%	35.66%	15.36%	1.57%	7.12%	10.62%	21.49%	26.66%	25.23%	16.23%
RN50-DeepAugment	92.88%	1.50%	6.62%	12.37%	32.40%	13.32%	1.36%	7.27%	10.62%	21.28%	26.28%	21.29%	15.28%
RN50-T	94.55%	1.05%	5.65%	7.38%	21.89%	10.42%	2.11%	4.74%	7.83%	17.46%	21.12%	19.60%	11.82%

training data. We find that the robustness is improved. For example, the dropped accuracy decreases from 10.62% to 9.05% on average compared to vanilla ViT-B.

5.3. Failure case analysis

To explore the reason why some robust trained models may fail, we leverage the LayerCAM [28] to generate the heat map for different models including vanilla RN50, RN50+SIN and RN50+Debiased for comprehensive studies. As shown in Fig. 6, the heat map of the Debiased model aligns better with the objects in the image than that of the original model. It is interesting to find that the SIN model sometimes makes wrong predictions even with its attention on the main object. We suspect that the SIN model relies too much on the shape. for example, the 'sea urchin' looks like the 'acron' with the shadow. However, its texture clearly indicates that it is the 'sea urchin'. In contrast, the Debiased model which is trained to focus on both the shape and texture can recognize it correctly. More studies can be found in Appendix C.4.

5.4. Model repairing

To validate that the evaluation on ImageNet (IN)-E can help to provide some insights for model's applicability and enhancement, we conduct a toy example for model repair-



Figure 6. Heat maps for explaining which parts of the image dominate the model decision through LayerCAM [28].

ing. Previous evaluation shows that the ResNet50 is vulnerable to background changes. Based on this observation, we randomly replace the backgrounds of objects with others during training and get a validation accuracy boost from 77.48% to 79.00%. Note that the promotion is not small as only 8781 training images with mask annotations are available in ImageNet. We also step further to find out if the improved model can get a boost the OOD robustness, as shown in the Tab. 3. It can be observed that with the insights provided by the evaluation on ImageNet-E, one can explore the model's attribute vulnerabilities and manage to repair the model and get a performance boost accordingly.

Table 3. Model repairing results. Top-1 accuracy (%) is reported except for IN-C, which is mCE (mean Corruption Error). Higher top-1 accuracy and lower mCE indicate better performance. IN-E reports the average accuracy on ImageNet-E.

Models	IN	IN-v2	IN-A	IN-C↓	IN-R	IN-Sketch	IN-E
RN50	77.5	65.7	6.5	68.6	39.6	27.5	83.7
RN50-repaired	79.0	67.2	9.4	65.8	40.7	29.4	85.0

6. Conclusion and Future work

In this paper, we put forward an image editing toolkit that can take control of object attributes smoothly. With this tool, we create a new dataset called ImageNet-E that can serve as a general dataset for benchmarking robustness against different object attributes. Extensive evaluations conducted on different state-of-the-art models show that most models are vulnerable to attribute changes, especially the adversarially trained ones. Meanwhile, other robust trained models can show worse results than vanilla models even when they have achieved a great robustness boost on other robustness benchmarks. We further discover ways for robustness enhancement from both preprocessing, network designing and training strategies.

Limitations and future work. This paper proposes to edit the object attributes in terms of backgrounds, sizes, positions and directions. Therefore, the annotated mask of the interest object is required, resulting in a limitation of our method. Besides, since our editing toolkit is developed based on diffusion models, the generalization ability is determined by DDPMs. For example, we find synthesizing high-quality person images is difficult for DDPMs. Under the consideration of both the annotated mask and data quality, our ImageNet-E is a compact test set. In our future work, we would like to explore how to leverage the edited data to enhance the model's performance, including both the validation accuracy and robustness.

References

- Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended diffusion for text-driven editing of natural images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18208–18218, 2022. 3, 4
- [2] Tao Bai, Jinqi Luo, Jun Zhao, Bihan Wen, and Qian Wang. Recent advances in adversarial training for adversarial robustness. arXiv preprint arXiv:2102.01356, 2021. 7
- [3] Yutong Bai, Jieru Mei, Alan L Yuille, and Cihang Xie. Are transformers more robust than cnns? *Advances in Neural Information Processing Systems*, 34:26831–26843, 2021. 7, 16
- [4] Hangbo Bao, Li Dong, and Furu Wei. Beit: Bert pre-training of image transformers. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2022. 7, 16
- [5] Andrei Barbu, David Mayo, Julian Alverio, William Luo, Christopher Wang, Dan Gutfreund, Josh Tenenbaum, and Boris Katz. Objectnet: A large-scale bias-controlled dataset for pushing the limits of object recognition models. *Advances in neural information processing systems*, 32, 2019.
- [6] Nicholas Carlini and David Wagner. Adversarial examples are not easily detected: Bypassing ten detection methods. In *Proceedings of the 10th ACM workshop on artificial intelli*gence and security, pages 3–14, 2017. 2, 3
- [7] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9640–9649, 2021. 16
- [8] Peng Cui and Susan Athey. Stable learning establishes some common ground between causal inference and machine learning. *Nature Machine Intelligence*, 4(2):110–115, 2022. 2
- [9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929, 2020. 6
- [10] Noam Eshed. Novelty detection and analysis in convolutional neural networks. Master's thesis, Cornell University, 2020. 13
- [11] Shanghua Gao, Zhong-Yu Li, Ming-Hsuan Yang, Ming-Ming Cheng, Junwei Han, and Philip Torr. Large-scale unsupervised semantic segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–20, 2022. 5, 13
- [12] Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673, 2020. 4
- [13] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. *arXiv preprint arXiv:1811.12231*, 2018. 2, 7, 19

- [14] Charles J Geyer. Practical markov chain monte carlo. Statistical science, pages 473–483, 1992. 3
- [15] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. arXiv preprint arXiv:1412.6572, 2014. 1
- [16] Robert M Haralick, Karthikeyan Shanmugam, and Its' Hak Dinstein. Textural features for image classification. *IEEE Transactions on systems, man, and cybernetics*, 11(6):610–621, 1973. 4
- [17] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition, pages 16000– 16009, 2022. 7
- [18] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition, pages 16000– 16009, 2022. 16
- [19] 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. 6
- [20] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8340–8349, 2021. 2, 7, 20
- [21] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *Proceedings of the International Conference on Learning Representations*, 2019. 1, 2, 3, 5
- [22] Dan Hendrycks, Norman Mu, Ekin D. Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshminarayanan. AugMix: A simple data processing method to improve robustness and uncertainty. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020. 7, 20
- [23] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. Advances in neural information processing systems, 30, 2017. 6
- [24] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. Advances in Neural Information Processing Systems, 33:6840–6851, 2020. 2, 3
- [25] 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. 6
- [26] Rui Huang, Andrew Geng, and Yixuan Li. On the importance of gradients for detecting distributional shifts in the wild. Advances in Neural Information Processing Systems, 34:677–689, 2021. 5, 17
- [27] Xiaowei Huang, Daniel Kroening, Wenjie Ruan, James Sharp, Youcheng Sun, Emese Thamo, Min Wu, and Xinping Yi. A survey of safety and trustworthiness of deep neural networks: Verification, testing, adversarial attack and defence,

and interpretability. *Computer Science Review*, 37:100270, 2020. 1

- [28] Peng-Tao Jiang, Chang-Bin Zhang, Qibin Hou, Ming-Ming Cheng, and Yunchao Wei. Layercam: Exploring hierarchical class activation maps for localization. *IEEE Transactions on Image Processing*, 30:5875–5888, 2021.
- [29] Oguzhan Fatih Kar, Teresa Yeo, and Amir Zamir. 3d common corruptions for object recognition. In *ICML 2022 Shift Happens Workshop*, 2022. 2
- [30] Ananya Kumar, Aditi Raghunathan, Robbie Jones, Tengyu Ma, and Percy Liang. Fine-tuning can distort pretrained features and underperform out-of-distribution. *arXiv preprint arXiv:2202.10054*, 2022. 6, 18
- [31] Yingwei Li, Qihang Yu, Mingxing Tan, Jieru Mei, Peng Tang, Wei Shen, Alan Yuille, and Cihang Xie. Shapetexture debiased neural network training. *arXiv preprint arXiv:2010.05981*, 2020. 7, 20
- [32] Geert Litjens, Thijs Kooi, Babak Ehteshami Bejnordi, Arnaud Arindra Adiyoso Setio, Francesco Ciompi, Mohsen Ghafoorian, Jeroen Awm Van Der Laak, Bram Van Ginneken, and Clara I Sánchez. A survey on deep learning in medical image analysis. *Medical image analysis*, 42:60–88, 2017. 1
- [33] Weitang Liu, Xiaoyun Wang, John Owens, and Yixuan Li. Energy-based out-of-distribution detection. Advances in Neural Information Processing Systems, 33:21464–21475, 2020. 5, 17
- [34] 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. 6
- [35] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2022. 6
- [36] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. In *International Conference on Learning Representations*, 2021. 3
- [37] Ahmet Murat Ozbayoglu, Mehmet Ugur Gudelek, and Omer Berat Sezer. Deep learning for financial applications: A survey. *Applied Soft Computing*, 93:106384, 2020. 1
- [38] 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. 16
- [39] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021. 2, 6, 18
- [40] Evgenia Rusak, Lukas Schott, Roland S Zimmermann, Julian Bitterwolf, Oliver Bringmann, Matthias Bethge, and

Wieland Brendel. A simple way to make neural networks robust against diverse image corruptions. In *European Conference on Computer Vision*, pages 53–69. Springer, 2020. 7, 20

- [41] 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(3):211–252, 2015. 2, 5
- [42] Hadi Salman, Andrew Ilyas, Logan Engstrom, Ashish Kapoor, and Aleksander Madry. Do adversarially robust imagenet models transfer better? Advances in Neural Information Processing Systems, 33:3533–3545, 2020. 7, 19
- [43] Steffen Schneider, Evgenia Rusak, Luisa Eck, Oliver Bringmann, Wieland Brendel, and Matthias Bethge. Improving robustness against common corruptions by covariate shift adaptation. Advances in Neural Information Processing Systems, 33:11539–11551, 2020. 17
- [44] Vikash Sehwag, Caner Hazirbas, Albert Gordo, Firat Ozgenel, and Cristian Canton. Generating high fidelity data from low-density regions using diffusion models. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 11492–11501, 2022. 3
- [45] Ian Naismith Sneddon. Fourier transforms. Courier Corporation, 1995. 4
- [46] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265. PMLR, 2015. 2, 3
- [47] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*, pages 6105–6114. PMLR, 2019. 6
- [48] Ross Wightman, Hugo Touvron, and Hervé Jégou. Resnet strikes back: An improved training procedure in timm. arXiv preprint arXiv:2110.00476, 2021. 7
- [49] Mitchell Wortsman, Gabriel Ilharco, Jong Wook Kim, Mike Li, Simon Kornblith, Rebecca Roelofs, Raphael Gontijo Lopes, Hannaneh Hajishirzi, Ali Farhadi, Hongseok Namkoong, et al. Robust fine-tuning of zero-shot models. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 7959–7971, 2022. 18
- [50] Kai Xiao, Logan Engstrom, Andrew Ilyas, and Aleksander Madry. Noise or signal: The role of image backgrounds in object recognition. *Proceedings of the International Conference on Learning Representations*, 2021. 2, 4
- [51] Qizhe Xie, Minh-Thang Luong, Eduard Hovy, and Quoc V Le. Self-training with noisy student improves imagenet classification. In *Proceedings of the IEEE/CVF conference* on computer vision and pattern recognition, pages 10687– 10698, 2020. 6, 18
- [52] Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, and Han Hu. Simmim: A simple framework for masked image modeling. In *Proceedings of* the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 9653–9663, 2022. 16

- [53] Hang Zhang, Chongruo Wu, Zhongyue Zhang, Yi Zhu, Haibin Lin, Zhi Zhang, Yue Sun, Tong He, Jonas Mueller, R Manmatha, et al. Resnest: Split-attention networks. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 2736–2746, 2022. 6, 7
- [54] Chuanxia Zheng, Tat-Jen Cham, Jianfei Cai, and Dinh Phung. Bridging global context interactions for high-fidelity image completion. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 11512–11522, June 2022. 4
- [55] Yao Zhu, Yuefeng Chen, Xiaodan Li, Kejiang Chen, Yuan He, Xiang Tian, Bolun Zheng, Yaowu Chen, and Qingming Huang. Toward understanding and boosting adversarial transferability from a distribution perspective. *IEEE Transactions on Image Processing*, 31:6487–6501, 2022. 17

A. Details for ImageNet-E

To guarantee the visual quality of the generated examples, we choose the animal classes from ImageNet since they appear more in nature without messy backgrounds. Specifically, images whose coarse labels in [fish, shark, bird, salamander, frog, turtle, lizard, crocodile, dinosaur, snake, trilobite, arachnid, ungulate, monotreme, marsupial, coral, mollusk, crustacean, marine mammals, dog, wild dog, cat, wild cat, bear, mongoose, butterfly, echinoderms, rabbit, rodent, hog, ferret, armadillo, primate] are picked. The corresponding coarse labels of each class we refer to can be found in $[10]^1$. Finally, our ImageNet-E consists of 373 classes. Since the number of masks provided in ImageNet-S [11] in these classes is 4352, thus the number of images in each edited kind is 4352. The ImageNet-E contains 11 kinds of attributes editing, including 5 kinds of background editing and 4 kinds of size editing, as well as one kind of position editing and one kind of direction editing. Finally, our ImageNet-E contains 47872 images. Experiments on more images can be found in section C.3. The comprehensive comparisons with the stateof-the-art robustness benchmarks are shown in Figure 7. In contrast to other benchmarks that investigate new out-ofdistribution corruptions or perturbations deep models may encounter, w conduct model debugging with in-distribution data to explore which object attributes a model may be sensitive to. The examples in ImageNet-E are shown in Figure 9. A demo video for our editing toolkit can be found at this url:https://drive.google.com/file/d/ 1h5EV3MHPGqkBww9qrhlvrl--kSIrD5Lp/view? usp=sharing. Our code can be found at an anonymous url: https://huggingface.co/spaces/ Anonymous-123/ImageNet-Editing.

Benchmarks	Description	Classes	Samples
ImageNet-A	Challenging examples collected by-hand	200	A * 🚫 🖗
ImageNet-C	Corruptions added on images	1000	
ImageNet-R	Various renditions of ImageNet object classes	200	A 🚳 🧖 🎶
ImageNet-3DCC	3D common corruptions	1000	
ImageNet-9	Images whose objects and backgrounds are disentangled with bbox	370	
ImageNet-E	Images with attribute- edited objects	373	1111

Figure 7. Benchmark comparison.



Figure 8. Comparisons between the amplitude supervision on high-frequency components (HF) and amplitude supervision on all frequency components (All).

B. Background editing

Intuitively, an image with complicated background tends to contain more high-frequency components, such as edges. Therefore, a straight-forward way is to define the background complexity as the amplitude of high-frequency components. However, this operation can result in noisy backgrounds, instead of the ones with complicated textures. Therefore, we directly define complexity as the amplitude of all frequency components. The compared results are shown in Figure 8. It can be observed that the amplitude supervision on high-frequency components tends to make the model generate images with more noise. In contrast, amplitude supervision on all frequency components can help to generate images with texture-complex backgrounds. To edit the background adversarially, we set $\mathcal{L}_c = CE(f(\mathbf{x}), y)$ where 'CE' is the cross entropy loss. f and y are the classifier and label of \mathbf{x} respectively. We adopt the classifier ffrom guided-diffusion².

C. Experimental details

C.1. Details for metrics

In this paper, we care more about how different attributes impact different models. Therefore, we choose the drop of top-1 accuracy as our evaluation metric. A lower dropped accuracy indicates higher robustness against our attribute

¹https://github.com/noameshed/noveltydetection/blob/master/imagenet_categories_synset.csv

²https://github.com/openai/guided-diffusion









Figure 9. Samples from ImageNet-E. From left to right, top to bottom, the images stand for background editing with $\lambda = -20$, $\lambda = 20$, $\lambda = 20$ -adv, randomly shuffled backgrounds, size editing with rate 0.1 and 0.05, randomly rotate, random position, randomly rotate based on images with object pixel rate 0.05 respectively.



Figure 10. Attribute editing with DDPMs. Give an input image and its corresponding object mask, the object is firstly removed with inpainting operation to get the pure background image. Then, we leverage the diffusion process to edit the background image x_0 and object image coherently. \odot denotes the element-wise blending of these two images using the object mask. For background editing, the background complexity objective function is added during the diffusion process (Alg. 1, line 5). For other object attributes editing, the object image needs to be transformed first (Alg. 2, line 1).

changes. The dropped accuracy is defined as:

$$DA = acc_{original} - acc.$$
 (8)

The detailed top-1 accuracy (Top-1) and dropped accuracy (DA)on our ImageNet-E are listed in Table 4, Table 5 and Table 6, Table 7. All the experiments are conducted for 5 runs and we report the mean value in the tables.

C.2. Classes whose accuracy drops the greatest

To find out which class gets the worst robustness against attribute changes, we plot the dropped accuracy in Figure 11. The evaluated models are vanilla RN50 and its Debiased model. It can be observed that objects that have tentacles with simple backgrounds are more easily to be attacked. For example, the dropped accuracy of the 'black widow' class reaches 47% for both vanilla and Debiased models. In contrast, the impact is smaller for images with complicated backgrounds such as pictures from 'squirrel monkey'.

C.3. Experiments on more data

To explore the model robustness against object attributes on large-scale datasets, we step further to conduct the image editing on all the images in the ImageNet-S validation set. Finally, the edited dataset ImageNet-E-L shares the same size as ImageNet-S, which consists of 919 classes and 10919 images. We conduct both background editing and size editing to them. The evaluation results are shown in Table 8. The same conclusion can also be observed. For instance, most models show vulnerability against attribute changing since the average dropped accuracies reach 12.22% and 22.21% in background and size changes respectively. When the model gets larger, the robustness is improved. The consistency implies that using our ImageNet-E can already reflect the model robustness against object attribute changes.

C.4. Bad case analysis

To make a comprehensive study of how the model behaves, we step further to make a comparison of the heat maps of the originals and edited ones. We choose the images that are recognized correctly at first but misclassified after editing. All the attributes editing including background, size, directions are explored. The heat maps are visualized in Figure 12. It can be observed that compared to the SIN and Debiased models, the vanilla RN50 is more likely to lose its focus on the interest area, especially in the size change scenario. For example, in the second row, as it puts his focus on the background, it returns a result with the 'nail' label. The same fashion is also observed in the background change scenario. The predicted label of 'night snake' turns into 'spider web' as the complex background has attracted its attention. In contrast, the SIN and Debiased models have robust attention mechanisms. The quantitative results in Table 5 also validate this. The dropped accuracy of RN50 (13.35%) is higher than SIN (12.19%) and Debiased (11.45%) even though the original accuracy of SIN (0.9157) is lower than vanilla RN50 (0.9269). However, the SIN also has its weakness. We find that though the SIN pays attention to the desired region, it can also make wrong predictions. As shown in the second row of Figure 12, when the object size gets smaller, the shape-based SIN model tends to make wrong predictions, e.g., mistaking the 'sea urchin' as 'acorn' due to the lack of texture analvsis. As a result, the dropped accuracy in the size change scenario is 24.23% for SIN, even lower than vanilla RN50, whose dropped accuracy is 21.26%. On the contrary, the Debiased model can recognize it correctly, profiting from its shape and texture-biased module. From the above observation, we can conclude that the texture matters in the small object scenario.

Table 4. Evaluations under different backgrounds.

Models	Ori	Inv	er	$\lambda = 1$	-20	$\lambda =$	= 20	$\lambda = 2$	0-Adv	Ran	dom
widdels	Top-1	Top-1	DA	Top-1	DA	Top-1	DA	Top-1	DA	Top-1	DA
RN50	92.69%	90.72%	1.97%	85.39%	7.30%	79.34%	13.35%	62.77%	29.92%	79.35%	13.34%
DenseNet121	92.10%	90.61%	1.49%	85.81%	6.29%	83.10%	9.00%	62.90%	29.20%	79.67%	12.43%
EF-B0	92.85%	91.78%	1.07%	85.75%	7.10%	82.14%	10.71%	57.97%	34.88%	77.21%	15.64%
ResNest50	95.38%	93.94%	1.44%	89.05%	6.33%	86.40%	8.98%	68.76%	26.62%	84.10%	11.28%
ViT-S	94.14%	93.32%	0.82%	87.72%	6.42%	85.16%	8.98%	63.02%	31.12%	81.08%	13.06%
Swin-S	96.21%	95.08%	1.13%	91.03%	5.18%	88.88%	7.33%	72.71%	23.50%	86.90%	9.31%
ConvNeXt-T	96.07%	94.64%	1.43%	91.38%	4.69%	89.81%	6.26%	76.24%	19.83%	88.14%	7.93%
RN101	94.00%	91.89%	2.11%	86.95%	7.05%	82.38%	11.62%	64.53%	29.47%	80.43%	13.57%
DenseNet169	92.37%	91.25%	1.12%	86.56%	5.81%	83.94%	8.43%	64.86%	27.51%	80.76%	11.61%
EF-B3	94.97%	93.10%	1.87%	87.20%	7.77%	86.57%	8.40%	65.07%	29.90%	82.05%	12.92%
ResNest101	95.54%	94.44%	1.10%	89.96%	5.58%	88.89%	6.65%	72.51%	23.03%	85.14%	10.40%
ViT-B	95.38%	94.55%	0.83%	90.06%	5.32%	86.95%	8.43%	68.78%	26.60%	84.40%	10.98%
Swin-B	95.96%	95.17%	0.79%	91.50%	4.46%	89.73%	6.23%	74.52%	21.44%	87.71%	8.25%
ConvNeXt-B	96.42%	95.73%	0.69%	92.67%	3.75%	91.56%	4.86%	79.93%	16.49%	90.38%	6.04%

Table 5. Evaluations with different robust models under different backgrounds.

Models	Ori	Inv	er	$\lambda =$	-20	$\lambda =$	= 20	$\lambda = 2$	0-Adv	Ran	dom
widdels	Top-1	Top-1	DA	Top-1	DA	Top-1	DA	Top-1	DA	Top-1	DA
RN50	92.69%	90.72%	1.97%	85.39%	7.30%	79.34%	13.35%	62.77%	29.92%	79.35%	13.34%
RN50-A	81.96%	81.30%	0.66%	77.21%	4.75%	68.34%	13.62%	44.09%	37.87%	66.71%	15.25%
RN50-SIN	91.57%	89.34%	2.23%	83.96%	7.61%	79.38%	12.19%	58.41%	33.16%	77.99%	13.58%
RN50-debiasd	93.34%	91.91%	1.43%	87.25%	6.09%	81.89%	11.45%	65.35%	27.99%	81.22%	12.12%
RN50-Augmix	93.50%	92.52%	0.98%	87.24%	6.26%	85.12%	8.38%	63.01%	30.49%	80.56%	12.94%
RN50-ANT	91.87%	90.19%	1.68%	85.25%	6.62%	79.93%	11.94%	56.21%	35.66%	76.51%	15.36%
RN50-DeepAugment	92.88%	91.38%	1.50%	86.26%	6.62%	80.51%	12.37%	60.48%	32.40%	79.56%	13.32%
RN50-T	94.55%	93.50%	1.05%	88.90%	5.65%	87.17%	7.38%	72.66%	21.89%	84.13%	10.42%

C.5. Details for robustness enhancements

Network design—-self-attention-like architecture. The results in Table 1 show that most vision transformers show better robustness than CNNs in our scenario. Previous study has shown that the self-attention-like architecture may be the key to robustness boost [3]. Therefore, to ablate whether incorporating this module can help attribute robustness generalization, we create a hybrid architecture (RN50d-hybrid) by directly feeding the output of res_3 block in RN50d into ViT-S as the input feature. The results are shown in Table 9. As we can find that while the added module maintains the robustness on background changes, it can help to boost the robustness against size changes. Moreover, the RN50-hybrid can also boost the overall performance compared to ViT-S.

Training strategy—-Masked image modeling. Considering that masked image modeling has demonstrated impressive results in self-supervised representation learning by recovering corrupted image patches [4], it may be robust to the attribute changes. Thus, we test the Masked AutoEncoder (MAE) [18] and SimMIM [52] training strategy based on ViT-B backbone. As shown in Table 10, the dropped ac-

curacies decrease a lot compared to vanilla ViT-B, validating the effectiveness of the masked image modeling strategy. Motivated by this success, we also test another kind of self-supervised-learning strategy. To be specific, we choose the representative method MoCo-V3 [7] in the contrastive learning family. After the end-to-end finetuning, it achieves top-1 83.0% accuracy on ImageNet. It can also improve the attribute robustness when compared to the vanilla ViT-B, showing the effectiveness of contrastive learning.

C.6. Hardware

Our experiments are implemented by PyTorch [38] and runs on RTX-3090TI.

D. Further exploration on backgrounds

Motivated by the models' vulnerability against background changes, especially for those complicated backgrounds. Apart from randomly picking the backgrounds from the ImageNet dataset as final backgrounds (random_bg), we also collect background templates with abundant textures, including leopard, eight diagrams, checker and stripe to explore the performance on out-of-distribution

Table 6. Evaluations under different object sizes.

Models	Ori	Fu	11	0.1	0	0.0	08	0.	05	0.0	5-rp	r	d
widdels	Top-1	Top-1	DA	Top-1	DA	Top-1	DA	Top-1	DA	Top-1	DA	Top-1	DA
RN50	92.69%	89.98%	2.71%	85.44%	7.25%	82.18%	10.51%	71.43%	21.26%	66.23%	26.46%	67.57%	25.12%
DenseNet121	92.10%	88.60%	3.50%	85.10%	7.00%	81.42%	10.68%	70.55%	21.55%	65.57%	26.53%	68.46%	23.64%
EF-B0	92.85%	89.82%	3.03%	84.85%	8.00%	81.28%	11.57%	69.57%	23.28%	64.94%	27.91%	73.74%	19.11%
ResNest50	95.38%	92.85%	2.53%	90.11%	5.27%	87.37%	8.01%	77.35%	18.03%	74.01%	21.37%	78.06%	17.32%
ViT-S	94.14%	93.34%	0.80%	88.77%	5.37%	85.55%	8.59%	76.77%	17.37%	71.28%	22.86%	77.01%	17.13%
Swin-S	96.21%	94.94%	1.27%	92.00%	4.21%	89.92%	6.29%	82.05%	14.16%	78.86%	17.35%	82.79%	13.42%
ConvNeXt-T	96.07%	94.32%	1.75%	92.79%	3.28%	90.89%	5.18%	83.31%	12.76%	80.36%	15.71%	80.29%	15.78%
RN101	94.00%	91.43%	2.57%	87.19%	6.81%	83.88%	10.12%	73.35%	20.65%	68.15%	25.85%	69.58%	24.42%
DenseNet169	92.37%	90.12%	2.25%	85.47%	6.90%	81.96%	10.41%	71.78%	20.59%	67.44%	24.93%	71.69%	20.68%
EF-B3	94.97%	93.61%	1.36%	88.17%	6.80%	84.81%	10.16%	73.61%	21.36%	69.99%	24.98%	77.73%	17.24%
ResNest101	95.54%	94.19%	1.35%	91.57%	3.97%	89.01%	6.53%	80.10%	15.44%	76.43%	19.11%	81.23%	14.31%
ViT-B	95.38%	94.76%	0.62%	91.38%	4.00%	89.08%	6.30%	80.87%	14.51%	76.56%	18.82%	80.43%	14.95%
Swin-B	95.96%	94.97%	0.99%	92.80%	3.16%	90.92%	5.04%	83.62%	12.34%	80.58%	15.38%	83.36%	12.60%
ConvNeXt-B	96.42%	95.43%	0.99%	94.17%	2.25%	93.06%	3.36%	86.95%	9.47%	84.02%	12.40%	83.41%	13.01%

Table 7. Evaluations with different robust models under different object sizes.

Models	Ori	Fu	11	0.10		0.08		0.05		0.05-rp		rd	
widdels	Top-1	Top-1	DA	Top-1	DA	Top-1	DA	Top-1	DA	Top-1	DA	Top-1	DA
RN50	92.69%	89.98%	2.71%	85.44%	7.25%	82.18%	10.51%	71.43%	21.26%	66.23%	26.46%	67.57%	25.12%
RN50-A	81.96%	77.09%	4.87%	72.34%	9.62%	68.02%	13.94%	56.45%	25.51%	49.45%	32.51%	50.00%	31.96%
RN50-SIN	91.57%	89.89%	1.68%	83.27%	8.30%	78.97%	12.60%	67.34%	24.23%	62.41%	29.16%	64.33%	27.24%
RN50-debiasd	93.34%	91.36%	1.98%	87.81%	5.53%	84.58%	8.76%	74.07%	19.27%	69.33%	24.01%	68.37%	24.97%
RN50-Augmix	93.50%	91.89%	1.61%	87.10%	6.40%	83.53%	9.97%	72.08%	21.42%	66.36%	27.14%	71.08%	22.42%
RN50-ANT	91.87%	90.30%	1.57%	84.75%	7.12%	81.25%	10.62%	70.38%	21.49%	65.21%	26.66%	66.64%	25.23%
RN50-DeepAugment	92.88%	91.52%	1.36%	85.61%	7.27%	82.26%	10.62%	71.60%	21.28%	66.60%	26.28%	71.59%	21.29%
RN50-T	94.55%	92.44%	2.11%	89.81%	4.74%	86.72%	7.83%	77.09%	17.46%	73.43%	21.12%	74.95%	19.60%

backgrounds. The evaluation results are shown in Table 12. It can be observed that the background changes can lead to a 13.34% accuracy drop. When the background is set to be a leopard or other images, the dropped accuracy can even reach 35.52%. Sometimes the robust models even show worse robustness. For example, when the background is eight diagrams, all the robust models show worse results than the vanilla RN50, which is quite unexpected. To comprehend the behaviour behind it, we visualize the heat maps of the different models in Figure 7. An interesting finding is that deep models tend to make decisions with dependency on the backgrounds, especially when the background is complicated and can attract some attention. For example, when the background is the eight diagrams, the SIN takes the goldfish as a dishwasher. We suspect it has mistaken the background as dishes. In the same fashion, the Debiased model and ANT take the 'sea slug' with eight diagrams as a 'shopping basket', which seems to make sense since the 'sea slug' looks like a vegetable.

E. Further discussion on the distribution

In this paper, our effort aims to give an editable image tool that can edit the object's attribute in the given image while maintaining it in the original distribution for model debugging. Thus, we choose the outof-distribution (OOD) detection methods including Energy [33] and GradNorm [26] following DRA [55] as the evaluation methods to find out whether our editing tool will move the edited image out of its original distribution. In contrast to FID which indicates the divergence of two datasets, the OOD detection is used to indicate the extent of the deviance of a single input image from the in-distribution dataset.

Covariate shift adaptation(*a.k.a* batch-norm adaptation, BNA) is a way for improving robustness against common corruptions [43]. Thus, it can help to get a top-1 accuracy performance boost in OOD data. One can easily find out if the provided dataset is OOD by checking whether the BNA can get a performance boost on its data. We have tested the full adaptation results using BNA on ResNet50. In contrast to the promotion on other out-of-distribution dataset, we find that this operation induces little changes to top-1 accuracy on both ImageNet validation set $(0.7615 \rightarrow 0.7613)$ and our ImageNet-E $(0.7934 \rightarrow 0.7933 \text{ under } \lambda = 20,$ $0.6521 \rightarrow 0.6514$ under random position scenario, mean accuracy of 5 runs). This similar tendency implies that our ImageNet-E shares a similar property with ImageNet.



Figure 11. Dropped accuracy (%) in each class. Classes whose number of images is less than 15 or dropped accuracy is zero are removed.

Models	Original	al Background Size-0.05		-0.05	Models	Original	Backg	round	Size-0.05		
widdels	Top-1	Top-1	DA	Top-1	DA	Widdels	Top-1	Top-1	DA	Top-1	DA
DenseNet121	86.60%	74.73%	11.87%	61.48%	25.12%	DenseNet169	87.66%	76.26%	11.40%	63.57%	24.09%
RN50	88.12%	71.64%	16.48%	63.13%	24.99%	RN101	89.52%	75.33%	14.19%	65.11%	24.41%
EF-B0	88.54%	75.64%	12.90%	62.16%	26.38%	EF-B3	92.12%	80.81%	11.31%	66.18%	25.96%
ResNest50	92.12%	80.61%	11.51%	70.05%	22.07%	ResNest101	92.78%	83.46%	9.32%	72.67%	20.11%
ViT-S	92.15%	78.94%	13.21%	69.30%	22.85%	ViT-B	94.12%	83.04%	11.08%	75.65%	18.47%
Swin-S	93.11%	82.98%	10.13%	75.36%	17.75%	Swin-B	93.18%	84.11%	9.07%	76.99%	16.19%
ConvNeXt-T	92.75%	84.00%	9.43%	76.41%	16.34%	ConvNeXt-B	94.05%	86.41%	7.64%	80.34%	13.71%

Table 8. Evaluations with more data.

F. Further evaluation on more state-of-the-art models

To provide evaluations on more state-of-the-art models, we step further to evaluate the CLIP [39] and EfficientNet-L2-Noisy-Student [51]. The average dropped accuracy in terms of different models can be found in Figure 13. CLIP shows a good robustness to out-of-distribution data [30]. Therefore, to find out whether the CLIP can also show a good robustness against attribute editing, we evaluate the CLIP model (Backbone ViT-B) with both the zero-shot and end-to-end finetuned version. To achieve this, we finetune the pretrained CLIP on the ImageNet training dataset based on prompt-initialized model following [49]. It acquires a 81.2% top-1 accuracy on ImageNet validation set while it is 68.3% for zero-shot version. The evaluation on ImageNet-E is shown in Table 11 and Table 13. Though previous studies have shown that the zero-shot CLIP model exhibits better out-of-distribution robustness than the finetuned ones, the finetuned CLIP shows better attribute robustness on ImageNet-E, as shown in Table 11 and Table 13. The tendency on ImageNet-E is the same with ImageNet (IN) validation set and ImageNet-V2 (IN-V2). This implies that the ImageNet-E shows a better proximity to ImageNet than other out-of-distribution benchmarks such as ImageNet-C (IN-C), ImageNet-A (IN-A). Another finding is that the CLIP model fails to show better robustness than ViT-B while they share the same architectures. We suspect that this is caused by that CLIP may have spared some capacity for out-of-distribution robustness. As the network gets larger, its attribute robustness gets better.

While EfficientNet-L2-Noisy-Student is one of the top models on ImageNet-A benchmark [51], it also shows superiority on ImageNet-E. To delve into the reason behind this, we test EfficientNet-L2-Noisy-Student-475 (EF-L2-NT-475) and EfficientNet-B0-Noisy-Student (EF-B0-NT). The EF-L2-NT-475 differs from EF-L2-NT in terms of input size, which former is 475 while it is 800 for the latter. It can be found that the input size can induce little improvement to the attribute robustness. In contrast, larger networks can benefit a lot to attribute robustness, which is consistent with the finding in Section 5.1.

Evaluations on 91 state-of-the-art models can be found in Figure 14. All the evaluated models in this figure are all



Figure 12. The heat map comparisons between original images and edited ones.

Table 9. Ablation study of the self-attention-like architecture.

Models	Ori		Ba	ckground o	changes			Size o	changes		Position	Direction	Ava
Widdels	UII	Inver	$\lambda = -20$	$\lambda = 20$	$\lambda = 20$ -Adv	Random	Full	0.1	0.08	0.05	rp	rd	Avg.
R50d	93.77%	1.23%	4.80%	6.48%	19.39%	8.28%	2.82%	4.36%	7.07%	16.95%	20.49%	19.31%	11.00%
ViT-S	94.74%	1.66%	7.32%	10.64%	32.17%	14.39%	1.22%	7.10%	10.64%	20.29%	25.08%	17.22%	14.61%
R50-hybrid	95.40%	1.04%	5.64%	7.16%	21.54%	9.19%	1.37%	3.53%	5.92%	13.92%	17.23%	14.12%	9.96%



Figure 13. The average accuracy drop of different models. The x-axis is the model's top-1 accuracy on ImageNet.

provided by the timm library, except for the MoCo-V3-FT and CLIP-FT, which are finetuned by us.

G. Failure cases of generated images

The failure cases of generated images are shown in Figure 16. The diffusion model fails to generate high-quality person images. Though the object is reserved, the whole image looks quite wired. Therefore, we only keep the animal classes, resulting a compact set of ImageNet-E. However, extensive evaluations to 919 in Section C.3 have witnessed a same conclusion with evaluations on 373 classes. This implies that using our ImageNet-E can already reflect the model robustness against object attribute changes.

H. Related literature to robustness enhancements

Adversarial training. [42] focus on adversarially robust ImageNet classifiers and show that they yield improved accuracy on a standard suite of downstream classification tasks. It provides a strong baseline for adversarial training. Therefore, we choose their officially released adversarially trained models³ as the evaluation model. Models with different architectures are adopted here⁴.

SIN [13] provides evidence that machine recognition today overly relies on object textures rather than global object shapes, as commonly assumed. It demonstrates the advantages of a shape-based representation for robust inference

³https://github.com/microsoft/robust-models-transfer

⁴https://github.com/alibaba/easyrobust

Background changes Size changes Position Direction Models Ori Avg. Inver $\lambda = -20$ $\lambda = 20$ $\lambda = 20$ -Adv Random Full 0.1 0.08 0.05 rd rp ViT-B 95.38% 0.83% 5.32% 8.43% 26.60% 10.98% 0.62% 4.00% 6.30% 14.51% 18.82% 14.95% 11.05% CLIP_finetune 93.68% 2.17% 9.82% 11.83% 38.33% 18.19% 9.06% 9.25% 12.67% 23.32% 28.56% 22.00% 18.30% MoCo-v3 95.70% 0.55% 4.91% 7.33% 24.33% 9.92% 0.92% 3.76% 5.62% 13.61% 17.85% 15.20% 10.35% MAE-ViT-B 96 12% 078% 4.77% 6.21% 21.09% 8.18% 0.78% 3.01% 4.86% 12.10% 15.47% 14 00% 9.05% SimMIM 96.14% 13.62% 9.82% 0.75% 5.13% 6.76% 23.58% 9.33% 0.97% 3.22% 5.33% 13.18% 17.12% original BG: λ=0 BG: λ=-20 BG: $\lambda = 20$ BG: random Size: 0.1 Size: 0.08 Size: 0.05 Position Direction Supervised ViTs Standard supervised CNNs Self-Supervised ViTs Data-rich models

Table 10. Ablation study of the self-supervised models. All the compared models are end-to-end finetuned on ImageNet except for ViT-B, which is supervised trained from the early start.

Figure 14. The top-1 accuracy performance under different editing scenarios of 91 state-of-the-art models.

Table 11. Evaluations on different robustness benchmarks. All results are top-1 accuracies(%) on corresponding datasets except for ImageNet-C, which is mCE (mean Corruption Error). Higher top-1 accuracy and lower mCE indicate better performance.

Models	IN	IN-V2	IN-A	IN-C	IN-R	IN-Sketch	IN-E
CLIP-zero-shot	68.3	61.9	50.1	43.1	77.6	48.3	62.1
CLIP-FT	81.2	70.7	35.3	479	65.0	44.9	77.2

(using their Stylized-ImageNet dataset to induce such a representation in neural networks)

Debiased [31] shows that convolutional neural networks are often biased towards either texture or shape, depending on the training dataset, and such bias degenerates model performance. Motivated by this observation, it develops a simple algorithm for shape-texture Debiased learning. To prevent models from exclusively attending to a single cue in representation learning, it augments training data with images with conflicting shape and texture information (*e.g.*, an image of chimpanzee shape but with lemon texture) and provides the corresponding supervision from shape and texture simultaneously. It empirically demonstrates the advantages of the shape-texture Debiased neural network training on boosting both accuracy and robustness.

Augmix [22] focuses on the robustness improvement to

unforeseen data shifts encountered during deployment. It proposes a data processing technique named Augmix that helps to improve robustness and uncertainty measures on challenging image classification benchmarks.

ANT [40] demonstrates that a simple but properly tuned training with additive Gaussian and Speckle noise generalizes surprisingly well to unseen corruptions, easily reaching the previous state of the art on the corruption benchmark ImageNet-C and on MNIST-C.

DeepAugment [20]. Motivated by the observation that using larger models and artificial data augmentations can improve robustness on real-world distribution shifts, contrary to claims in prior work. It introduces a new data augmentation method named DeepAugment, which uses image-to-image neural networks for data augmentation. It improves robustness on their newly introduced ImageNet-R benchmark and can also be combined with other augmentation methods to outperform a model pretrained on 1000× more labeled data.

There are some more tables and figures in the next pages.

Models	Original	Random_bg		Leopard		Eight diagrams		Checker		Stripe	
		Top-1	DA	Top-1	DA	Top-1	DA	Top-1	DA	Top-1	DA
RN50	92.69%	79.35%	13.34%	57.17%	35.52%	64.32%	28.37%	65.13%	27.56%	62.90%	29.79%
RN50-A	81.96%	66.71%	15.25%	25.05%	56.91%	37.21%	44.75%	32.47%	49.49%	46.96%	35.00%
RN50-SIN	91.57%	77.99%	13.58%	62.74%	28.83%	48.74%	42.83%	51.15%	40.42%	52.65%	38.92%
RN50-debiasd	93.34%	81.22%	12.12%	68.58%	24.76%	62.68%	30.66%	67.10%	26.24%	63.16%	30.18%
RN50-Augmix	93.50%	80.56%	12.94%	57.35%	36.15%	56.20%	37.30%	68.78%	24.72%	65.68%	27.82%
RN50-ANT	91.87%	76.51%	15.36%	58.11%	33.76%	59.04%	32.83%	51.91%	39.96%	54.69%	37.18%
RN50-DeepAugment	92.88%	79.56%	13.32%	62.83%	30.05%	57.71%	35.17%	59.46%	33.42%	61.80%	31.08%
R50-T	94.55%	84.13%	10.42%	72.93%	21.62%	73.98%	20.57%	79.42%	15.13%	76.43%	18.12%

Table 12. Evaluation of images generated with different backgrounds.



Figure 15. Heat maps under different backgrounds.

Table 13. More evaluations on state-of-the-art models including CLIP and EfficientNet-L2-Noisy-Student.

Models	Ori	Background changes						Size changes				Direction	1	
		Inver	$\lambda = -20$	$\lambda = 20$	$\lambda = 20$ -Adv	Random	Full	0.1	0.08	0.05	rD	rd	Avg.	
ViT-B/16	95 38%	0.83%	5 32%	8 43%	26.60%	10.98%	0.62%	4 00%	6 30%	14 51%	18.82%	14.95%	11.05%	
Trenshit											11100 /0			
CLIP_RN50 72 38% 6.03% 11.64% 16.72% 35.07% 21.82% 8.78% 14.39% 17.69% 26.48%										20 70%	25 31%	20.77%		
CLIP RN101	73 35%	4 51%	10.77%	14 42%	33.42%	19.63%	6 30%	14.53%	18 19%	26.58%	30.08%	24.51%	19.85%	
CLIP PN50x4	77 18%	4.51%	10.77%	13 27%	31 30%	18 51%	7 16%	12 37%	15.66%	20.30%	27 10%	24.51%	18 /8%	
CLIP RN50x16	82 10%	4 39%	10.10%	12 41%	27 14%	16.62%	6.62%	11 10%	13.53%	27.09%	25.27%	23.13%	16.80%	
CLIP PN50×64	85 660L	4.3770	8 800%	10 70 %	27.1470	12 44 07-	6 200%	0.200%	11.020	10 170%	23.27%	20.57%	14 57 0%	
CLIP_KN50x04	74 0907	4.1170	0.07 % 12 2407	10.79 70	42.1570	13.44 % 26.200	2.000	9.20%	10.740	20.050	21.02%	20.37%	14.37%	
CLIP_VII-B/32	74.08%	5.55%	15.24%	18.04%	45.20%	20.39%	2.99%	13.39%	19.74%	29.03%	33.37%	24.89%	22.72%	
CLIP_ViT-B/16	80.01%	4.88%	11.56%	15.28%	36.14%	20.09%	4.88%	12.67%	15.77%	25.31%	28.87%	21.57%	19.21%	
CLIP_ViT-L/14	87.61%	4.35%	11.04%	14.46%	33.69%	18.35%	1.81%	11.67%	15.09%	23.66%	27.19%	18.05%	17.50%	
CLIP_ViT-L/14-336	88.01%	3.16%	9.07%	12.25%	29.69%	16.08%	3.16%	9.20%	11.78%	19.94%	22.89%	16.15%	15.02%	
CLIP_ViT-L/14-336	88.01%	3.16%	9.07%	12.25%	29.69%	16.08%	3.16%	9.20%	11.78%	19.94%	22.89%	16.15%	15.02%	
Finetune														
CLIP_ViT-B/16-FT	93.68%	2.17%	9.82%	11.83%	38.33%	18.19%	4.66%	9.25%	12.67%	23.32%	28.56%	22.00%	17.86%	
CLIP_ViT-L/14-336-FT	96.97%	1.29%	5.16%	6.18%	19.93%	8.09%	1.29%	3.47%	4.90%	10.98%	13.74%	10.96%	8.47%	
EF-B0	92.85%	1.07%	7.10%	10.71%	34.88%	15.64%	3.03%	8.00%	11.57%	23.28%	27.91%	19.11%	16.12%	
EF-B0-NT	94.30%	1.97%	8.43%	10.51%	34.93%	15.99%	1.79%	7.91%	11.50%	22.96%	27.62%	19.07%	16.07%	
EF-B7	97.10%	1.80%	6.37%	7.20%	23.36%	10.78%	1.65%	4.16%	6.25%	14.13%	17.12%	10.56%	10.16%	
EF-B7-NT	97.38%	1.30%	5.26%	6.10%	19.96%	9.15%	0.55%	3.31%	4.75%	10.67%	12.87%	7.98%	8.06%	
EF-L2-NT-475	97.84%	1.08%	3.60%	4.51%	14.88%	7.14%	0.51%	2.21%	2.71%	5.50%	7.35%	4.58%	5.30%	
EF-L2-NT	97.63%	1.26%	3.50%	4.06%	12.73%	6.90%	0.71%	2.27%	2.79%	5.01%	6.03%	4.55%	4.85%	



Figure 16. The failure cases of attribute editing.